



Politecnico di Bari

Repository Istituzionale dei Prodotti della Ricerca del Politecnico di Bari

Knowledge-Enabled Recommender Systems in the Linked Data Era

This is a PhD Thesis

Original Citation:

Knowledge-Enabled Recommender Systems in the Linked Data Era / Anelli, Vito Walter. - ELETTRONICO. - (2020).
[10.60576/poliba/iris/anelli-vito-walter_phd2020]

Availability:

This version is available at <http://hdl.handle.net/11589/191260> since: 2020-02-23

Published version

DOI:10.60576/poliba/iris/anelli-vito-walter_phd2020

Publisher: Politecnico di Bari

Terms of use:

(Article begins on next page)



Politecnico
di Bari

Department of Electrical and Information Engineering
Electrical and Information Engineering

Ph.D. Program

SSD: ING-INF/05-INFORMATION PROCESSING SYSTEMS

Final Dissertation

Knowledge-Enabled Recommender Systems in the Linked Data Era

by

Vito Walter Anelli

Supervisors:

Prof. Eng. Tommaso Di Noia

Prof. Eng. Eugenio Di Sciascio

Coordinator of Ph.D. Program:

Prof. Eng. Alfredo Grieco

Course n° 32, 01/11/2016-31/10/2019



Politecnico
di Bari

LIBERATORIA PER L'ARCHIVIAZIONE DELLA TESI DI DOTTORATO

Al Magnifico Rettore
del Politecnico di Bari

Il/la sottoscritto VITO WALTER ANELLI nato a BARI (BA) il 13/10/1984 residente a BARI in via DALLA CHIESA 24/A e-mail vitowalteranelli@gmail.com , vitowalter.anelli@poliba.it

iscritto al 3° anno di Corso di Dottorato di Ricerca in INGEGNERIA ELETTRICA E DELL'INFORMAZIONE ciclo XXXII

ed essendo stato ammesso a sostenere l'esame finale con la prevista discussione della tesi dal titolo:

KNOWLEDGE-ENABLED RECOMMENDER SYSTEMS IN THE LINKED DATA ERA

DICHIARA

- 1) di essere consapevole che, ai sensi del D.P.R. n. 445 del 28.12.2000, le dichiarazioni mendaci, la falsità negli atti e l'uso di atti falsi sono puniti ai sensi del codice penale e delle Leggi speciali in materia, e che nel caso ricorressero dette ipotesi, decade fin dall'inizio e senza necessità di nessuna formalità dai benefici conseguenti al provvedimento emanato sulla base di tali dichiarazioni;
- 2) di essere iscritto al Corso di Dottorato di ricerca INGEGNERIA ELETTRICA E DELL'INFORMAZIONE ciclo XXXII, corso attivato ai sensi del "Regolamento dei Corsi di Dottorato di ricerca del Politecnico di Bari", emanato con D.R. n.286 del 01.07.2013;
- 3) di essere pienamente a conoscenza delle disposizioni contenute nel predetto Regolamento in merito alla procedura di deposito, pubblicazione e autoarchiviazione della tesi di dottorato nell'Archivio Istituzionale ad accesso aperto alla letteratura scientifica;
- 4) di essere consapevole che attraverso l'autoarchiviazione delle tesi nell'Archivio Istituzionale ad accesso aperto alla letteratura scientifica del Politecnico di Bari (IRIS-POLIBA), l'Ateneo archiverà e renderà consultabile in rete (nel rispetto della Policy di Ateneo di cui al D.R. 642 del 13.11.2015) il testo completo della tesi di dottorato, fatta salva la possibilità di sottoscrizione di apposite licenze per le relative condizioni di utilizzo (di cui al sito <http://www.creativecommons.it/Licenze>), e fatte salve, altresì, le eventuali esigenze di "embargo", legate a strette considerazioni sulla tutelabilità e sfruttamento industriale/commerciale dei contenuti della tesi, da rappresentarsi mediante compilazione e sottoscrizione del modulo in calce (Richiesta di embargo);
- 5) che la tesi da depositare in IRIS-POLIBA, in formato digitale (PDF/A) sarà del tutto identica a quelle **consegnate**/inviata/da inviarsi ai componenti della commissione per l'esame finale e a qualsiasi altra copia depositata presso gli Uffici del Politecnico di Bari in forma cartacea o digitale, ovvero a quella da discutere in sede di esame finale, a quella da depositare, a cura dell'Ateneo, presso le Biblioteche Nazionali Centrali di Roma e Firenze e presso tutti gli Uffici competenti per legge al momento del deposito stesso, e che di conseguenza va esclusa qualsiasi responsabilità del Politecnico di Bari per quanto riguarda eventuali errori, imprecisioni o omissioni nei contenuti della tesi;
- 6) che il contenuto e l'organizzazione della tesi è opera originale realizzata dal sottoscritto e non compromette in alcun modo i diritti di terzi, ivi compresi quelli relativi alla sicurezza dei dati personali; che pertanto il Politecnico di Bari ed i suoi funzionari sono in ogni caso esenti da responsabilità di qualsivoglia natura: civile, amministrativa e penale e saranno dal sottoscritto tenuti indenni da qualsiasi richiesta o rivendicazione da parte di terzi;
- 7) che il contenuto della tesi non infrange in alcun modo il diritto d'Autore né gli obblighi connessi alla salvaguardia di diritti morali ed economici di altri autori o di altri aventi diritto, sia per testi, immagini, foto, tabelle, o altre parti di cui la tesi è composta.

Luogo e data BARI, 17/02/2020

Firma Vito Walter Anelli

Il/La sottoscritto, con l'autoarchiviazione della propria tesi di dottorato nell'Archivio Istituzionale ad accesso aperto del Politecnico di Bari (POLIBA-IRIS), pur mantenendo su di essa tutti i diritti d'autore, morali ed economici, ai sensi della normativa vigente (Legge 633/1941 e ss.mm.ii.),

CONCEDE

- al Politecnico di Bari il permesso di trasferire l'opera su qualsiasi supporto e di convertirla in qualsiasi formato al fine di una corretta conservazione nel tempo. Il Politecnico di Bari garantisce che non verrà effettuata alcuna modifica al contenuto e alla struttura dell'opera.
- al Politecnico di Bari la possibilità di riprodurre l'opera in più di una copia per fini di sicurezza, back-up e conservazione.

Luogo e data BARI, 17/02/2020

Firma Vito Walter Anelli



Politecnico
di Bari

Department of Electrical and Information Engineering

Electrical and Information Engineering

Ph.D. Program

SSD: ING-INF/05-INFORMATION PROCESSING SYSTEMS

Final Dissertation

Knowledge-Enabled Recommender Systems in the Linked Data Era

by

Vito Walter Anelli

Referees:

Prof. Alejandro Bellogín

Prof. Mathieu D'Aquin

Supervisors:

Prof. Eng. Tommaso Di Noia

Prof. Eng. Eugenio Di Sciascio

Coordinator of Ph.D. Program:

Prof. Eng. Alfredo Grieco

Course n°32, 01/11/2016-31/10/2019

*If your daily life
seems poor,
do not blame it.*

*Blame yourself,
tell yourself that
you are not poet enough
to call forth its riches.*

R. M. R.

Contents

| | |
|---|-----------|
| Abstract | i |
| 1 Introduction | 1 |
| 1.1 Motivation | 2 |
| 1.2 Research Questions | 3 |
| 1.3 Structure of the work | 9 |
| I Background | 11 |
| 2 Knowledge Representation | 13 |
| 2.1 Semantic Web | 14 |
| 2.1.1 Resource Description Framework | 15 |
| 2.1.2 RDF Schema | 17 |
| 2.1.3 Linked Data | 18 |
| 2.1.4 SPARQL | 19 |
| 2.2 Preference representation | 20 |
| 2.2.1 CP-theories | 22 |
| 3 Recommender Systems | 33 |
| 3.1 Introduction | 33 |
| 3.1.1 The recommendation problem | 33 |
| 3.1.2 From Rating Prediction to Ranking | 37 |
| 3.2 Collaborative-Filtering Recommender Systems | 38 |

| | | |
|-----------|---|-----------|
| 3.2.1 | Neighborhood-Based Models | 39 |
| 3.2.2 | Matrix Factorization - SVD & SVD++ | 42 |
| 3.2.3 | Advances in Matrix factorization | 48 |
| 3.2.4 | Matrix Factorization - Funk MF | 49 |
| 3.2.5 | Matrix Factorization for positive-only feedback | 50 |
| 3.3 | Optimization techniques | 51 |
| 3.3.1 | ALS - Alternating Least Squares | 51 |
| 3.3.2 | BPR - Bayesian Personalized Ranking Criterion | 53 |
| 3.4 | Content-Based Recommender Systems | 54 |
| 3.4.1 | Vector Space Model | 58 |
| 3.4.2 | From Vector space Model to Knowledge-aware Recommender Systems | 60 |
| 3.5 | Hybrid Recommenders | 62 |
| 3.6 | Recommender Systems Evaluation | 65 |
| 3.6.1 | Protocols | 66 |
| 3.6.2 | Accuracy | 67 |
| 3.6.3 | Diversity | 69 |
| 3.6.4 | Novelty | 72 |
| II | Feeding RSs with explicit knowledge | 75 |
| 4 | Introduction | 77 |
| 5 | From semantic graphs to matrices | 81 |
| 5.1 | Introduction | 81 |
| 5.2 | Related Work | 84 |
| 5.3 | Approach | 89 |
| 5.3.1 | Motivation | 89 |
| 5.3.2 | Data Model | 89 |
| 5.3.3 | Problem Formulation | 90 |
| 5.3.4 | Pure FF | 92 |
| 5.3.5 | Jaccard- f_{sr} | 94 |

| | | |
|----------|---|------------|
| 5.3.6 | Content Based- and Hybrid- f_{sr} | 95 |
| 5.3.7 | f_{rsCBF} | 97 |
| 5.4 | Experimental Evaluation | 98 |
| 5.4.1 | Experiments for Pure FF | 98 |
| 5.4.2 | Experimental Evaluation for FF extensions | 101 |
| 5.5 | Conclusion | 116 |
| 6 | Metadata to address Cold-start problem | 119 |
| 6.1 | Introduction | 119 |
| 6.2 | Related work | 122 |
| 6.2.1 | Cross-domain recommender systems | 122 |
| 6.2.2 | Matrix factorization-based cross-domain recommender systems | 125 |
| 6.3 | Matrix factorization models for cross-domain recommendation | 129 |
| 6.3.1 | Regularization through similarity prediction | 130 |
| 6.3.2 | Regularization based on item neighborhoods | 132 |
| 6.3.3 | Regularization based on item centroids | 135 |
| 6.4 | Experiments | 138 |
| 6.4.1 | Dataset | 138 |
| 6.4.2 | Evaluation methodology and metrics | 142 |
| 6.4.3 | Evaluated methods | 146 |
| 6.4.4 | Results | 150 |
| 6.4.5 | Inter-domain item semantic similarity | 150 |
| 6.4.6 | Item ranking accuracy | 151 |
| 6.4.7 | Recommendation diversity | 157 |
| 6.5 | Conclusions and future work | 160 |
| 7 | Interpretability of Factorization Machines | 165 |
| 7.1 | Introduction | 165 |
| 7.2 | Related Work | 169 |
| 7.3 | Approach | 172 |
| 7.3.1 | Knowledge Graphs and Linked Data | 173 |

| | | |
|----------|---|------------|
| 7.3.2 | Formal Model | 174 |
| 7.3.3 | Optimization | 178 |
| 7.3.4 | Personalized Recommendation | 180 |
| 7.3.5 | Semantic Accuracy | 181 |
| 7.3.6 | Robustness | 181 |
| 7.4 | Experimental Evaluation | 182 |
| 7.4.1 | Experimental Setup | 182 |
| 7.4.2 | Features extraction | 185 |
| 7.4.3 | Accuracy, Diversity and Novelty with kaHFM | 187 |
| 7.4.4 | Semantic Accuracy | 190 |
| 7.4.5 | Generative Robustness | 191 |
| 7.5 | Conclusion | 192 |
| 8 | Reasoning about Preferences | 195 |
| 8.1 | Introduction | 195 |
| 8.2 | Related work | 198 |
| 8.3 | Motivating scenario | 203 |
| 8.4 | CP-theories and Linked Open Data | 205 |
| 8.4.1 | The special case of CP-nets | 214 |
| 8.4.2 | Ordering SPARQL results via CP-theories | 215 |
| 8.5 | Query formulation algorithm for CP-theories | 217 |
| 8.6 | Ordering Query for the Book Example | 220 |
| 8.6.1 | Instantiation of the framework | 226 |
| 8.7 | Application | 229 |
| 8.8 | Experiments | 233 |
| 8.8.1 | Test on Real Users | 234 |
| 8.8.2 | Test on Simulated Users | 235 |
| 8.9 | Conclusion and future work | 236 |
| 9 | Hybrid Relevance: How to enhance traditional relevance weighting schemes | 243 |
| 9.1 | Introduction | 243 |
| 9.2 | CReL-FM: Collaborative-aware relevance for recommendation | 245 |

| | | |
|---|---|------------|
| 9.3 | Experimental Evaluation | 249 |
| 9.3.1 | Results | 252 |
| 9.3.2 | Discussion | 254 |
| 9.4 | Related Work | 256 |
| 9.5 | Conclusions and Future Work | 259 |
| III Modeling and evaluating without an explicit knowledge representation | | 261 |
| 10 | Introduction | 263 |
| 11 | A re-ranking algorithm exploiting Time and Diversity | 267 |
| 11.1 | Introduction | 267 |
| 11.2 | Related Work | 269 |
| 11.3 | Intent-aware diversification for recommendations | 270 |
| 11.4 | Intent modeling with Temporal Dynamics | 271 |
| 11.4.1 | Time-Based Intent Modeling | 271 |
| 11.4.2 | Session-Based Intent Modeling | 272 |
| 11.5 | Experimental setting | 273 |
| 11.6 | Conclusions and future work | 276 |
| 12 | A study over the dimensions of Time, and Popularity | 279 |
| 12.1 | Introduction | 279 |
| 12.2 | Time-aware Local Popularity | 281 |
| 12.3 | Experimental Evaluation | 285 |
| 12.4 | Conclusion | 287 |
| 13 | Beyond similarities: Dissimilarity and Asymmetric Similarities | 289 |
| 13.1 | Introduction | 289 |
| 13.2 | Dissimilarity in Recommendation | 291 |
| 13.2.1 | Motivation | 291 |
| 13.2.2 | Metrics | 294 |

| | | |
|-----------|---|------------|
| 13.3 | Experimental Evaluation | 296 |
| 13.4 | Conclusion and Future Work | 302 |
| 14 | An investigation over hyperparameters tuning: A Discriminative Power perspective | 305 |
| 14.1 | Introduction | 305 |
| 14.2 | Experimental Settings | 308 |
| 14.3 | Discriminative Power of metrics on Hyperparameters | 310 |
| 14.4 | Metrics Confidence | 313 |
| 14.5 | Dominant Hyperparameter | 314 |
| 14.6 | Conclusions and Future Work | 316 |
| 15 | Generalized Cross-Entropy for Fairness Evaluation | 317 |
| 15.1 | Introduction | 317 |
| 15.2 | Background and Prior Art | 323 |
| 15.2.1 | Fairness notions | 323 |
| 15.2.2 | Fairness and accuracy trade-off | 324 |
| 15.2.3 | From reciprocal recommendation to multiple stakeholders | 326 |
| 15.2.4 | Evaluating fairness in recommender systems | 328 |
| 15.3 | A probabilistic framework to evaluate fairness | 330 |
| 15.3.1 | Using Generalized Cross Entropy to measure user and item fairness | 330 |
| 15.3.2 | Toy example | 335 |
| 15.4 | Experimental settings | 337 |
| 15.4.1 | Datasets | 337 |
| 15.4.2 | Evaluation Protocol and Temporal Splitting | 338 |
| 15.4.3 | Attribute selection and discretization | 339 |
| 15.4.4 | Baseline recommenders | 342 |
| 15.4.5 | Evaluation metrics | 347 |
| 15.5 | Results | 350 |
| 15.5.1 | Analysis of item fairness results | 350 |
| 15.5.2 | Analysis of user fairness results | 353 |

CONTENTS

vii

15.5.3 Discussion 354
15.6 Conclusions and future work 356

IV Conclusion and Future Work 365

Bibliography 369

Acknowledgements 435

Abstract

Recommender Systems are unfamiliar to ordinary people. However, they are almost everywhere. In a world that overwhelms us with relevant and irrelevant information, they make the difference. They process catalogs from thousands to millions of items to return us only the relevant and personalized information. Otherwise, we would be as castaways in the ocean of information that try to drink it all. On the other side, they are constantly whispering in our ear, suggesting to enjoy news, movies, songs. If they are Jiminy or Lamp-Wick, it is our responsibility. At the same time, the Web is evolving, providing us rich and semantic information. The so-called Semantic Web lets us feed Recommender Systems with high-quality knowledge. This knowledge lets Recommender Systems understand the domain, provide explanations, improve the quality of recommendations. In this research journey, we have faced different aspects of the recommendation and the multiple ways semantic knowledge can be beneficial. We have first focused on Matrix Factorization, a recommendation technique, and we have proposed several ways to exploit knowledge. Feature Factorization, Graph Spreading Relevance, and interpretable Factorization Machines are a few examples. We have developed a recommender that takes into account conditional pair-wise preferences to lower the human-machine barrier. We have also faced the semi-structured knowledge, proposing models that consider temporal diversification, personalized popularity, and dissimilarity. Finally, we have focused on Recommender Systems evaluation, proposing new techniques for tuning hyperparameters, and a new notion of fairness.

We hope you will enjoy the journey.

Chapter 1

Introduction

In the last decades, we have observed a quick growth of the Web. In the beginning, it was often used only by insiders. Now, it has become a powerful tool used by everyone. Concurrently with its spread, also the role of users has changed. Before, they were only consumers of the available information. Nowadays, they also produce a large amount of information available on the Web. This evolution, named "Web 2.0" consisted of a mass injection of documents and websites on the Web.

When we type a single keyword on a search engine, it could even return millions of documents. Our parents and grandparents have left us many different senses for the term Freedom, the possibility of choosing among alternatives. However, when we have a million alternatives, a human is not able to process them. In practical terms, having one million alternatives is like having no one. This effect takes the name of "misery-inducing tyranny".

To avoid this problem, a new need has emerged of processing, annotating, and indexing these documents to make their search easier. The development of these techniques has led to the search engines era. These agents can process a considerable amount of data to return a relevance value for each item. All the issues related to the search have given birth to a new research field: Information Retrieval. However, only a small fraction of search engines consider some degrees of personalization of search results. With the explicit aim of personalization, a new class of agents has emerged: Recommender systems. The rationale of these tools is filter-

ing out irrelevant information on a per user basis. They are designed to analyze big collections of objects, products, or services and extract personalized relevant shortlists.

To make a Recommender Systems work properly, the system must understand the domain knowledge. Frequently, the problem is worked around by the designers encoding the domain knowledge in the Recommender Systems implementation. However, it is reasonable that all the needed knowledge is already available on the Web. This claim leads directly to the second problem. Until now, we have discussed information and documents on the Web irrespective of their nature. If we consider the Web as a collection of documents, we are claiming that some of these documents contain the desired knowledge. Here, a problem arises. Usually, the interpretation of a non-structured document (i.e., plain text) is a typically human task. Unfortunately, since we have observed that the number of documents on the Web is overwhelming, we obtain a new overloading problem.

This is basically due to the documental nature of the Web, in which documents enclose non-structured information. However, even though each document contained only one structured fact, human interpretation would be unfeasible. These considerations lead to two new needs. The first is these facts should be machine-readable. If we had these kinds of data, a software agent could automatically analyze them. The second is the meaning of data should be explicitly defined. If the semantics comes with data, the software agent can understand it. This new Web is named "Semantic Web".

1.1 Motivation

The emerging of the new knowledge bases based on Semantic Web technologies enabled the creation of a new class of Recommender Systems. The exploitation of this kind of knowledge can improve Recommender Systems under many different aspects and alleviate some known issues. For instance, some traditional Recommender Systems approaches suffer the Cold-Start and the over-specialization problem. Cold-Start is a problem that a designer has to face when a new item or a new

user enters the platform. If the suggestions provided to the user are produced based on the platform's past transactions, a new item will not ever be recommended. In this case, we are talking about the Cold-Start item problem. Analogously, if the suggestions are produced based on the user's past transactions, we can not provide any recommendations to the new user. Instead, the over-specialization problem is the tendency of some Recommender Systems to suggest always items that are very similar to the user's history.

These problems indicate that there are many different dimensions to consider to evaluate Recommender Systems. Beyond the classic accuracy dimension, nowadays, we usually also evaluate Coverage, Novelty, Diversity, and Serendipity of results. This research work was born with the idea of exploiting the knowledge enclosed in Linked Open Data datasets to alleviate the Recommender Systems problems.

Almost immediately, we have faced many other problems to inject this knowledge into Recommender Systems. From the choice of the preferred knowledge source to the evaluation of the knowledge source quality, this research work has become a journey through the grey areas of the Semantic Web and the Recommender Systems.

Despite the newborn research field, we have not been alone during this travel. Other researches have perceived the importance of leveraging Linked Data knowledge in Recommender Systems and Information Retrieval. Several technical tracks on Knowledge-based Recommender Systems have begun appearing in top-notch conferences. Although each researcher mainly focuses on a few specific topics, a community has begun to take shape.

1.2 Research Questions

In this dissertation, we have focused on some specific aspects of knowledge-aware Recommender Systems. We have explored if the exploitation of semantics in these aspects could be beneficial. In detail, we have focused on two different kinds of side information: structured and semi-structured information.

Regarding the exploitation of structured information, we have explored the possibility of using this information to deal with sparse matrices. Our idea has been to exploit the graph-based nature of Linked Data to compensate for the lack of connections in the original users-items matrix. These matrices frequently have a density of less than 0.1%. If the representation of user preferences moved from items to semantic features, we could exploit the information delivered by a knowledge graph. Another common problem generated by the lack of data is the Cold-Start problem. Again, we have found interesting exploring if structured and semi-structured meta-data can help to produce meaningful recommendations even in extreme scenarios like Cross-Domain recommendations. Another common problem we have wanted to face is the interpretability of the recommendations. The most effective models exploit latent factors that can not be associated with any explicit features. For this reason, interpret a particular recommendation list is not possible. Eventually, this leads to a gradual decrease in trust. We have tried to adapt Latent Factors models to work with Linked Data knowledge. We have substituted latent factors with explicit knowledge to make these models interpretable. This has been a particularly exciting research line. The most challenging aspect has been measuring if the system was preserving the original semantics.

Based on these considerations, we have formulated the following research questions:

- RQ1 – Can the injection of Linked Data alleviate the sparsity problem in Collaborative-Filtering Recommender systems?
- RQ2 – Can we exploit Linked Data to face the Cold-Start Problem?
- RQ3 – Is semantic knowledge a way to make Recommender Systems interpretable?
- RQ4 – Can we define a recommendation approach based on pair-wise preferences elicitation?
- RQ5 – Can an Information Retrieval weighting scheme also consider collaborative information?

Regarding the semi-structured information, we have focused on different aspects. The first aspect is related to the exploitation of the temporal dimension. The platform transactions usually contain information about the time items are enjoyed by users. Unfortunately, Recommender Systems often ignore that information. It happens because there are multiple ways to consider Time, and the choice is highly dependant on the specific scenario. On the other hand, Knowledge Representation research has widely explored the Time dimension. We have taken advantage of this privileged point of view to enclose Time in Recommender Systems.

In particular, we have been interested in the diversification process and the notion of Popularity. In literature, several works proposed methods to increase the diversification of the recommendation lists. However, even in this context, the temporal dimension is usually ignored. We have found interesting considering the temporal aspect to enhance the diversification. Since a possible representation of data over time is considering temporal snapshots, we have considered the idea of re-defining the notion of Popularity over time. Moreover, if Popularity can change over time, maybe it could also change from one person to another.

From these considerations, we have formulated a slightly more complex notion of Popularity, proposing a time-aware personalized variant. To exploit this new notion, we have taken advantage of neighborhood models, with a specialized similarity. However, similarities are only a drop in the ocean of relatedness and vectors' closeness. We have found that the current definitions of Similarity were quite limiting. Keeping in mind the semantic relatedness measures, we have proposed to consider also a dissimilarity signal in the classic formulations. Moreover, when we deal with relatedness between resources in a knowledge graph, it is straightforward that relatedness may be different in the two ways. For this reason, we have generalized the formulation of asymmetric similarity. The experiments have shown an interesting misalignment regarding the behavior of the metrics during evaluation. Hence, we have decided to explore more the importance of the choice of the metric to tune and evaluate our recommender systems. Our idea has been measuring the degree of significance and robustness of the metrics during the tuning of hyperparameters. We have found some interesting findings that can drive and improve the

tuning process. Finally, the formulation of the metrics has piqued our interest. Inspired by the exploitation of metadata in Diversity evaluation, we have proposed a generalized formulation to measure the fairness of the proposed recommendation lists.

Based on these observations, we defined the following research questions:

- RQ6 – Can we formulate a diversification scheme that takes into account temporal aspects?
- RQ7 – Is that possible to formulate a personalized notion of the Popularity of items that also considers Time?
- RQ8 – What happens in Recommender Systems’ performance if we exploit a more complex notion of similarity that also considers dissimilarity and asymmetric similarity?
- RQ9 – Are the evaluation metrics equally valid for tuning Recommender Systems hyperparameters? Dually, do the different parameters have the same degree of importance?
- RQ10 – Is that possible measuring a generalized notion of fairness exploiting semi-structured information?

Almost all the ideas proposed in this research work have passed a peer-review, and they are published in Specialistic Journals or International Conferences Proceedings. To provide a more holistic overview, we have reported the contributions connecting them to the related research questions:

- RQ1 – Can the injection of Linked Data alleviate the sparsity problem in Collaborative-Filtering Recommender systems?
 - **Feature Factorization for Top-N Recommendation: From Item Rating to Features Relevance**, [27] RecSysKTL @RecSys 2017
 - **Feature Augmentation in Top-n Recommendation Scenarios via Linked Data**, TKDE 2020 (submitted)

- RQ2 – Can we exploit Linked Data to face the Cold-Start Problem?
 - **Addressing the user cold start with cross-domain collaborative filtering: exploiting item metadata in matrix factorization**, [144] User Modeling and User-Adapted Interaction Journal 2019
- RQ3 – Is semantic knowledge a way to make Recommender Systems interpretable?
 - **How to make latent factors interpretable by feeding Factorization machines with knowledge graphs**, [30] ISWC 2019. **This work has been awarded as Best Student Research Paper.**
 - **Semantic Interpretation of Top-N Recommendations**, TKDE 2019 (Major Revision)
- RQ4 – Can we define a recommendation approach based on pair-wise preferences elicitation?
 - **Combining RDF and SPARQL with CP-theories to reason about preferences**, [386] Semantic Web Journal 2019
- RQ5 – Can an Information Retrieval weighting scheme also consider collaborative information?
 - **A Principled Approach to Hybrid Relevance in Top-N Recommendation**. SIGIR 2020 (submitted)
- RQ6 – Can we formulate a diversification scheme that takes into account temporal aspects?
 - **An Analysis on Time- and Session-aware Diversification in Recommender Systems**. [23] UMAP 2017
- RQ7 – Is that possible to formulate a personalized notion of the Popularity of items that also considers Time?

- **Time-aware Personalized Popularity in top-N Recommendation**, [29] ComplexRec @RecSys 2018
- **Local Popularity and Time in top-N Recommendation**, 41st European Conference on Information Retrieval, [32] ECIR 2019
- RQ8 – What happens in Recommender Systems performance if we exploit a more complex notion of similarity that also considers dissimilarity and asymmetric similarity?
 - **The Importance of being Dissimilar in Recommendation**, [31] 34rd ACM/SIGAPP Symposium On Applied Computing, SAC 2019
- RQ9 – Are the evaluation metrics equally valid for tuning Recommender Systems hyperparameters? Dually, do the different parameters have the same degree of importance?
 - **On the discriminative power of Hyperparameters in Cross-Validation and how to choose them**, [28] RecSys 2019
- RQ10 – Is that possible measuring a generalized notion of fairness exploiting semi-structured information?
 - **A Flexible Framework for Evaluating User and Item Fairness in Recommender Systems**, UMUAI 2019 (Major Revision)

As a careful reader can imagine, this research work also covers some other topics like Semantic Knowledge extraction and Semantic Data management. For the sake of space, and to avoid dispersive storytelling, we have excluded them from this research summary. To provide a reader with a more comprehensive overview, we report below the reference to these works:

- Semantic Knowledge extraction
 - **LOSM: a SPARQL endpoint to query Open Street Map**, [275] 14th International Semantic Web Conference, ISWC 2015

- **Exposing Open Street Map in the Linked Data cloud**, [24] IEA/AIE 2016: The 29th International Conference on Industrial, Engineering & Other Applications of Applied Intelligent Systems, 2016
- **Querying deep web data sources as linked data**. [22] WIMS 2017
- Semantic knowledge management
 - **Etytree: A Graphical and Interactive Etymology Dictionary based on Wiktionary**. [294] WWW (Companion Volume) 2017

Additionally, we have organized two workshops in top-notch conferences, to connect researchers with the same research interests:

- Knowledge-aware and Conversational Recommender Systems Workshop, [21, 26] RecSys 2018
- Second Knowledge-aware and Conversational Recommender Systems Workshop, CIKM 2019 [25]

1.3 Structure of the work

The remainder of this dissertation is structured as follows. Part I is devoted to providing a general overview of the necessary technologies useful to understand the more advanced approaches presented in the following parts. We have first given an idea of what Semantic Web is. We start from a historical viewpoint to motivate the necessary advent of Semantic Web. Then, we briefly describe the technologies we widely use in this work. After these sections, we provide a brief introduction of preference representation.

The following chapter focuses on Recommender Systems. An introduction to the research field opens the chapter. Next, we describe some of the most known approaches for the recommendation. Sections devoted to optimization techniques and evaluation metrics close the chapter.

The following two parts focus on the main research lines of this work. Part II focuses on the exploitation of explicitly formalized knowledge. This part covers

different research lines from the leveraging of semantic features in matrix factorization to model interpretability, to the generation of recommendations based on pair-wise preferences. Part III focuses on different kinds of knowledge representation, in which data semantics is not fully structured. We focus on the notions of Time, Popularity, and Similarity, proposing approaches in line with our background. Part III continues with an analysis of evaluation metrics to assess their discriminative power when tuning a Recommender System. A proposal of a new framework to measure fairness in recommendation lists closes the part.

Part IV is devoted to drawing some conclusions about this research work.

Part I

Background

Chapter 2

Knowledge Representation

This section is devoted to a brief overview of the basic concepts regarding knowledge representation. In detail, the chapter describes the birth and development of the Semantic Web and the Linked Open Data initiative. In particular, the aim is to provide the crucial ideas and technologies that stand behind the broad term of the Semantic Web.

We introduce to the staple language, Resource Description Framework (RDF). Then we briefly describe the query language SPARQL. The following chapter focuses on the representation of Preferences. We motivate its importance in this work, and we provide a brief overview of the "ceteris paribus" theories.

This section covers only the most basic concepts about them. A more exhaustive literature review about advanced theories and applications is realized on a per research line basis in the following sections. This choice was made necessary because the corresponding literature reviews could seem a bit confusing if the reader is not already involved in the specific research field. We hope this choice lets the reader appreciate these fascinating concepts properly.

2.1 Semantic Web

The extensive and pervasive growth of the Web has produced the most voluminous information archive ever. The amount of available information is so large that there is no chance that a single person would be able to analyze and index it. The reason for it is also due to the nature itself of the Web, which is document-based and fragmented. Most of the available information is unstructured. It misses the optional additional data necessary to analyze it before accessing the document itself. Moreover, entire sections of the Web are, de facto, inaccessible due to linguistic barriers. These are only some examples of the necessity of a total rebuilding (and rethinking) of the Web. It was needed to move the focus of the Web from documents to the data. The single facts, the assertions, the data should be the heart element of the Web instead of documents. It is straightforward that these data should have a unique interpretation. These data should not be ambiguous as the natural language usually is. Even though we would have built this so-called Web of Data, a single person still would not be able to analyze the whole Web. This consideration has two consequences. The first one is the awareness that the Web can be analyzed only through an automatic process that excludes the necessity of human interpretation of data. Hence, data should be machine-interpretable. The second consequence is that the automatic interpretation is only possible if data comes with an explicit semantics attached to it. This new Web, composed of interpretable data, is called Semantic Web. The birth of the Semantic Web is commonly associated with the publication of an article by Tim Berners-Lee on Scientific American in 2001 [60]. The article is focused on what the Web could become, and in particular, a sentence could serve as an "Ante-literam" definition: "The Semantic Web is not a separate Web but an extension of the current one, in which information is given well-defined meaning, better enabling computers and people to work in cooperation". Hence, the Semantic Web is an information network that lets machines understand the contained information. Human users can thus prepare complex queries that machines can understand and process. The underlying needed comprehension tacitly imposes to exploit a well-formalized knowledge representation. The knowledge representation would have been just the first step in the process of development of the Se-

semantic Web. Although, since the organization of knowledge and its interpretability is crucial, it has quickly become the keystone of the development of the Semantic Web. This branch of the Semantic Web, specifically devoted to the representation and organization of knowledge archives (datasets), has obtained a specific name, Linked Data. In 2006, Berners-Lee published a new article [338] partially rethinking the idea of the Semantic Web under a Linked Data perspective. The community has spent a big effort to define and refine the languages and technologies involved. Nowadays, these technologies are exploited to expose data on the Web and connect data. This has resulted in a brief but summarizing assertion: "Linked Data is the Semantic Web did well".

2.1.1 Resource Description Framework

Resource Description Framework (RDF) is a framework designed to represent knowledge. It is one of the pillars of Semantic Web, proposed by the World Wide Web Consortium (W3C) in 1998. The knowledge encoded in RDF can be easily be represented in the form of a graph, queried, and processed by automatic reasoning techniques. Along with the framework, different serialization formats have been proposed, that avoid redundancy and allow wide interoperability. RDF is a general model useful to describe the resources on the Web. The skeleton of the RDF syntax, its Data model, is a set of triples. Each of them is composed of a subject, a predicate, and an object. An RDF triple could express an assertion, a logic expression, or a generic statement. An example of an RDF triple is:

```
<http://example.org/#spiderman>  
<http://www.perceive.net/schemas/relationship/enemyOf>  
<http://example.org/#green-goblin>.
```

A set of RDF triples is named RDF graph. An RDF graph is a directed graph, in which each triple is a connection Vertex-Arc-Vertex. Formally, an RDF graph is the conjunctive set of its triples. A triple is composed of three distinct elements: the subject expressed as Internationalized Resource Identifier (IRI) or as a Blank Node; the predicate expressed as IRI; the object expressed as IRI, as a Literal, or

as a Blank Node. Different triples can share the same IRI. IRIs, Literals, and Blank Nodes since they compose RDF triples are called RDF terms.

IRIs [133] are a generalization of Uniform Resource Identifiers (URIs) [59] that enable the usage of a more wide range of Unicode characters. This means that all URIs can be considered IRIs, but the vice-versa is not necessarily true. The definition and the meaning of IRIs are apart from RDF. On the one hand, this implies absolute freedom in the IRIs creation process. On the other hand, the same IRIs can be used by third-party organizations in their knowledge bases.

Literals are used to express values like strings, numbers, and dates. Moreover, in the case of strings of type

`http://www.w3.org/1999/02/22-rdf-syntax-ns#langString`, it is allowed to add a language tag. Literals are composed of three optional elements: value, type, and language. For this reason, two Literals are considered equals only if the three elements are equals. Datatypes are used with Literals to represent the different kinds of data. The adopted abstraction in RDF is XML Schema compatible. Thus, it is possible to introduce new custom data types and use them in RDF. A data type is composed of: lexical space, a value space, and a lexical-to-value mapping. The lexical space is a set of unicode strings. The lexical-to-value mapping represents a function which has as domain the lexical space and as range the value space.

Blank Nodes are local identifiers that denote a generic RDF term that is present within the considered collection, with a local scope. They usually used if we are not interested in the values of the variable associated with the Blank Node.

Resources, or entities, can be defined using IRIs or strings. Resources can indicate everything: physical objects, documents, concepts, numbers, strings. An example of the description of the "Pulp Fiction" RDF resource is:

```
<http://dbpedia.org/resource/Pulp_Fiction>  
<http://www.w3.org/2000/01/rdf-schema#label>  
"Pulp Fiction"@en  
<http://dbpedia.org/resource/Pulp_Fiction>  
<http://dbpedia.org/ontology/director>
```

<http://dbpedia.org/resource/Quentin_Tarantino>

which states that the resource has an English name, "Pulp Fiction", and the resource has been directed by Quentin Tarantino.

RDF vocabulary is a collection of IRIs that can be used in an RDF graph. RDF schema, for instance, is an RDF vocabulary. Moreover, RDF vocabularies can be further used to define other vocabularies. If the initial part of the path of the IRIs is common among more IRIs, this substring is usually defined as a namespace. These namespaces can also be associated with a shorter and intuitive string named prefix. Some RDF serialization formats exploit prefixes to reduce the verbosity of IRIs.

RDF Dataset is a collection of RDF graphs. RDF Datasets have some specific features: i) one RDF graph has to be defined as default graph; ii) they contain zero or more named graph. A named graph is a pair composed of an IRI or a Blank Node, and an RDF graph. To conclude this brief overview of RDF, we mention **rdf:type** predicate

(<http://www.w3.org/1999/02/22-rdf-syntax-ns#type>). This predicate is particularly important under an ontological perspective since it states the belonging of the subject to a certain object class.

2.1.2 RDF Schema

RDF model is a flexible graph-based model that can be easily used for knowledge representation. Nevertheless, it contains only a few raw tools to define the semantics of the resources. For this reason, RDF Schema (RDFS), has been proposed. RDFS provides essential tools to define ontologies. Ontologies, or RDF vocabularies, are a useful tool to structure domain knowledge. RDFS eases the definition of classes and their properties. It provides ways to define which properties should be used with instances of specific classes. RDFS is completely compatible with RDF. Thus an RDFS vocabulary still is a valid RDF graph. In particular, RDFS is formalized as an RDF vocabulary with the namespace www.w3.org/2000/01/rdf-schema# (with the common prefix **rdfs:**). To define the kind of resource, RDFS provides the term *Class*. Classes are described by the RDFS resources **rdfs:Class** and

rdfs:Resource, **rdfs:subClassOf**, and through the well-known RDF IRI **rdf:type**. As the reader could imagine, these entities let create a hierarchical organization of the entities described in the RDF graph. A single entity can still be defined as belonging to more than one class. The predicate **rdfs:domain** and **rdfs:range** are intended to describe other properties. In detail, **rdfs:domain** denotes that the subject of the triple that makes use of the target property should belong to the class indicated in the object. **rdfs:range** denotes that the object of the triple that makes use of the target property should belong to the class indicated in the object. **rdfs:subClassOf** is a particular predicate with domain and range entities of the type `rdfs:Class`. Furthermore, the properties can be hierarchically organized, making use of **rdfs:subPropertyOf**. Finally, RDFS provides `rdfs:comment` and `rdfs:label` that lets insert comments and labels in natural language.

2.1.3 Linked Data

One of the core points of the documental Web was to connect different documents through specific links, hyperlinks. Similarly, in the Semantic Web, the connections between data are fundamental. The datasets that provide links with other datasets and the technologies involved are consequently named Linked Data. To enable modern applications to work with these knowledge bases, data should be exposed using standard formats, should provide an endpoint to query the knowledge base. When a dataset also provides connections with other datasets, we can define it as a Linked Dataset. Two common examples of Linked Datasets are DBpedia and Wiki-Data, which basically expose the content of Wikipedia pages following the Linked Data best practices. Since DBpedia has been used extensively during this thesis, I have chosen it as an example in the following paragraphs. I hope this will make the reader more familiar with this dataset, enabling her to understand the remainder of the thesis better. To better explain the importance of inter-datasets links, let us consider an example. DBpedia provides, for instance, links to the Geonames dataset. Providing this information, expressed as triples, enables the creation of new applications that can access the information about a resource, retrieve the Geonames

link, access Geonames, and retrieve additional information. It is important to underline that all the additional information that is present in Geonames is not present in DBpedia. These connections let developers build applications that can integrate different data sources easily. The Linking Data Initiative started in 2007 intending to extend the actual Web using the standards proposed in the Semantic Web. In this while, a huge number of datasets have been created and connected together, composing the so-called Linked Open Data (LOD) cloud. The LOD cloud clearly shows that this new Web is not substituting the documental Web, but it is enriching it.

2.1.4 SPARQL

SPARQL (Simple Protocol and RDF Query Language) is the de-facto standard query language for RDF data. The language was formalized by Data Access Working Group of W3C (World Wide Web Consortium) in 2008. A newer formalization, SPARQL 1.1, was released in 2013. SPARQL 1.1 shares most of the RDF specifics, but it is stricter in some cases. An example is the identification of an RDF graph that, in SPARQL 1.1, does not allow the usage of Blank Nodes. SPARQL is a key element of the Semantic Web and it lets extract information from remote RDF knowledge bases. The fundamental role of SPARQL is really simple: it processes the RDF graph finding one or more specific subgraphs that correspond to the informative need expressed by the query. The SPARQL processing is made possible by providing two graphs: the data graph (usually already loaded in a remote RDF store) and the query graph (even here described through triples). SPARQL adopts the Turtle syntax, an intuitive and non-verbose RDF serialization format.

Assertions are expressed employing subject-predicate-object sequences, and they are terminated by a full stop. PREFIX introduces prefixes and namespaces, whereas URIs are enclosed within angle brackets and Literals within quotation marks. With only these elements we can already provide an example:

```
@prefix cd: <http://example.org/cd/>
@prefix: <http://example.org/esempio/>
```

```
:Permutation cd:autore "Amon Tobin".  
:Bricolage cd:autore "Amon Tobin".  
:Amber cd:autore "Autechre".  
:Amber cd:anno 1994.
```

The results of the SPARQL query can be returned in more or less machine-understandable formats: XML, JSON, RDF, HTML.

The query language reflects the graph-based nature of the data model. Indeed, the query-answering mechanism is done via pattern-matching over the knowledge graph (graph-matching). Usually, a SPARQL query is composed of at least one graph pattern. A graph pattern is a set of triples, also called triple-pattern. These patterns are similar to RDF, thus they consist of a subject-predicate-object sequence. The main difference is that each of these components could be a variable. The query is solved by looking for the subgraphs that correspond to the pattern. In particular, the search for the subgraphs aims to find all the possible instantiations of RDF terms with variables that correspond to existing triples in the knowledge graph. The basic syntax is inspired by SQL. For instance, SPARQL uses the SELECT clause to define the variables that will appear in the results, whereas the WHERE clause encloses the graph pattern. As an example, a graph pattern can be: `?title cd:author ?author .` where the two variables are denoted by the question mark. Similarly to SQL, we can exploit the GROUP BY clause to aggregate results or FILTER to filter results based on some conditions. Instead of SELECT, we can also use other clauses as ASK, which returns a boolean value, DESCRIBE, which returns the complete description of a resource, and CONSTRUCT, which returns graphs.

2.2 Preference representation

Dealing with user preferences is an important aspect of every application designed to provide personalized information to the end-user. The original interest in preferences can be found in decision theory, as a way to support complex, multifactorial decision processes [147], and nowadays every personalized system needs a preference model to capture what the user likes or dislikes. Once the user model has been

represented, it is then exploited to filter information coming from a data source, e.g., a database, in order to provide a ranked list of results matching the order encoded in the preferences of the user.

Query languages usually let us specify the information that we want to be returned (hard constraints), where the result set contains elements with no specific order with reference to user preferences. It contains all those resources that exactly match the constraints represented by the query. As a matter of fact, if just one of the requirements representing the query is not fulfilled, the result set can be empty. At the same time, returning huge and unordered sets of answers could be useless and even counter-productive. A possible way to bypass these issues is to allow the language to represent both hard constraints—used to return only relevant results—and soft ones, i.e., preferences—to rank the results by fulfilling a user’s tastes. Approaches to preference representation can be either quantitative or qualitative [123]. The formers are based on a total ordering of the outcomes, given by a scoring function, while the latters enable the representation of partial orders since preferences are treated as independent dimensions. From a user perspective, a qualitative approach is more natural than the quantitative one [284]. Indeed, in the first case, the user has just to provide pairwise qualitative comparisons, while in the second case, she has to assign a value to many alternatives, which very often are represented in a multi-attribute setting. Regarding the Linked (Open) Data world, the notion of qualitative preferences in SPARQL queries was introduced in [347], by Siberski et al., whose *preference-based querying* language extends SPARQL through the introduction of solution modifiers (the PREFERRING clause). Their query formulation retrieves only items that are the most preferred ones, or equivalently *undominated*. The work [162] builds on the earlier approach of [347], but adds preferences at the level of filters, rather than as a solution modifier. The *PrefSPARQL* syntax of [162] needs no additional solution modifier to express qualitative preferences, as it leverages the expressive power of SPARQL 1.1. However, the approaches proposed in [347] and [162] both have an important limitation: they are not able to provide an order of all the available outcomes that reflects user preferences. That is, they return only the undominated (a.k.a. Pareto-optimal) results, i.e., those outcomes

best satisfying user preferences. Unfortunately, the size of the resulting answer set could be too small to be of practical use. This fosters moving beyond the Pareto-optimal set identification to a top- k scenario [284], where firstly available outcomes are ordered, even if with the ties implied by a qualitative approach, from best (most preferred) to worst (less preferred) according to a given user's preferences, and then the first k results are returned.

2.2.1 CP-theories

Utility functions can be considered as the ideal tool for representing and reasoning with preferences, but the total order that they allow representing does not always reflect the actual user model. Partial orders among preferences are a more natural way to represent a user's tastes. Qualitative statements, e.g., "*given u , I prefer x_i over x'_i* ", permit a system to encode partial orders among user preferences, thus granting the representation of a more realistic user model. Let us consider the following example. The previous example is a typical conditional statement where the core notion of a CP-statement is explicitly represented. We have that when a particular condition u is true, a user prefers to enjoy items where also x_i is true, rather than items where x'_i is true, given that x_i and x'_i cannot be true at the same time.

Example 1 (Books)

"Giorgio has just finished his exams and wants to relax with a book. Giorgio can read both English and French, but he would like to improve his French to enrich his curriculum and so he prefers to read French books. Giorgio prefers reading crime books over autobiographical ones for French books since he believes that crime plots are more captivating and, therefore, more useful while learning a foreign language. The reverse order holds for English books. Giorgio is a good reader and, therefore, given an English book, he prefers those being part of a saga. The literary genre and the presence of a subsequent work not have the same importance to Giorgio: in case of English books, he considers the choice on the genre more important than the one dependant on sequels, while the opposite happens for French books. Finally, for books characterized by a sequel, Giorgio regards the presence of

a cinematographic version positively”.

■

By looking at Example 1, we find it quite hard (even impossible) to directly represent Giorgio’s preferences by means of a score assigned to his preferential statements. In fact, Giorgio’s preferences can be better expressed as qualitative (pairwise) comparisons.

Relevant frameworks to represent and reason with qualitative preferences are built according to the *ceteris paribus* semantics [166] and specifically consist of conditional preference networks (or CP-nets) [71] and a formalism along similar lines to CP-nets, but with a richer language of preference statements, namely, conditional preference theories (or CP-theories) [395].

Syntax. Formally, given a set of variables V , a CP-theory Γ is a set of preference statements φ of the general form:

$$u_\varphi : x_\varphi > x'_\varphi [W_\varphi],$$

where u_φ is an assignment to a set of variables $U_\varphi \subset V$, x_φ and x'_φ are different assignments to some variable $X_\varphi \notin U_\varphi$, and W_φ is some subset of $V - U_\varphi - \{X_\varphi\}$.

Semantics. The interpretation of φ is that, given u_φ , x_φ is strictly preferred to x'_φ , all else being equal, but irrespective of the values of variables in W_φ .

That is, φ compactly states that for all assignments w, w' to W_φ and assignments t to $T_\varphi = V - U_\varphi - \{X_\varphi\} - W_\varphi$, $tu_\varphi x_\varphi w$ is preferred to $tu_\varphi x'_\varphi w'$.

In what follows, we will use the word *outcome* to indicate a complete assignment to all the variables in V and denote the set of all outcomes as \mathcal{O} . For the statement φ , we denote by φ^* the set of pairs of outcomes $(tu_\varphi x_\varphi w, tu_\varphi x'_\varphi w')$, where t is an assignment to T_φ , and w, w' are assignments to W_φ . Further defining $\Gamma^* = \cup_{\varphi \in \Gamma} \varphi^*$, it is then natural to define for the CP-theory Γ a strict partial order $>_\Gamma$, induced by Γ on the set of outcomes \mathcal{O} , as the transitive closure of Γ^* . The CP-theory formalism allows to express the usual CP-net *ceteris paribus* statements by simply considering $W_\varphi = \emptyset$ and identifying U_φ with $Pa(X)$, the parents of

| | |
|-------------|--|
| φ_1 | $\top : C_F > C_{UK}[\emptyset]$ |
| φ_2 | $C_{UK} : LG_A > LG_C[SW]$ |
| φ_3 | $C_F : SW_{No} > SW_{Yes}[LG]$ |
| φ_4 | $C_{UK} : SW_{Yes} > SW_{No}[\emptyset]$ |
| φ_5 | $C_F : LG_C > LG_A[\emptyset]$ |
| φ_6 | $SW_{Yes} : F_{Yes} > F_{No}[\emptyset]$ |
| φ_7 | $SW_{No} : F_{No} > F_{Yes}[\emptyset]$ |

Table 2.1: The CP-theory $\Gamma_{C-LG-SW-F}$.

a variable X , that is, variables which the preferences on X depend on. However, as anticipated in Section 8.1, under the stricter CP-net formalism, for each variable X , a parent set $Pa(X)$ must be defined and instantiated when the preference order over values of X is established. This is not required in CP-theories, where you can find more statements related to the same variable X , but with different sets U . In addition, CP-theories allow stronger conditional preference statements than CP-nets, which are natural for users to express. For example, they represent a formalism even more general than TCP-nets [72], an enhancement of CP-nets where (conditional) relative importance between variables can be expressed. TCP-nets can be represented through statements with W containing at most one variable. Moreover, there are statements, such as *I prefer x_i over x'_i irrespective of the values of all other variables*, that cannot be expressed in CP-nets or TCP-nets, but correspond in the new formalism of CP-theories to $\top : x > x' [V - \{X\}]$, where \top is the assignment to an empty set of variables.

Example 2 (Books cont'd)

The overall profile of Giorgio may be modelled by means of the CP-theory $\Gamma_{C-LG-SW-F}$, that is, the set of statements given in Table 2.1. There, a set of binary variables $V = \{Country, LiteraryGenre, SubsequentWork, FilmVersion\}$ (abbreviated as C, LG, SW, and F, respectively) is considered. Their domains are given by:

- $dom(\text{Country}) = \{C_F, C_{UK}\}$ (for *France* and *United Kingdom*),
- $dom(\text{LiteraryGenre}) = \{LG_C, LG_A\}$ (for *Cri-me_fiction* and *Autobiographical_novel*),
- $dom(\text{SubsequentWork}) = \{SW_{Yes}, SW_{No}\}$ (indicating if a book has a sequel or not),
- $dom(\text{FilmVersion}) = \{F_{Yes}, F_{No}\}$ (indicating whether there is a cinematographic version of the book or not). ■

For a CP-theory Γ , the preference ranking over outcomes $>_\Gamma$ introduced above, can be equivalently induced under the *worsening swap* semantics. Hereafter, we use the notation $o(X_i) = x_i$ to indicate that the variable X_i is assigned the value x_i in o , and analogously $o(\{X_j, \dots, X_{j+k}\}) = \{x_j, \dots, x_{j+k}\}$ to state that $X_j = x_j, \dots, X_{j+k} = x_{j+k}$ in o .

Given two outcomes o and o' of \mathcal{O} , there is a worsening swap from o to o' , if there exist a variable

$$X_i \in V - \{X_j, \dots, X_{j+k}\} - W, x_i, x'_i \in dom(X_i)$$

and an assignment x_j, \dots, x_{j+k} to the variables set $\{X_j, \dots, X_{j+k}\}$ such that:

- (i) $o(X_i) = x_i$ and $o'(X_i) = x'_i$,
- (ii) $o(\{X_j, \dots, X_{j+k}\}) = o'(\{X_j, \dots, X_{j+k}\}) = \{x_j, \dots, x_{j+k}\}$,
- (iii) $o(V - \{X_i\} - \{X_j, \dots, X_{j+k}\} - W) = o'(V - \{X_i\} - \{X_j, \dots, X_{j+k}\} - W)$, and
- (iv) $x_j \dots x_{j+k} : x_i > x'_i[W] \in \Gamma$.

The preference relation $>_\Gamma$ over \mathcal{O} is therefore the transitive closure of worsening swaps.

A CP-theory Γ is *consistent*, if it has a model, i.e., if there exists a strict total order $>$ that *satisfies* Γ , which is equivalent to $> \supseteq \Gamma^*$, that is, $>$ extends $>_\Gamma$. In [394], it is proved that the irreflexivity of $>_\Gamma$ (or equivalently the acyclicity of Γ^*) is a necessary and sufficient condition for consistency.

The consistency of a CP-theory implies that there are no cycles generated by $>_{\Gamma}$. Avoiding cycles is of paramount importance, as they introduce conflicting information while ordering the outcomes in \mathcal{O} . Let us consider the ordering represented in Figure 2.1, where an edge from o_i to o_j represents $o_i >_{\Gamma} o_j$. As we cannot establish what is the correct ordering of the outcomes due to the cycle created by o_2 and o_3 , we could even have situations where we cannot compute the most preferred (undominated) outcome. For example, suppose in Figure 2.1, we do not have o_1 . What would then be the best solution for the user in this case?

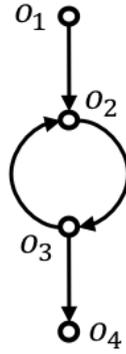


Figure 2.1: An example of cycle among outcomes.

A necessary condition for consistency is local consistency. Consider a CP-theory Γ , a variable $X \in V$, and an assignment a to a set of variables $A \subseteq V$. An ordered pair (x, x') of X values is *validated* by a , if there exists a statement φ of the form $u : x > x' [W]$ in Γ , such that a extends u , that is, a projected to U_{φ} gives u .

The *local ordering* $\succ_a^X(\Gamma)$ (abbreviated as \succ_a^X) on X values is defined as the transitive closure of the set of pairs (x, x') validated by a . Γ is *locally consistent*, if \succ_a^X is irreflexive for all variables X and outcomes α . Local consistency is a necessary condition for consistency, since if Γ is not locally consistent, then there exist an outcome α , a variable X , and a sequence x_1, \dots, x_k of values of X with associated statements $u_i : x_i > x_{i+1} [W_i] \in \Gamma$ such that α extends u_i and $\alpha(X) = x_1 = x_k$.

This gives a worsening swapping sequence from α to α (only involving changing variable X), thus implying that $>_{\Gamma}$ is not irreflexive, or equivalently that Γ is

not consistent. In general, deciding whether a CP-theory is locally consistent is coNP-complete, but it can be shown that if the size of the parent sets and the size of the domain sets are bounded by a constant, then deciding local consistency is polynomial [395]. Moreover, for CP-nets and TCP-nets, local consistency is always guaranteed [395].

Given a CP-theory Γ , there are several kinds of directed graphs that can be defined on the set of variables V . For $S, T \subset V$, we indicate with $S \rightarrow T$ the set of edges $\{(X, Y) : X \in S, Y \in T\}$, omitting the set brackets, if S or T is a singleton set, e.g., abbreviating $S \rightarrow \{Y\}$ with $S \rightarrow Y$.

The *dependency graph* $H(\Gamma)$ consists of edges $U_\varphi \rightarrow X_\varphi$ for all φ in Γ , that is, all the pairs of the form (Y, X_φ) , with $\varphi \in \Gamma$ and $Y \in U_\varphi$. That is, the edge (Y, X) belongs to $H(\Gamma)$ iff there is some conditional preference statement $\varphi \in \Gamma$ that makes the preferences for X conditional on Y . Relative importance is not encoded in $H(\Gamma)$.

On the other side, we define $G(\Gamma)$ to contain $U_\varphi \rightarrow X_\varphi$ and $X_\varphi \rightarrow W_\varphi$ for all φ in Γ , i.e., $G(\Gamma) = H(\Gamma) \cup \{X_\varphi \rightarrow W_\varphi \mid \varphi \in \Gamma\}$. $G(\Gamma)$ contains both dependency and relative importance information: it is $H(\Gamma)$ with the addition of edges (X, Z) , if there is any statements φ representing a preference on values of X irrespective of the value of Z (then, X is *more important* than Z , with importance meant as in the TCP-net formalism [72]).

A CP-theory Γ is *fully acyclic*, if $G(\Gamma)$ is acyclic. For fully acyclic CP-theories, consistency and local consistency are equivalent [395].

For a CP-theory Γ and assignment a to a set of variables $A \subseteq V$, we can define another directed graph $J_a(\Gamma)$ on V made of the set of edges $U_\varphi \rightarrow \{X_\varphi\} \cup W_\varphi$ for all $\varphi \in \Gamma$ and also the set $\{X_\varphi\} \rightarrow W_\varphi$ for all $\varphi \in \Gamma$ such that $U_\varphi \subset A$ and a extends u_φ .

A CP-theory Γ is *context-uniformly conditionally acyclic* (or *cuc-acyclic*), if it is locally consistent, and for each outcome $o \in \mathcal{O}$, $J_o(\Gamma)$ is acyclic. It can be proved that a cuc-acyclic CP-theory is always consistent [394]. The condition of cuc-acyclicity is weaker than the full acyclic one, and it requires only $H(\Gamma)$ to be acyclic [395] instead of $G(\Gamma)$.

Cuc-acyclicity implies a less expressive power in terms of preferences that the user may express, but, as we will see in the following, it grants a nicer computa-

tional complexity when ranking a set of outcomes (along with guaranteeing always consistency, while deciding the consistency of general CP-nets and CP-theories is PSPACE-complete [158]).

Example 3

To better clarify the expressive limits of *cuc-acyclic* CP-theories, we consider the following CP-theory $\hat{\Gamma}$ whose $H(\hat{\Gamma})$ is shown in Figure 2.2.

- (1) *Crime Fiction* : *Agatha Christie* > *Andrea Camilleri* [\emptyset]
- (2) *Isaac Asimov* : *Science Fiction* > *Crime Fiction* [\emptyset]

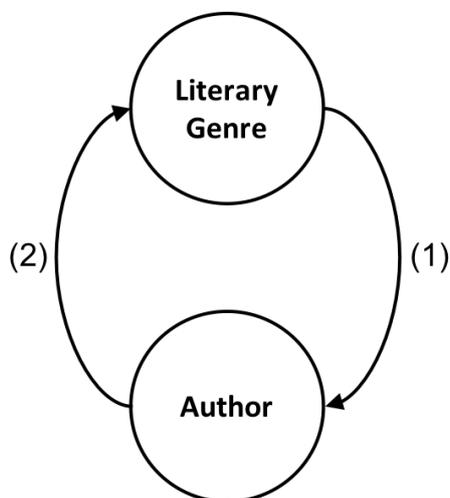


Figure 2.2: The graph $H(\hat{\Gamma})$ representing preferences in Example 3.

Due to the interrelation between the two preferences, $H(\hat{\Gamma})$ is cyclic, and so $\hat{\Gamma}$ cannot be *cuc-acyclic*.

Although there might be cases where *cuc-acyclicity* is a strong limitation in the representation of user preferences, it represents a limited subset of all possible CP-theories that can be used to model a user profile. ■

In what follows, we show two results, both proved in [395], which determine a strict partial order extending $>_{\Gamma}$. The approaches proposed to compare outcomes are strongly related to the ordering queries defined in [71] and already discussed in [320] for CP-nets, and can be seen as a generalization of Corollary 4 and Theorem

5 of [71]. The first result is applicable for a fully acyclic CP-theory Γ . It proposes to compare two outcomes by looking (using the appropriate local ordering) at their value on each of the most important variables on which they differ, where importance is defined according to the graph $G(\Gamma)$.

In Theorem 1, we denote with $\Delta(\alpha, \beta)$ the set of variables of V on which outcomes α and β differ, i.e., $\Delta(\alpha, \beta) = \{Y \in V \mid \alpha(Y) \neq \beta(Y)\}$. If $\alpha \neq \beta$, we build $\Theta(\alpha, \beta)$ as the set of G' -maximal elements of $\Delta(\alpha, \beta)$, being G' the transitive closure of $G(\Gamma)$. $\Theta(\alpha, \beta)$ is therefore the set of variables $Y \in \Delta(\alpha, \beta)$ such that there exists no $Z \in \Delta(\alpha, \beta)$ with $(Z, Y) \in G'$.

Theorem 1

Let Γ be a locally consistent and fully acyclic CP-theory, and let the binary relation $\succ_{p(\Gamma)}$ on \mathcal{O} be defined as follows, for any pair of outcomes α and β : $\alpha \succ_{p(\Gamma)} \beta$ iff $\alpha \neq \beta$ and $\alpha(Y) \succ_{\alpha}^Y \beta(Y)$ for all $Y \in \Theta(\alpha, \beta)$. Then, $\succ_{p(\Gamma)}$ is a strict partial order extending $>_{\Gamma}$, and the comparison between any pair of outcomes requires polynomial time.

The second result deals with cuc-acyclic CP-theories [395]. In Theorem 2, $\Delta(\alpha, \beta)$ is still used to denote the set of different variables in outcomes α and β . If $\alpha \neq \beta$, $\Theta'(\alpha, \beta)$ is defined as the set of \triangleright_{α} -undominated elements of $\Delta(\alpha, \beta)$, being \triangleright_{α} the transitive closure of $J_{\alpha}(\Gamma)$: $Y \in \Theta'(\alpha, \beta)$ iff there exists no Z in $\Theta'(\alpha, \beta)$ with $(Z, Y) \in \triangleright_{\alpha}$.

Theorem 2

Let Γ be a cuc-acyclic CP-theory, and let the binary relation \gg_{Γ} on \mathcal{O} be defined as follows, for any pair of outcomes α and β : $\alpha \gg_{\Gamma} \beta$ iff $\alpha \neq \beta$ and $\alpha(Y) \succ_{\alpha}^Y \beta(Y)$ for all $Y \in \Theta'(\alpha, \beta)$. Then, \gg_{Γ} is a strict partial order extending $>_{\Gamma}$, and the comparison between any pair of outcomes requires polynomial time.

Theorem 2 proposes a more general approach to generate a strict partial order on \mathcal{O} , which although requiring a pretty strong condition on the CP-theory, i.e., cuc-acyclicity, does not need full acyclicity. In particular, if Γ is fully acyclic, and α and β are two outcomes to compare, then $J_{\alpha}(\Gamma) \subseteq G(\Gamma)$ and so $\triangleright_{\alpha} \subseteq G'$. Therefore,

$\Theta'(\alpha, \beta) \supseteq \Theta(\alpha, \beta)$. This implies that if $\alpha \gg_{\Gamma} \beta$ then $\alpha \succ_{p(\Gamma)} \beta$, so that \gg_{Γ} is a closer approximation of \succ_{Γ} than $\succ_{p(\Gamma)}$.

Section 8.4.2 will ground on this more general Theorem 2 to formulate a SPARQL query able to rank outcomes according to user preferences encoded in a CP-theory model.

Example 4 (Books cont'd)

Relative to the CP-theory $\Gamma_{C-LG-SW-F}$ with Giorgio's preferences (see Table 2.1), the use of Theorem 2 produces the following sound ranking solution:

$$\begin{aligned} & \langle C_F L G_C S W_{No} F_{No}, \\ & (C_F L G_C S W_{No} F_{Yes}, C_F L G_A S W_{No} F_{No}), \\ & C_F L G_A S W_{No} F_{Yes}, C_F L G_C S W_{Yes} F_{Yes}, \\ & (C_F L G_C S W_{Yes} F_{No}, C_F L G_A S W_{Yes} F_{Yes}), \\ & C_F L G_A S W_{Yes} F_{No}, C_{UK} L G_A S W_{Yes} F_{Yes}, \\ & C_{UK} L G_A S W_{Yes} F_{No}, C_{UK} L G_A S W_{No} F_{No}, \\ & C_{UK} L G_A S W_{No} F_{Yes}, C_{UK} L G_C S W_{Yes} F_{Yes}, \\ & C_{UK} L G_C S W_{Yes} F_{No}, C_{UK} L G_C S W_{No} F_{No}, \\ & C_{UK} L G_C S W_{No} F_{Yes} \rangle, \end{aligned}$$

where outcomes within round parentheses are not comparable.

As an example, we may consider the outcome $C_F L G_C S W_{No} F_{No}$ which is favoured in the comparisons made according to the order \gg_{Γ} over every other outcome. In fact, the set

$$\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{No} F_{Yes})$$

coincides with the set of distinct variables in the compared outcomes, namely,

$$\Delta(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{No} F_{Yes}),$$

being both equal to $\{F\}$, and the preference φ_7 can be used to locally order the first outcome over the second one. Analogously,

$$\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{No} F_{No})$$

is equal to

$$\begin{aligned} & \Delta(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{No} F_{No}) \\ &= \{LG\}, \end{aligned}$$

and the preference φ_5 can be exploited. The set

$$\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{No} F_{Yes})$$

coincides with

$$\begin{aligned} & \Delta(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{No} F_{Yes}) \\ &= \{LG, F\}, \end{aligned}$$

and the preferences φ_5 and φ_7 can be exploited for the variables LG and F , respectively. While

$$\begin{aligned} & \Delta(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{Yes} F_{Yes}) \\ &= \{SW, F\}, \end{aligned}$$

it holds that

$$\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{Yes} F_{Yes})$$

is composed of the only variable SW , and the preference φ_3 can be used for this comparison. The preference φ_3 can also be used when dealing with

$$\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{Yes} F_{No}),$$

which coincides with the set

$$\begin{aligned} & \Delta(C_F L G_C S W_{No} F_{No}, C_F L G_C S W_{Yes} F_{No}) \\ &= \{SW\}. \end{aligned}$$

As for Giorgio, *SubsequentWork* takes priority over *LiteraryGenre* for French books, according to φ_3 , and takes priority over *FilmVersion*, because of the preference φ_6 (or φ_7), it follows that the set of distinct variables

$$\Delta(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{Yes} F_{Yes})$$

has three elements $\{LG, SW, F\}$, while

$$\begin{aligned} & \Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{Yes} F_{Yes}) \\ &= \{SW\}, \end{aligned}$$

and the preference φ_3 can be used in this case. For the last comparison involving France, the preference φ_3 is still determinant, as

$$\begin{aligned} & \Delta(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{Yes} F_{No}) \\ &= \{LG, SW\}, \end{aligned}$$

but $\Theta'(C_F L G_C S W_{No} F_{No}, C_F L G_A S W_{Yes} F_{No}) = \{SW\}$. When comparing the outcome $C_F L G_C S W_{No} F_{No}$ with any outcome o' in which C_{UK} appears, Θ' contains only the undominated variable *Country*, and the preference φ_1 can be used to advantage $C_F L G_C S W_{No} F_{No}$ over o' . ■

Chapter 3

Recommender Systems

3.1 Introduction

In this section, we provide an overview of the different Recommender Systems approaches. We present the most prominent families of Recommender Systems to draw the necessary background to face the remainder of the dissertation. As a representative of Collaborative Filtering techniques, we focus first on Neighborhood-based and Matrix Factorization models. Then, Vector Space Model is detailed as a representative of Content-based models. The section is closed by an overview of the main Recommender Systems evaluation methodologies with particular emphasis on evaluation protocols and metrics.

3.1.1 The recommendation problem

Recommender Systems are software tools and techniques devoted to providing suggestions to users. The idea is to provide the user with a personalized shortlist of items or services she could appreciate. The most common recommendation strate-

gies rely on the content (Content-Based Recommender Systems) or collaborative information (Collaborative-Filtering Recommender Systems). Beyond these families, there is still a wide range of other approaches. Graph-based, context-aware, Semantic-aware, and more recently Deep Neural Network Recommender Systems are only some examples. Taxonomies of Recommender Systems are often defined considering the source data or the recommendation technique. First, we focus on Collaborative-Filtering algorithms (CFs). On the one hand, CFs have shown the highest accuracy performance. On the other hand, we do not need to introduce any other data source. Indeed, most of them are entirely based on the transactions (or ratings) matrix. The rating matrix, usually denoted by a capital R , is a matrix that contains all the transactions that happened on the recommendation platform. These transactions could be ratings, purchases, clicks, or other feedback provided explicitly or implicitly by the users. The core idea of CFs is to exploit this information to recommend items similar users like or items that usually are consumed along with the already experienced ones. Content-Based (CBF) Recommender Systems usually do not use collaborative information. CBFs take advantage of a certain representation of items to compute similarities among them or building the user profile or both. Recommender systems, in their simplest form, can, therefore, be defined as tools that provide ranked lists of items of a specific category, based on explicit or implicit information that relies on user tastes and by exploiting additional information on items, other users, context and/or the past history of the user to whom we want to provide a recommendation list.

Up to the Netflix prize [52], the research community defined the recommendation problem as a rating prediction task. In a few words, the problem was to predict a likeness value for an unexperienced item based on the user's history. However, in real recommender systems applications, only a small subset of relevant items is provided to users [178]. Indeed, several studies pointed out that rating prediction optimization was not able to produce the optimal top-N recommendation lists [256]. The recommendation problem was hence re-defined as a top-N recommendation task [108], in which the optimization goal shifted to items ranking. In general, a good formalization of the recommendation problem comes from [334]: Let

$U = \{u_1, u_2, \dots, u_m\}$ be the set of all users, and let $I = \{i_1, i_2, \dots, i_n\}$ be the set of all possible items that can be recommended. Each user u_i has a list of items I_u , which the user has expressed her opinions about (explicitly as a rating score or implicitly derived from purchase records). Hence, we define an utility function of item i for user u , as $f : U \times I \rightarrow R$, where R is a completely ordered set. The recommendation problem corresponds to finding the item $i^{max,u} \in I$ that maximizes the utility function f for each user $u \in U$. Formally:

$$\forall u \in U, i^{max,u} = \arg \max_{i \in I} f(u, i)$$

The fundamental problem is that the utility value is not known for each item $i \in I$. We only know the utility values for the items the user has already rated in his history. Thus, the recommender system goal is predicting the utility function value on unknown data. Once we learn the utility function parameters, it is possible to realize a system that deems items' scores that are not yet been rated by the user.

It is straightforward that RSs need different kinds of data to produce their recommendations. In particular, three fundamental elements underpin any RS: **users**, **items**, and **ratings**. These components are usually represented using a rating matrix $\mathbf{R} : Users \times Items$ where each row corresponds to a user, whereas each column denotes an item. It is straightforward to infer that each matrix element is a rating given by a specific user to a specific item. In most of the cases, these matrices are very sparse because each user experienced only a small portion of the platform's items. Figure 3.1 shows a typical User-Item matrix in which question marks represent the lack of a rating. Let us draw, in detail, the three kinds of data:

- **Items** are an abstract representation of the objects or services proposed by the platform. In [264], the authors proposed an interesting taxonomy to categorize them. Atomic items can be considered as low complexity items. A few examples are news, webpages, books, songs, and movies. On the other side, some items are more complex since they can be a composition of other entities. Moreover, the complexity of an item can also vary because of the recommendation technique. In a common CF-RS, items are usually defined using only a numeric ID. Instead, in a CBF-RS, items can be extensively described by exploiting a language devoted to Knowledge Representation.

| | Item 1 | Item 2 | Item 3 | ... | Item n |
|--------|--------|--------|--------|-----|--------|
| User 1 | 2 | 3 | ? | ... | 5 |
| User 2 | ? | 4 | 3 | ... | ? |
| User 3 | 3 | 2 | ? | ... | 3 |
| ... | ... | ... | ... | ... | ... |
| User m | 1 | ? | 5 | ... | 4 |

Figure 3.1: Example of a rating matrix

- **Users** are the entities that receive recommendations and consume them. Each recommendation technique has to face at least two important problems regarding users. The first is how to deduce the behavior of users. On some platforms, where users provide an explicit rating to the items, the answer could seem obvious. In other scenarios, the behavior of users is deduced by clicks, listenings, plays, or other implicit feedback. The second problem is how to represent the user. This representation is dependant on the recommendation technique. For instance, a Neighborhood-based RS represents a user exploiting her rating vector. More, a matrix factorization algorithm uses a vector of latent factors. Finally, a CBF-RS usually represents a user with a weighted features vector.
- **Ratings**, or in general transactions, represent the relation between users and items. As mentioned, this relation can be established exploiting explicit feedback or implicit signals. A user can express explicit feedback in multiple ways. One example is the binary like/dislike rating. In this case, the user expresses her likeness of the item in a very strong and polarized way. Another example is the exploitation of a scale of values (e.g., from 1 to 5) in which the relation between the user and the item is weighted by a likeness degree. Finally, in the videos or music domain, explicit feedback is the number of playings of the video or the song. Explicit feedback comes with the ad-

vantage of avoiding misinterpretations of user behavior [296]. Indeed, with implicit feedback, this could occur at any moments, since the transformation of feedback in likeness is on the system designer. However, explicit feedback comes with an important drawback. Users are less prone to spend their time rating an item. Thus, only a small fraction of consumed items are effectively rated [199]. On the other hand, implicit feedback usually can catch only positive feedback (purchases, plays, listenings), without any information about what user dislikes. For this reason, this kind of ratings is usually also called 'unary' ratings. To sum up, explicit feedback can be considered more robust, even though they usually show inconsistencies due to users' behavior. Finally, in some cases, explicit and implicit information are combined for more accurate results [221].

3.1.2 From Rating Prediction to Ranking

As mentioned above, the recommendation task can be defined typically in two different ways. In rating prediction task, the aim is estimating the items rating for a user. Instead, Learning to Rank is the problem of recommending to users a shortlist [108] of Top-N recommendations [356]. These differences lead to different optimization goals. The former goal is modeled as a maximization of the prediction accuracy, whereas the latter is focused on generating the best list to provide to the user. Top-N recommendation task is also called **Item recommendation task** [312] since the optimization focus shifts from ratings to items. Learning to Rank [88] algorithms can be further categorized in Point-wise [225], Pair-wise [312, 247] and List-wise [344, 343]. It is important to underline that ratings are crucial in both tasks. However, one main difference is that, in a Rating Prediction task, the value associated with the user-item pair is an estimation of the rating. On the other side, in a Top-N task, the same value is used only to sort the recommendation list. We could say that the value is an estimation of a value associated with the position of the item, with no direct relation with the corresponding rating.

3.2 Collaborative-Filtering Recommender Systems

We have stressed that Recommender Systems techniques are usually classified into two main categories based on the kind of data they use to compute recommendations: Collaborative-Filtering and Content-Based. Beyond these classical models, a more exhaustive dissertation on recommendation techniques is provided in [317]. Since many of the proposed approaches are hybrid, Burke's classification [81] of hybrid Recommender Systems closes the section.

The basic idea behind a Collaborative-Filtering Recommender System is to take advantage of the transactions similarities among users or items [336]. Since the popularity of items among similar users is a very strong signal, this class of recommendation algorithms has become one of the most influential and successful. In general, these algorithms only need the user-item matrix with the corresponding feedback information. The underlying assumption is that similar users show similar preferences. Since the algorithms take advantage of users transactions, the performance of all the Collaborative-Filtering techniques is heavily affected by the number of users. The rationale behind this class of algorithms is to replicate a classic human behavior: asking similar people for reliable suggestions. Collaborative-Filtering algorithms can be further categorized in Memory-Based, and Model-Based Recommender Systems [73]. Memory-Based Recommender Systems directly exploit the user-item transactions to compute recommendations. The most extensively used approach for Memory-Based Recommender Systems is the K Nearest Neighbors (k-NN). In particular, they compute a similarity matrix to identify the more similar entities (users or items). The choice of the entities for which the similarity is computed defines the 'schema' of the k-NN algorithm. Instead, in Model-Based Recommender Systems, the algorithm trains a model exploiting the stored transactions. Usually, this model contains a latent representation of items and users. This representation is then used to produce recommendation lists.

3.2.1 Neighborhood-Based Models

User-Based k - Nearest Neighbors

This technique has been proposed by GroupLens Usenet article recommender [316] and resumed later by Ringo music recommender [342] and BellCore video recommender [181]. It computes the similarity between two customers and uses the preferences of the most similar users to estimate ratings to be assigned to items. In order to determine which user are “similar”, similarity functions need to be defined. They will be detailed in Section 3.2.1. There are several ways for computing rating predictions in user-based approaches. A popular one is shown in the equation below where $r_{u,i} \in \{1, \dots, 5\}$ is the known (five stars) rating of user u for the item i and $\hat{r}_{u,i}$ is the system prediction of the rating of u for i :

$$\hat{r}_{u,i} = \bar{r}_u + \frac{\sum_{u_j \in N_i(u)} (r_{u_j,i} - \bar{r}_u) \cdot w_{u,u_j}}{\sum_{u_j \in N_i(u)} |w_{u,u_j}|}$$

Here \bar{r}_u denotes the average rating of user u , $N_i(u)$ is the set of neighbors who rated item i and w_{u,u_j} is the similarity of the user u and u_j . The user to user similarity can be computed using different approaches. A popular one is the *Pearson correlation* [316]:

$$w_{u,u_j} = \frac{\sum_i (r_{u,i} - \bar{r}_u) \cdot (r_{u_j,i} - \bar{r}_{u_j})}{\sqrt{\sum_{i=1} (r_{u,i} - \bar{r}_u)^2} \cdot \sqrt{\sum_{i=1} (r_{u_j,i} - \bar{r}_{u_j})^2}}$$

Item-Based k - Nearest Neighbors

As the number of users increases, user-based kNN suffers from scalability problem. To overcome to this drawback, item to item kNN was introduced by Sarwar et al. in [334] and [209]. Thus, when the number of users exceeds the number of items, as is it most often the case, item-based recommendation approaches require much less memory and time to compute the similarity weights than user-based ones, making them more scalable. They are based on the similarity computation between items instead of users but the metrics used to compute similarity are the same used in user-user case. If two items shares same ratings, positive and negative ones, by the same users than they can be considered similar and hence can be assumed that a users

rates in a similar way similar items. It should be noticed that, when a new rating is inserted in the platform, it impacts only peripherally on item-item similarity while in a user-based approach, users that may not have been similar since now, could become a few moments later since ratings are constantly entered in the platform. Moreover, the increasing number of users will not affect real time operations since the system will use the item-item similarity matrix to make estimations.

Similarity measures

Various similarities have been suggested, namely the Euclidean Distance, the Cosine Distance and the Jaccard similarity, among others. Which one works best in Collaborative Filtering recommender systems is openly argued, having some authors claim that the Pearson Correlation is the best suited algorithm for User based approaches [73, 177], and others identifying the Cosine Distance on item based approaches as the best one [197]. The similarity weights have two main tasks in kNN approaches: they allow to select the best and trusted neighbors and they are used to weight the candidate items from these neighbors in a different way. Among all the similarity functions, two of them are particularly important for this dissertation: *Cosine similarity* and *Jaccard similarity*. The first is computed as:

$$w_{u,u_j} = \frac{\sum_{i \in I_{u,u_j}} r_{u,i} r_{u_j,i}}{\sqrt{\sum_{i \in I_u} r_{u,i}^2 \sum_{i \in I_{u_j}} r_{u_j,i}^2}}$$

where w_{u,u_j} represents the similarity of users u and u_j , across all items commonly rated. The second is:

$$w_{u,u_j} = \frac{\sum_i^n b_{i,u} b_{i,u_j}}{\sum_i^n b_{i,u} - \sum_i^n b_{i,u_j} + \sum_i^n b_{i,u} b_{i,u_j}}$$

where $b_{i,u}$ is a binary value that is 1 when item i has been rated by the user u and n is the number of items in the overall platform.

The two previous formulas refer to a user-based case, but it is straightforward to derive the item based counterpart.

Advantages

Advantages of these methods should be highlighted:

- **Simplicity:** Neighborhood-based methods are intuitive and relatively simple to implement. In their simplest form, only one parameter (the number of neighbors) has to be chosen.
- **Efficiency:** Compared to other recommender systems that need to be re-trained entirely, memory-based kNN may give an immediate feedback to users. In fact, in these systems not only similarity can be precomputed but also the K-neighbors for each item. In such a way the search operation and the prediction can be estimated in a rapid manner.

Drawbacks

However the neighborhood approaches have several flaws too:

- **Limited coverage** because of sparsity: The similarities between two users are computed by comparing their ratings for the same items. This can lead to an information loss. For example if user's u neighbors never rated item i this item will be never recommended to u despite her preferences.
- **Overfitting:** When no real neighbors can be computed, the recommender system will find them anyway leading to wrong prediction.
- For some set of neighbors the notion of rating notation in the referring space could be a mistaken assumption and the presence of neighbors could lead to erroneous estimations.
- No relationships between neighbors are considered. For example, movie's sequels will be considered and computed independently giving more weight to that kind of movie.
- **Cold Start:** at the heart of the kNN algorithm the user rating lays. Without ratings, no recommendation can take place for any user of the system. This problem named cold start, describes the inability of a recommender system to start offering recommendations when new users or items newly added to the system may have no ratings at all.

Neighborhood models summary

Building a kNN recommender system have three main steps: **normalization** of data, neighbors **selection** and interpolation weight **definition**. For each of these components, several different alternatives are available. It has to be noticed that the best approach may differ from one recommendation problem to another.

The **normalization** of data is important for each CF approach. Each user has personal interpretation of the rating scale. While one rater might tend to give high marks to films he likes, another rater might keep the highest grades for exceptional movies. There are two widely used systems to compensate for variations in the rating approach by different users, one is called mean-centering and the other is called Z-score.

The neighbors **definition** has a key role for recommendations as it emphasizes the relations between users and items. For this task, it is of paramount importance the similarity function choice. Last component is the interpolation weights **definition** because it affects how to consider item or user neighbors.

3.2.2 Matrix Factorization - SVD & SVD++

From memory-based to model-based Recommenders

Model Based algorithms received a significant push in the research community after the million dollar prize competition opened by Netflix in October 2006. During the three years that took the competition to be won, more and more contenders adopted Model Based approaches as part of their strategy. For the first time in history, the research community gained access to 100 millions movie ratings imparting a boost in new recommender systems approaches partly thanks to the nature of the competition. In fact, 1 milion dollars were offered to anyone who could improve over Netflix existent system in term of accuracy. During this competition model-based CF have proven to be more accurate than memory-based one [221]. These models create an offline model and apply it online to compute recommendations, which usually leads to greater accuracy and stability in recommendations. Although a great variety of models belong to this group, such as graph-based approaches and

probabilistic recommendation approaches, the most popular are the matrix factorization ones, also called latent factor models. These models analyze the user-item matrix in order to find latent factors, which can be described as latent features that characterize the user-item relationships. Since latent factors are computed in a non-supervised way, their interpretation is not trivial, and in some cases it is not possible at all (this is the reason why they are called latent). The computation of latent factors is performed through matrix factorization techniques: the user-item matrix is decomposed into two smaller matrices, the product of which is an approximation of the original matrix. Items are then recommended to users if they are close in the latent factor space and ratings are computed through element-wise products between user and item vectors. The first approaches used the Singular Value Decomposition (SVD) to decompose the user-item matrix. On the one hand, the obvious drawback of the Model Based method is that the model, not only takes substantially longer to be computed but needs to be computed anew if the matrix of data changes, which happens every time a new user enters a rating. Generally, small changes are left unprocessed, but when they become substantial, the model needs to be retrained. On the other hand these techniques tries to overcome to two major problems: sparsity of matrix and cold start. The first one is the sparsity of matrix caused by the insufficient number of the transactions and feedback data. The user-item matrix is usually almost entirely empty because only a small part of items, w.r.t. the entire item dataset, are rated by each user. The second one concerns the personalized recommendations for users with no or few past history(new users). Providing recommendations to these users becomes a difficult problem because their learning and predictive ability is limited. Multiple research have been conducted in this direction using hybrid models. These models use auxiliary information(side information) to overcome the cold start problem. The dual problem also occurs when a new item is introduced. Since no user have previously rated this item, it will not be recommended to anyone.

Baseline predictors

Since typical CF data exhibit large user and item biases, these need to be encapsulated within the *Baseline predictors* that depends only on user and item without involving interactions between them. For instance, some users have systematic tendencies to give very high/low ratings w.r.t. to others and some items to receive very high/low ratings. It is really important to model these biases meticulously because it helps to isolate exactly interactions between users and items to compare with other users' ones. We denote the baseline of a user u on an item i as b_{ui} . The simplest baseline that can be computed is the overall average rating: $b_{ui} = \mu$. The previous formula can be improved by computing the average of the specific user or towards that specific item. By also considering the observed deviation of user u and of item i from the average, the baseline could be improved. In general a baseline predictor for an unknown rating r_{ui} is computed as in the equation below:

$$b_{ui} = \mu + b_u + b_i$$

where b_u and b_i are the aforementioned deviations for user u and i respectively. These can be estimated by solving a least squares problem but an easier way to compute them is by decoupling the calculation of the b_i 's from the calculation of the b_u 's; it should be noticed that this can lead to less accurate results. They can be computed as follows:

$$b_u = \frac{1}{|I_u|} \sum_{i \in I_u} (r_{u,i} - \mu)$$
$$b_i = \frac{1}{|U_i|} \sum_{u \in U_i} (r_{u,i} - b_u - \mu)$$

A regularized version of the baseline can be obtained by using the regularization parameters β_u e β_i :

$$b_u = \frac{1}{|I_u| + \beta_u} \sum_{i \in I_u} (r_{u,i} - \mu)$$
$$b_i = \frac{1}{|U_i| + \beta_i} \sum_{u \in U_i} (r_{u,i} - b_u - \mu)$$

Considering the temporal dynamics in baseline predictors could be a successful approach. So just to finalize, these baseline predictors may catch the real user bias,

useful to obtain a new user-item matrix containing only the effective interactions between users and items.

Matrix Factorization Models

The most successful latent factor models are based on matrix factorization(MF)[224]. In its basic form MF characterizes both items and users by vectors of factors inferred from items rating patterns. These methods have recently become popular by combining good scalability with predictive accuracy. They offer much flexibility for modeling various real-life applications. The classic SVD model, widely used in Information Retrieval[118], is not well suited for CF domain. In fact, applying SVD to explicit ratings shows critical issues related to the high portion of missing values. One possible way to overcome to this problem is to fill the matrix with some estimations but this can lead to an inaccurate system[215, 332]. For such a reason, systems that directly model only the observed ratings have been deployed and regularization has been introduced to prevent overfitting.

SVD - Singular Value Decomposition

Matrix factorization models map both users and items to a joint latent factor space of dimensionality f , such that user-item interactions are modeled as inner products in that space. The latent space tries to explain ratings by characterizing both products and users on factors automatically inferred from user feedback[317]. These factors could represent both existing properties such has the genre of a movie or as is customary, completely unknown and that may not be interpreted by a human being. As a consequence, each item i is represented with a vector $q_i \in R^f$ while each user with a vector $p_u \in R^f$. Therefore, each item i will have a measure that characterizes it in each f dimensions. Same for user u in the same f dimensional space. In particular, for the user, the latent feature will represent her interest towards that property while for the item it will measure how much that property belongs to the item. The dot product $q_i^T p_u$ will capture the user's estimated interest in the item. The Figure 3.4

represents the basic form of matrix factorization. Finally, the estimated rating can be computed by combining it with the baseline predictor previously defined:

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T p_u$$

Our cost function to minimize will be the regularized squared error in the equation below:

$$\min_{b_*, q_*, p_*} = \sum_{(u,i) \in K} (\mu + b_i + b_u + q_i^T p_u)^2 + \lambda_4 (b_i^2 + b_u^2 + \|q_i\|^2 + \|p_u\|^2)$$

The regularization constant λ_4 is usually computed by cross validation while for the learning parameters b_u, b_i, p_u, q_i stochastic gradient descent or alternating least squares is performed. It is important to notice that alternating least square techniques fixes p_u to minimize q_i or viceversa. In this way the resulting cost function is convex and can be optimally solved[46, 45].

SVD++

Prediction accuracy could be improved through implicit feedback, useful to exploit additional information about user interest especially when few explicit ratings are provided. Moreover, in the event that no implicit feedback is available, it can be captured by considering items that user rated without considering the rating value. This lead to the development of several techniques [221, 298, 327]. Here we focus on the SVD++ method that modifies the previous equation of the estimated rating as follows:

$$\hat{r}_{ui} = \mu + b_i + b_u + q_i^T \left(p_u + |R(u)|^{-\frac{1}{2}} \sum_{j \in R(u)} y_j \right)$$

where y_j is a generic factor vector related to the item i based on implicit feedback. Moreover the user component has been modified in $p_u + |R(u)|^{-\frac{1}{2}} \sum_{j \in R(u)} y_j$. The sum of factor vectors y_j is normalized by $|R(u)|^{-\frac{1}{2}}$ in order to bring the variance across the range of observed values of $|R(u)|$. The cost function can now be minimized as mentioned in the SVD case. Several kind of implicit feedback may be considered and they will be automatically weighted during the minimization step.

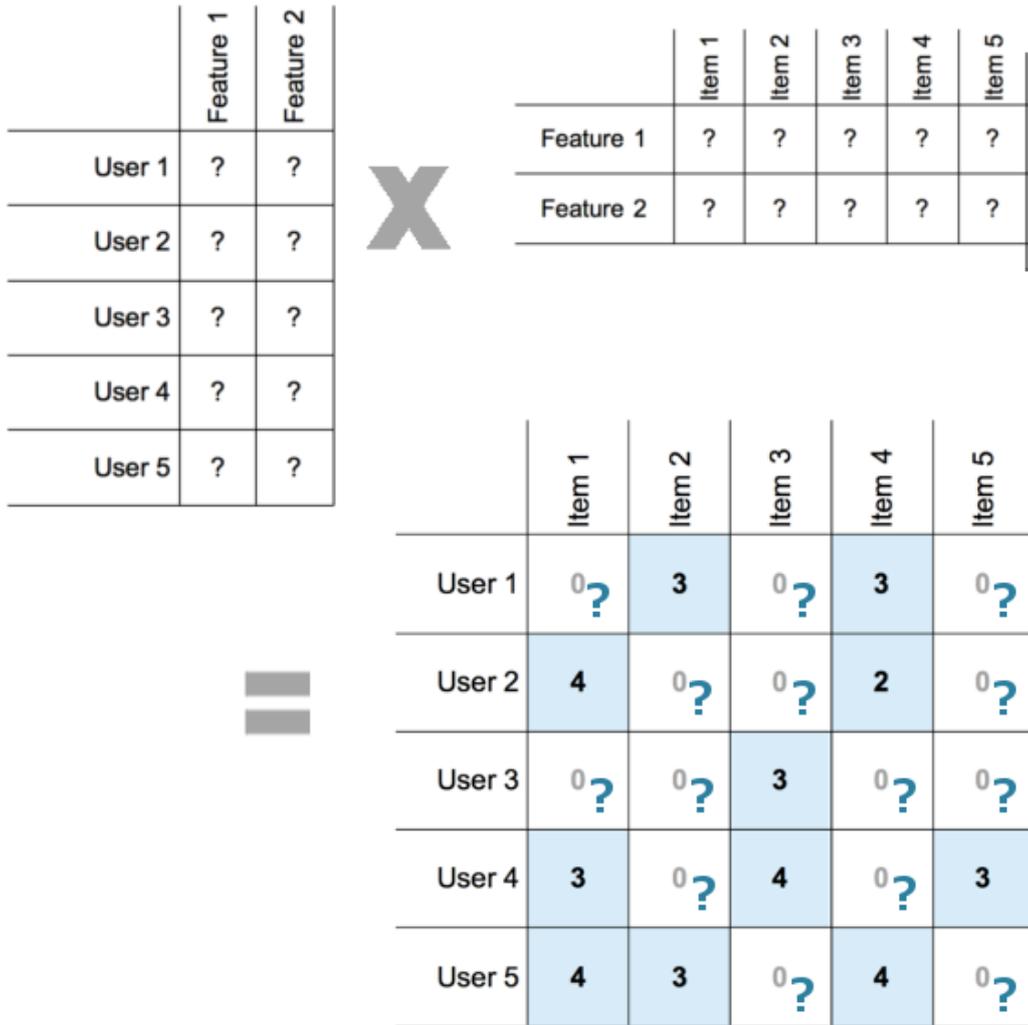


Figure 3.2: Basic representation of matrix factorization

3.2.3 Advances in Matrix factorization

Matrix Factorization (MF) models are among the most popular approaches to collaborative filtering, and have been actively investigated since they were introduced in the context of the Netflix prize competition [44]. As opposed to classic user- and item-based collaborative filtering heuristics [175, 239], MF methods train a statistical model from the available data using machine learning techniques. Specifically, they perform a *dimensionality reduction* of the highly sparse rating matrix into a subspace of latent factors, which aim to capture implicit properties of users and items. In order for MF to be effective, the dimension k of the latent subspace is assumed to be much smaller than the number of users and items, $k \ll \min(|U|, |I|)$, essentially acting as a *bottleneck* that compresses the sparse input while retaining enough information to explain the observed user-item interactions.

Matrix factorization models for rating prediction

Recommendation models based on MF have their roots on the *Latent Semantic Analysis* (LSA) technique [118], widely used in Natural Language Processing and Information Retrieval. LSA attempts to automatically infer concepts implicit in text documents by approximating the term-document matrix with a truncated Singular Value Decomposition (SVD) of lower rank. The first MF approaches for recommendation borrowed the same idea, and applied it to the user-item matrix in the rating prediction task [332]. In contrast to LSA, the SVD is not well defined for sparse matrices as those commonly found in recommender systems, and hence the above approaches relied on imputation techniques to fill the missing matrix entries before applying SVD.

Rather than filling the rating matrix, which may introduce inaccurate information, subsequent approaches aimed to only factorize observed ratings instead of the whole matrix. One of the first and most popular methods in this line is the model proposed by [151], in which each user u is assigned a vector $\vec{p}_u \in \mathbb{R}^k$ of latent features automatically inferred from the data, and similarly each item i is assigned a vector $\vec{q}_i \in \mathbb{R}^k$ in the same subspace. Intuitively, latent features aim to capture properties implicit in the data —such as the amount of *comedy* or *action* in the case of

movies—, but does not need to be interpretable at all, as this is not enforced in the model [223]. Ratings are then estimated as the dot product of latent feature vectors:

$$\hat{r}(u, i) = \langle \vec{p}_u, \vec{q}_i \rangle \quad (3.1)$$

Equivalently, the rating matrix \mathbf{R} is factorized as $\mathbf{R} \approx \mathbf{P}\mathbf{Q}^\top$, where \mathbf{P} is a $|U| \times k$ matrix with the user vectors \vec{p}_u as rows, and respectively \mathbf{Q} is $|I| \times k$ contains the \vec{q}_i as rows. The values of these matrices are automatically estimated from the data, by minimizing the Mean Squared Error of the ratings predicted against the ratings observed in a training set. That is, \mathbf{P} and \mathbf{Q} are chosen to minimize to following loss function:

$$\mathcal{L}(\mathbf{P}, \mathbf{Q}) = \sum_{(u,i) \in \mathcal{R}} (r_{ui} - \langle \vec{p}_u, \vec{q}_i \rangle)^2 + \lambda (\|\vec{p}_u\|^2 + \|\vec{q}_i\|^2) \quad (3.2)$$

where \mathcal{R} is the set of observed ratings, i.e. the set of non-zero entries of the rating matrix \mathbf{R} , and $\lambda > 0$ is a regularization hyper-parameter used to prevent overfitting.

3.2.4 Matrix Factorization - Funk MF

In [151] this function is minimized using *Stochastic Gradient Descent*, a widely used optimization technique that iteratively updates the parameters in the opposite direction of the gradient. When applied to 3.2, this technique yields the following update rules for the parameters \vec{p}_u and \vec{q}_i for each rating r_{ui} in the training set:

$$\vec{p}_u \leftarrow \vec{p}_u - \eta (e_{ui}\vec{q}_i + \lambda\vec{p}_u) \quad (3.3)$$

$$\vec{q}_i \leftarrow \vec{q}_i - \eta (e_{ui}\vec{p}_u + \lambda\vec{q}_i) \quad (3.4)$$

The *learning rate* η is a hyper-parameter that controls the extent to which the model parameters are updated in each iteration, and is carefully chosen; too large values may make the algorithm fail to converge, while too small values may make its convergence very slow. e_{ui} is the prediction error, and is defined as $e_{ui} \triangleq r_{ui} - \hat{r}(u, i)$.

In addition to Stochastic Gradient Descent, other optimization techniques have been explored in the literature, such as *Alternating Least Squares* [45], which is the standard technique followed in MF models for positive-only feedback (3.2.5).

The basic SVD model by [151] is easily extensible, and has served as a building block for more complex matrix factorization models. For instance, [221] proposed the SVD++ model, which includes additional parameters to account for implicit feedback in rating predictions. Further extensions of SVD introduce temporal variables to capture the evolution of user preferences through time [223].

3.2.5 Matrix Factorization for positive-only feedback

The core ideas behind the standard Matrix Factorization model for collaborative filtering have also been applied to the item ranking task when positive-only feedback is available instead of numeric ratings. Recommendation models designed for this type of data must take into account its particular characteristics, most notably the absence of negative feedback, but also the possible uncertainty in the positive feedback, as an observed user-item interaction may not necessarily indicate a preference of the user towards the item.

In one of the most representative works in this direction, [184] proposed an adaptation of the rating-based MF model described previously to deal with positive-only feedback. As opposed to the rating-based SVD, which only considers the observed ratings, Hu et al.’s method models the full set of $|U| \cdot |I|$ interactions. Since negative feedback is not available in this scenario, the authors argue that the algorithm has also to model the missing information as an indirect source of negative user preferences. For such purpose, they introduce a parameter c_{ui} for each possible user-item pair that measures the confidence on the corresponding interaction, whether observed or not:

$$c_{ui} = 1 + \alpha k_{ui} \quad (3.5)$$

where k_{ui} is the count of implicitly collected interactions between user u and item i –such as number of clicks on a product web page on an e-commerce site, and the number of listening records of a given song in an online music provider–, and $\alpha > 0$ is a scaling parameter. When no interaction is observed, $k_{ui} = 0$ and the model assigns minimum confidence to the user-item pair, as it is unknown whether the lack of interaction is because the user really does not like the item, or just because the user does not know the item. Likewise, the more interactions are collected and the

greater k_{ui} , the larger is the confidence on that observation. Moreover, focusing on the item ranking task, Hu et al.'s approach only aims to predict if the user will interact with the item, rather than the actual number of observations k_{ui} . Hence, a new set of variables is introduced so that $x_{ui} = 1$ if $k_{ui} > 0$, and $x_{ui} = 0$ otherwise.

Similarly to the SVD model for ratings, the recommendation score of item i for user u is estimated as the dot product of their corresponding latent feature vectors:

$$s(u, i) = \langle \vec{p}_u, \vec{q}_i \rangle \quad (3.6)$$

The model parameters \vec{p}_u and \vec{q}_i are again automatically learned by minimizing the mean squared error for the score predictions, but now accounting for the different confidence levels and the full set of possible user-item pairs:

$$\mathcal{L}(\mathbf{P}, \mathbf{Q}) = \sum_u \sum_i c_{ui} (x_{ui} - \langle \vec{p}_u, \vec{q}_i \rangle)^2 + \lambda (\|\mathbf{P}\|^2 + \|\mathbf{Q}\|^2) \quad (3.7)$$

3.3 Optimization techniques

3.3.1 ALS - Alternating Least Squares

The loss function can be minimized with different numerical optimization techniques such as Stochastic Gradient Descent, but in [184] the authors propose an *Alternating Least Squares* (ALS) procedure that efficiently handles the greater cost of accounting for the missing values. Clearly, the loss function in 3.7 involves many more terms than that of 3.2, as the number of observed entries in the user-item matrix is usually very small due to the data sparsity.

The key observation behind ALS is that when one set of parameters is fixed, the optimization problem in 3.7 is convex and analytically solvable using ordinary least-squares estimation. In particular, fixing the \vec{q}_i parameters and setting the gradient with respect to \vec{p}_u to zero yields the solution

$$\vec{p}_u = \left(\mathbf{Q}^\top \mathbf{C}^u \mathbf{Q} + \lambda \mathbf{I} \right)^{-1} \mathbf{Q}^\top \mathbf{C}^u \vec{x}_u. \quad (3.8)$$

where \mathbf{I} is the $k \times k$ identity matrix, \mathbf{C}^u is a $|I| \times |I|$ diagonal matrix with the c_{ui} values, and \vec{x}_u is a column vector of length $|I|$ with the x_{ui} values. The same procedure

procedure ALS-TRAINInitialize \mathbf{P}, \mathbf{Q} at random **repeat****P step**Fix \mathbf{Q} and optimize all \vec{p}_u in parallel using 3.8**Q step**Fix \mathbf{P} and optimize all \vec{q}_i in parallel using 3.9**until** *convergence*;**Algorithm 1:** Alternating Least Squares training algorithm.

can be applied by fixing the user factors, and optimizing the item factors, leading to the solution

$$\vec{q}_i = \left(\mathbf{P}^\top \mathbf{C}^i \mathbf{P} + \lambda \mathbf{I} \right)^{-1} \mathbf{P}^\top \mathbf{C}^i \vec{x}_i. \quad (3.9)$$

Similarly, \mathbf{C}^i is a $|U| \times |U|$ diagonal matrix with the c_{ui} confidence values, and \vec{x}_i is a column vector of length $|U|$ containing the binary values of x_{ui} .

As pointed out by the authors, the products $\mathbf{Q}^\top \mathbf{C}^u \mathbf{Q}$ and $\mathbf{P}^\top \mathbf{C}^i \mathbf{P}$ require time $\mathcal{O}(k^2|U|)$ and $\mathcal{O}(k^2|I|)$ for each user and item, respectively, and represent a computational bottleneck during the training phase. However, these terms can be computed more efficiently noting that $\mathbf{Q}^\top \mathbf{C}^u \mathbf{Q} = \mathbf{Q}^\top \mathbf{Q} + \mathbf{Q}^\top (\mathbf{C}^u - \mathbf{I}) \mathbf{Q}$, where $\mathbf{Q}^\top \mathbf{Q}$ is independent of the user and thus can be precomputed, and $\mathbf{C}^u - \mathbf{I}$ only has non-zero entries in the diagonal for the $|I(u)|$ items with $k_{ui} > 0$, which is much smaller than $|I|$. Considering the computation of the matrix inverse, the total complexity of 3.8 for a single user is $\mathcal{O}(k^2|I(u)| + k^3)$. Likewise, the complexity for 3.9 is $\mathcal{O}(k^2|U(i)| + k^3)$.

The main advantage of ALS is that the optimal factors for each user in Equation (3.8) can be computed in parallel once the item factors are fixed (P step). Symmetrically, once the user factors are obtained and fixed, the item factors in 3.9 can be found for each item in parallel (Q step). This observation leads to the alternating nature of ALS, respectively fixing one set of parameters and optimizing the other until convergence is reached or for a given number of iterations, as illustrated in 1.

The ALS-based method by [184] became the standard baseline for matrix factorization models with positive-only feedback, and has been extended and improved in

subsequent works since it was first proposed. One notable paper by [304] presents a new training procedure to boost the time complexity of the P step of each user to $\mathcal{O}(k^2 + k|I(u)|)$, and analogously the Q step. In order to achieve this significant improvement, the authors propose an approximate solution to the least-squares problem in each step. Rather than directly finding the k -dimensional solution as in Equations (3.8) and (3.9), which involves the costly computation of a matrix inverse, their approach iteratively solves k one-dimensional least squares problems, one for each latent dimension, much less expensive to solve. As reported in the paper, the loss of accuracy due to the approximate algorithm is small compared to the saved time for training. In subsequent work, [365] extended ALS to a ranking-based MF approach that learns to predict the relative ordering of items instead of individual point-wise scores. More recently, [295] proposed a graph-based Bayesian model that is able to capture the meaning of missing values, distinguishing between a user disliking an item or being unaware of it.

3.3.2 BPR - Bayesian Personalized Ranking Criterion

Matrix Factorization models can be easily trained to reduce the prediction error via gradient descent methods, alternating least-squares (ALS) and MCMC. However, in a *top-N* recommendation task, MF models can be trained using a learning to rank approach like Bayesian Personalized Ranking Criterion (BPR) [312]. The BPR criterion is optimized using a stochastic gradient descent algorithm on a set D_S of triples (u, i, j) , with $i \in I^u$ and $j \notin I^u$, where I^u is the set of items for the user u . Each triple is selected through a random sampling from a uniform distribution. The BPR optimization criterion can thus be formulated as:

$$\begin{aligned} \text{BPR-OPT} &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{x}_{uij}) - \lambda_{\Theta} \|\Theta\|^2 \\ &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{y}(\mathbf{x}_{ui}) - \hat{y}(\mathbf{x}_{uj})) - \lambda_{\Theta} \|\Theta\|^2 \end{aligned} \quad (3.10)$$

where \hat{x}_{uij} is the estimated rating difference, and $\hat{y}(\mathbf{x}_{ui})$ is the prediction the user-item pair $u - i$. In this formulation, $\sigma(\cdot)$ is a sigmoid function, and the update step

is defined as:

$$\Theta \leftarrow \Theta + \alpha \left(\frac{e^{-\hat{x}_{uij}}}{1 + e^{-\hat{x}_{uij}}} \cdot \frac{\partial}{\partial \Theta} \hat{x}_{uij} + \lambda \Theta \right) \quad (3.11)$$

where α is the chosen learning rate. In an implicit feedback setting, we may assume that there is only an instance for the pair user-item. Hence, in the model we can derive \hat{x}_{uij} as:

$$\begin{aligned} \hat{x}_{uij} &= \hat{y}(\mathbf{x}_{ui}) - \hat{y}(\mathbf{x}_{uj}) = w_i - w_j + \\ &+ \sum_{f \in F} v_{(u)} v_{(i)} - v_{(u)} v_{(j)} \end{aligned} \quad (3.12)$$

where w_i is the bias for the given item i , and $v_{(u)}$ is the vector of latent factors for the user u . This computation can be performed in an efficient way computing the partial derivatives (to update the factorized parameters) for the only active entities involved in the transactions, w_i , w_j , v_u , v_i , and v_j :

$$\frac{\partial}{\partial \Theta} \hat{x}_{uij} = \begin{cases} 1, & \text{if } \theta = w_i, \\ -1, & \text{if } \theta = w_j, \\ v_{(u)}, & \text{if } \theta = v_{(i)}, \\ -v_{(u)}, & \text{if } \theta = v_{(j)}, \\ (v_{(i)} - v_{(j)}), & \text{if } \theta = v_{(u)}, \\ 0, & \text{otherwise} \end{cases} \quad (3.13)$$

Using Equation (3.13) in Equation (3.11) the model parameters can be iteratively updated to maximize the BPR-OPT criterion. The algorithm updates sequentially each sampled triple and continues the training until it reaches the provided number of iterations.

3.4 Content-Based Recommender Systems

Content-based filtering Recommender Systems exploit item content to provide recommendations based on what the users rated in the past. The recommendation pro-

cess starts with the definition of a user profile. While items often have a technical description in terms of features, the challenging task is to extract significant features that can drive the recommendation process. In content-based methods, feedback of users is combined with the content information available in the items. In most cases, the items' attributes are simple keywords that are extracted from their description. Semantic indexing techniques are also used to represent the item and user profiles using concepts instead of keywords. The recommendation process basically consists of matching up the attributes of the user profile against the attributes of a new item. The result is a relevance estimation that represents the user's level of interest in that object. If a profile accurately reflects user preferences, the effectiveness of a content-based approach is considerably improved. The most common technique to define a user profile from the descriptions of the items she liked is to extract a list of words that describes such items. In this case, the main issue is the definition of the criteria to apply for selecting the most important words. Therefore, thanks to the availability of several open knowledge sources such as DBpedia¹ and to the increasing interest in semantic technologies, lot of research works moved from the classic keyword-based approach to a concept-based one.

The architecture of a content-based recommender system is composed of three main components: **Content Analyzer**, **Profile Learner**, and **Filtering Component**. Figure 3.3 shows the high-level architecture of a content-based recommender system.

The *Content Analyzer* is a component devoted to extracting the relevant features of an item. It exploits feature extraction techniques to build a new item representation that could be processed in the simplest way by a recommender system.

However, this item could be of any kind and this is the reason why this task could lead to a very difficult issue to solve. For example, analyzing a document or a web page may be a simple task while in the case of a movie or an image it might not be that easy. In the first case, for instance, an approach that takes into account most relevant words could be adopted by using some metrics that define the significance of a word in that document or web page. It is straightforward that the choice of

¹<http://wiki.dbpedia.org/>

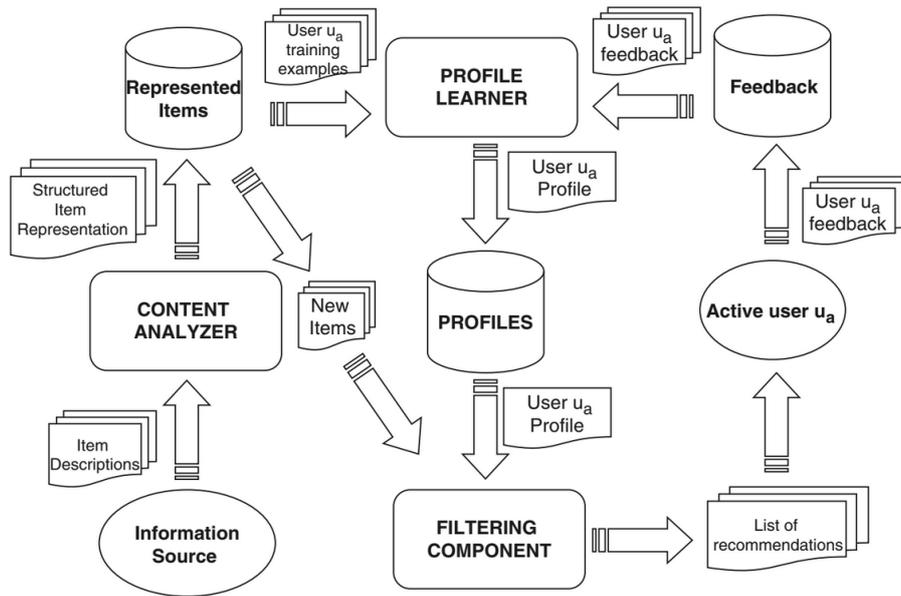


Figure 3.3: High level architecture of a content-based recommender [317]

approach affects the performance of the entire recommender system.

For movies or images, the analysis becomes more complex. Besides visual features extraction, the problem can be alleviated if some description of those items is available in a knowledge base. Then, we can exploit this structured information to feed the recommender system. However, once a description is retrieved, it is still important to state what information should be considered.

The new representation of the item becomes the input of the *Profile Learner* and *Filtering Component*

The *Profile Learner* exploits the descriptions to match with the abstract user model. Its goal consists in inferring user preferences and interests to build a user model. Usually, Machine learning techniques are employed to train the user model. Finally, the *Filtering Component* is responsible for estimating a score for a generic item exploiting the user profile representation. This is the component that produces the list of recommendations.

Conclusion

The adoption of the content-based recommendation paradigm has several advantages:

- Content-based recommenders are based only on ratings provided by the history of the target user and on their ability to learn a user profile representation that approximates user preferences as well as possible. This means that no collaborative information is needed to feed the recommender. Instead, a Collaborative-Filtering method needs ratings from other users to provide recommendations. Consequently, Content-based recommenders do not suffer the so-called Cold-start items problem.
- They can provide an understandable and immediate *explanation* of the recommended items thanks to the nature of the recommendation process. Indeed, *Explanations* can be provided, in their simplest form, by explicitly listing content features.

Content-based methods do have several shortcomings as well:

- Even though content-based methods are effective at providing recommendations for new items, they are not effective at providing recommendations for new users. The reason is the dependency of the user profile from the history of her ratings. In general, a large number of ratings implies a robust representation of the user profile and hence robust predictions.
- As mentioned before, content-based recommender systems are based on the past of the user to suggest new items that are consistent with it. This approach will tend to produce always the same kind of recommendations over time. As an example, a person that has not experienced romantic movies would never receive a recommendation for a romantic movie. This behavior is called *over-specialization* problem or lack of *Serendipity*. These terms highlight the tendency of the content-based systems to not produce recommendations with unexpected items.

- Best Content-based techniques require something more than factual information, they require domain knowledge. As an example, the description of a movie can distinguish among actors and the director, and this can lead to accuracy improvements. Typically, the elected tool to let a recommender understand structured information is providing domain ontologies. Unfortunately, domain ontologies are not always available. Consequently, feeding the recommender with only factual information can lead to losing some important information for the recommendation task.
- The vast size of the item set is a problem that content-based methods have to face. Since we need to find items that correlate the most with the user's interests, we are forced to examine all the items. Moreover, we must examine the content of every item to estimate a score, whereas collaborative filtering systems only need to examine users ratings. Therefore, since the number of items rises very quickly, a content-based suffers scalability issues. As a result, the solution is simplifying the item representation. Commonly, this leads to a performance decrease.

3.4.1 Vector Space Model

Many content-based recommender systems represent documents using a simple spatial representation. Probably the most know representation is the Vector Space Model (VSM). In a VSM, each user or item (document) is represented in an n -dimensional space, where n is the cardinality of the considered features (or keywords). In practical terms, users and items are represented through Boolean or weighted vectors. Their respective positions and the distance, or better the proximity, between them, provides a measure of how these two entities are related or similar. The choice of features may substantially differ depending on their availability and application scenario: crowd-sourced tags, categorical, ontological, or textual knowledge are just some of the most exploited ones. All in all, in a CB approach we need:

1. to get reliable items descriptions;

2. a way to measure the strength of each feature for each item;
3. to represent users;
4. to measure similarities.

Item descriptions The set of available features is usually generated processing the textual description of items. Some of the most common operations in this sense are: tokenization, stopwords removal, and stemming, or even more advanced Natural Language Processing methods [35].

Feature strength The different terms within a document usually deserve various degrees of importance for categorizing the document. What we do need is a weighting procedure, a scheme, that lets us assign the correct weight to different terms. Some of the most adopted schemes are TF-IDF and BM25. TF-IDF, in particular, can be derived from the probabilistic distribution of terms, and it reflects some common observations [328]. First of all, rare terms should not be considered less relevant than frequent terms. Second, if a term is frequent in a document it is not less relevant than occurring once. Last, the contribution of each document in the weighting should be independent of the length of the description.

TF-IDF Given a set of items $I = \{i_1, i_2, \dots, i_N\}$ in a catalog and their associated features f_i in a Collection \mathcal{C} we may build the set of all possible features as $F = \{f \mid f_i \in \mathcal{C} \text{ with } i \in I\}$. In the following, we use f to denote a feature in F irrespectively of the item. Let us denote with $freq(f, i)$ the frequency of the feature f in the document i . Consequently, $\max_z freq(z, i)$ denotes the maximum frequency considering all the features of item i . Let us $I(f)$ represents the subset of I which considers only those items that contain f . Each item can be then represented as a vector of weights $\omega_i = [\omega_{(i,1)}, \dots, \omega_{(i,f)}, \dots, \omega_{(i,|F|)}]$ where $\omega_{(i,f)}$ is computed as the normalized TF-IDF value for f :

$$\omega_{(i,f)} = \frac{\frac{freq(i,f)}{\max_z freq(z,i)}}{\sqrt{\underbrace{\sum_{k \in F} \left(\frac{freq(i,k)}{\max_z freq(z,i)} \cdot \log \frac{|I|}{|I(k)|} \right)^2}_{TF}}} \cdot \underbrace{\log \frac{|I|}{|I(f)|}}_{IDF}$$

User Representation Analogously, when we have a set U of users, we may represent them using the features describing the items they enjoyed in the past. Given a user u , if we denote with I^u the set of the items enjoyed by u . Let us denote with $I^u(f)$ the set of items experienced by u which contain f . Hence, we have $\omega_u = [\omega_{(u,1)}, \dots, \omega_{(u,f)}, \dots, \omega_{(u,|F|)}]$ with

$$\omega_{(u,f)} = \frac{\sum_{i \in I^u} \omega_{(i,f)}}{|I^u(f)|}$$

Similarity As mentioned, a similarity measure is required to measure the closeness between users and items in the features space. The proximity of two vectors can be computed in multiple ways. However, a very common strategy to evaluate similarities between the vectors i and j is evaluating the cosine vector similarity (CSV) of their corresponding vectors in \mathbf{F} :

$$CSV(i, j) = \frac{\omega_i \cdot \omega_j}{\|\omega_i\| \cdot \|\omega_j\|}$$

3.4.2 From Vector space Model to Knowledge-aware Recommender Systems

Regarding the descriptions of items, nowadays we can easily retrieve them from the Web. In particular, thanks to the Linked Open Data initiative a lot of semantically structured knowledge is publicly available in the form of Linked Data datasets. In [279], the authors originally proposed to encode a Linked Data knowledge graph in a vector space model to develop a CB recommender system. Let us recall that, in a Linked Data dataset, a resource is described using triples in the shape *subject-predicate-object*, where the subject is the resource itself. Accordingly to the graph

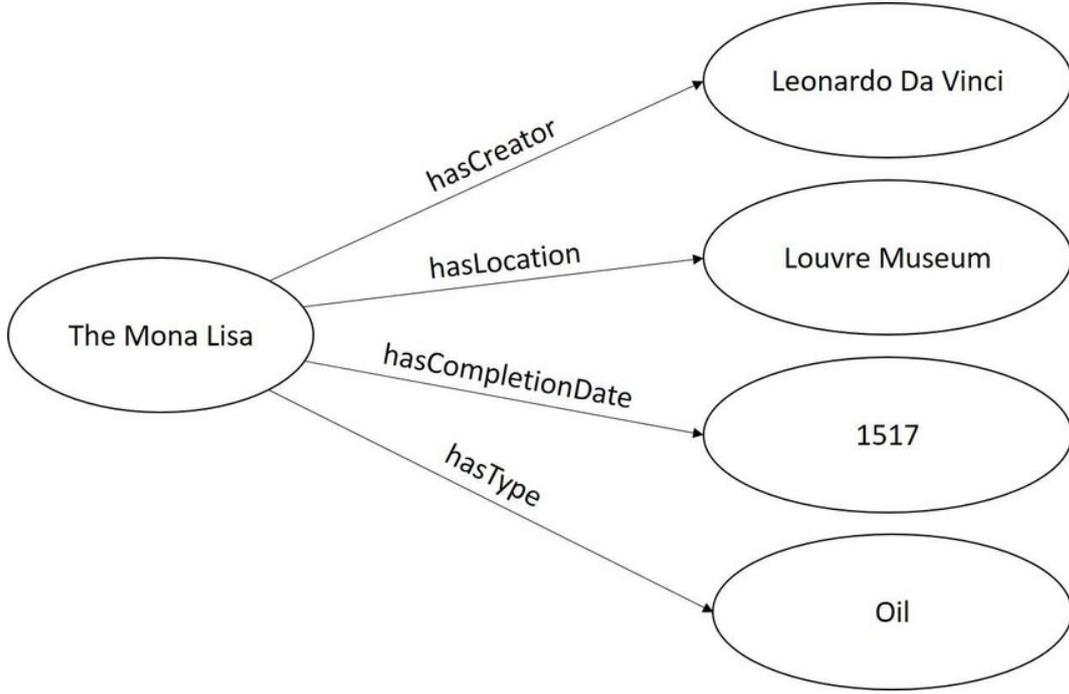


Figure 3.4: Basic example of item semantic description

nature of data, the underlying problem is building the vectors ω_i , and ω_u which represent, respectively the item and the user vectors. This difference requires to redefine the features representation. Given a set of items $I = \{i_1, i_2, \dots, i_N\}$ in a catalog and their associated triples $\langle i, \rho, \omega \rangle$ in a knowledge graph \mathcal{KG} we may build the set of all possible features as $F = \{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG} \text{ with } i \in I\}$. In the following, when no confusion arises, we use f to denote a feature $\langle \rho, \omega \rangle$ in F . Each item can be then represented as a vector of weights $\mathbf{i} = [v_{(i,1)}, \dots, v_{(i,\langle \rho, \omega \rangle)}, \dots, v_{(i,|F|)}]$ where $v_{(i,\langle \rho, \omega \rangle)}$ is computed as the normalized TF-IDF value for $\langle \rho, \omega \rangle$:

$$v_{(i,\langle \rho, \omega \rangle)} = \underbrace{\frac{|\{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG}\}|}{\sqrt{\sum_{f \in F} |\{f \mid \langle i, f \rangle \in \mathcal{KG}\}|^2}}}_{TF^{\mathcal{KG}}} \cdot \log \underbrace{\frac{|I|}{|\{j \mid \langle j, \rho, \omega \rangle \in \mathcal{KG} \text{ and } j \in I\}|}}_{IDF^{\mathcal{KG}}}$$

Analogously, when we have a set U of users, we may represent them using the features describing the items they enjoyed in the past. Given a user u , if we denote

with I^u the set of such items and we have $\mathbf{u} = [v_{(u,1)}, \dots, v_{(u,f)} \dots, v_{(u,|F|)}]$ with

$$v_{(u,f)} = \frac{\sum_{i \in I^u} v_{(i,f)}}{|\{i \mid i \in I^u \text{ and } v_{(i,f)} \neq 0\}|}$$

Once we have generated the new item and user representation, the computation of the user-item similarity (and hence the generation of recommendations) can proceed as aforementioned.

3.5 Hybrid Recommenders

Hybrid Filtering is a combination of different Filtering approaches [287]. The rationale is proposing a filtering method able to go beyond the limitations of traditional approaches. Some of these limitations are the Cold-Start problem, the over-specialization of Content-Based approaches and the sparsity of the Users-Items matrix. Moreover, some hybrids recommenders aim directly to increase the accuracy of recommendations and performance of the whole system. The goal of exploiting a hybrid recommender is to compensate for the limitations of an approach using the advantages of another. There exist multiple ways to implement this kind of recommender [17]:

- Implementing the two approaches separately and then use an aggregation function to merge the results;
- Injecting some Content-Based approach characteristics into a Collaborative-Filtering approach;
- Injecting some Collaborative-Filtering approach characteristics into a Content-Based approach;
- Building a single system that exploits at the same time collaborative and content information.

Typical examples of hybrid recommenders are streaming platforms like Netflix or Amazon Prime. On one side, they provide recommendations based on collaborative information. On the other side, they also exploit movie features to propose

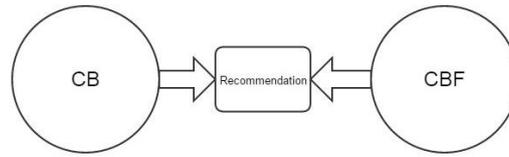


Figure 3.5: Hybrid system that implements the two approaches separately and aggregates the results

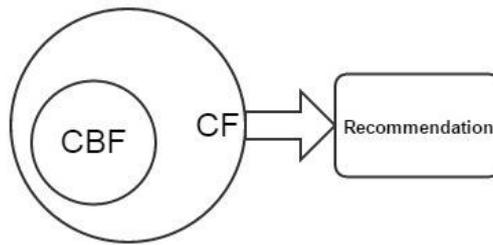


Figure 3.6: Hybrid system that injects some Content-Based approach characteristics into a Collaborative-Filtering approach

more accurate recommendations. In detail, in 2006 Netflix released an anonymized dataset regarding one million users challenging researchers to overcome their recommendation algorithm. Among the best approaches, most of them were simple hybrids that had taken advantage of ensembles of different methods.

It is straightforward that the choice of the integration scheme affects the resulting architecture. Figure 3.5 shows that the two recommenders are completely distinct and the integration takes place only during the generation of the recommendations. Figures 3.6 and 3.8 show that we can make one of the recommenders completely dependant on the other, which usually becomes the new knowledge source. Figure 3.7 depicts a single system that is fed by two different information. In this classification scheme, there are no constraints about the exploitation of that information.

Another interesting classification for hybrid models has been proposed by Burke [80]:

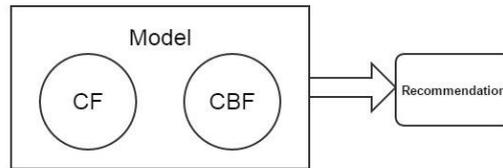


Figure 3.7: Hybrid system that exploits at the same time collaborative and content information

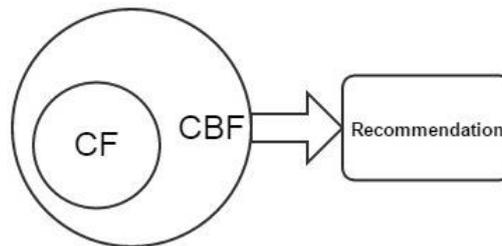


Figure 3.8: Hybrid system that injects some Collaborative-Filtering approach characteristics into a Content-Based approach

- **Weighted:** The ratings returned by the different approaches are then aggregated through a weighting function;
- **Switching:** The system evaluates the context and based on that it selects the best algorithm to propose recommendations;
- **Mixed:** The different algorithms generate different recommendation lists. These lists are then merge using an arbitrary mechanism like the mean or a combination of parts of the lists;
- **Feature combination:** The weights of the features obtained after the training are used as the new knowledge source. Another recommender receives this input and proposes the recommendations;
- **Cascade:** The recommendation process is considered as a sequential process in which the output of a recommender feeds another;

- **Feature augmentation:** The output of a recommender is used to enrich items and users' descriptions of another recommender;
- **Meta-level:** The overall system exploits one recommendation algorithm to generate the model and this model is used as input for another recommender.

3.6 Recommender Systems Evaluation

In the last decades, we observed a flourishing of Recommender Systems algorithms. Most of them are very specialized approaches, designed to perfectly fit a specific scenario. The same recommender fed with different data would probably perform poorly. This happens because identifying what makes Recommender Systems behave differently is a hard task. For instance, it could be due to operating in different settings. However, even the degree of sparsity of the Users-Items matrix can affect the performance of the system. Moreover, even though two systems operate in the same scenario, with the same data, identifying the causes of the different behaviors would be hard. Indeed, there is no total convergence on which is the most important dimension of evaluation, neither of which is the best metric to measure [178].

Evaluating a recommender can be hard for many reasons. As an example, given the same set of data, an algorithm can show oscillating performance based on the evaluation target. Moreover, different evaluation settings can affect the evaluation result. Finally, even if we consider the recommendation accuracy, measuring the rating prediction error is a completely different task from evaluating the ranking of recommended items.

In [256], the authors have pointed out that accuracy is not the only aspect to consider during evaluation. Others consider the online evaluation as the only effective evaluation. However, online evaluation with real users means bearing the costs of the platform and involving users. Even not considering the necessary evaluation time, it is clear that this kind of evaluation is not affordable for every research group. For these reasons, we focus on offline evaluation. This section is devoted to providing a brief overview of the evaluation protocol, metrics, and methods [178].

3.6.1 Protocols

How to assess the quality of a recommendation list, even though we are limiting our analysis to offline evaluation, is not an easy task. It could be argued that the main aspect to consider is the accuracy of the recommendations. However, the definition of accuracy is not unique. For instance, if we change the recommendation task from rating prediction to items' recommendation, also the notion of accuracy should be re-formulated. Before going into the different definitions, we need to give some Recommenders Evaluation background. Bellogin [227] has shown that there are many aspects in common between the Recommender Systems and Information Retrieval evaluation. This is reasonable since the research field was born as a branch of Information Retrieval. Moreover, if we focus on a Top-N recommendation task, the similarities are more evident. Consequently, most of the evaluation techniques and metrics we use were originally designed for Information Retrieval tasks. However, between the two evaluation paradigms, there are important differences. First of all, in Information Retrieval domain experts can objectively evaluate the relevance of a document. This is not possible in a recommendation scenario.

The relevance of a recommended item can be established only on a per-user basis. This has dramatic effects on the evaluation. In theory, we should evaluate a Recommender System only online, proposing recommendations and waiting for the feedback. However, this means that a study conducted by a small/medium-sized laboratory can rely only on a few users. Moreover, the number of items that are effectively rated is usually very small. How a researcher should deal with the remaining unrated items is an open problem. Indeed, the choice of ignoring them or considering them as negative examples affects either the training phase and the evaluation procedure. Moreover, the splitting of data in a training set and a test set can be realized in multiple ways. Here too, the choice has significant effects on the evaluation results [356]. The authors underline that, even though the test set contains only a small number of items, the ranking affects the entire collection. For this reason, they propose two different evaluation protocols: **all unrated items**, and **rated test-items**.

The **all unrated items** protocol imposes to use as candidate items all the items

that have never been rated by the user. A candidate item is an item that can be potentially recommended. Consequently, the items that are present in the test set and the recommendation list are considered as relevant [48][325]. The rationale is that in a real case scenario the recommender analyzes all the items in the collection. Instead, in the **rated test-items** protocol, the candidate items correspond to test items. The difference among algorithms is measured by ranking ability or prediction errors.

3.6.2 Accuracy

As mentioned, the accuracy metrics are broadly categorized in error metrics, and ranking metrics. According to the historical development of metrics, first we focus on error metrics.

Rating prediction metrics

In the past, the metrics based on prediction error have been the most widely used. Among them, the most known ones are the **Root Mean Squared Error (RMSE)** and **Mean Absolute Error (MAE)**. The goal of a **rating prediction task** is to estimate the rating for an item by a given user. Once the recommender is trained, we use its prediction function to estimate the rating value \hat{r}_{ui} , for the test set (TS) item. For this user-item pair (u, i) , the actual value of rating is known but hidden to the recommender.

$$\text{MAE} = \frac{1}{|TS|} \sum_{(u,i) \in TS} |\hat{r}_{ui} - r_{ui}| \quad (3.14)$$

$$\text{RMSE} = \sqrt{\frac{1}{|TS|} \sum_{(u,i) \in TS} (\hat{r}_{ui} - r_{ui})^2} \quad (3.15)$$

Ranking Metrics

Since the research community has moved toward the **Top-N recommendation task**, the previous metrics have fallen out of use. The reason is that a low rating error does

not imply a correct ranking for Top-N recommendations. Indeed, these metrics evaluate errors of all relevant items without considering any sort order. Consequently, the error of the last relevant item is equally important to the error on the most relevant one. To sum up, rating prediction metrics may estimate the same error for top-N items and bottom-N items, without taking into account that an error in top-N items should be more relevant compared to an error for lower ranked items.

The new evaluation need has pushed researchers to propose new metrics. Among them, we there are **Precision**, **Recall**, and **Normalized Discounted Cumulative Gain**. Since these metrics check the presence of items in the first N recommendations, the metric value is coupled with a number which denotes the list length. Common values are: 1, 5, 10, 25, 50, 100. Precision and Recall require binary values. In practical terms, they assign 1 to relevant items and 0 otherwise.

Given a user u , Precision ($P_u@n$) is computed as the ratio of the number of relevant items in the recommendation list over N . Recall ($R_u@N$) is computed as the ratio of the number of relevant items in the recommendation list over the number of relevant items in the Test Set.

Precision is defined as the proportion of retrieved items that are relevant to the user.

$$P_u@N = \frac{|L_u(N) \cap TS_u^+|}{n} \quad (3.16)$$

where $L_u(N)$ is the recommendation list up to the N -th element and TS_u^+ is the set of relevant test items for u . Precision measures the system's ability to reject any non-relevant documents in the retrieved set. Recall is defined as the proportion of relevant items that are retrieved.

$$Recall@N = \frac{|L_u(N) \cap TS_u^+|}{TS_u^+}$$

Recall measures the system's ability to find all the relevant documents.

where TS_u^+ is the set of relevant items in the Test Set for user u and $L_u(N)$ is the recommendation list truncated at N . Precision is heavily affected by the evaluation protocol. On the other side, Recall is affected by the threshold chosen to identify a test item as relevant. The items with an actual rating lower than the threshold are considered non-relevant. For this reason, this threshold is usually called the

relevance threshold. When a relevance threshold is used to compute a metric, it is adopted the *threshold-based relevant items* condition [83]

Precision and recall can be combined with each other in the F1 measure computed as the harmonic mean between precision and recall.

$$F1@N = 2 \cdot \frac{Precision@N \cdot Recall@N}{Precision@N + Recall@N}$$

In information retrieval, Discounted cumulative gain (DCG) is a metric of ranking quality that measures the usefulness of a document based on its position in the result list. Recommended results may vary in length depending on the user, therefore is not possible to compare performance among different users, so the cumulative gain at each position should be normalized across users. Let us define r_{uk} as the rating assigned by the user u to the item at the position k of the recommendation list. Hence, normalized discounted cumulative gain, or nDCG, is computed as:

$$nDCG_u@N = \frac{1}{IDCG@N} \sum_{k=1}^N \frac{2^{r_{uk}} - 1}{\log_2(1+k)} \quad (3.17)$$

where k is the position of an item in the recommendation list and $IDCG@N$ indicates the score obtained by an ideal ranking of the recommendation list $L_u(N)$ that contains only relevant items. $IDCG@N$ is used as a normalization factor. However, for all unrated items, there is no rating for the majority of items. A common adopted solution [356] is defining a single fix value instead.

3.6.3 Diversity

Accuracy metrics are needed to measure performance of a top-N recommender. Nevertheless, recommendation quality also depends on other factors that could improve the user experience, such as Novelty and Diversity.

Diversity, as an alternative dimension of the performance of recommender systems, can be measured through common measures as *item coverage*, *catalog coverage*, *Gini index*, and *Shannon entropy*.

Some recommenders could be designed to produce recommendations only for some users or recommend only a small number of items. This problem is known as

the long tail or heavy tail problem. Nevertheless, if recommendations include very popular items, the accuracy of recommendations can still be high. On the other side, the interest of users toward only very popular items could be limited. These considerations have pushed researchers to define new metrics that take into account these problems. They have proposed two metrics: **User Coverage**, and **Item (or Catalog) Coverage**.

Coverage

Catalog Coverage computes the number of items that are effectively inserted in recommendation lists. It denotes the propensity of a system to recommend always the same items. An example of measuring this quantity is **aggregate diversity**. This metric, defined as **diversity-in-top-n** in [14], measures the overall number of distinct items recommended in all recommendation lists.

$ADiv@N$ is computed as:

$$ADiv@N = \frac{|\bigcup_u L_u(N)|}{|I|}$$

It is worth to note that in this formula it is present a normalization for the length of the catalog.

Moreover, this metric provides some information about the personalization of recommendation. A high value of Item Coverage implies that most of the users received completely different lists. Dually, a low value of Item Coverage suggests that the recommendation lists are not diversified.

Similarly, the **User Coverage** [163] denotes the overall number of different users the recommender is able to produce a recommendation for. It denotes the ability of a system to produce recommendation for all the users.

Distributional inequality

Another measure of diversity is **sales diversity** [148], which measures how the items are unequally distributed in the recommendation lists. The *Gini index* (*Gini*) and Shannon entropy (*SE*) are two different metrics used to measure the distribu-

tional inequality [91]. Distributional inequality measures how unequally different items are chosen by users when a particular recommender system is used. Considering L_I the list of items in the collection ordered in increasing order of the probability to be put in a recommendation list $p(i)$, the *Gini index* is defined as:

$$Gini = \frac{1}{n-1} \sum_{j=1}^n (2j-n-1)p(i_j)$$

where i_1, \dots, i_n is the list of items sorted by increasing $p(i)$.

The index goes to 0 when all items are equally chosen. On the other side, it goes to 1 when the recommender always selects the same item. According to the evaluation target, it is possible to adopt another formulation that considers the number of times an item is recommended. Usually, the Gini index is represented through a reversed scale, obtained as $1 - \text{Index}$. This choice makes the index more readable because, with this formulation, higher is better. In this case, low values correspond to a scenario in which items are not equally chosen.

Shannon entropy (*SE*), on the other hand, takes the same probability and computes a different value:

$$SE = - \sum_{i=1}^n p(i) \log p(i)$$

Entropy goes to 0 when a single item is recommended. When a recommender recommends n items, the entropy value reaches $\log(n)$.

Intra-List Diversity

Another interesting characteristic is Intra-List diversity. In general, we could say that a recommendation list is diversified if the items it contains are different. Usually, a recommender does not care about the similarity or dissimilarity of recommended items. However, let us suppose to design a recommender for a travel agency. If the different alternatives are similar, the recommender is not suggesting anything interesting to users. To improve user experience [352], we should provide a varied recommendation list. In this scenario, to measure diversity, we have to measure a degree of diversity between items in the same list [377].

Since diversity is an opposite signal to similarity, we can choose a similarity metric for this task. In literature, several strategies have been proposed: means, summations, minimum and maximum distances. An example is Intra-List Diversity (ILD) [429].

Expected Intra-List Diversity is a diversity metric which measure how much is diversificate the recommendation.

$$ILD_u@N = \frac{1}{2} \sum_{x_i \in L_u(n)} \sum_{x_j \in L_u(n)} 1 - sim(x_i, x_j) \quad (3.18)$$

$$ILD@N = \frac{1}{|U|} \sum_{u \in U} ILD_u@N \quad (3.19)$$

Finally, the diversification can also be computed considering the different intents of the user. In order to evaluate the diversification power in this case we could measure ERR-IA[95].

$$ERR - IA = \sum_{r=1}^n \frac{1}{r} \sum_t P(t|q) \prod_{i=1}^{r-1} (1 - R_i^t) R_r^t$$

where r is the position of an item i , t is the topic, $P(t|q)$ is the conditional probability of the topic given the query (user profiles in recommendation scenario), R_i is the probability of the relevance of the item and R_r is the probability of the relevance of the list of items from 1 to r . With this metric, the contribution of each item in the recommendation list is based on the relevance of documents ranked above it. The discount function then also depends on the relevance of previously ranked documents.

3.6.4 Novelty

The **novelty** of recommendations [377] could be defined as the presence in the recommendation list of completely unknown items for the user [220]. There is not a unified definition of novelty. Indeed recommender systems use different definitions of novelty. However, in general, it is considered as a measure of how many new items are recommended. A common approach is measuring the number of recommended items that come from the long tail. It is a signal of the ability of a system of

recommending items that never the user could have discovered. Hence, measuring novelty is crucial for a successful recommender.

A metric that measures the ability of the system of proposing long-tail items is Entropy-Based Novelty (EBN) [47]. Let us define again the recommendation list as $L_u(N)$. Entropy-Based Novelty can be computed as:

$$EBN_u@N = - \sum_{i \in L_u} p_i \cdot \log_2 p_i$$

in cui:

$$p_i = \frac{|\{u \in U \mid i \text{ è rilevante per l'utente } u\}|}{|U|}$$

With this *EBN* formulation, when $EBN_u@N$ decreases the novelty increases.

A metric that measures the ability of the system of proposing long-tail items is Entropy-Based Novelty (EBN) [47]. Let us define again the recommendation list as $L_u(N)$. Entropy-Based Novelty can be computed as: Another novelty metric is Expected Popularity Complement which corresponds to the expected number of relevant items that come from the long-tail. In this case, the binary relevance formulation of *EPC* [91] can be computed as:

$$EPC = C \sum_{i_k \in \mathbb{R}} disc(k) p(rel|i_k, u) (1 - p(seen|i_k)) \quad (3.20)$$

where $disc(k)$ is a discount function, $p(rel|i_k, u)$ is the relevancy of the item in the recommendation list and $(1 - p(seen|i_k))$ reflects a factor of item novelty. C is a normalizing constant, which stabilize the metric against unwanted biases. Two approaches are used in information retrieval to define $\frac{1}{C}$. The former defines $\frac{1}{C}$ as the maximum metric value obtainable by an ideal recommendation ranking as in nDCG or a-nDCG. The latter defines it as the expected browsing depth.

EFD, on the other hand, is a measure of the expected inverse collection frequency of relevant and seen items:

$$EFD_u@N = - \frac{1}{|Rec_u^N|} \sum_{i \in Rec_u^N} \log_2 p(i|seen)$$

The aforementioned metrics provide a measure of the ability of a system to recommend relevant long-tail items.

Part II

Feeding RSs with explicit knowledge

Chapter 4

Introduction

This chapter is focused on the exploitation of a formal representation of knowledge to feed Recommender Systems. We have already provided a broad overview of Semantic Web technologies and Recommender Systems techniques. Hidden in the overview, a careful reader may have read some hints on the approaches we are proposing. External knowledge can help propose more accurate recommendations. However, a fine representation of knowledge can highlight other similarities. This usually leads to improvements in terms of other dimensions like diversity and novelty. Moreover, the collaborative filtering algorithms usually deal with very sparse matrices. This lack of data can be compensated by external knowledge.

In the last decade, collaborative filtering approaches have shown their effectiveness in computing accurate recommendations starting from the user-item matrix. Unfortunately, due to their inner nature, collaborative algorithms show their limits when they deal with sparse matrices and, in these cases, encoding user preferences only through past ratings may lead to unsatisfactory recommendations. Hybrid approaches have been proposed to cope with this issue by exploiting side information about the items within the catalog. In the first line of research, we propose to

inject knowledge **from semantic graphs to matrices**. In practical terms, we propose to exploit past user ratings, and Linked Open Data to evaluate the relevance of every single feature within each user profile thus moving from a user-item to a user-feature matrix. Here, each value is a pair representing both the popularity of the feature in the user profile and its estimated rating. We then propose two computationally efficient content-based approaches and two hybrids, that make use of matrix factorization techniques to compute recommendations. The evaluation has been performed on three datasets referring to different domains (movies, music, and books) and experimental results show that the proposed methods outperform state of the art approaches in terms of accuracy, novelty, and diversity of results.

Providing relevant personalized recommendations for new users is one of the major challenges in recommender systems. This problem, known as the user *cold start* has been approached from different perspectives. In particular, cross-domain recommendation methods exploit data from source domains to address the lack of user preferences in a target domain. Most of the cross-domain approaches proposed so far follow the paradigm of collaborative filtering and avoid analyzing the contents of the items, which are usually highly heterogeneous in the cross-domain setting. Content-based filtering, however, has been successfully applied in domains where item content and metadata play a key role. Such domains are not limited to scenarios where items do have text contents (e.g., books, news articles, scientific papers, and web pages), and where text mining and information retrieval techniques are often used. Potential application domains include those where items have associated metadata, e.g., genres, directors and actors for movies, and music styles, composers and themes for songs. With the advent of the Semantic Web and its reference implementation Linked Data, a plethora of structured, interlinked metadata is available on the Web. These metadata represent a potential source of information to be exploited by content-based and hybrid filtering approaches. Motivated by the use of Linked Data for recommendation purposes, in the second line of research we present and evaluate a number of **matrix factorization models for cross-domain collaborative filtering that leverage metadata** as a bridge between items liked by users in different domains. We show that in case the underlying knowledge graph

connects items from different domains and then in situations that benefit from cross-domain information, our models can provide better recommendations to new users while keeping a good trade-off between recommendation accuracy and diversity.

Model-based approaches to recommendation can recommend items with a very high level of accuracy. Unfortunately, even when the model embeds content-based information, if we move to a latent space we miss references to the actual semantics of recommended items. Consequently, this makes non-trivial the interpretation of a recommendation process. In the third line of research, **we show how to initialize latent factors in Factorization Machines by using semantic features coming from a knowledge graph in order to train an interpretable model**. With our model, semantic features are injected into the learning process to retain the original informativeness of the items available in the dataset. The accuracy and effectiveness of the trained model have been tested using two well-known recommender systems datasets. By relying on the information encoded in the original knowledge graph, we have also evaluated the semantic accuracy and robustness for the knowledge-aware interpretability of the final model.

Preference representation and reasoning play a central role in supporting users with complex and multi-factorial decision processes. In fact, user tastes can be used to filter information and data in a personalized way, thus maximizing their expected utility. Over the years, many frameworks and languages have been proposed to deal with user preferences. Among them, one of the most prominent formalism to represent and reason with (qualitative) conditional preferences (CPs) are *conditional preference theories (CP-theories)*. In the fourth line of research, **we show how to combine CP-theories with Semantic Web technologies in order to encode in a standard SPARQL 1.1 query the semantics of a set of CP statements representing user preferences** by means of RDF triples that refer to a “preference” OWL ontology. In particular, here we focus on context-uniform conditional (cuc) acyclic CP-theories [395]. The framework that we propose allows a standard SPARQL client to query Linked Data datasets, and to order the results of such queries relative to a set of user preferences.

In *Recommender Systems* and, more broadly, in *Information Retrieval* scenarios,

the notion of relevance for the attributes of an item (or a document) plays a crucial role. As an example, all the items belonging to a recommended list are supposed to be of interest to the user because there is a similarity between the relevance of their attributes and those belonging to the items already enjoyed by the user in the past. Relevance measures such as *TF-IDF* or *BM25* are a representation of the informativeness of attributes in item description and are based exclusively on content-based information. In this investigation, we propose to enhance pure content-based relevance values for item attributes by exploiting collaborative information. **The idea is that of representing a vector of attributes whose weights encode both content-based and collaborative knowledge in a principled way.** To show the effectiveness of our proposal, we have tested it on three different datasets with respect to *state-of-the-art* algorithms in *Recommender Systems*.

Chapter 5

From semantic graphs to matrices

5.1 Introduction

Recent years have seen the flourishing of many and diverse recommendation techniques based on collaborative information encoded in the user-rating matrix. Factorization techniques have proven their effectiveness in improving the performance of recommendation engines and are implemented in many industrial and commercial systems [203, 44]. More recently, deep learning arrived as a new player in the field to develop powerful collaborative and hybrid approaches [103]. The core idea of collaborative filtering is to exploit the user-user connections through items and ratings to estimate user tastes and then predicting a list of items the user may be interested in. State-of-the-art algorithms can capture complex non-linear or latent factors-based relationships between users and items thus resulting more effective in all those scenarios where several users partially overlap their ratings or, in other words, the user-rating matrix is less sparse. To overcome the limits of pure collaborative approaches, hybrid ones [80] have been proposed that also encode side information, typically content-based, about the items. Indeed, whenever avail-

able, descriptions of the items can be used as a valuable source of information to augment the knowledge injected in and exploited by the system to compute a recommendation list of items. In this direction, an interesting class of recommender systems is the so called semantics-aware [116] where the information describing items goes beyond text and keywords and is represented by categorical/ontological data. Semantics-aware (SA) approaches have been widely adopted to integrate domain knowledge in a recommender system. SA approaches make use of ontologies or encyclopedic sources to encode and exploit such knowledge and in the last years many approaches have been proposed [260, 65, 230]. In fact, the domain-specific knowledge eases the process of interpreting documents and extracting relevant information from them. More recently, thanks to the Linking Open Data initiative, many structured data have become freely available to represent the content of items in different knowledge domains and they have been used to feed recommendation engines [279].

As a general remark, we can say that most of the recommendation algorithms available in the literature focus on computing the relevance of a set of items with reference to the user profile. Recommendation algorithms are designed around the computation of a relevance score for an item by evaluating its similarity with reference to other items. Features composing the description of an item, whatever the source, are not considered per se in the recommendation process but are usually exploited to evaluate the similarity between items or users. We believe that more attention might be paid to modeling the recommendation problem with a focus on recommending features rather than items. Expanding an item in its features give us a new set of explicit connections between items to be exploited with collaborative filtering algorithms. Finally, recommending items via feature recommendation may lead to an easier generation of explanation for the recommended list of items.

Unfortunately, moving from items to features is not that straight as in a forest of many features, most of them may result not relevant to a user. Moreover, once we design an algorithm able to compute a recommendation list of features, we have to go back to the items space, as the ultimate goal of a recommender systems is to suggest items to a user.

In this research line we present `FF` (for `Feature Factorization`), a *top-N* recommendation algorithm relying on user’s feature preferences and collaborative filtering information in the features space. The main goal of `FF` is to compute an ordered list of features preferred by the user and, starting from such list, to reassemble the relevance values of each returned feature to produce a *top-N* list of items to recommend. All the side information adopted by `FF` is retrieved from `DBpedia`, the cornerstone dataset of the Linked Data cloud. For each item in the user profile we retrieve its features by querying `DBpedia` thus having them as a set of entities. This avoids all problems related to synonymy and polysemy which usually occur when dealing with keyword-based features. By combining the popularity of a feature in the user profile and the ratings assigned to the items it is part of, for each user we compute a pair containing the relevance of the feature and its inferred rating. The resulting matrix in the user-feature space can be then manipulated via factorization techniques to compute, for each user, a ranked list of features which is in turn post-processed to produce the final list of recommendations.

Experimental evaluations of `FF` on two datasets related to the domains of books and music show its effectiveness in terms of accuracy, diversity and novelty of results in very sparse settings.

Please note that, although our approach shares some points with multi-criteria recommender systems [15], `Feature Factorization` differs from them because of the following main reasons:

- In `FF`, we do not assume any explicit rating to a specific feature of an item but we rather try to infer it starting from the global rating of the item;
- We target a top-N recommendation task and not a rating prediction one;

First, we formalize an initial version of `FF`, purely based on feature factorization. It is straightforward this version deserves its own experimental evaluation to assess if it is a feasible research direction. Then, starting from the initial version of `FF` we have worked on the following research questions:

RQ1 Are all the retrieved features needed to get results comparable with state-of-the-art algorithms?

- RQ2** How much does $\mathbb{F}\mathbb{F}$ performance depend on the number of items we are able to find a mapping for?
- RQ3** Is it possible to reduce the computational effort of $\mathbb{F}\mathbb{F}$ and keep at least the same results in terms of accuracy?
- RQ4** What is the influence of the quality for Linked Data information on the recommendation results?
- RQ5** How does $\mathbb{F}\mathbb{F}$ perform with reference to diversity and novelty of recommendations?

The remainder of the chapter is structured as follows. In the next section we report some related work on LOD-based and feature-based approaches to recommendation. We continue in Section 5.3 by introducing and describing $\mathbb{F}\mathbb{F}$ and its extensions. Experimental evaluations are presented in Section 5.4 while in Sections 5.4.1 and 5.4.2 we present and discuss the corresponding results. Conclusions and future work close the chapter.

5.2 Related Work

Several works have tried to build recommender systems by exploiting Linked Open Data (LOD) as side information for representing users or items, in addition to the user preferences usually collected through their ratings. Such approaches usually rely on *DBpedia*, a dataset which acts as a hub for most of the knowledge in the so-called LOD cloud. In the following, we review the recent literature on LOD-based recommender systems and, besides, since we propose an approach that leverages the relevance of single features in the user profile, we present related work in feature-based recommender systems.

LOD-based RS. Recommender Systems can exploit the knowledge coming from the LOD cloud for different tasks and in several ways. A detailed review of the literature, up to 2015, on Recommender Systems leveraging Linked Open Data is presented in [116]. Properties gathered from *DBpedia* may be used, e.g.,

to produce cross-domain recommendations [145], to build a multirelational graph for a graph-based recommender [282], or to generate effective natural-language recommendation explanations [269]. In [144] the authors propose three matrix factorization models for cross-domain recommendation. LOD are exploited to compute inter-domain item similarities, addressing the lack of data in the cold-start scenario. In [386] a web tool is provided, to encode the semantics of *Conditional preferences theories* [395]. Users preferences are expressed as RDF triples using a specific OWL ontology. These preferences are then used to produce meaningful recommendation lists via SPARQL queries. *ExpLOD* [269] is a novel tool able to generate natural language explanations exploiting information encoded in Linked Open Data encyclopedic datasets. In [266, 267] an extensive evaluation is performed, to establish whether LOD features are beneficial or not in using a graph-based recommender system, specifically a *PageRank with Priors* [168] algorithm. *Linked Data Semantic Distance (LDS)* measure has been proposed by [302] to compute a new distance metric to feed LOD-enabled recommender systems. In [280] authors summarize the main adopted techniques to feed a recommender system with LOD, explaining all the crucial steps to extract this kind of knowledge and to exploit it. In [268] popularity, collaborative, and content-based information is exploited with and without Linked Data to evaluate the impact on the usage of LOD features to produce recommendations.

When working with a Linked Data dataset, e.g., DBpedia, its properties may be used in very different ways:

1. to define semantic similarity measures for providing more accurate recommendations [297, 270, 259, 302, 281];
2. to deal with classical problems of recommender systems, as the *limited content analysis* or *cold-start*, e.g., by introducing new relevant features to improve item representations [68, 337], or to cope with the increasing data sparsity [266];
3. to improve the overall accuracy of a recommender system [288, 265], or to provide a good balance between different recommendation objectives, such

as accuracy and diversity [212, 266, 286].

Feature-based RS. The most widely adopted recommendation algorithms rely on the assumption that there are sufficient historical data for measuring similarity between items or users. Unfortunately, this assumption does not hold in several domains, where new items often lack ratings or comments, and where products that are less often purchased have fewer records of ratings. Several works attempt to analyze the user purchasing behavior based on item features and, consequently, user preferences towards specific features of items are then analyzed and exploited to provide potential accurate recommendations. In [393], products are represented using vectors of features, and a customer profile module computes the level of interest of the customer in product features as the ratio of features among the products purchased, and the product quantity purchased by that customer. Euclidean distance is then used to calculate the similarity between the customer and product profiles to recommend the *top-N* products. Similarly, in [165] a feature-based recommender system is presented for domains without enough historical data to effectively measure user or item similarities. The authors build the system based on the idea that *users who bought items with specific features also buy items with the same or similar features*. A similar approach is proposed in [274], in which effective strategies to incorporate item features for *top-N* recommender systems are developed. In graph-based recommender systems, an interesting work was proposed in [170], where recommendations are produced inferring user preferences, evaluating item-preferences and attribute-preferences. The paper points out the importance of the feature evaluation and a method is proposed which exploits explicit feature ratings named attributes. Another approach called Feature Preferences Matrix Factorization (FPMF) has been proposed in [272]. FPMF incorporates user feature preferences in a matrix factorization to predict user likes. Preferences on features are interpreted as user's attributes, which are taken into account when computing predictions on items by introducing additional latent factor vectors corresponding to attributes. It is worth noticing that none of the previous mentioned approaches exploits features coming from the Linked Open Data cloud.

Yet, approaches on feature-based RS need suitable datasets containing the ex-

explicit opinions of users on features of items, i.e., the ground truth for the preferences of users on items' features, in order to evaluate user models in terms of matching between estimated and true preferences on features. In [373], a Synthetic Data Generator is presented able to produce user-item and user/item-attribute datasets, while in [271], the authors used crowdsourcing to collect a dataset in the movie domain containing the explicit preferences of users on both items and their attributes. In [420] authors face the problem of splitting sparse data in time dimension. Instead of using pre-assumed distributions, they exploit item features extracted from textual reviews to realize a feature-level user model. In another interesting work, User-specific Feature-based item-Similarity Models (UFSM) [136] are proposed, in which multiple global item similarity functions are computed, together with user-specific weights for each function. The personalized weighted function is then used to produce recommendation lists. In [403] authors propose ReMF, a matrix factorization algorithm able to deal with hierarchically structured features. They also propose a new regularization method, named recursive regularization. Collaborative Knowledge Base Embedding framework (CKE) is proposed in [414]. The framework is designed to collaboratively learn latent representation starting from visual, textual, and structured knowledge. Knowledge base embedding is performed using three source-specific embedding techniques: TransR, SDAE, and SCAE respectively for structured, textual, and visual sources. The three representation are then combined in a single pair-wise ranking optimization function, to learn users and items collaborative representations. In [179] parallel Recurrent Neural Network (p-RNN) architecture is used to model sessions based on features extracted from pictures and text descriptions. Another interesting work is [374], in which 3-dimensional convolutional neural networks are used to produce recommendations. Authors chose this architecture to cope with item factual and categorical features within a single session. In details, this architecture is exploited to catch spatiotemporal patterns. In Boosted Factorization Machines (BoostFM) [406] a re-weighting scheme is proposed to inject boosting technique into a Factorization-Machine-based recommender system. Authors also propose two methods for pair-wise and list-wise learning to rank optimization. In [206] explicit items and features pair-wise pref-

erences are mapped to item comparisons in order to provide more accurate recommendations. DeepLIFT [346] technique is proposed to highlight the most important features in a neural network model. A different perspective is proposed by [114], in which the collaborative information is extracted and encoded in the form of new features to feed a content-based recommender system. Discrete Factorization Machines (DFMs) [240] are a recent proposal in which the high computational cost of feature-based and factorization machines algorithms is dramatically reduced by exploiting binarization of real-valued parameters. In order to avoid the quantization loss, the authors propose a different updating rule. The implicit short-term interests of the user are modeled through attentive neural networks in [322] with the aim of maximizing the time spent by users on the platform. Instead of polyphonic timbres of the song or topics, authors make use of embeddings of tags for the item occurred before the song to be predicted. In [318] a new graph embedding technique, RDF2Vec, is introduced, that make use of Weisfeiler-Lehman Subtree RDF Graph Kernels and graph walks to generate latent representations. More recently a Multi-Task Feature Learning approach (MKR) for Knowledge Graph Enhanced Recommendation [388] has been proposed that learn high-order interactions between items and entities in the knowledge graph. This approach takes advantage of knowledge graph embeddings and associates each item in a catalog to one or more entities in the knowledge graph. Other researchers build their approaches by leveraging user-defined features, such as tags, in order to produce better recommendations. We will not survey this line of research since, although tags seem to convey preferences over features, tagging an item does not necessarily mean that the user liked the attribute denoted by that tag. Finally, in our opinion feature-based Rs research line may have an overlap with multi-criteria RSs [15], that studies innovative approaches in collaborative recommendation by attempting to capture and model user preferences in a more comprehensive and nuanced manner by engaging multi-criteria ratings, in order to represent users' subjective preferences for various components of individual items.

5.3 Approach

5.3.1 Motivation

This line of research aims at investigating the role of feature rating and relevance in the item rating process. The main intuition behind FF is that items can be considered as a collection of features. Hence, when users rate an item, they are actually expressing their preference over the whole collection. In the evaluation of a movie, the user implicitly evaluates the director, the actors, the producer, the country in which the movie is set. Each feature has its own rating and a relevance degree, hence a Recommender System should consider these factors. The item rating action can be then interpreted as an attempt to choose an overall rate for the entire set of features. Our assumption is that if we want to discover the contribution of each single feature in the evaluation, first of all, we need to unpack each item in its composing features. Then, by combining the overall popularity of each feature in the user profile (feature relevance) and the rating assigned to items containing that feature we may estimate the implicit rating the user is giving to that specific feature.

In our model the user profile is not just a set of $\langle item, rating \rangle$ pairs but it contains information about the relevance (popularity) of each feature composing the rated items and its estimated rating. It is then represented as a set of triples $\langle feature, relevance, rating \rangle$. In the following, we see principled methods to estimate both the user-feature rating and the user-feature relevance, and then we move to the recommendation problem using the features that compose the user profile.

5.3.2 Data Model

For a better understanding of the data we use to reshape the user profile in terms of knowledge-aware features, we first introduce the multidimensional graph we have used to build them. As we can see from Fig. 5.1 the user profile is built by considering information coming from both the `Users-Items` matrix and from `DBpedia` as external knowledge source modeled as `Linked Data`. Now we see how the graph-based nature of this latter is exploited to identify features used to represent

items. Linked Data are represented as RDF labeled oriented graphs and their data model is based on the notion of triple $\langle \textit{subject}, \textit{predicate}, \textit{object} \rangle$ where *predicate* represents the relation connecting the two entities *subject* and *object*. With reference to Figure 5.1, we have that an item i in the catalog of a recommendation engine may represent the subject of a triple $\langle i, p, e \rangle \in \text{DBpedia}$. As an example, in the movie domain we have triples like $\langle \text{Matrix}, \text{starring}, \text{Keanu_Reeves} \rangle$, $\langle \text{Matrix}, \text{director}, \text{The_Wachowskis} \rangle$. It is noteworthy that the same pair *predicate*, *object* may occur for different items as for $\langle \text{Point_Break}, \text{starring}, \text{Keanu_Reeves} \rangle$ and, at the same time, the same entity may appear as object in diverse triples as in the case of $\langle \text{V_for_Vendetta}, \text{producer}, \text{The_Wachowskis} \rangle$ where the entity `The_Wachowskis` is connected to the subject via `producer` instead of `director` as for the previous triple. In order to catch the different knowledge encoded in the use of the same entity as object in triples with diverse predicates, in our model, we consider the chain $p \circ e$ as a feature associated to the item i which in turn represents the subject of the corresponding triple. In the model we propose, each item in the user profile is associated with a *relevance function* denoted with $\rho^{ui}(\cdot)$. Its value represents an estimation of how important is a particular item to the user u . Analogously, we have a value associated to each feature in the profile computed via the function $\rho^{uf}(\cdot)$ that for each feature represented by the chain $p \circ e$ returns a value representing the relevance of that feature in the user profile. As we have pointed out in Section 5.3.1, we assume there is a relation between the two relevance values which should be reflected in the mathematical formulation of the two functions $\rho^{ui}(\cdot)$ and $\rho^{uf}(\cdot)$. Actually, each feature is also associated with a *rating* $r^{uf}(\cdot)$ which is inferred by considering the rating of all the items containing the specific feature.

5.3.3 Problem Formulation

By considering the data associated to the user profile we can move from a rating matrix connecting user and items to a user-feature matrix where each value is represented by the pair $\langle \rho^{uf}(\cdot), r^{uf}(\cdot) \rangle$. In other words, we may consider two user-feature matrices: the one \mathcal{P} containing relevance values $\rho^{uf}(\cdot)$, the other \mathcal{R} the inferred

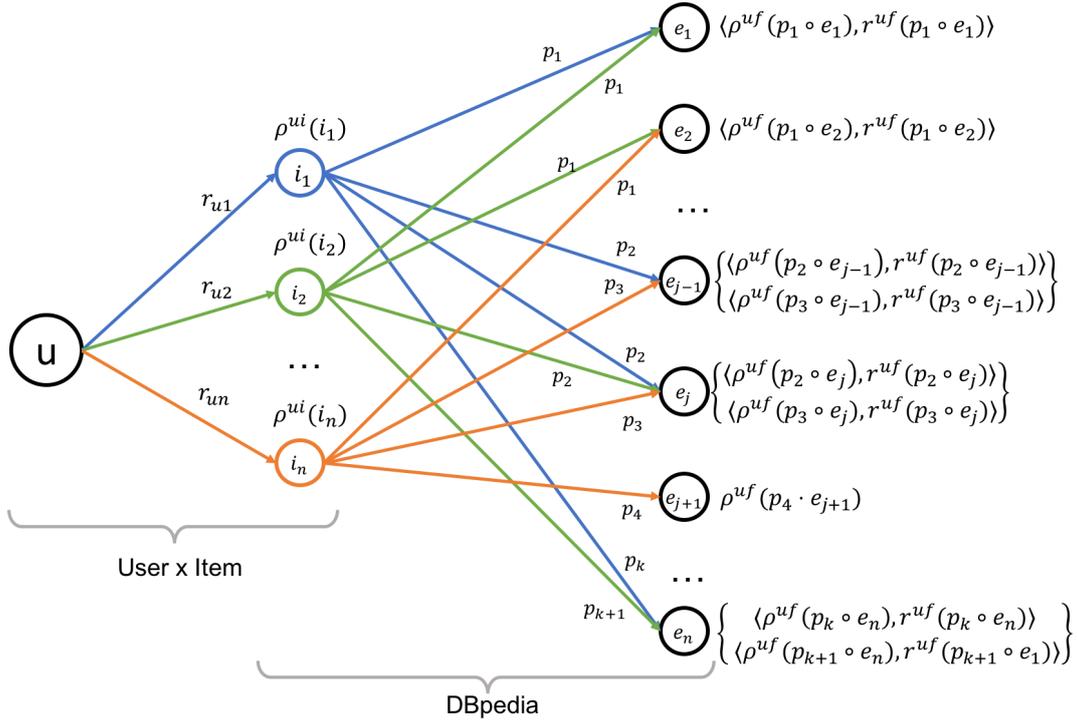


Figure 5.1: A graph-based representation of the data behind the computation of the user profile.

ratings $r^{uf}(\cdot)$. Given I representing the catalog of items and I^u containing the items experienced by the user u , for the sake of compactness, in the formulation of the problem we now introduce the following sets:

$$F^i = \{p \circ e \mid \langle i, p, e \rangle \in \text{DBpedia} \wedge i \in I\}$$

$$F^u = \{p \circ e \mid \langle i, p, e \rangle \in \text{DBpedia} \wedge i \in I^u\}$$

$$F = \bigcup_{i \in I} F^i$$

$$I^{uf}(p \circ e) = \{i \mid i \in I^u \wedge p \circ e \in F^i\}$$

In FF, the relevance of a feature $p \circ e$ is computed as its probability of belonging to I^u . More formally we have:

$$\rho^{uf}(p \circ e) = \frac{|I^{uf}(p \circ e)|}{|I^u|}$$

The idea behind this computation is quite straight: the more a feature is connected to the items in the user profile (i.e., the more a feature is popular), the higher is its relevance for the user. Once we have computed the relevance of all the features in the user profile, we can move to the computation of the relevance for the items $i \in I^u$. It can be computed as the normalized summation of the relevance for all the features it is composed by. In formulas, we have:

$$\rho^{ui}(i) = \frac{\sum_{p \circ e \in F^i} \rho^{uf}(p \circ e)}{|F^i|}$$

Please note that $\rho^{ui}(i)$ is not defined for $i \notin I^u$. As discussed before, we see that the relevance of an item is influenced by the relevance of the features it is composed by. In terms of the final user-feature matrices we want to compute, the relevance of an item is used to estimate the rating associated to the features occurring in the items already rated by the user. Given a feature $p \circ e$, the computation of $r^{uf}(p \circ e)$ exploits both the rating and the relevance of each item $i \in I^u$ containing $p \circ e$.

$$r^{uf}(p \circ e) = \frac{\sum_{i \in I^{uf}(p \circ e)} r^{ui} \cdot \rho^{ui}(i)}{\sum_{i \in I^{uf}(p \circ e)} \rho^{ui}(i)}$$

To completely move from a user-item space to a user-feature one, we combine \mathcal{P} and \mathcal{R} by introducing a new matrix \mathcal{S} obtained by element-wise multiplication of \mathcal{P} and \mathcal{R} . Each element of \mathcal{S} is then computed as:

$$s^{uf}(p \circ e) = \rho^{uf}(p \circ e) \cdot r^{uf}(p \circ e)$$

Because of the spreading mechanism exploited by the method, we name $s^{uf}(\cdot)$ **feature spreading relevance** (f_{sr}). Our intuition is that features with the higher values of f_{sr} are the most representative for the users and are those influencing their ratings. To test our intuition, we have defined different similarity measures based on $s^{uf}(\cdot)$ and we then have tested them in a recommendation scenario.

5.3.4 Pure FF

The profiles we have built contain only the features the user has met before, but usually the number of those features is dramatically smaller than the overall number of

features and this results in \mathcal{P} and \mathcal{R} being very sparse. To complete the information they contain, we compute, via Biased Matrix Factorization, the missing values $\hat{\rho}^{uf}(p \circ e)$ for \mathcal{P} and $\hat{r}^{uf}(p \circ e)$ for \mathcal{R} . We run matrix factorization independently on \mathcal{P} and \mathcal{R} . Biased Matrix Factorization is a Matrix Factorization model that minimizes RMSE using stochastic gradient descent [224]. It computes user's and item's biases to improve the estimation of the predicted value. Biased Matrix Factorization represents a state-of-the-art algorithm in rating prediction task. $\hat{\rho}^{uf}(p \circ e)$ and $\hat{r}^{uf}(p \circ e)$ represent the predicted relevance and the predicted rating for all those features not belonging to any of the items in I_u . As the resulting matrices contain both content-based and collaborative information (due to the Matrix Factorization), we refer to them as *hybrid profiles*.

With the *hybrid profile* we can estimate a ranked list for all the remaining items within the collection. In fact, the ranking of an item in the list is computed by considering the rating of the features belonging to the item and their relevance.

$$\begin{aligned} \hat{r}^{ui}(i) = & \sum_{((i,p,e) \in \text{DBpedia}) \wedge (i \in I_u)} \rho^{uf}(p \circ e) \cdot r^{uf}(p \circ e) + \\ & + \sum_{((i,p,e) \in \text{DBpedia}) \wedge (i \notin I_u)} \hat{\rho}^{uf}(p \circ e) \cdot \hat{r}^{uf}(p \circ e) \end{aligned} \quad (5.1)$$

It is important to point out that these estimations do not correspond to an actual rating. Instead, the goal is trying to preserve a correct item ranking.

Post-filtering

To improve the results of the final recommendation process, we propose a post-filtering step for reducing the number of features considered while computing the final rank. The proposed filtering springs from the following observations:

- Not all the features items are relevant in the computation of the ranking for an item. All those features whose rating results low just introduce noise in the final values we compute.
- Feature ranking and relevance values evaluated via pure content-based approaches, i.e., before the Matrix Factorization, have a different influence if

compared with the collaborative ones representing latent factors computed after the Matrix Factorization.

To lower the number of features involved in the computation, and produce recommendations based only on the best ratings of the estimated features, we propose a filter that operates on directly estimated features (content-based), and estimated features coming from collaborative computation. Then, we introduce two thresholds α and β that act as filters on the feature rating values respectively in the content-based and the collaborative cases. Hence, Equation (5.1) is slightly modified:

$$\begin{aligned} \hat{r}^{ui}(i) = & \sum_{((i,p,e) \in \text{DBpedia}) \wedge (i \in I_u) \wedge r^{uf}(p \circ e) > \alpha} \rho^{uf}(p \circ e) \cdot r^{uf}(p \circ e) + \\ & + \sum_{((i,p,e) \in \text{DBpedia}) \wedge (i \notin I_u) \wedge \hat{r}^{uf}(p \circ e) > \beta} \hat{\rho}^{uf}(p \circ e) \cdot \hat{r}^{uf}(p \circ e) \end{aligned} \quad (5.2)$$

5.3.5 Jaccard-fsr

One of the simplest way to compute the similarity between two objects (items or users) is by Jaccard similarity. It is computed by taking into account the common elements between the two objects. In our case, if we consider a pure content-based setting, we may compute the similarity between a user and an item by looking at the common features available within the user profile u and the item description i . Actually, as the features with a higher value of `fsr` should be more relevant for u , then we may select only those having the highest value of $s^{uf}(\cdot)$ for u and check only these latter against i instead of the whole user profile. To identify these valuable features, we have first computed $\mu(u)$ representing the average value of $s^{uf}(\cdot)$ for u and we eventually select only those features with a value of `fsr` higher than $\mu(u)$. More formally, based on the following definitions:

$$\mu(u) = \frac{\sum_{p \circ e \in F^u} s^{uf}(p \circ e)}{|F^u|} \quad (5.3)$$

$$Rel^{u+} = \{p \circ e \mid p \circ e \in F^u \wedge s^{uf}(p \circ e) > \mu(u)\}$$

Then, we define a Jaccard-`fsr` similarity value between u and $i \notin I^u$ as:

$$J\text{-fsr}(u, i) = \frac{|Rel^{u+} \cap F^i|}{|Rel^{u+} \cup F^i|} \quad (5.4)$$

5.3.6 Content Based- and Hybrid-fsr

Given the values of $s^{uf}(\cdot)$, we may compute a score associated to an item i not belonging to the original user profile by summing the values of relevant fsr in Rel^{u+} for those features belonging to the description of i . As we are adopting a pure content based approach we refer to this measure as Content-based feature spreading relevance ($CB\text{-fsr}$).

$$CB\text{-fsr}(u, i) = \sum_{p \circ e \in Rel^{u+} \cap F^i} s^{uf}(p \circ e) \quad (5.5)$$

In the previous equation, as well as in Equation (5.4), we see that we consider only the relevant features in the users profile, or, in other words, the features they like.

In fact, this model makes use only of positive feedback. Even though this is fine-tuned feature-specific feedback, it still is positive feedback.

Instead, many other approaches take advantage of both positive and negative feedback. Some of them can rely on explicit negative feedback provided by users. Others adopt the *missing not at random* principle to set as negative all the unseen items. In this scenario, the situation is dramatically different, because the estimated rating is related to a feature, instead of an item.

To evaluate negative feedback for specific features, our idea is to exploit a Matrix Factorization (MF) model. These kinds of algorithms are often designed as iterative algorithms, with the aim of minimizing the rating prediction error. In a Matrix Factorization model, entities (users, and items) are represented by means of a bias and a vector of latent factors of predefined size D . In our scenario, we want to substitute items with features. For each user $u \in U$, and each feature $f \in F$ we build a binary vector $\mathbf{r}^{uf}(p \circ e) \in \mathbb{R}^{1 \times |F|}$, representing the estimated interaction between u and f . In this modeling, $\mathbf{r}^{uf}(p \circ e)$ contains only two 1 values corresponding to u and f if $f \in F^u$, while all the other values are set to 0. The prediction formula can be defined for each uf pair as:

$$r_H^{uf}(p \circ e) = w_0 + w_u + w_f + \sum_{k=1}^D v_{(u,k)} \cdot v_{(f,k)} \quad (5.6)$$

where the parameters to learn are: w_0 representing the global bias; w_u and w_f repre-

senting the individual biases for u , and f ; the pair $v_{(u,k)}$ and $v_{(f,k)}$ in $\sum_{k=1}^D v_{(u,k)} \cdot v_{(f,k)}$ measuring the strength of the interaction between each pair u and f . The number of latent factors is represented by D . This value is usually selected at design time when implementing the MF algorithm.

Matrix Factorization can be easily trained in order to minimize the rating prediction error via gradient descent methods, alternating least-squares (ALS) and MCMC. Usually, when MF is optimized via gradient descent methods, the prediction formula can be defined as:

$$\delta^{uf} = w_0 + w_u + w_f + \sum_{k=1}^D v_{(u,k)} \cdot v_{(f,k)} \quad (5.7)$$

$$\hat{r}_H^{uf} = r_{min} + \frac{1}{1 + e^{-\delta^{uf}}} \cdot range \quad (5.8)$$

where r_{min} is the minimum rating of all the considered transactions, whereas $range$ is the rating range size. The rating prediction problem is bound to the $(0, 1)$ range of the sigmoid function. The consequent cost function is:

$$J(\theta) = \sum_{u \in U} \sum_{f \in F^u} (r^{uf} - \hat{r}^{uf}) + \lambda_{\Theta} \cdot \|\Theta\|^2 \quad (5.9)$$

where r_H^{uf} is the feature rating estimated by our method, whereas \hat{r}_H^{uf} is the rating estimated by the matrix factorization model. We can introduce $\sigma(\cdot)$ as a sigmoid function, and the update step can be defined as:

$$\begin{aligned} \Theta \leftarrow & \Theta + \alpha \cdot ((r^{uf} - \hat{r}^{uf}) \cdot \sigma(\delta^{uf}) \cdot (1 - \sigma(\delta^{uf}))) \cdot \\ & \cdot range \cdot \frac{\partial}{\partial \Theta} \hat{r}^{uf} + \lambda \cdot \Theta \end{aligned} \quad (5.10)$$

To update the factorized parameters, partial derivatives can be computed as:

$$\frac{\partial}{\partial \Theta} \hat{r}^{uf} = \begin{cases} 1, & \text{if } \theta = w_f, \\ v_{(u,k)}, & \text{if } \theta = v_{(f,k)}, \\ v_{(f,k)}, & \text{if } \theta = v_{(u,k)}, \\ 0, & \text{otherwise} \end{cases} \quad (5.11)$$

Using Equation (5.11) in Equation (5.10) the model parameters can be iteratively updated to minimize the rating prediction error.

As the resulting matrix contains both content-based and collaborative information (due to the Matrix Factorization), we refer to it as *hybrid profile*. To exploit it, we set the lowest K features (Rel^{u-}) as negative examples with a value \hat{b}^{uf} :

$$\hat{b}^{uf}(p \circ e) = -1 + \hat{s}^{uf}(p \circ e) \quad \{p \circ e \in Rel^{u-}\} \quad (5.12)$$

where \hat{s}^{uf} is the estimated feature f value for user u with the Matrix Factorization model. With both content-based, and collaborative filtering information we may compute a score associated to item i . We refer to this measure as *Hybrid-fsr*:

$$\begin{aligned} Hybrid\text{-}fsr(u, i) &= \sum_{p \circ e \in Rel^{u+} \cap F^i} s^{uf}(p \circ e) \\ &+ \sum_{p \circ e \in Rel^{u-} \cap F^i} \hat{b}^{uf}(p \circ e) \end{aligned} \quad (5.13)$$

5.3.7 frsCBF

Values in \mathcal{S} represent an indicator on the importance of a feature for a specific user. Based on this information we may compute recommendation lists for items $\bar{i} \notin I^u$ by exploiting $s^{uf}(p \circ e)$ values. We may compare features in F^u with those in $F^{\bar{i}}$ in different ways. The most intuitive way is that of computing a score for \bar{i} is by summing of the values associated to the features it is composed by with a normalization factor.

$$\hat{r}^u(\bar{i}) = \frac{\sum_{p \circ e \in F^u \cap F^{\bar{i}}} s^{uf}(p \circ e)}{\sqrt{\sum_{p \circ e \in F^u \cap F^{\bar{i}}} (s^{uf}(p \circ e))^2}} \quad (5.14)$$

We can see that the value computed by Equation (5.14) corresponds to the cosine similarity between the user profile vector containing $s^{uf}(p \circ e)$ if $p \circ e \in F^u$ and 0 vice versa, and the corresponding binary vector for \bar{i} containing 1 if $p \circ e \in F^{\bar{i}}$ and 0 vice versa.

5.4 Experimental Evaluation

5.4.1 Experiments for Pure FF

In this section the experimental evaluation settings and the metrics used to evaluate the proposed algorithm are presented. We have evaluated the algorithms in terms of ranking accuracy for *top-N* recommendations. The evaluation has been carried out on two datasets, `LibraryThing` and `Last.fm` belonging respectively to the domains of books and music.

Datasets description and Pre-Processing

To alleviate the popularity bias from the evaluation results we have removed the 1% most popular items [108]. Moreover, we have removed users with a number of ratings smaller than five as we want to evaluate the algorithms in a non cold-start setting. The `LibraryThing` dataset contains 7,564 users, 39,515 items, and 797,299 ratings. The minimum, mean and maximum number of ratings for user in the dataset are 20, 63, 3,018, respectively. `Last.fm` contains 1,892 users, 17,632 items and 92,834 ratings. In `LibraryThing`, ratings are distributed over a 1 – 10 scale. In `Last.fm` the rating is the number of times a song has been played, hence that number has been rescaled for each user in a 1 – 10 scale. Table 5.1 shows some statistics of the datasets subsets considering only the items mapped to `DBpedia` (using publicly available mappings [288]) after the pre-processing step. In case a mapping does not exist, a simple placeholder feature is used, that inherits the corresponding item values in terms of rating and relevance.

Table 5.1 also reports the sparsity values both for users-items and users-features matrices.

Evaluation protocol and experiment setting

To evaluate Pure FF we use the *all unrated items* [356] evaluation protocol, in which the ability to choose the correct set of items to propose to the users is favorite despite of the local ranking ability (*rated test-items* evaluation protocol). In

| | | | | |
|---------------------|---------|------------|-----------|--------------|
| LibraryThing | # users | # items | # ratings | sparsity (%) |
| user-item space | 6,909 | 12,656 | 248,589 | 99.7157 |
| | # users | # features | # ratings | sparsity (%) |
| user-feature space | 6,909 | 141,531 | 8,680,619 | 99.11226 |
| Last.fm | # users | # items | # ratings | sparsity (%) |
| user-item space | 1,866 | 8,502 | 39,557 | 99.75066 |
| | # users | # features | # ratings | sparsity (%) |
| user-feature space | 1,866 | 274,523 | 4,989,281 | 99.02603 |

Table 5.1: Datasets Statistics.

all unrated items the recommendation list is produced using as candidate list the Cartesian product between users and item minus the items the user experimented in the training set. Evaluation has been conducted using a hold-out 80 – 20 splitting, in which 20% of the ratings are retained as test set. We have evaluated the accuracy of our approach by computing Precision ($P@N$), Recall ($R@N$), and nDCG ($nDCG@N$).

Baselines. In the experimental evaluation we have compared FF with the popularity baseline (PopRank) and, as we rely on Matrix Factorization, the well-known matrix factorization algorithm BPRMF [312] both in its pure collaborative version and in the hybrid one considering side information BPRMF+SI. We have also included PopRank as it is acknowledged that popularity ranking can show good performance and it is an important baseline to compare against [108]. To produce recommendation lists from these well-known algorithms we have used the *MyMediaLite*¹ implementation [154]. As for the selection of α and β parameters needed in Equation (5.2), in these experiments we have kept a conservative approach and we have set respectively α to the mean μ of the rated items, and β to the mean μ plus the standard deviation σ . Clearly, these values are not the optimal ones and the performance could be improved by a cross-validation setting of these parameters.

¹<http://www.mymedialite.net/>

Experimental Results

Tables 5.2 and 5.3 show the performance of FF compared with the competing algorithms described in Section 5.4.1. In **bold** we mark the best result for each metric. All the evaluations have been performed by using the same protocols as implemented in RankSys² library [91].

In Table 5.2 we show the evaluation results on `LibraryThing` dataset with a threshold set to 7/10 in a *Top* – 10 recommendation list. The ranking accuracy performance, measured through nDCG, Precision and Recall shows that `Feature Factorization` performs better than the competing algorithms. In details, `Pure FF` performs 4 to 6 times better than `BPRMF`, the second best accurate algorithm.

| Alg | P@N | R@N | nDCG@N |
|-----------------|----------------|----------------|----------------|
| Pure FF | 0.03251 | 0.06576 | 0.06129 |
| BPRMF | 0.00837 | 0.01280 | 0.01020 |
| BPRMF+SI | 0.00777 | 0.01325 | 0.01007 |
| PopRank | 0.00023 | 0.00095 | 0.00044 |

Table 5.2: Comparative results on `LibraryThing` dataset, Top-10 recommendation list and relevance threshold of 7/10.

As the rescaling operation in `Last.fm` affects the values of the items in the test set, we have decided to evaluate considering all the items in test set as relevant (i.e., setting the relevance threshold to 0). Table 5.3 shows ranking accuracy evaluation results on `Last.fm` dataset with a threshold of 0/10 for a *Top* – 10 recommendation list. For precision metric the best performing algorithm is `Pure FF` that performs 4 times better than `BPRMF`. For nDCG, `Feature Factorization` performs at least 5 times better than the competing algorithms. The differences about accuracy metrics between `Pure FF` and the other algorithms are statistically significant according to the Student’s paired *t-test* with $p < 0.001$ for every cases.

²<https://github.com/RankSys/RankSys>

| Alg | P@N | R@N | nDCG@N |
|-----------------|----------------|----------------|----------------|
| Pure FF | 0.01543 | 0.02701 | 0.02330 |
| BPRMF | 0.00348 | 0.00902 | 0.00495 |
| BPRMF+SI | 0.00032 | 0.00073 | 0.00028 |
| PopRank | 0.00027 | 0.00089 | 0.00021 |

Table 5.3: Comparative results on Last.fm dataset, Top-10 recommendation list and no relevance threshold.

5.4.2 Experimental Evaluation for FF extensions

The evaluation has been carried out on three well-known datasets, LibraryThing, MovieLens, and Last.fm belonging respectively to the domains of books, movies, and music.

Datasets description and Pre-Processing

Again, to alleviate the popularity bias from the evaluation results, we have removed the 1% most popular items [108]. Moreover, we have removed users with a number of ratings smaller than five as we want to evaluate the algorithms in a non-cold-start setting. The MovieLens dataset contains 1,000,209 ratings over of approximately 3,900 movies made by 6,040 users on the MovieLens platform. Each rating is expressed on a 1 – 5 scale. The LibraryThing dataset contains 7,564 users, 39,515 items and 797,299 ratings. The minimum, mean and maximum number of ratings per user in the dataset are 20, 63 and 3,018, respectively. Last.fm contains 1,892 users, 17,632 items and 92,834 ratings. In LibraryThing, ratings are distributed over a 1 – 10 scale. In Last.fm the rating is the number of times a song has been played, hence that number has been rescaled for each user on a 1 – 10 scale. Table 5.4 shows some statistics of the datasets subsets considering only the items mapped to DBpedia (using publicly available mappings [288]) after the pre-processing step. Table 5.4 also reports the sparsity values both for Users-Items and Users-Features matrices.

| | | | | |
|--|---------|------------|------------|--------------|
| LibraryThing user-item space | # users | # items | # ratings | sparsity (%) |
| | 6,909 | 12,656 | 248,589 | 99.7157 |
| user-feature space | # users | # features | # ratings | sparsity (%) |
| | 6,909 | 141,531 | 8,680,619 | 99.11226 |
| Last.fm user-item space | # users | # items | # ratings | sparsity (%) |
| | 1,866 | 8,502 | 39,557 | 99.75066 |
| user-feature space | # users | # features | # ratings | sparsity (%) |
| | 1,866 | 274,523 | 4,989,281 | 99.02603 |
| MovieLens user-item space | # users | # items | # ratings | sparsity (%) |
| | 6,040 | 3,171 | 689,867 | 96.3981 |
| user-feature space | # users | # features | # ratings | sparsity (%) |
| | 6,040 | 100,845 | 56,732,878 | 90.68584 |

Table 5.4: Datasets Statistics.

Evaluation protocol and experiment setting

To evaluate \mathbb{FF} , even here, we adopt the *all unrated items* [356] evaluation protocol. In *all unrated items* the recommendation list is produced using as candidates list the Cartesian product between users and item minus the items the user experimented in the training set. The evaluation has been conducted using a hold-out 80 – 20 splitting, in which 20% of the ratings are retained as the test set.

We chose to evaluate the approaches through accuracy, novelty and aggregate diversity metrics. The *top-N* recommendation accuracy metrics we used are Precision ($P@N$), Recall ($R@N$) and nDCG ($nDCG@N$). As for novelty, in the last years, several metrics have been proposed, which measure the ability of a system to recommend items that are not very popular within the catalog and that usually belong to the long tail [220]. In this investigation, *EPC* (Expected Popularity Complement) is used for novelty evaluation. *EPC* measures how much a system is capable to propose long tail items to the users. In this case, the binary relevance formulation of *EPC* has been adopted [91]. Diversity has been measured through *catalog coverage* (aggregate diversity in *top-N* list), *Gini index* and *Shannon entropy*. The *catalog coverage*, also called aggregate diversity or the diversity-in-top-N ($ADiv@N$), formulated in [14], denotes the overall number of different items recommended within

all recommendation lists. It denotes the propensity of a system to recommend always the same items. *Gini index (Gini)* and Shannon entropy (*SE*) are two different metrics used to measure the distributional inequality [91]. The evaluation has been performed considering *Top – 5* and *Top – 10* recommendations for `MovieLens`, `Last.fm` and `LibraryThing` datasets and two different thresholds for deeming items as relevant: 7 and 9. Accuracy and novelty metrics have been computed on a per-user basis and the results have been averaged.

Baselines

In the experimental evaluation, we have compared `FF` with the popularity baseline and three well-known Matrix Factorization algorithms. These latter have been tested both in their pure-collaborative variant and in their hybrid ones by also considering side information coming with the dataset (`Tag`), i.e., tags associated to the items, or considering `Linked Data` side information extracted from `DBpedia` (`LOD`). Hybrid variants of matrix factorization algorithms have been realized feeding them with new transactions corresponding to items-features associations, as suggested in [274]. The competing algorithms we selected are then:

- `PopRank`, a non-personalized algorithm that produces the same recommendation list for all the users. This list is computed measuring the items' popularity and ordering them in descending order. It is acknowledged that popularity ranking can show good performance and it is an important baseline to compare against [108].
- `Bayesian personalized ranking – Matrix Factorization (BPRMF)`, a matrix factorization algorithm that exploits the `Bayesian Personalized Ranking` criterion[312] to minimize the ranking errors.
- `Soft Margin Ranking – Matrix Factorization (SMRMF)`, a matrix factorization model that uses stochastic gradient descent to optimize a soft margin [392].

- **Weighted Matrix Factorization (WRMF)** a weighted matrix factorization model that makes use of alternate least squares [184] for optimization.

To produce recommendation lists from these well-known algorithms we have used the *MyMediaLite*³ implementation [154]. The experimental setting is designed to evaluate the performance of the proposed approach against different matrix factorization algorithms, varying datasets, and side-information fed to the recommender, keeping fixed the models' hyperparameters. For this reason, the parameters are set considering the values suggested by the authors in their original papers.

Experimental Results

Tables 5.5 to 5.18 show the performance of FF^4 compared with the competing algorithms described in Section 5.4.2. In **bold** we mark the best result for each metric while we underline the second best result. The differences about accuracy metrics between FF and the other algorithms are statistically significant according to the Student's paired *t-test* [351] with $p \ll 0.001$.

Ranking accuracy evaluation The first three metrics in Tables 5.5 to 5.18 are devoted to accuracy evaluation of FF with respect to competing baselines. As it could be appreciated in each experiment at least one variant of FF outperforms the baselines. However, it should be noticed that, with respect to results from the literature, the removal of the 1% heavily affects the collaborative algorithms. For this reason, it is important to underline that *Hybrid-fsr* also suffers from this removal. It is worth to notice that this performance drop is alleviated by its hybrid nature. Despite this consideration, collaborative algorithms in *MovieLens* experiments still show high performance, with *BPRMF* as the winner among the competitors. Nevertheless, *LibraryThing* and *Last.fm* showed that the algorithm that best

³<http://www.mymedialite.net/>

⁴Implementation available at: <https://github.com/sisinflab/Features-Factorization>

deals with increasing sparsity is WRMF. Among the proposed variants the discussion is more complex because from experiments it is clear that they behave as they were completely different algorithms. In details, there could be noticed two different dimensions for the results analysis: quality of meta-data, and statistics of datasets. The latter seems to impose the magnitude of the metrics value, which is comparable to competing algorithms. The former heavily affects the results and we consider it so important to be detailed in the next subsection to explain the details of the measured differences.

RQ1: Are all the retrieved features needed to get results comparable with state-of-the-art algorithms?

This question is motivated by a rich literature in feature selection research field. In our previous experiment (regarding Pure FF), we have decided to take advantage of all available features. However, since it was a collaborative filtering algorithm, we have found really hard to understand if all the features have been effectively beneficial. There have been too many parameters that could affect the results, and too many assumptions to be made. To establish whether they have been beneficial or not, we have decided to define a simple a clear content-based variant of FF based on cosine similarity and compare it against a Jaccard-based variant. In this variant, a threshold been defined to consider only relevant features. The choice of the threshold could be unfair, thus we have decided to adopt the mean of features values as a reasonable threshold (see Equation (5.3)). The experiments clearly show that the performance of the two variants depends on the specific dataset and items' descriptions. In `MovieLens` we can observe rich descriptions for almost all the considered items, it is therefore not surprising that the winner is `CB-fsr`, in which all features are involved. `LibraryThing` shows that the winner is `Jaccard-fsr Tag`. This leads to two considerations: 1) item descriptions in `KB` are not as effective (for recommendation task) as dataset tags are; 2) Among the considered tags, performing a feature selection is beneficial. Even in `Last.fm`, `LOD` features are worse than tags, and it could be noticed that items descriptions may vary from very rich ones to items with less than 5 features.

RQ2: How much does FF performance depend on the number of items we are able to find a mapping for?

This is a very common question that is usually posed when considering the injection of `Linked Open Data` in a recommendation scenario. Despite the obvious fact that collecting high-quality users-based tags using crowdsourcing may result in a difficult task, we have decided to perform a series of specific experiments to answer the question. To test how much the lack of some items affects the overall recommendation results, we have dropped half of the `Hybrid-fsr` items repeatedly until we have reached one eighth of the original number of items. It is possible to appreciate the number of considered items in the first column of each table, near each `Hybrid-fsr` experiment (e.g., `Hybrid-fsr 3, 171, 1,586, 793, 396` for `MovieLens`). For each case, the items have been dropped with uniform distribution, and the model has been completely retrained. In `MovieLens`, the performance of `Hybrid-fsr` with one-half items shows results that are still comparable against the competing collaborative filtering algorithms. Moreover, the results of one-eighth items experiments show that the decrease w.r.t. original performance is only of two-thirds of the original value. In `LibraryThing`, the decrease is more evident, but its hybrid nature makes `Hybrid-fsr` still comparable with other algorithms even though the number of considered items is reduced to one eighth of the original number. In `Last.fm`, the decrease in performance is even stronger. However, this may be due to the specific dataset, because the same behavior can be observed on all the baselines, but `WRMF`.

RQ3: Is it possible to reduce the computational effort of FF and keep at least the same results in terms of accuracy?

The adoption of a model that takes into account matrix factorization in a much bigger matrix (`Users-Features` matrix) could be considered a questionable option in very big datasets. For this reason, we have decided to create the pure content-based variants of `FF`. Since we have created two different content-based variants for the aforementioned reasons, we consider useful to discuss the results of all the three variants together. Moreover, since `Hybrid-fsr` has been designed to work

mainly with LOD side-information, we consider the LOD versions of the variants. In `MovieLens` we may notice rich descriptions and highly collaboratively-connected users and items, thus it is reasonable that the winner is `CB-fsr`, followed by `Hybrid-fsr`. In `LibraryThing` and `Last.fm`, the ranking order is different, and `Hybrid-fsr` is the winner in these cases. These behaviors are reasonably derived from the quality of descriptions: with worse descriptions, `CB-fsr` is the last of the ranked list. However, it is noteworthy that with less-described items, `Hybrid-fsr` is able to exploit collaborative information and it performs better than the others.

RQ4: What is the influence of the quality of Linked Data information on the recommendation results?

We have mentioned the differences in terms of performance in previous discussions, however, this is another common question related to the adoption of `Linked Open Data`. In order to answer it, we have decided to feed all the available algorithms with both: LOD, and Tag side-information. In `MovieLens` all the FF variants work better with LOD information than tags. For the competing approaches, `WRMF` seems to be the only one in which performance with LOD is better than its other variants. In `LibraryThing`, we can observe the opposite behavior, with `Jaccard-fsr Tag` as a winner, and this means that tags are better than LOD for item recommendation. However, we can observe that the winning of `Jaccard` variant means that almost half of these features have limited importance for the recommendation task. The consideration we have made before about not homogeneous `Last.fm` LOD description is confirmed by the performance of Tag versions against LOD ones.

RQ5: How does FF performs with reference to diversity and novelty of recommendations?

Even though accuracy is really important to establish the quality of a recommender system, when different variants are proposed, an evaluation of the impact on novelty and diversity of recommendation is mandatory. Novelty and diversity metrics have

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|--|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-f_{sr} LOD | 0.06328 | <u>0.02982</u> | 0.05783 | 0.05970 | 804 | 0.03454 | 6.82679 |
| CB-f_{sr} LOD | 0.08563 | 0.03593 | 0.07435 | 0.07903 | 552 | 0.01935 | 5.98567 |
| Jaccard-f_{sr} Tag | 0.02238 | 0.00939 | 0.01578 | 0.02016 | 591 | 0.03404 | 7.06170 |
| CB-f_{sr} Tag | 0.02377 | 0.00790 | 0.01563 | 0.02127 | 599 | 0.05414 | 7.79185 |
| Hybrid-f_{sr} 3171 | <u>0.06411</u> | 0.02965 | <u>0.05831</u> | <u>0.06061</u> | 823 | 0.03519 | 6.84044 |
| Hybrid-f_{sr} 1586 | 0.05593 | 0.01997 | 0.04292 | 0.04913 | 718 | 0.03970 | 7.08804 |
| Hybrid-f_{sr} 793 | 0.05066 | 0.01547 | 0.03817 | 0.04651 | 527 | 0.03629 | 7.13862 |
| Hybrid-f_{sr} 396 | 0.02321 | 0.00466 | 0.01759 | 0.02641 | 298 | 0.02500 | 6.66303 |
| BPRMF | 0.05808 | 0.02494 | 0.04548 | 0.04934 | 1109 | 0.09286 | 8.54816 |
| BPRMF Tag | 0.05801 | 0.02564 | 0.04669 | 0.04944 | <u>1069</u> | <u>0.08744</u> | <u>8.43765</u> |
| BPRMF LOD | 0.04394 | 0.00972 | 0.03588 | 0.04050 | 772 | 0.03234 | 6.70888 |
| SMRMF | 0.03828 | 0.01635 | 0.02838 | 0.03375 | 612 | 0.03670 | 7.22453 |
| SMRMF Tag | 0.04152 | 0.01830 | 0.03222 | 0.03802 | 621 | 0.04212 | 7.41911 |
| SMRMF LOD | 0.02454 | 0.00877 | 0.01852 | 0.02270 | 481 | 0.02262 | 6.49055 |
| WRMF | 0.02950 | 0.02250 | 0.02786 | 0.02344 | 535 | 0.04192 | 7.47313 |
| WRMF Tag | 0.03411 | 0.02447 | 0.03119 | 0.02727 | 532 | 0.04000 | 7.39674 |
| WRMF LOD | 0.03404 | 0.02473 | 0.03139 | 0.02695 | 513 | 0.03841 | 7.34936 |
| MostPopular | 0.00970 | 0.00431 | 0.00764 | 0.00683 | 51 | 0.00240 | 3.47028 |

Table 5.5: Comparative results on MovieLens dataset, Top-5 recommendation list and relevance threshold of 4/5.

been measured for all the considered experiments, and we suddenly may notice that recommendation lists of Hybrid- f_{sr} show to be more diversified w.r.t. CB- f_{sr} . This may due to the nature of the proposed hybrid algorithm, that makes use of novel features to perform recommendations. Moreover, it could be noticed that about CB- f_{sr} the decrease in diversity seems to be related to the increase in accuracy metrics. If we look at the literature, this is probably due to the overspecialization issue that affects content-based algorithms. In this evaluation, these considerations are coherent to results on LibraryThing, and Last.fm datasets. For the sake of completeness, it is important to underline that the absolute winner in terms of aggregate diversity and distributional inequality on MovieLens, and Last.fm is BPRMF. However, its low performance in terms of accuracy makes us not consider it in this discussion. Finally, Novelty results behave in a coherent way along with accuracy metrics, and this denotes that performance registered for popular items are confirmed on the long tail.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.05801 | 0.04868 | 0.06291 | 0.05509 | 1065 | 0.04348 | 7.25795 |
| CB-fsr LOD | 0.07363 | 0.05875 | 0.07810 | 0.07044 | 736 | 0.02571 | 6.47117 |
| Jaccard-fsr Tag | 0.02071 | 0.01725 | 0.01820 | 0.01936 | 853 | 0.04949 | 7.61902 |
| CB-fsr Tag | 0.02174 | 0.01510 | 0.01774 | 0.02016 | 845 | 0.07750 | 8.32482 |
| Hybrid-fsr 3171 | <u>0.05838</u> | <u>0.04882</u> | <u>0.06334</u> | <u>0.05568</u> | 1093 | 0.04457 | 7.28407 |
| Hybrid-fsr 1586 | 0.05240 | 0.03500 | 0.04681 | 0.04723 | 890 | 0.05024 | 7.53621 |
| Hybrid-fsr 793 | 0.04627 | 0.02575 | 0.03969 | 0.04368 | 620 | 0.04512 | 7.53465 |
| Hybrid-fsr 396 | 0.01732 | 0.00684 | 0.01511 | 0.02088 | 339 | 0.03069 | 7.02743 |
| BPRMF | 0.05356 | 0.04640 | 0.05169 | 0.04698 | 1387 | 0.11563 | 8.88650 |
| BPRMF Tag | 0.05430 | 0.04829 | 0.05347 | 0.04733 | <u>1358</u> | <u>0.10926</u> | <u>8.78973</u> |
| BPRMF LOD | 0.03672 | 0.01830 | 0.03413 | 0.03511 | 1113 | 0.04512 | 7.30863 |
| SMRMF | 0.03619 | 0.02958 | 0.03269 | 0.03238 | 905 | 0.05203 | 7.72057 |
| SMRMF Tag | 0.03854 | 0.03349 | 0.03678 | 0.03568 | 857 | 0.05725 | 7.88452 |
| SMRMF LOD | 0.02169 | 0.01487 | 0.01952 | 0.02067 | 731 | 0.03263 | 7.00065 |
| WRMF | 0.03015 | 0.04300 | 0.03639 | 0.02416 | 663 | 0.05506 | 7.88662 |
| WRMF Tag | 0.03276 | 0.04592 | 0.03937 | 0.02684 | 670 | 0.05327 | 7.82256 |
| WRMF LOD | 0.03369 | 0.04630 | 0.04015 | 0.02719 | 635 | 0.05073 | 7.76383 |
| MostPopular | 0.01235 | 0.01399 | 0.01189 | 0.00838 | 83 | 0.00475 | 4.35932 |

Table 5.6: Comparative results on MovieLens dataset, Top-10 recommendation list and relevance threshold of 4/5.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.07430 | 0.02487 | 0.06133 | 0.06975 | 804 | 0.03454 | 6.82679 |
| CB-fsr LOD | 0.10020 | 0.02993 | 0.07913 | 0.09172 | 552 | 0.01935 | 5.98567 |
| Jaccard-fsr Tag | 0.03513 | 0.00958 | 0.01897 | 0.03140 | 591 | 0.03404 | 7.06170 |
| CB-fsr Tag | 0.03225 | 0.00753 | 0.01821 | 0.02873 | 599 | 0.05414 | 7.79185 |
| Hybrid-fsr 3171 | <u>0.07560</u> | <u>0.02487</u> | <u>0.06194</u> | <u>0.07099</u> | 823 | 0.03519 | 6.84044 |
| Hybrid-fsr 1586 | 0.06868 | 0.01753 | 0.04700 | 0.06030 | 718 | 0.03970 | 7.08804 |
| Hybrid-fsr 793 | 0.06497 | 0.01379 | 0.04296 | 0.05964 | 527 | 0.03629 | 7.13862 |
| Hybrid-fsr 396 | 0.03341 | 0.00441 | 0.02093 | 0.03728 | 298 | 0.02500 | 6.66303 |
| BPRMF | 0.07295 | 0.02167 | 0.05050 | 0.06216 | 1109 | 0.09286 | 8.54816 |
| BPRMF Tag | 0.07215 | 0.02233 | 0.05126 | 0.06176 | <u>1069</u> | <u>0.08744</u> | <u>8.43765</u> |
| BPRMF LOD | 0.05384 | 0.00857 | 0.03938 | 0.04944 | 772 | 0.03234 | 6.70888 |
| SMRMF | 0.05179 | 0.01512 | 0.03253 | 0.04582 | 612 | 0.03670 | 7.22453 |
| SMRMF Tag | 0.05632 | 0.01767 | 0.03703 | 0.05182 | 621 | 0.04212 | 7.41911 |
| SMRMF LOD | 0.03152 | 0.00773 | 0.02082 | 0.02906 | 481 | 0.02262 | 6.49055 |
| WRMF | 0.03675 | 0.01998 | 0.02977 | 0.02926 | 535 | 0.04192 | 7.47313 |
| WRMF Tag | 0.04185 | 0.02130 | 0.03328 | 0.03323 | 532 | 0.04000 | 7.39674 |
| WRMF LOD | 0.04152 | 0.02194 | 0.03351 | 0.03284 | 513 | 0.03841 | 7.34936 |
| MostPopular | 0.01179 | 0.00355 | 0.00814 | 0.00825 | 51 | 0.00240 | 3.47028 |

Table 5.7: Comparative results on MovieLens dataset, Top-5 recommendation list and relevance threshold of 3/5.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.06980 | 0.04184 | 0.06590 | 0.06554 | 1065 | 0.04348 | 7.25795 |
| CB-fsr LOD | 0.08791 | 0.04943 | 0.08197 | 0.08312 | 736 | 0.02571 | 6.47117 |
| Jaccard-fsr Tag | 0.03180 | 0.01664 | 0.02117 | 0.02980 | 853 | 0.04949 | 7.61902 |
| CB-fsr Tag | 0.02957 | 0.01426 | 0.02016 | 0.02736 | 845 | 0.07750 | 8.32482 |
| Hybrid-fsr 3171 | <u>0.07020</u> | 0.04179 | <u>0.06630</u> | <u>0.06621</u> | 1093 | 0.04457 | 7.28407 |
| Hybrid-fsr 1586 | 0.06573 | 0.03091 | 0.05061 | 0.05891 | 890 | 0.05024 | 7.53621 |
| Hybrid-fsr 793 | 0.06078 | 0.02312 | 0.04415 | 0.05695 | 620 | 0.04512 | 7.53465 |
| Hybrid-fsr 396 | 0.02533 | 0.00654 | 0.01787 | 0.02986 | 339 | 0.03069 | 7.02743 |
| BPRMF | 0.06770 | 0.04054 | 0.05604 | 0.05944 | 1387 | 0.11563 | 8.88650 |
| BPRMF Tag | 0.06950 | <u>0.04302</u> | 0.05773 | 0.06031 | <u>1358</u> | <u>0.10926</u> | <u>8.78973</u> |
| BPRMF LOD | 0.04551 | 0.01605 | 0.03723 | 0.04325 | 1113 | 0.04512 | 7.30863 |
| SMRMF | 0.04940 | 0.02766 | 0.03659 | 0.04426 | 905 | 0.05203 | 7.72057 |
| SMRMF Tag | 0.05224 | 0.03213 | 0.04122 | 0.04859 | 857 | 0.05725 | 7.88452 |
| SMRMF LOD | 0.02803 | 0.01341 | 0.02167 | 0.02659 | 731 | 0.03263 | 7.00065 |
| WRMF | 0.03795 | 0.03892 | 0.03812 | 0.03040 | 663 | 0.05506 | 7.88662 |
| WRMF Tag | 0.04066 | 0.04082 | 0.04118 | 0.03306 | 670 | 0.05327 | 7.82256 |
| WRMF LOD | 0.04190 | 0.04124 | 0.04201 | 0.03365 | 635 | 0.05073 | 7.76383 |
| MostPopular | 0.01581 | 0.01206 | 0.01247 | 0.01058 | 83 | 0.00475 | 4.35932 |

Table 5.8: Comparative results on MovieLens dataset, Top-10 recommendation list and relevance threshold of 3/5.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|-----------------|
| Jaccard-fsr LOD | 0.02822 | 0.04786 | 0.04560 | 0.03203 | <u>4043</u> | <u>0.09912</u> | <u>10.47816</u> |
| CB-fsr LOD | 0.02921 | 0.04786 | 0.04908 | 0.03478 | 3825 | 0.07815 | 9.88738 |
| Jaccard-fsr Tag | 0.04009 | 0.07187 | 0.07124 | 0.04863 | 2989 | 0.05186 | 9.30556 |
| CB-fsr Tag | <u>0.03763</u> | <u>0.06848</u> | <u>0.06936</u> | <u>0.04570</u> | 1519 | 0.01962 | 7.92150 |
| Hybrid-fsr 13073 | 0.02866 | 0.04836 | 0.04593 | 0.03237 | 4098 | 0.10247 | 10.53733 |
| Hybrid-fsr 6851 | 0.01719 | 0.02409 | 0.02720 | 0.02072 | 3127 | 0.06968 | 9.87870 |
| Hybrid-fsr 3425 | 0.00978 | 0.01306 | 0.01623 | 0.01255 | 2017 | 0.04219 | 8.99890 |
| Hybrid-fsr 1712 | 0.00423 | 0.00487 | 0.00687 | 0.00557 | 1156 | 0.02221 | 8.17881 |
| BPRMF | 0.00481 | 0.00686 | 0.00643 | 0.00490 | 1604 | 0.03418 | 9.11855 |
| BPRMF Tag | 0.00417 | 0.00680 | 0.00607 | 0.00430 | 2317 | 0.04500 | 9.35826 |
| BPRMF LOD | 0.00090 | 0.00102 | 0.00132 | 0.00103 | 1310 | 0.01386 | 7.14172 |
| SMRMF | 0.00347 | 0.00506 | 0.00468 | 0.00369 | 1202 | 0.01937 | 8.08369 |
| SMRMF Tag | 0.00249 | 0.00345 | 0.00313 | 0.00253 | 1089 | 0.01164 | 7.10852 |
| SMRMF LOD | 0.00127 | 0.00153 | 0.00185 | 0.00146 | 1261 | 0.01503 | 7.36249 |
| WRMF | 0.00625 | 0.01163 | 0.01001 | 0.00643 | 370 | 0.00771 | 7.06538 |
| WRMF Tag | 0.00454 | 0.00837 | 0.00686 | 0.00451 | 619 | 0.00864 | 7.11308 |
| WRMF LOD | 0.00321 | 0.00626 | 0.00485 | 0.00314 | 359 | 0.00482 | 6.38107 |
| MostPopular | 0.00012 | 0.00025 | 0.00011 | 0.00008 | 12 | 0.00033 | 2.46979 |

Table 5.9: Comparative results on LibraryThing dataset, Top-5 recommendation list and relevance threshold of 9/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|-----------------|
| Jaccard-fsr LOD | 0.01966 | 0.06244 | 0.05009 | 0.02468 | <u>5213</u> | <u>0.11331</u> | <u>10.69759</u> |
| CB-fsr LOD | 0.01944 | 0.06220 | 0.05273 | 0.02598 | 5192 | 0.09369 | 10.18603 |
| Jaccard-fsr Tag | 0.02723 | 0.09652 | 0.07852 | 0.03661 | 4038 | 0.05882 | 9.51908 |
| CB-fsr Tag | <u>0.02435</u> | <u>0.08563</u> | <u>0.07396</u> | <u>0.03352</u> | 2011 | 0.02213 | 8.20660 |
| Hybrid-fsr 13073 | 0.01987 | 0.06265 | 0.05035 | 0.02491 | 5274 | 0.11708 | 10.75585 |
| Hybrid-fsr 6851 | 0.01087 | 0.02992 | 0.02781 | 0.01504 | 3907 | 0.07767 | 10.09293 |
| Hybrid-fsr 3425 | 0.00575 | 0.01512 | 0.01573 | 0.00875 | 2357 | 0.04635 | 9.24221 |
| Hybrid-fsr 1712 | 0.00265 | 0.00587 | 0.00660 | 0.00399 | 1337 | 0.02543 | 8.50527 |
| BPRMF | 0.00408 | 0.01163 | 0.00802 | 0.00435 | 2187 | 0.04422 | 9.50059 |
| BPRMF Tag | 0.00389 | 0.01247 | 0.00822 | 0.00405 | 3263 | 0.05861 | 9.76858 |
| BPRMF LOD | 0.00069 | 0.00184 | 0.00152 | 0.00084 | 1878 | 0.01721 | 7.46591 |
| SMRMF | 0.00310 | 0.00937 | 0.00622 | 0.00333 | 1809 | 0.02768 | 8.65574 |
| SMRMF Tag | 0.00285 | 0.00837 | 0.00512 | 0.00276 | 1686 | 0.01784 | 7.79521 |
| SMRMF LOD | 0.00096 | 0.00248 | 0.00210 | 0.00117 | 1942 | 0.02260 | 8.04500 |
| WRMF | 0.00576 | 0.02056 | 0.01351 | 0.00599 | 514 | 0.01181 | 7.67982 |
| WRMF Tag | 0.00434 | 0.01583 | 0.00977 | 0.00437 | 935 | 0.01274 | 7.69802 |
| WRMF LOD | 0.00358 | 0.01295 | 0.00754 | 0.00340 | 584 | 0.00753 | 7.02718 |
| MostPopular | 0.00019 | 0.00101 | 0.00040 | 0.00014 | 20 | 0.00075 | 3.57568 |

Table 5.10: Comparative results on LibraryThing dataset, Top-10 recommendation list and relevance threshold of 9/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|-----------------|
| Jaccard-fsr LOD | 0.05853 | 0.06017 | 0.06312 | 0.06646 | <u>4043</u> | <u>0.09912</u> | <u>10.47816</u> |
| CB-fsr LOD | 0.05182 | 0.05403 | 0.06314 | 0.06052 | 3825 | 0.07815 | 9.88738 |
| Jaccard-fsr Tag | 0.07735 | 0.08662 | 0.09597 | 0.09279 | 2989 | 0.05186 | 9.30556 |
| CB-fsr Tag | <u>0.06195</u> | <u>0.06839</u> | <u>0.08517</u> | <u>0.07385</u> | 1519 | 0.01962 | 7.92150 |
| Hybrid-fsr 13073 | 0.05931 | 0.06053 | 0.06332 | 0.06722 | 4098 | 0.10247 | 10.53733 |
| Hybrid-fsr 6851 | 0.03671 | 0.02982 | 0.03702 | 0.04334 | 3127 | 0.06968 | 9.87870 |
| Hybrid-fsr 3425 | 0.02009 | 0.01583 | 0.02194 | 0.02549 | 2017 | 0.04219 | 8.99890 |
| Hybrid-fsr 1712 | 0.00932 | 0.00579 | 0.00937 | 0.01198 | 1156 | 0.02221 | 8.17881 |
| BPRMF | 0.00926 | 0.00702 | 0.00814 | 0.00938 | 1604 | 0.03418 | 9.11855 |
| BPRMF Tag | 0.00857 | 0.00723 | 0.00760 | 0.00849 | 2317 | 0.04500 | 9.35826 |
| BPRMF LOD | 0.00165 | 0.00104 | 0.00158 | 0.00175 | 1310 | 0.01386 | 7.14172 |
| SMRMF | 0.00674 | 0.00455 | 0.00569 | 0.00713 | 1202 | 0.01937 | 8.08369 |
| SMRMF Tag | 0.00524 | 0.00374 | 0.00440 | 0.00559 | 1089 | 0.01164 | 7.10852 |
| SMRMF LOD | 0.00240 | 0.00198 | 0.00259 | 0.00267 | 1261 | 0.01503 | 7.36249 |
| WRMF | 0.01366 | 0.01679 | 0.01515 | 0.01421 | 370 | 0.00771 | 7.06538 |
| WRMF Tag | 0.00964 | 0.01116 | 0.00958 | 0.00963 | 619 | 0.00864 | 7.11308 |
| WRMF LOD | 0.00701 | 0.00846 | 0.00726 | 0.00686 | 359 | 0.00482 | 6.38107 |
| MostPopular | 0.00012 | 0.00016 | 0.00011 | 0.00008 | 12 | 0.00033 | 2.46979 |

Table 5.11: Comparative results on LibraryThing dataset, Top-5 recommendation list and relevance threshold of 7/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|-----------------|
| Jaccard-fsr LOD | 0.04183 | 0.07904 | 0.06855 | 0.05195 | <u>5213</u> | <u>0.11331</u> | <u>10.69759</u> |
| CB-fsr LOD | 0.03484 | 0.06922 | 0.06678 | 0.04559 | 5192 | 0.09369 | 10.18603 |
| Jaccard-fsr Tag | 0.05303 | 0.11633 | 0.10434 | 0.07030 | 4038 | 0.05882 | 9.51908 |
| CB-fsr Tag | <u>0.04148</u> | <u>0.08711</u> | <u>0.08966</u> | <u>0.05526</u> | 2011 | 0.02213 | 8.20660 |
| Hybrid-fsr 13073 | 0.04247 | 0.07893 | 0.06871 | 0.05262 | 5274 | 0.11708 | 10.75585 |
| Hybrid-fsr 6851 | 0.02390 | 0.03677 | 0.03742 | 0.03204 | 3907 | 0.07767 | 10.09293 |
| Hybrid-fsr 3425 | 0.01236 | 0.01851 | 0.02105 | 0.01818 | 2357 | 0.04635 | 9.24221 |
| Hybrid-fsr 1712 | 0.00592 | 0.00710 | 0.00883 | 0.00867 | 1337 | 0.02543 | 8.50527 |
| BPRMF | 0.00837 | 0.01280 | 0.01020 | 0.00869 | 2187 | 0.04422 | 9.50059 |
| BPRMF Tag | 0.00777 | 0.01325 | 0.01007 | 0.00793 | 3263 | 0.05861 | 9.76858 |
| BPRMF LOD | 0.00127 | 0.00166 | 0.00172 | 0.00145 | 1878 | 0.01721 | 7.46591 |
| SMRMF | 0.00592 | 0.00836 | 0.00718 | 0.00638 | 1809 | 0.02768 | 8.65574 |
| SMRMF Tag | 0.00556 | 0.00835 | 0.00653 | 0.00568 | 1686 | 0.01784 | 7.79521 |
| SMRMF LOD | 0.00191 | 0.00315 | 0.00287 | 0.00223 | 1942 | 0.02260 | 8.04500 |
| WRMF | 0.01253 | 0.02829 | 0.01982 | 0.01315 | 514 | 0.01181 | 7.67982 |
| WRMF Tag | 0.00909 | 0.02006 | 0.01333 | 0.00921 | 935 | 0.01274 | 7.69802 |
| WRMF LOD | 0.00763 | 0.01679 | 0.01080 | 0.00729 | 584 | 0.00753 | 7.02718 |
| MostPopular | 0.00023 | 0.00095 | 0.00044 | 0.00017 | 20 | 0.00075 | 3.57568 |

Table 5.12: Comparative results on LibraryThing dataset, Top-10 recommendation list and relevance threshold of 7/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.00482 | 0.01943 | 0.01561 | 0.00568 | 899 | 0.02222 | 7.88347 |
| CB-fsr LOD | 0.00397 | 0.01679 | 0.01185 | 0.00451 | 772 | 0.01347 | 6.77499 |
| Jaccard-fsr Tag | <u>0.00547</u> | <u>0.02206</u> | <u>0.01705</u> | <u>0.00653</u> | <u>1058</u> | <u>0.03078</u> | <u>8.45781</u> |
| CB-fsr Tag | 0.00740 | 0.03135 | 0.02356 | 0.00851 | 1008 | 0.02621 | 8.07217 |
| Hybrid-fsr 10004 | 0.00514 | 0.02041 | 0.01612 | 0.00602 | 946 | 0.02401 | 8.01508 |
| Hybrid-fsr 5002 | 0.00311 | 0.01215 | 0.00841 | 0.00341 | 1028 | 0.03232 | 8.57696 |
| Hybrid-fsr 2051 | 0.00118 | 0.00482 | 0.00320 | 0.00143 | 907 | 0.02905 | 8.44240 |
| Hybrid-fsr 1250 | 0.00032 | 0.00134 | 0.00101 | 0.00034 | 656 | 0.02387 | 8.22125 |
| BPRMF | 0.00032 | 0.00080 | 0.00047 | 0.00025 | 253 | 0.00599 | 6.28848 |
| BPRMF Tag | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 151 | 0.00099 | 3.43110 |
| BPRMF LOD | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 1112 | 0.02934 | 8.14561 |
| SMRMF | 0.00011 | 0.00054 | 0.00034 | 0.00011 | 364 | 0.00647 | 6.24558 |
| SMRMF Tag | 0.00011 | 0.00018 | 0.00000 | 0.00007 | 411 | 0.00648 | 5.99848 |
| SMRMF LOD | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 777 | 0.01832 | 7.51538 |
| WRMF | 0.00236 | 0.00857 | 0.00491 | 0.00208 | 186 | 0.00839 | 6.79522 |
| WRMF Tag | 0.00171 | 0.00589 | 0.00393 | 0.00160 | 300 | 0.00820 | 6.71812 |
| WRMF LOD | 0.00171 | 0.00679 | 0.00351 | 0.00170 | 334 | 0.00829 | 6.75719 |
| MostPopular | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 9 | 0.00043 | 2.50907 |

Table 5.13: Comparative results on Last.fm dataset, Top-5 recommendation list and relevance threshold of 9/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.00311 | 0.02479 | 0.01750 | 0.00417 | 1254 | 0.02848 | 8.24498 |
| CB-fsr LOD | 0.00284 | 0.02291 | 0.01440 | 0.00354 | 1196 | 0.01932 | 7.34914 |
| Jaccard-fsr Tag | <u>0.00402</u> | <u>0.03224</u> | <u>0.02012</u> | <u>0.00512</u> | <u>1544</u> | 0.03931 | <u>8.80789</u> |
| CB-fsr Tag | 0.00456 | 0.03760 | 0.02570 | 0.00613 | 1519 | 0.03539 | 8.52017 |
| Hybrid-fsr 10004 | 0.00327 | 0.02568 | 0.01803 | 0.00440 | 1325 | 0.03072 | 8.38414 |
| Hybrid-fsr 5002 | 0.00193 | 0.01563 | 0.00955 | 0.00246 | 1421 | <u>0.03900</u> | 8.85578 |
| Hybrid-fsr 2051 | 0.00086 | 0.00674 | 0.00384 | 0.00112 | 1191 | 0.03486 | 8.71436 |
| Hybrid-fsr 1250 | 0.00048 | 0.00402 | 0.00190 | 0.00045 | 807 | 0.02910 | 8.54335 |
| BPRMF | 0.00027 | 0.00134 | 0.00061 | 0.00023 | 334 | 0.00768 | 6.66212 |
| BPRMF Tag | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 242 | 0.00226 | 4.75911 |
| BPRMF LOD | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 1591 | 0.03678 | 8.45905 |
| SMRMF | 0.00016 | 0.00161 | 0.00069 | 0.00015 | 554 | 0.00952 | 6.75383 |
| SMRMF Tag | 0.00016 | 0.00098 | 0.00021 | 0.00012 | 652 | 0.01040 | 6.72105 |
| SMRMF LOD | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 1234 | 0.02770 | 8.17090 |
| WRMF | 0.00188 | 0.01358 | 0.00654 | 0.00183 | 257 | 0.01195 | 7.28991 |
| WRMF Tag | 0.00150 | 0.01063 | 0.00567 | 0.00149 | 452 | 0.01346 | 7.45752 |
| WRMF LOD | 0.00150 | 0.01179 | 0.00529 | 0.00155 | 552 | 0.01376 | 7.48693 |
| MostPopular | 0.00005 | 0.00018 | 0.00000 | 0.00004 | 14 | 0.00097 | 3.48598 |

Table 5.14: Comparative results on Last.fm dataset, Top-10 recommendation list and relevance threshold of 9/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.00675 | 0.02364 | 0.01996 | 0.00792 | 899 | 0.02222 | 7.88347 |
| CB-fsr LOD | 0.00547 | 0.02108 | 0.01464 | 0.00603 | 772 | 0.01347 | 6.77499 |
| Jaccard-fsr Tag | <u>0.00868</u> | <u>0.03131</u> | <u>0.02330</u> | <u>0.00971</u> | <u>1058</u> | <u>0.03078</u> | <u>8.45781</u> |
| CB-fsr Tag | 0.01115 | 0.04118 | 0.03115 | 0.01257 | 1008 | 0.02621 | 8.07217 |
| Hybrid-fsr 10004 | 0.00707 | 0.02449 | 0.02057 | 0.00829 | 946 | 0.02401 | 8.01508 |
| Hybrid-fsr 5002 | 0.00482 | 0.01657 | 0.01258 | 0.00571 | 1028 | 0.03232 | 8.57696 |
| Hybrid-fsr 2051 | 0.00161 | 0.00634 | 0.00427 | 0.00195 | 907 | 0.02905 | 8.44240 |
| Hybrid-fsr 1250 | 0.00075 | 0.00259 | 0.00181 | 0.00082 | 656 | 0.02387 | 8.22125 |
| BPRMF | 0.00054 | 0.00143 | 0.00081 | 0.00045 | 253 | 0.00599 | 6.28848 |
| BPRMF Tag | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 151 | 0.00099 | 3.43110 |
| BPRMF LOD | 0.00011 | 0.00018 | 0.00004 | 0.00008 | 1112 | 0.02934 | 8.14561 |
| SMRMF | 0.00032 | 0.00161 | 0.00078 | 0.00025 | 364 | 0.00647 | 6.24558 |
| SMRMF Tag | 0.00011 | 0.00013 | 0.00000 | 0.00007 | 411 | 0.00648 | 5.99848 |
| SMRMF LOD | 0.00011 | 0.00018 | 0.00004 | 0.00009 | 777 | 0.01832 | 7.51538 |
| WRMF | 0.00397 | 0.01456 | 0.00920 | 0.00385 | 186 | 0.00839 | 6.79522 |
| WRMF Tag | 0.00279 | 0.00987 | 0.00627 | 0.00259 | 300 | 0.00820 | 6.71812 |
| WRMF LOD | 0.00236 | 0.00831 | 0.00493 | 0.00230 | 334 | 0.00829 | 6.75719 |
| MostPopular | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 9 | 0.00043 | 2.50907 |

Table 5.15: Comparative results on Last.fm dataset, Top-5 recommendation list and relevance threshold of 7/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.00477 | 0.03391 | 0.02333 | 0.00610 | 1254 | 0.02848 | 8.24498 |
| CB-fsr LOD | 0.00391 | 0.02873 | 0.01786 | 0.00476 | 1196 | 0.01932 | 7.34914 |
| Jaccard-fsr Tag | <u>0.00665</u> | <u>0.04586</u> | <u>0.02774</u> | <u>0.00791</u> | <u>1544</u> | 0.03931 | <u>8.80789</u> |
| CB-fsr Tag | 0.00734 | 0.05332 | 0.03524 | 0.00943 | 1519 | 0.03539 | 8.52017 |
| Hybrid-fsr 10004 | 0.00482 | 0.03356 | 0.02361 | 0.00629 | 1325 | 0.03072 | 8.38414 |
| Hybrid-fsr 5002 | 0.00327 | 0.02282 | 0.01465 | 0.00430 | 1421 | <u>0.03900</u> | 8.85578 |
| Hybrid-fsr 2051 | 0.00123 | 0.00984 | 0.00543 | 0.00158 | 1191 | 0.03486 | 8.71436 |
| Hybrid-fsr 1250 | 0.00070 | 0.00482 | 0.00270 | 0.00076 | 807 | 0.02910 | 8.54335 |
| BPRMF | 0.00043 | 0.00214 | 0.00095 | 0.00040 | 334 | 0.00768 | 6.66212 |
| BPRMF Tag | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 242 | 0.00226 | 4.75911 |
| BPRMF LOD | 0.00011 | 0.00045 | 0.00012 | 0.00008 | 1591 | 0.03678 | 8.45905 |
| SMRMF | 0.00027 | 0.00268 | 0.00112 | 0.00024 | 554 | 0.00952 | 6.75383 |
| SMRMF Tag | 0.00021 | 0.00112 | 0.00022 | 0.00015 | 652 | 0.01040 | 6.72105 |
| SMRMF LOD | 0.00011 | 0.00071 | 0.00022 | 0.00010 | 1234 | 0.02770 | 8.17090 |
| WRMF | 0.00295 | 0.02086 | 0.01151 | 0.00316 | 257 | 0.01195 | 7.28991 |
| WRMF Tag | 0.00252 | 0.01791 | 0.00898 | 0.00246 | 452 | 0.01346 | 7.45752 |
| WRMF LOD | 0.00209 | 0.01518 | 0.00721 | 0.00212 | 552 | 0.01376 | 7.48693 |
| MostPopular | 0.00005 | 0.00018 | 0.00000 | 0.00004 | 14 | 0.00097 | 3.48598 |

Table 5.16: Comparative results on Last.fm dataset, Top-10 recommendation list and relevance threshold of 7/10.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.03355 | 0.03254 | 0.04071 | 0.03832 | 899 | 0.02222 | 7.88347 |
| CB-fsr LOD | 0.02304 | 0.02246 | 0.02695 | 0.02458 | 772 | 0.01347 | 6.77499 |
| Jaccard-fsr Tag | <u>0.04812</u> | <u>0.04328</u> | <u>0.04865</u> | <u>0.05145</u> | <u>1058</u> | <u>0.03078</u> | <u>8.45781</u> |
| CB-fsr Tag | 0.04930 | 0.04582 | 0.06167 | 0.05507 | 1008 | 0.02621 | 8.07217 |
| Hybrid-fsr 10004 | 0.03376 | 0.03236 | 0.04133 | 0.03912 | 946 | 0.02401 | 8.01508 |
| Hybrid-fsr 5002 | 0.02583 | 0.02294 | 0.02731 | 0.02841 | 1028 | 0.03232 | 8.57696 |
| Hybrid-fsr 2051 | 0.00911 | 0.00762 | 0.00796 | 0.01011 | 907 | 0.02905 | 8.44240 |
| Hybrid-fsr 1250 | 0.00493 | 0.00384 | 0.00470 | 0.00568 | 656 | 0.02387 | 8.22125 |
| BPRMF | 0.00407 | 0.00553 | 0.00372 | 0.00370 | 253 | 0.00599 | 6.28848 |
| BPRMF Tag | 0.00000 | 0.00000 | 0.00000 | 0.00000 | 151 | 0.00099 | 3.43110 |
| BPRMF LOD | 0.00107 | 0.00091 | 0.00057 | 0.00091 | 1112 | 0.02934 | 8.14561 |
| SMRMF | 0.00482 | 0.00403 | 0.00351 | 0.00512 | 364 | 0.00647 | 6.24558 |
| SMRMF Tag | 0.00182 | 0.00155 | 0.00074 | 0.00136 | 411 | 0.00648 | 5.99848 |
| SMRMF LOD | 0.00086 | 0.00085 | 0.00035 | 0.00064 | 777 | 0.01832 | 7.51538 |
| WRMF | 0.01983 | 0.02150 | 0.02291 | 0.01936 | 186 | 0.00839 | 6.79522 |
| WRMF Tag | 0.01426 | 0.01281 | 0.01567 | 0.01406 | 300 | 0.00820 | 6.71812 |
| WRMF LOD | 0.01050 | 0.01013 | 0.00997 | 0.00917 | 334 | 0.00829 | 6.75719 |
| MostPopular | 0.00011 | 0.00013 | 0.00002 | 0.00007 | 9 | 0.00043 | 2.50907 |

Table 5.17: Comparative results on Last.fm dataset, Top-5 recommendation list and no relevance threshold.

| Alg | P@N | R@N | nDCG@N | EPC | ADiv@N | Gini | SE |
|-------------------------|----------------|----------------|----------------|----------------|-------------|----------------|----------------|
| Jaccard-fsr LOD | 0.02556 | 0.04863 | 0.04784 | 0.03099 | 1254 | 0.02848 | 8.24498 |
| CB-fsr LOD | 0.01929 | 0.03838 | 0.03477 | 0.02138 | 1196 | 0.01932 | 7.34914 |
| Jaccard-fsr Tag | <u>0.03735</u> | <u>0.06844</u> | <u>0.05942</u> | <u>0.04263</u> | <u>1544</u> | 0.03931 | <u>8.80789</u> |
| CB-fsr Tag | 0.03794 | 0.07064 | 0.07165 | 0.04510 | 1519 | 0.03539 | 8.52017 |
| Hybrid-fsr 10004 | 0.02653 | 0.05034 | 0.04874 | 0.03214 | 1325 | 0.03072 | 8.38414 |
| Hybrid-fsr 5002 | 0.01902 | 0.03343 | 0.03220 | 0.02273 | 1421 | <u>0.03900</u> | 8.85578 |
| Hybrid-fsr 2051 | 0.00702 | 0.01214 | 0.01071 | 0.00832 | 1191 | 0.03486 | 8.71436 |
| Hybrid-fsr 1250 | 0.00370 | 0.00617 | 0.00600 | 0.00457 | 807 | 0.02910 | 8.54335 |
| BPRMF | 0.00348 | 0.00902 | 0.00495 | 0.00338 | 334 | 0.00768 | 6.66212 |
| BPRMF Tag | 0.00032 | 0.00073 | 0.00028 | 0.00021 | 242 | 0.00226 | 4.75911 |
| BPRMF LOD | 0.00096 | 0.00188 | 0.00084 | 0.00088 | 1591 | 0.03678 | 8.45905 |
| SMRMF | 0.00407 | 0.00723 | 0.00471 | 0.00447 | 554 | 0.00952 | 6.75383 |
| SMRMF Tag | 0.00225 | 0.00402 | 0.00155 | 0.00181 | 652 | 0.01040 | 6.72105 |
| SMRMF LOD | 0.00080 | 0.00156 | 0.00058 | 0.00067 | 1234 | 0.02770 | 8.17090 |
| WRMF | 0.01656 | 0.03566 | 0.02847 | 0.01711 | 257 | 0.01195 | 7.28991 |
| WRMF Tag | 0.01372 | 0.02467 | 0.02217 | 0.01369 | 452 | 0.01346 | 7.45752 |
| WRMF LOD | 0.01099 | 0.02072 | 0.01522 | 0.00991 | 552 | 0.01376 | 7.48693 |
| MostPopular | 0.00027 | 0.00089 | 0.00021 | 0.00019 | 14 | 0.00097 | 3.48598 |

Table 5.18: Comparative results on Last . fm dataset, Top-10 recommendation list and no relevance threshold.

Results Discussion

FF shows interesting performance in terms of accuracy, diversity, and novelty on all the experimental settings and it results as a highly competitive approach compared to other pure-collaborative and hybrid variants of state-of-the-art algorithms. In particular, we see that on the MovieLens, LibraryThing, and Last . fm datasets, at least a variant of FF gets the best results for the evaluated accuracy metrics. On the other hand, our Feature Factorization is not the absolute champion when compared on the MovieLens, and Last . fm datasets considering diversity metrics. However, if we consider the trade-off between accuracy, novelty, and diversity, we see that FF is the best performing algorithm. The differences between the different behaviors on the three datasets have been detailed in the previous discussion, and summing up they can be explained by looking at different dimensions of the datasets (LOD, Tag side information; quality of descriptions, sparsity of the datasets, and popularity bias). In Last . fm, we have rescaled the users' feedback represented as the number of times they played a song and normalized it on a 1-10

scale. This could have affected the final results especially in terms of accuracy for all the algorithms. If we consider the feature-augmented dataset, by looking at the data represented in Table 5.4 the first observation we make is that the number of features in `Last.fm` is two orders of magnitude higher than the number of items while in `LibraryThing` it is just one. Then, the decrease in performance of `FF` may be also attributed to the curse of dimensionality problem. Moreover, a deeper investigation of the quality of the adopted `LOD` dataset is needed. A few papers have been published on this topic [411, 89], however, still there is not a community-endorsed metric to evaluate the quality of the knowledge encoded in a Linked Data dataset for recommendation tasks.

5.5 Conclusion

In this line of research, we have introduced `FF`, a novel algorithm that bases on feature recommendation as an intermediate step for computing *top-N* items recommendation lists. The main idea behind `FF` is that feature relevance in a user profile plays a key role in the selection and rating of an item in a collection. Based on this observation we have developed an algorithm that shifts the recommendation problem from a user-item space to a user-feature one. In this new space, we have introduced the notion of feature relevance and feature rating. We have combined them with well-known factorization techniques computing rating and relevance for each feature unknown to the user. Then, by combining the values associated with the features composing an item we have predicted a *top-N* recommendation list of items. We have compared `FF` with well-known factorization techniques (both pure collaborative and hybrid variants with side information) on three datasets in the domains of movie, books, and music. In all the datasets `FF` results as the best algorithm in terms of a trade-off between accuracy, diversity, and novelty of results. This can be considered as a strong clue to confirm our intuition that recommending items via feature ranking is a feasible way to develop content-aware recommendation engines. As future work, we are investigating the behavior of `FF` with different factorization techniques in the item-feature space. Moreover, since we have col-

lected content-based data from Linked Open Data datasets, an analysis of the influence of such datasets on the recommendation results is also in progress. Another aspect we are willing to deepen is related to results explanation. Indeed, very interestingly, item recommendation via feature ranking paves the way to new proposals for explanation services.

Chapter 6

Metadata to address Cold-start problem

6.1 Introduction

Cross-domain recommendation has recently emerged as a potential solution to the cold start problem in recommender systems [86], aiming to mitigate the lack of data by exploiting user preferences and item attributes in domains distinct but related to the target domain. In this line, most of the cross-domain approaches proposed so far are based on collaborative filtering [110], exploiting user preferences as a bridge to relate source and target domains, and ignoring the content of the items. Hence, they benefit from the fact that they do not need to perform any kind of analysis of item contents, which are in general highly heterogeneous across domains, and whose inter-domain relationships may be difficult to be establish.

These difficulties, however, can be addressed nowadays thanks to the Semantic Web initiative [338], and more specifically to its reference implementation the Linked Open Data (LOD) project [64], which has originated a large number of

inter-linked knowledge repositories publicly available in the Web, following the Semantic Web standards for data representation and access. Hence, in the current Web there is a wide array of structured data sources with information of items belonging to a variety of domains, such as history, arts, science, industry, media and sports, to name a few. This information not only consists of particular multimedia contents and associated metadata, but also explicit, semantic relations between items and metadata.

Motivated by the availability of large amounts of item metadata and semantic relations in the Linked Data cloud, we aim to address the cross-domain recommendation problem not only focusing on user preferences and item attributes, but also exploiting content-based relations between items from different domains. More specifically, we propose to use the set of LOD semantic features and relations as inter-domain links for supporting knowledge transfer across domains, enabling cross-domain item similarities, and providing recommendations for cold start users in the target domain.

Previous work has proposed graph-based algorithms to address the recommendation problem in heterogeneous datasets [405, 282], analyzing the topology of semantic networks to jointly exploit user preferences and item metadata. These approaches have been shown to be effective for recommendation, but suffer from computational issues caused by the size of the semantic networks, which are in general very large. Differently, we avoid these issues by working in two steps. First, we exploit the semantic networks to compute inter-domain similarities that link items from different domains. Then, we leverage the computed similarities in hybrid Matrix Factorization (MF) models for recommendation, which no longer need to deal with the whole networks.

Therefore, this research line has led to the development of novel, effective hybrid matrix factorization models that jointly exploit user preferences and item metadata for cross-domain recommendation. Moreover, we adapt a fast learning algorithm by [304] for efficiently building our models, and evaluate them in cold start scenarios on several domains, in terms of both precision and diversity.

We evaluate the performance of the proposed models using a dataset of Face-

book¹ *likes* about books, movies and music. In order to obtain semantic metadata for the different items, we first mapped the items in our dataset to entities in LOD by means of SPARQL queries, and then extracted their attributes and relations to enhance the item profiles.

In a first experiment, we compared several state-of-the-art semantic similarity metrics for content-based recommendation, aiming to understand which is more suitable for later injecting in our cross-domain MF models, and achieved the best results using the link-based approach by [262]. Second, we evaluated the ranking precision and diversity of the recommendations computed by the proposed models. We show that, depending on the involved source and target domains, our models generate more accurate suggestions than the baselines in severe cold start situations. Moreover, the proposed approaches provide a better trade-off between accuracy and diversity, which are in general difficult to balance.

We point out that the presented approaches can be effectively used if the underlying LOD knowledge graph encodes direct or indirect connections between items in different domains. In fact, we need to compute semantic similarity values between items not belonging to the same domain. These connections are quite common for knowledge domains with some degree of information overlapping such as in the case of music, movies, and books but, in case they are missing or rare, this may result as a limitation for the performances of the approaches we introduce here.

The remainder of the chapter is structured as follows. In 6.2, we revise related work on cross-domain recommender systems, focusing on those approaches that are based on Matrix Factorization. In 6.3, we present the developed cross-domain hybrid matrix factorization models. Next, in 6.4, we report and analyze the empirical results achieved in the experiments conducted to analyze user cold start situations. Finally, in 6.5 we end with some conclusions and future research lines.

¹Facebook online social networking, <https://www.facebook.com>

6.2 Related work

In this section, we survey the state of the art on cross-domain recommender systems. First, in 6.2.1 we describe the cross-domain recommendation problem and present a categorization of the approaches, giving representative examples of each category. Next, in 6.2.2 we focus on those cross-domain recommendation approaches that use the matrix factorization technique to bridge the source and target domains.

6.2.1 Cross-domain recommender systems

Nowadays, the majority of recommender systems offer recommendations for items belonging to a single domain. For instance, Netflix² recommends movies and TV shows, Spotify³ recommends songs and music albums, and Barnes & Noble⁴ recommends books. These domain-specific systems have been successfully deployed by numerous web platforms, and the single-domain recommendation functionality is not perceived as a limitation, but rather pitched as a focus on a certain market.

Nonetheless, in large e-commerce sites such as Amazon.com⁵ and eBay⁶ users often provide feedback for items from multiple domains, and in social networks like Facebook⁷ and Twitter⁸ users express their tastes and interests for a variety of topics. It may, therefore, be beneficial to leverage all the available user data provided in various systems and domains, in order to generate more encompassing user models and better recommendations. Instead of treating each domain (e.g., movies, music and books) independently, knowledge acquired in a *source* domain could be transferred to and exploited in another *target* domain. The research challenge of transferring knowledge, and the business potential of delivering recommendations spanning across multiple domains, have triggered an increasing interest in *cross-domain recommendations*.

²Netflix streaming media and video provider, <https://www.netflix.com>

³Spotify digital music service, <https://www.spotify.com>

⁴Barnes & Noble online bookseller, <http://www.barnesandnoble.com>

⁵Amazon electronic commerce site, <https://www.amazon.com>

⁶eBay consumer-to-consumer and business-to-consumer sales, <http://www.ebay.com>

⁷Facebook social network, <https://www.facebook.com>

⁸Twitter online news and social networking service, <https://twitter.com>

The cross-domain recommendation problem has been addressed from various perspectives in different research areas. It has been handled by means of user preference aggregation and mediation strategies for cross-system personalization in User Modeling [12, 58, 341], as a potential solution to mitigate the cold start and sparsity problems in Recommender Systems [110, 345, 369], and as a practical application of knowledge transfer in Machine Learning [155, 232, 292]. Focusing on how knowledge is exploited by cross-domain recommender systems, in [86] we categorized existing works according a two-level taxonomy.

- *Aggregating knowledge.* Knowledge from various source domains is aggregated to perform recommendations in a target domain. Depending on the stage in the recommendation process where the aggregation is performed we can further distinguish three cases. First, we find approaches that *merge user preferences* e.g., ratings, tags, transaction logs, and click-through data. The aggregation can be done by means of a multi-domain rating matrix [57, 323], using a common representation for user preferences such as social tags [361, 12, 143] or semantic concepts [207], linking the preferences via a multi-domain graph [110, 369], or mapping user preferences to domain-independent features such as personality traits [85] or user-item interaction features [250]. In the second case, user modeling data from various recommender systems is *mediated* to improve target recommendations. For instance, [57, 370, 341] import user neighborhoods and user-user similarities computed in the source domain into the target. Finally, some approaches directly *combine* single-domain recommendations, e.g., rating estimations [57, 156] and rating probability distributions [428].
- *Linking and transferring knowledge.* Knowledge linkage or transfer between domains is established to support recommendations. In this case, we find methods that (i) *link domains* by a common knowledge such as item attributes [101], association rules [85], semantic networks [142, 207], and inter-domain correlations [417, 345, 324]; methods that (ii) *share latent features* between source and target domains factor models, either by using same model parameters [291, 182, 172] in both factorizations, or by introducing new parameters

that extend the factorizations [137, 141]; and methods that (iii) *transfer rating patterns* extracted by co-clustering the source domain rating matrix and exploit them in the target domain [233, 155, 109]. After defining the problem, in [290] three different knowledge transfer strategies for collaborative recommendation with auxiliary data (TL-CRAD) are introduced: (i) adaptive knowledge transfer, (ii) collective knowledge transfer and (iii) integrative knowledge transfer. Then, for each of them the author surveys related work with reference to different knowledge strategies with an emphasis on: transfer via prediction rule, transfer via regularization and transfer via constraint.

In terms of the goals addressed by cross-domain recommenders, we find great diversity among the reviewed approaches. Most proposals focus on improving accuracy by reducing data sparsity [232, 345, 87, 417, 292, 369, 250, 426]. In many domains, the average number of ratings per user and item is low, which may negatively affect the quality of the recommendations. Data collected outside the target domain can increase the rating density, and thus may upgrade the recommendation quality. Others seek to enhance user models, which may have personalization-oriented benefits such as (i) discovering new user preferences for the target domain [357, 362], (ii) enhancing similarities between users and items [11, 58], and (iii) measuring vulnerability in social networks [157, 190]. Cross-domain methods have also been applied to bootstrap recommender systems by importing preferences from another source outside the target domain [341], and have been proposed to improve the diversity of recommendations by providing better coverage of the range of user preferences [396]. Finally, a few approaches have dealt with the new user problem [182, 323, 370, 137]. When a user starts using a recommender system, this has no knowledge of the user's tastes and interests, and cannot produce personalized recommendations. This may be solved by exploiting the user's preferences collected in a different source domain.

We observe that addressing the cold-start has been barely investigated as in [242] where the authors present a neighborhood-based algorithm for the dual cold-start problem. The generalization of items and users into a cluster level to obtain high-quality relations also in cold start scenario is the focus of [263]. They first

employ biased matrix factorization to map rating matrix into lower-dimension latent spaces. After this step they apply the K-means clustering algorithm to categorize users and items. Cold start is also the main topic of [398] where the authors propose a novel approach to cross-system personalization based on two assumptions: the existence of a user model that could be shared among platforms and that a specific system can maintain (and provide) the user models built by its system.

As we shall present in 6.3, we aim to deal with the cold-start problem by means of novel matrix factorization models that jointly exploit user ratings and item metadata. Before, in 6.2.2, we revise state of the art cross-domain recommender systems based on matrix factorization.

6.2.2 Matrix factorization-based cross-domain recommender systems

Although matrix factorization models can be applied in cross-domain approaches based on knowledge aggregation –essentially as a standard recommendation problem once the user preferences from both domains are combined– they have been mostly used in knowledge linkage or transfer approaches. In these settings, latent factors from source and target domains are either shared or related in order to establish the bridge between the domains.

One way of linking domains explored in previous works exploits inter-domain similarities by integrating them into the probabilistic matrix factorization method [327]. Specifically, such similarities are imposed as constraints over user or item latent factors when jointly factorizing rating matrices. For instance, [87] proposed an approach in which inter-domain similarities are implicitly learned from data, as model parameters in a non-parametric Bayesian framework. Since user feedback is used to estimate the similarities, user overlap between the domains is required. Addressing the sparsity problem, [417] adapted the probabilistic matrix factorization method to include a probability distribution of user latent factors that encodes inter-domain correlations. One strength of this approach is that user latent factors shared across domains are not needed, allowing more flexibility in capturing the

heterogeneity of domains. Instead of automatically learning implicit correlations in the data, [345] argued that explicit common information is more effective, and relied on shared social tags to compute cross-domain user-to-user and item-to-item similarities. Similarly to previous approaches, rating matrices from the source and target domains are jointly factorized; but in this case user and item latent factors from each domain are restricted, so that their product is consistent with the tag-based similarities.

Latent factors shared between domains can be exploited to support cross-domain recommendations. In this context, two types of approaches have been studied to perform the actual transfer of knowledge; namely, *adaptive* and *collective* models. In the former, latent factors are learned in the source domain, and are integrated into a recommendation model in the target domain, while in the latter, latent factors are learned simultaneously optimizing an objective function that involves both domains. [292] addressed the sparsity problem in the target domain following the adaptive approach, proposing to exploit user and item information from auxiliary domains where user feedback may be represented differently. In particular, they studied the case in which users express binary like/dislike preferences in the source domain, and utilize 1-5 ratings in the target domain. Their approach performs singular value decomposition (SVD) in each auxiliary domain, in order to separately compute user and item latent factors, which are then shared with the target domain. Specifically, transferred factors are integrated into the factorization of the rating matrix in the target domain and added as regularization terms so that specific characteristics of the target domain can be captured. Latent factors can also be shared in a collective way, as studied by [291]. In this case, instead of learning latent features from the source domains and transferring them to the target domain, the authors proposed to learn the latent features simultaneously in all the domains. Both user and item factors are assumed to generate the observed ratings in every domain, and, thus, their corresponding random variables are shared between the probabilistic factorization models of each rating matrix. Moreover, the factorization method is further extended by incorporating another set of factors that capture domain-dependent information, resulting in a tri-factorization scheme. A limitation of the

proposed approach is that the users and items from the source and target domains have to be identical. Instead of focusing on sharing latent factors, [137], and [141] studied the influence of social tags on rating prediction, as a knowledge transfer approach for cross-domain recommendations. The authors presented a number of models based on the SVD++ algorithm [221] to incorporate the effect of tag assignments into rating estimation. The underlying hypothesis is that information about item annotation in a source domain can be exploited to improve rating prediction in a target domain, as long as a set of common tags between the domains exists. In the proposed models, tag factors are added to the latent item vectors, and are combined with user latent features to compute rating estimations. The difference between these models is in the set of tags considered for rating prediction. In all the models knowledge transfer is performed through the shared tag factors in a collective way, since these are computed jointly for the source and the target domains. [182] presented a more complex approach that takes domain factors into account. There, the authors argue that user-item dyadic data cannot fully capture the heterogeneity of items, and that modeling domain-specific information is essential to make accurate predictions in a setting where users typically express their preferences in a single domain. They referred to this problem as the *unacquainted world*, and proposed a tensor factorization algorithm to exploit the triadic user-item-domain data. In that method, rating matrices from several domains are simultaneously decomposed into shared user, item, and domain latent factors, and a genetic algorithm automatically estimates optimal weights of the domains. In a recent work, [426], the authors propose a two-step approach where the latent factors learned via MF for both the source and target domains are linked by training a deep neural network (DNN) representing their connections. Interestingly, the training process of the DNN is driven by the sparsity degrees of individual users and items in the source and target domains. Contextual and content-based information is exploited in [366] to cluster users in the source domain prior to a tensor factorization. The proposed Cross Domain- Multi Dimensional Tensor Factorization (CD-MDTF) mitigates the sparsity and cold-start problem by transferring the aggregated knowledge from the source domain to target domain. An approach based on linking and transferring

knowledge is proposed in [421] where the main assumption is that correspondences among entities in different domains are unknown but can be computed with a cost. Starting from this assumption, the authors propose a unified framework aimed at actively mapping entities in different domain and then transferring knowledge via collaborative filtering. This latter step leverages partial mappings among entities for knowledge transfer. The authors also show how to integrate in their framework various extended matrix factorization techniques in a transfer learning manner. An emphasis on the meaningfulness of the knowledge extracted from the source domain to the target domain is the main topic in [415]. A clustering step among users and items is performed both in the source and target before a matrix factorization. Then, by comparing the resulting matrices, it is possible to evaluate the consistency of the information transfer.

Rather than sharing user or item latent factors for knowledge transfer, a different set of approaches analyzes the structure of rating data at the community level. These methods are based on the hypothesis that even when their users and items are different, close domains are likely to have user preferences sampled with the same population. Therefore, latent correlations may exist between preferences of groups of users for groups of items, which are referred to as *rating patterns*. In this context, rating patterns can act as a bridge that relates the domains, such that knowledge transfer can be performed in either adaptive or collective manners. In the adaptive setting, rating patterns are extracted from a dense source domain [232, 155]. In the collective setting, data from all the domains are pulled together and jointly exploited, even though users and items do not overlap across domains [233]. In [172], the authors propose to alleviate the data sparsity problem in a target domain by transferring rating patterns from multiple incomplete source domains. The proposed approach extracts rating patterns from a sparse source domain that are eventually combined with collaborative filtering to approximate the target domain and predict missing values. In particular, they take into account the effects related to negative transfer to obtain a more robust recommendation.

6.3 Matrix factorization models for cross-domain recommendation

In this section we refer to state-of-art models and optimization techniques explained in Sections: 3.2.3, 3.2.3, 3.2.5, 3.2.4, 3.3.1. In contrast to previous works that rely on graph-based methods for exploiting semantic metadata, the proposed approach first computes inter-domain content-based similarities between the items, and then exploits these similarities to regularize the joint learning of matrix factorizations in the source and target domains. In particular, we present three alternative hybrid models that make different assumptions about the relationships between source and target domain item latent factors, simultaneously exploiting user preferences and item metadata.

Moreover, in the chapter we detail our adaptations of the fast alternating least squares training algorithm for matrix factorization proposed by [304], in order to deal with the increased complexity of our models, which not only learn the auxiliary source domain user preferences, but also the item metadata used to bridge the domains.

Items from different domains tend to have very diverse attributes that are not straightforward related. For instance, a book may be characterized by its *author* or by its *book genres*, and a movie can be described using its *cast*, *director* or *movie genres*. In fact, content-based features are often different between domains, and even when they refer to related concepts, such as *book genres* and *movie genres*, the features may not be directly aligned, e.g., *funny movies* vs. *comedy books*.

In order to overcome the heterogeneity of features of items from different domains, we propose to exploit Linked Data for linking entities from multiple and diverse domains. Specifically, we map the items in our datasets to entities in DBpedia, a multi-domain repository that provides a semantic-based, structured representation of knowledge in Wikipedia. In 6.4.1 we shall describe the process of mapping items to semantic entities from DBpedia. Once the items are mapped to their corresponding entities, we use the DBpedia graph to compute semantic similarities between such entities, mining both the attributes and the structure of the graph with seman-

tic relations. More specifically, we exploit the information in DBpedia to compute a semantic similarity matrix $\mathbf{S} \in \mathbb{R}^{|\mathcal{I}_S| \times |\mathcal{I}_T|}$ between the source domain items \mathcal{I}_S and the target domain items \mathcal{I}_T :

$$s_{ij} = \text{sim}(i, j), \quad i \in \mathcal{I}_S, j \in \mathcal{I}_T \quad (6.1)$$

In 6.4.4 we shall report recommendation performance results by using several semantic similarity metrics from the state of the art.

The computed inter-domain item similarities are then used to *link* the domains for cross-domain recommendation. In the cold start, when a user has rated a few (if any) items in the target domain, a recommender system could suggest the user with items in the target domain that are semantically similar to those the user liked in the source domain. Hence, the system could be effective only if there is an overlap of users between the domains. Moreover, even cold start users in the target domain should have some preferences in the source domain.

In the next subsections we present our three recommendation models based on the exploitation of semantic similarities to regularize item factors in MF, so that similar items from different domains tend to have similar parameters. In this way, even if the user’s preferences in the target domain are unknown, a recommender system could suggest the user with target items that are most similar to those she preferred in the source.

6.3.1 Regularization through similarity prediction

The first semantic-based matrix factorization cross-domain model we propose is based on the assumption that latent vectors of related items should explain the items semantic similarities, in addition to the users’ preferences. That is, we not only seek to predict the preferences $r_{ui} \approx \langle \vec{p}_u, \vec{q}_i \rangle$, but also the inter-domain similarities $s_{ij} \approx \langle \vec{q}_i, \vec{q}_j \rangle$, where $i \in \mathcal{I}_S$ and $j \in \mathcal{I}_T$.

Hence, our model jointly factorizes the rating and inter-domain item similarity matrices that link the source and target domains. Let $\mathcal{U} = \mathcal{U}_S \cup \mathcal{U}_T$ be the set of all users, which we assume overlaps between the domains, and let $\mathcal{I} = \mathcal{I}_S \cup \mathcal{I}_T$ be the set of all items, which we assume do not overlap. Our model learns a latent vector

$\vec{p}_u \in \mathbb{R}^k$ for each user $u \in \mathcal{U}$, but separately models source and target domain items \vec{q}_i and \vec{q}_j , with $i \in \mathcal{I}_S$ and $j \in \mathcal{I}_T$, as follows:

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \mathbf{Q}_S, \mathbf{Q}_T) = & \sum_{u \in \mathcal{U}} \sum_{a \in \mathcal{I}} c_{ua} (r_{ua} - \langle \vec{p}_u, \vec{q}_a \rangle)^2 \\ & + \lambda_C \sum_{i \in \mathcal{I}_S} \sum_{j \in \mathcal{I}_T} (s_{ij} - \langle \vec{q}_i, \vec{q}_j \rangle)^2 + \lambda \left(\|\mathbf{P}\|^2 + \|\mathbf{Q}_S\|^2 + \|\mathbf{Q}_T\|^2 \right) \end{aligned} \quad (6.2)$$

where \mathbf{Q}_S and \mathbf{Q}_T are matrices containing the item latent vectors as rows from the source and target domains, respectively. We note that the summation in the first term iterates over all items $a \in \mathcal{I}$ from both domains, as we want to factorize the source and target user-item preference matrices simultaneously. The cross-domain regularization parameter $\lambda_C > 0$ controls the contribution of the inter-domain semantic similarities; large values of the parameter will force items to have too similar latent vectors, whereas low values will result in limited transfer of knowledge between domains.

As in standard matrix factorization, we train our model using Alternating Least Squares. First, we fix \mathbf{Q}_S and \mathbf{Q}_T , and solve analytically for each \vec{p}_u by setting the gradient to zero. Since the user factors do not appear in the additional cross-domain regularization term, we obtain the same solution as for the baseline MF model (see 3.8):

$$\vec{p}_u = \left(\mathbf{Q}^\top \mathbf{C}^u \mathbf{Q} + \lambda \mathbf{I} \right)^{-1} \mathbf{Q}^\top \mathbf{C}^u \vec{r}_u \quad (6.3)$$

In order to simplify the notation, we have defined the matrix \mathbf{Q} as the row-wise concatenation of \mathbf{Q}_S and \mathbf{Q}_T . The matrix \mathbf{C}^u is a diagonal matrix with the confidence values c_{ua} for all $a \in \mathcal{I}$, and the vector \vec{r}_u contains the preferences of user u , again for all items $a \in \mathcal{I}$.

Next, we fix the user factors \mathbf{P} and the target domain item factors \mathbf{Q}_T , and compute the optimal values for the source domain item factors. Again, by setting the corresponding gradient to zero and solving analytically we obtain:

$$\vec{q}_i = \left(\mathbf{P}^\top \mathbf{C}^i \mathbf{P} + \lambda_C \mathbf{Q}_T^\top \mathbf{Q}_T + \lambda \mathbf{I} \right)^{-1} \left(\mathbf{P}^\top \mathbf{C}^i \vec{r}_i + \lambda_C \mathbf{Q}_T^\top \vec{s}_i \right) \quad (6.4)$$

As previously, the vector \vec{r}_i contains the preferences assigned to item i , and \vec{s}_i is the i -th row of the inter-domain semantic similarity matrix \mathbf{S} . Finally, we proceed as

before fixing \mathbf{P} and \mathbf{Q}_S to compute the optimal solution for the target domain item latent vectors:

$$\vec{q}_j = \left(\mathbf{P}^\top \mathbf{C}^j \mathbf{P} + \lambda_C \mathbf{Q}_S^\top \mathbf{Q}_S + \lambda \mathbf{I} \right)^{-1} \left(\mathbf{P}^\top \mathbf{C}^j \vec{r}_j + \lambda_C \mathbf{Q}_S^\top \vec{s}_j \right) \quad (6.5)$$

The computation of the optimal factors can be parallelized within each step, but the larger number of items to consider and the extra step required for the source domain greatly increase the training time with respect to the MF baseline. In order to address this issue, we adapt the fast RRI training algorithm for ALS proposed by [304]. Since the computation of the user factors is the same as in the original MF model, the procedure remains the same for the P-step. For the source domain Q-step, by inspecting 6.2 and 6.4, we note that the additional terms that arise from the inter-domain similarities can be treated just like user preferences as follows. For each source item i :

1. Generate examples for each rating r_{ui} as for baseline MF (see [304])
2. For each target item $j \in \mathcal{I}_T$:
 - Generate an input example $\vec{x}_j := \vec{q}_j$.
 - Use the similarity as the dependent variable, $y_j := s_{ij}$.
 - Use a constant confidence value $c_j := \lambda_C$.
 - The parameter to optimize is $\vec{w} := \vec{q}_j$.

The above procedure will produce the similarity terms of 6.4, which can be defined by means of the confidence matrix $\tilde{\mathbf{C}}^i = \lambda_C \mathbf{I}$. The procedure for the target domain Q-step is completely analogous.

6.3.2 Regularization based on item neighborhoods

Our second semantic-based matrix factorization cross-domain model exploits the item semantic similarities in a different fashion. Instead of forcing pairwise item interactions to reproduce the observed similarity values, the approach we present here leverages \mathbf{S} to regularize the item latent vectors, so that feature vectors of

similar items are pushed together in the latent space. Intuitively, items that are semantically similar should also have similar latent parameters.

As previously, let $\mathcal{U} = \mathcal{U}_S \cup \mathcal{U}_T$ and $\mathcal{I} = \mathcal{I}_S \cup \mathcal{I}_T$ be the sets of all users and items, respectively. Our approach jointly factorizes the source and target domain rating matrices, and regularizes similar item factors proportionally to the items similarity. However, instead of considering all the potentially similar source domain items, we limit the regularization of a target domain item $j \in \mathcal{I}_T$ to its neighborhood, i.e., to the set $N(j) \subseteq \mathcal{I}_S$ of the top- n most similar source domain items:

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \mathbf{Q}_S, \mathbf{Q}_T) = & \sum_{u \in \mathcal{U}} \sum_{a \in \mathcal{I}} c_{ua} (r_{ua} - \langle \vec{p}_u, \vec{q}_a \rangle)^2 \\ & + \lambda_C \sum_{j \in \mathcal{I}_T} \sum_{i \in N(j)} s_{ij} \|\vec{q}_j - \vec{q}_i\|^2 + \lambda \left(\sum_{u \in \mathcal{U}} \|\vec{p}_u\|^2 + \sum_{a \in \mathcal{I}} \|\vec{q}_a\|^2 \right) \end{aligned} \quad (6.6)$$

We note that items with greater similarity values are more heavily regularized, whereas items with values of $s_{ij} \approx 0$ in their neighborhoods are barely affected. However, it may still be convenient to regularize such items so that they benefit from cross-domain information, and thus may be eligible for recommendation to cold start users. Therefore, we also experiment normalizing the similarity scores in the item neighborhoods so that $\sum_{i \in N(j)} s_{ij} = 1$. In this way all target items are equally regularized, but each is affected by its source domain neighbors proportionally to their similarity scores.

By assigning latent vectors to target domain items close to those of similar source domain items, our model is able to generate recommendations in cold start settings. Specifically, let \vec{q}_j be the latent vector learned for target item $j \in \mathcal{I}_T$, and let \vec{q}_i be the latent vector of source item $i \in \mathcal{I}_S$, which we assume is semantically similar to j . Our model will regularize both factors so that their distance $\|\vec{q}_j - \vec{q}_i\|$ is small, or equivalently, $\vec{q}_j \approx \vec{q}_i$. Consider now a cold start user u who only provided preferences in the source domain, so that her corresponding latent vector \vec{p}_u is therefore only adjusted using source domain preferences. In standard MF, it is not guaranteed that \vec{p}_u will extrapolate to the target domain, and will provide an accurate prediction for \vec{q}_j . In contrast, our model ensures that $\langle \vec{p}_u, \vec{q}_j \rangle \approx \langle \vec{p}_u, \vec{q}_i \rangle$, i.e., target domain items yield relevance prediction scores close to that of similar

source domain items. Hence, u will be recommended with a target domain item j if the user liked the source domain item i , or if i would be recommended to u in the source domain.

Once more, we train our neighborhood-based matrix factorization model using Alternating Least Squares. As in the previous model, the user factors are not affected by the extra regularization, and can be computed again using 6.3, leaving the P-step unchanged. For the target domain item factors \vec{q}_j we proceed as usual, fixing the user and source item factors, and finding the values such that $\frac{\partial \mathcal{L}}{\partial \vec{q}_j} = 0$, which yields the solution:

$$\vec{q}_j = \left[\mathbf{P}^\top \mathbf{C}^j \mathbf{P} + \left(\lambda + \lambda_C \sum_{i \in N(j)} s_{ij} \right) \mathbf{I} \right]^{-1} \left(\mathbf{P}^\top \mathbf{C}^j \vec{r}_j + \lambda_C \sum_{i \in N(j)} s_{ij} \vec{q}_i \right) \quad (6.7)$$

Repeating the same procedure for the source item factors \vec{q}_i we obtain:

$$\vec{q}_i = \left[\mathbf{P}^\top \mathbf{C}^i \mathbf{P} + \left(\lambda + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij} \right) \mathbf{I} \right]^{-1} \left(\mathbf{P}^\top \mathbf{C}^i \vec{r}_i + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij} \vec{q}_j \right) \quad (6.8)$$

where $N^{-1}(i)$ is the *inverse neighborhood* of item i , i.e., the set of target domain items that have i among their neighbors: $N^{-1}(i) = \{j \in \mathcal{I}_T | i \in N(j)\}$.

Unlike the model presented in the previous section, we cannot apply RR1 directly by treating the new similarity terms as additional user preferences. Instead, we derive again the update rules for each component of the source and target domain item parameters. As mentioned before, user parameters remain unchanged. Let $j \in \mathcal{I}_T$ be a target item, and consider the optimization of the α -th component $q_{j\alpha}$ of its corresponding latent vector \vec{q}_j . We can rewrite the loss in 6.6 as a function only of $q_{j\alpha}$ as follows:

$$\begin{aligned} \mathcal{L}_\alpha(q_{j\alpha}) &= \sum_{u \in \mathcal{U}} c_{uj} (e_{uj} - p_{u\alpha} q_{j\alpha})^2 + \lambda q_{j\alpha}^2 \\ &\quad + \lambda_C \sum_{i \in N(j)} s_{ij} (q_{j\alpha} - q_{i\alpha})^2 + \text{constant} \end{aligned} \quad (6.9)$$

where $e_{uj} \triangleq r_{uj} - \sum_{\beta \neq \alpha} p_{u\beta} q_{j\beta}$, and the constant includes terms that do not depend

on $q_{j\alpha}$. If we set the derivative $\frac{d\mathcal{L}_\alpha}{dq_{j\alpha}} = 0$, we obtain:

$$q_{j\alpha} = \frac{\sum_{u \in \mathcal{U}} c_{uj} e_{uj} p_{u\alpha} + \lambda_C \sum_{i \in N(j)} s_{ij} q_{i\alpha}}{\sum_{u \in \mathcal{U}} c_{uj} p_{u\alpha}^2 + \lambda + \lambda_C \sum_{i \in N(j)} s_{ij}} \quad (6.10)$$

Using the optimizations described in [304], the computational cost of the above formula for all items is $\mathcal{O}(k^2|\mathcal{U}| + k|\mathcal{R}| + n|\mathcal{I}_T|)$, since all the neighborhoods are formed using the top n most similar items, $|N(j)| \leq n$. Applying the same procedure to the source domain item factor \vec{q}_i we obtain:

$$q_{i\alpha} = \frac{\sum_{u \in \mathcal{U}} c_{ui} e_{ui} p_{u\alpha} + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij} q_{j\alpha}}{\sum_{u \in \mathcal{U}} c_{ui} p_{u\alpha}^2 + \lambda + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij}} \quad (6.11)$$

The main difference with respect to 6.10 is that the sets $N^{-1}(i)$ are not bounded, as a source item can potentially be the neighbor of an arbitrary number of target items, so that $|N^{-1}(i)| \leq |\mathcal{I}_T|$, resulting in a theoretical worst-case cost of $\mathcal{O}(k^2|\mathcal{U}| + k|\mathcal{R}| + |\mathcal{I}_S||\mathcal{I}_T|)$. We observe, however, that in practice most of the source items appear only in a few neighborhoods and that the algorithm is still very efficient.

6.3.3 Regularization based on item centroids

When neighbor source domain items are mutually diverse, the neighborhood-based model presented in the previous section may struggle to regularize a target domain item that has to be simultaneously close to all its neighbors. The model we propose in this section works like the neighborhood-based model, but, instead of using the neighbor source domain items *individually* in the regularization, it uses their *centroid* (average) latent vector:

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \mathbf{Q}_S, \mathbf{Q}_T) = & \sum_{u \in \mathcal{U}} \sum_{a \in \mathcal{I}} c_{ua} (r_{ua} - \langle \vec{p}_u, \vec{q}_a \rangle)^2 \\ & + \lambda_C \sum_{j \in \mathcal{I}_T} \left\| \vec{q}_j - \sum_{i \in N(j)} s_{ij} \vec{q}_i \right\|^2 + \lambda \left(\sum_{u \in \mathcal{U}} \|\vec{p}_u\|^2 + \sum_{a \in \mathcal{I}} \|\vec{q}_a\|^2 \right) \end{aligned} \quad (6.12)$$

The same considerations regarding the neighborhood $N(j)$ and the normalization of the similarity scores also apply to this model. However, the effect on the item relevance predictions for cold start users is different. Let \vec{q}_j be an item in the target

domain, and let $N(j)$ be its neighborhood of most similar source domain items. The regularization scheme in our centroid-based approach aims to minimize the distance $\|\vec{q}_j - \sum_{i \in N(j)} s_{ij} \vec{q}_i\|$, so that the latent vector of item j is close, *on average*, to those of the source items in $N(j)$, i.e., $\vec{q}_j \approx \sum_{i \in N(j)} s_{ij} \vec{q}_i$. Let u be a cold start user in the target domain that has some preferences in the source domain. Again, her feature vector \vec{p}_u is only learned using the user's source preferences, and may not be reliable for computing relevance predictions for target domain items in standard MF. Our model, however, ensures that

$$\langle \vec{p}_u, \vec{q}_j \rangle \approx \left\langle \vec{p}_u, \sum_{i \in N(j)} s_{ij} \vec{q}_i \right\rangle = \sum_{i \in N(j)} s_{ij} \langle \vec{p}_u, \vec{q}_i \rangle$$

That is, the predicted relevance score is roughly the average of the relevance scores for the neighbor source domain items, weighted by their corresponding semantic similarity.

As in the previous models, the user parameters are not affected by the item regularization terms, and can be computed in the standard fashion using 6.3. For the target domain item factors $\vec{q}_j, j \in \mathcal{S}_T$, we set the gradient of 6.12 to zero to obtain:

$$\vec{q}_j = \left[\mathbf{P}^\top \mathbf{C}^j \mathbf{P} + (\lambda + \lambda_C) \mathbf{I} \right]^{-1} \left(\mathbf{P}^\top \mathbf{C}^j \vec{r}_j + \lambda_C \sum_{i \in N(j)} s_{ij} \vec{q}_i \right) \quad (6.13)$$

Comparing the above to 6.7 we observe that both are equivalent when $\sum_{i \in N(j)} s_{ij} = 1$, i.e., normalizing the similarity values has the same effect of than centroid-based regularization on the target domain item factors. The solution for source item factors \vec{q}_i , in contrast, has a different form:

$$\vec{q}_i = \left[\mathbf{P}^\top \mathbf{C}^i \mathbf{P} + \left(\lambda + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij}^2 \right) \mathbf{I} \right]^{-1} \cdot \left(\mathbf{P}^\top \mathbf{C}^i \vec{r}_i + \lambda_C \sum_{j \in N^{-1}(j)} s_{ij} (\vec{q}_j - \vec{z}_{j \setminus i}) \right) \quad (6.14)$$

where we have defined $\vec{z}_{j \setminus i} = \sum_{l \in N^{-1}(i), l \neq j} s_{lj} \vec{q}_l$ to simplify the notation. We note that, differently to the previous models, the computation of the source domain latent

vectors cannot be parallelized, as the value of $\vec{q}_i, i \in \mathcal{I}_S$ depends on the values of other $\vec{q}_l, l \in \mathcal{I}_S$ through the parameter $\vec{z}_{j \setminus i}$. As a result, the training process can be slow when the set of source domain items is large. In our experiments, however, we observed that the time penalty of computing the source factors sequentially is usually compensated by the faster RR1 algorithm, although we do not provide any quantitative analysis as it falls out of the scope of this work.

In order to apply RR1 to our centroid-based approach, we derive again the solutions for each α -th coordinate separately. Once more, the solution for the user factors remains the same as it is not affected by the regularization terms. For the target domain item factors \vec{q}_j , we consider the loss in 6.12 as a function only of the α -th component $q_{j\alpha}$:

$$\begin{aligned} \mathcal{L}_\alpha(q_{j\alpha}) &= \sum_{u \in \mathcal{U}} c_{uj} (e_{uj} - p_{u\alpha} q_{j\alpha})^2 + \lambda q_{j\alpha}^2 \\ &\quad + \lambda_C \left(q_{j\alpha} - \sum_{i \in N(j)} s_{ij} q_{i\alpha} \right)^2 + \text{constant} \end{aligned} \quad (6.15)$$

As previously, the constant includes terms that do not depend on $q_{j\alpha}$, and e_{uj} is defined as in 6.9. Setting the derivative $\frac{d\mathcal{L}_\alpha}{dq_{j\alpha}} = 0$ yields:

$$q_{j\alpha} = \frac{\sum_{u \in \mathcal{U}} c_{uj} e_{uj} p_{u\alpha} + \lambda_C \sum_{i \in N(j)} s_{ij} q_{i\alpha}}{\sum_{u \in \mathcal{U}} c_{uj} p_{u\alpha}^2 + \lambda + \lambda_C} \quad (6.16)$$

We note, once again, the similar form of the above solution with respect to the previous model in 6.10. If we apply the same procedure to the source domain item factors, we obtain:

$$q_{i\alpha} = \frac{\sum_{u \in \mathcal{U}} c_{ui} e_{ui} p_{u\alpha} + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij} (q_{j\alpha} - \vec{z}_{(j \setminus i)\alpha})}{\sum_{u \in \mathcal{U}} c_{ui} p_{u\alpha}^2 + \lambda + \lambda_C \sum_{j \in N^{-1}(i)} s_{ij}^2} \quad (6.17)$$

The computational complexity for the target domain factors is equivalent to the model from the previous section, whereas for the source domain factors it is $\mathcal{O}(k^2 |\mathcal{U}| + k |\mathcal{R}| + n |\mathcal{I}_S| |\mathcal{I}_T|)$ in the worst case, which is similar to the neighborhood-based model since the size of the neighborhoods n is in general small.

6.4 Experiments

In a first experiment, we compared several state-of-the-art semantic similarity metrics for content-based recommendation, aiming to understand which is more suitable for later injecting in our cross-domain MF models, and achieved the best results using the link-based approach by [262]. Second, we evaluated the ranking precision and diversity of the recommendations computed by the proposed models. We show that, depending on the involved source and target domains, our models generate more accurate suggestions than the baselines in severe cold start situations. Moreover, the proposed approaches provide a better trade-off between accuracy and diversity, which are in general difficult to balance.

6.4.1 Dataset

Our dataset initially consisted of a large set of *likes* assigned by users to items in Facebook. Using the Facebook Graph API, a user's *like* is retrieved in the form of a 4-tuple with the following information: the identifier, name and category of the liked item, and the timestamp of the like creation, e.g., `{id: "35481394342", name: "The Godfather", category: "Movie", created_time: "2015-05-14T12:35:08+0000"}`. The name of an item is given by the user who created the Facebook page of such item. In this context, distinct names may exist for a particular item, e.g., *The Godfather*, *The Godfather: The Movie*, *The Godfather - Film series*, etc. Users thus may express likes for different Facebook pages which actually refer to the same item. Aiming to unify and consolidate the items extracted from Facebook likes, we developed a method that automatically maps the items names with the unique URIs of the corresponding DBpedia entities, e.g., `http://dbpedia.org/resource/The_Godfather` for the identified names of *The Godfather* movie.

Linking items to DBpedia entities Given a particular item, we first identified DBpedia entities that are labeled with the name of the item. For such purpose, we

launched a SPARQL query targeted on the subjects of triples that have `rdfs:label`⁹ as property and the item title as object. The next query is an example for *The Matrix 2* title:

```
SELECT DISTINCT ?item WHERE {
  {
    ?item rdf:type dbo:Film .
    ?item rdfs:label ?name .
    FILTER regex(?name, "the.*matrix.*2", "i") .
  }
  UNION
  {
    ?item rdf:type dbo:Film .
    ?tmp dbo:wikiPageRedirects ?item .
    ?tmp rdfs:label ?name .
    FILTER regex(?name, "the.*matrix.*2", "i") .
  }
}
```

To resolve ambiguities in those names that correspond to multiple items belonging to different domains, we specify the type of the item we wanted to retrieve in each case. Specifically, the previous query includes a triple clause with `rdf:type`¹⁰ (or `dbo:type`¹¹) as property. Hence, in the given example, the subject *The Matrix 2* refers to the “movie” type, which is associated to the `dbo:Film` class in DBpedia. The item types were set from the item categories provided in Facebook, and their associated DBpedia and YAGO¹² classes¹³ were identified by manual inspection of the `rdf:type` values of several entities. 6.1 shows the list of item types and DBpedia/YAGO classes we considered for the three domains of our dataset.

⁹Namespace for rdfs, <http://www.w3.org/2000/01/rdf-schema>

¹⁰Namespace for rdf, <http://www.w3.org/1999/02/22-rdf-syntax-ns#>

¹¹Namespace for dbo, <http://dbpedia.org/ontology>

¹²The YAGO knowledge base, <http://www.mpi-inf.mpg.de/yago-naga/yago>

¹³Namespace for yago, <http://dbpedia.org/class/yago>

Table 6.1: Considered item types and their DBpedia and YAGO classes for the three domains of the dataset.

| | Item type | DBpedia/YAGO classes |
|--------|---------------------|--|
| Books | Book | dbo:Book, yago:Book102870092, yago:Book102870526 |
| | Genre | yago:LiteraryGenres |
| | Writer | dbo:Writer, yago:Writer110794014 |
| | Fictional character | dbo:FictionalCharacter, yago:FictionalCharacter109587565 |
| Movies | Movie | dbo:Film, yago:Movie106613686 |
| | Genre | dbo:MovieGenre, yago:FilmGenres |
| | Director | yago:FilmDirector110088200, yago:Director110014939 |
| | Actor | dbo:Actor, yago:Actor109765278 |
| | Fictional character | dbo:FictionalCharacter, yago:FictionalCharacter109587565 |
| Music | Composition | dbo:Song, dbo:MusicalWork, dbo:Single, dbo:ClassicalMusicComposition, dbo:Opera |
| | Genre | dbo:MusicGenre, yago:MusicGenres, yago:MusicGenre107071942 |
| | Album | dbo:Album, yago:Album106591815 |
| | Musician | dbo:MusicalArtist, yago:Musician110339966, yago:Musician110340312, yago:Composer109947232 |
| | Band | dbo:Band, yago:MusicalOrganization108246613 |

Moreover, running the previous query template we observed that a number of items were not linked to DBpedia entities because the labels corresponded to Wikipedia *redirection* webpages. In these cases, to reach the appropriate entities the query makes use of the `dbo:wikiPageRedirects` property. The result of the previous query for *The Matrix 2* is `http://dbpedia.org/resource/The_Matrix_Reloaded`, which actually is the DBpedia entity of the second movie in *The Matrix* saga. Here, it is important to note that thanks to the Wikipedia page redirect component we were able to link items whose names do not have a direct syntactic match with the label of its DBpedia entity, but with the label of a redirected entity, e.g., the *Matrix 2* title matches the `The Matrix Reloaded` entity.

Final semantically annotated dataset For every linked entity, we finally accessed DBpedia to retrieve the metadata that afterward will be used as input for the recommendation models. In this case, we launched a SPARQL query asking for all the properties and objects of the triples that have the target entity as subject. Following the example given before, such a query would be:

```
SELECT ?p ?o WHERE {  
  dbr:The_Matrix_Reloaded ?p ?o .  
}
```

This query returns all the DBpedia property-value pairs of the `dbr:The_Matrix_Reloaded`¹⁴ entity. However, since our ultimate goal is item recommendation, we should only exploit metadata that may be relevant to relate common preferences of different users. Thus, we filtered the query results by considering certain properties in each domain. Specifically, 6.2 shows the list of DBpedia properties selected for each of the three domains of our dataset. Hence, for example, for the movie items, we would have as metadata the movies genres, directors, and actors, among others.

The items and relations shown in the table thus represent a *semantic network* that is automatically obtained from DBpedia for each particular domain. 6.3 shows

¹⁴Namespace for `dbr`, `http://dbpedia.org/resource`

statistics of the dataset for the three domains of interest, namely books, movies, and music. Additionally, users may express preferences in more than one domain. 6.4 shows the number of users shared between each pair of domains.

Semantically enriched item profiles Fixing *books, movies, musicians* and *bands* as the target items to be recommended, we can distinguish the following three types of item metadata obtained:

- *attributes*, which correspond to item-attribute entities associated to the considered item types of 6.2, and are distinct to the entities of target items, e.g., the genre(s), director(s) and actors of a particular movie.
- *related items*, which correspond to the item-item properties in 6.2 that derive related entities, e.g., the novel a movie is based on (`dbo:basedOn` property), the prequel/sequel of a movie (`dbo:previousWork` / `dbo:subsequentWork` properties), or the musicians belonging to a band (`dbo:bandMember` property).
- *extended attributes*, which correspond to attribute-attribute properties that generate extended item attributes, originally not appearing as metadata, e.g., the subgenres of a particular music genre (`dbo:musicSubgenre` property).

The above three types of item metadata constitute the semantically enriched item profiles that we propose to use in our recommendation models. We note that they differ from the commonly used content-based item profiles composed of plain attributes. We also note that in the conducted experiments, the results achieved by exploiting the enriched profiles were better than those achieved by only using item attributes.

6.4.2 Evaluation methodology and metrics

The evaluation of the proposed models was conducted utilizing a modified user-based 5-fold cross-validation strategy, based on the methodology by [217] for cold

Table 6.2: DBpedia properties considered as item metadata; *item* can be book, movie and composition, musician and band.

| Relation | DBpedia properties |
|------------------------------------|--|
| item – genre | dct:subject, dbo:genre |
| book – genre | dbo:literaryGenre |
| music genre – music genre | dbo:musicSubgenre, dbo:musicFusionGenre, dbo:movement, dbo:derivative, dbo:stylisticOrigin |
| item – author | dbo:author, dbo:creator |
| book – writer | dbo:writer |
| movie – actor, character, director | dbo:starring, dbo:cinematography, dbo:director |
| composition – musician | dbo:artist, dbo:composer, dbo:musicComposer, dbo:musicalArtist, dbo:associatedMusicalArtist |
| music item – album | dbo:album |
| band – musician | dbo:bandMember, dbo:formerBandMember, dbo:musicalBand, dbo:associatedBand |
| item – item, character | dbo:series |
| item – character | dbo:portrayer |
| item – item | dbo:basedOn, dbo:previousWork, dbo:subsequentWork, dbo:notableWork |

Table 6.3: Statistics of the extracted dataset enriched with metadata.

| | Books | Movies | Music |
|-----------------|--------------|---------------|--------------|
| Users | 1876 | 26943 | 49369 |
| Items | 3557 | 3901 | 5748 |
| Likes | 42869 | 876501 | 2084462 |
| Sparsity (%) | 99.4 | 99.2 | 99.3 |
| Avg. items/user | 22.85 | 32.53 | 42.22 |
| Avg. users/item | 12.05 | 224.69 | 362.64 |

Table 6.4: User overlap between domains. Right to each target, the ratio of shared users relative to the source domain.

| Source | Target | | | | | |
|---------------|---------------|----------|---------------|----------|--------------|----------|
| | Books | % | Movies | % | Music | % |
| Books | 1876 | 100.0 | 1495 | 79.7 | 1519 | 81.0 |
| Movies | 1495 | 5.5 | 26943 | 100.0 | 21720 | 80.6 |
| Music | 1519 | 3.1 | 21720 | 44.0 | 49369 | 100.0 |

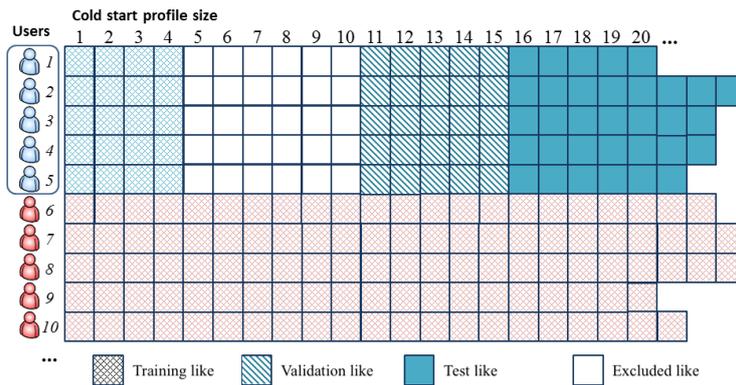


Figure 6.1: Overview of the cold start evaluation setting in a given cross-validation fold. The box indicates the test users in the current fold, whose profiles are split into training, validation, and testing sets. Different cold start profile sizes are simulated by sequentially adding *likes* to their training sets —four in the figure.

start evaluation. Our goal is to understand how the different approaches perform as the number of observed *likes* in the target domain increases. First, we divide the set of users into five subsets of roughly equal size. In each cross-validation stage, we keep all the data from four of the groups in the training set. Then, for each user u in the fifth group –the test users– we randomly split her *likes* into three subsets, as depicted in 6.1:

1. *Training data*, initially filled with u 's *likes* and iteratively downsampled discarding one by one to simulate different cold start profile sizes,
2. *Validation data* containing the set of *likes* used for tuning hyperparameters, and
3. *Testing data* used to compute the performance metrics.

The above procedure was modified for the cross-domain scenario by extending the training set with the full set of *likes* from the auxiliary domain, in order to obtain the actual training data for the predictive models. For each cold start profile size, we built the recommendation models using the data in the final training set. Then, for each test user, we generated a ranked list of the top 10 suggested items from

the set of target domain items in the training set that are not yet known to the user. The performance is estimated from the output of each model and the test set using rank-based metrics. We note that in our evaluation, any item ranked after position 10 by the model is considered not relevant when computing the metrics, as we are interested in the more realistic setting where the user only examines a limited subset of the recommendations.

Regarding the metrics, we used the Mean Reciprocal Rank (MRR) to evaluate the ranking accuracy of the recommendations, which computes the average reciprocal rank of the first relevant item in the recommendation list. Binomial Diversity Framework (BinomDiv) [376] was used to evaluate the individual diversity, namely the degree of diversity in the recommendation lists based on item genres extracted from DBpedia.

6.4.3 Evaluated methods

We compared the performance of our proposed methods against the following baseline algorithms:

- **POP.** Non personalized baseline that always recommends the most popular items not yet liked by the user. Popularity is measured as the number of users in the dataset that liked the item.
- **UNN.** User-based nearest neighbors with Jaccard similarity. The size of the neighborhood is tuned for each dataset using a validation set.
- **INN.** Item-based nearest neighbors with Jaccard similarity and indefinite neighborhood size.
- **iMF.** Matrix factorization method for positive-only feedback [184] trained using the fast ALS technique by [304].
- **BPR.** Bayesian personalized ranking from implicit feedback [312]. We used for our experiments the implementation available in LibRec [164].

- **FISM**. Factored item similarity model by [205]. We used the implementation of the FISMauc variant optimized for the item ranking problem available in LibRec [164].
- **HeteRec**. Graph-based recommender system proposed in [405], based on a diffusion method of user preferences following different meta-paths.
- **SPRank**. Originally proposed in [282], it implements a hybrid approach to compute recommendations with LOD datasets. We used a publicly available implementation of SPRank¹⁵.

With the exception of POP (which only uses target domain data) and SPRank, we considered the application of all the baselines to both single- and cross-domain scenarios. We were not able to compute meaningful results for SPRank by using DBpedia properties shown in Table 6.2 due to the structure of the connections between domains in the underlying knowledge graph. All the paths calculated by SPRank to link items in different domains resulted in being not very relevant thus bringing to shallow performances of the algorithm. Moreover, given the datasets adopted for the experimental evaluation, we were not able to generate all the meta-path needed to compute recommendations. We used machines with up to 3 TB of disk space but it was not sufficient.

Hereafter we use the prefix CD- to indicate that the algorithm is operating in cross-domain mode using the union of the rating matrices from the source and target domains. We did not consider for our evaluation the SemanticSVD++ method by [321], as it is designed for rating prediction rather than item ranking. Moreover, preliminary tests showed that its performance was much lower than the other methods, and that its training time was about one order of magnitude larger.

We evaluated the three methods presented in this study:

- **SimMF**. Our matrix factorization model regularized with similarity prediction described in 6.3.1.

¹⁵<https://github.com/sisinflab/lodreclib>

- **NeighborMF**. Our proposed matrix factorization model with neighborhood-based regularization from 6.3.2.
- **CentroidMF**. Our matrix factorization model from 6.3.3 that uses the neighbor’s centroid to regularize the target domain item factors.

We tuned the hyperparameters of the considered recommendation models using a held-out validation set of likes, as we explain in the next section. For UNN, we only had to select the size of the user neighborhoods. For the matrix factorization models, in contrast, the number of hyperparameters is larger, namely, the dimensionality of the latent factor space k , the amount of regularization λ , and the confidence parameter for positive-only feedback α . Moreover, the models proposed also include the cross-domain regularization rate λ_C , which controls the contribution of the inter-domain item similarities. Finally, for NeighborMF and CentroidMF, we tuned the size n of the item neighborhoods $N(j)$, and the possibility to normalize the neighbors’ similarities so that the sum to 1, as explained in 6.3.2.

The high number of parameters to tune rules out the possibility of performing a grid search for the best values. Hence, we used Bayesian Optimization techniques [353] that train Machine Learning models to predict candidate values that are likely to maximize a given function while simultaneously reducing the uncertainty of over unknown parameter values.

We tuned the parameters of the single-domain methods only on the target domain, and used the same values for their cross-domain variants. For UNN, the optimal number of neighbors was $n = 50$ for books, and $n = 100$ for movies and music. For iMF we obtained the optimal parameters $k = (10, 29, 21)$, $\lambda = (10^{-5}, 0.823, 1)$, and $\alpha = (6, 7, 10)$ for books, movies, and music, respectively. For BPR we used $\lambda = 0.01$ for regularization and $\eta = 0.01$ as learning rate. In the case of FISM, we used $\lambda = 0.001$ and $\eta = 10^{-5}$. The optimal values for our proposed cross-domain models are reported in 6.5.

Table 6.5: Optimal hyperparameters for SimMF, NeighborMF, and CentroidMF. The last column indicates whether the similarities in the neighborhood are normalized or not.

| | Source | Method | k | λ | α | λ_C | n | Norm. |
|---------------|---------------|---------------|-----|-----------|----------|----------------------|-----|--------------|
| Books | Movies | SimMF | 112 | 0 | 1 | 10^{-8} | | |
| | | NeighborMF | 134 | 1 | 1 | 9.125 | 49 | ✓ |
| | | CentroidMF | 153 | 0.999 | 1 | 8.778 | 100 | ✓ |
| | Music | SimMF | 10 | 1 | 16 | 10^{-8} | | |
| | | NeighborMF | 10 | 0 | 18 | 10 | 100 | ✓ |
| | | CentroidMF | 10 | 0 | 14 | 0.109 | 100 | |
| Movies | Books | SimMF | 12 | 1 | 1 | 0.002 | | |
| | | NeighborMF | 12 | 1 | 1 | 10 | 81 | ✓ |
| | | CentroidMF | 14 | 0.100 | 1 | 0.200 | 1 | ✓ |
| | Music | SimMF | 35 | 0 | 1 | 1.6×10^{-6} | | |
| | | NeighborMF | 51 | 1 | 1 | 10 | 100 | |
| | | CentroidMF | 29 | 1 | 1 | 9.494 | 99 | ✓ |
| Music | Books | SimMF | 10 | 1 | 1 | 0.039 | | |
| | | NeighborMF | 10 | 0.995 | 1 | 3.014 | 100 | ✓ |
| | | CentroidMF | 10 | 0.724 | 1 | 1.673 | 14 | |
| | Movies | SimMF | 11 | 0.571 | 4 | 0.641 | | |
| | | NeighborMF | 10 | 0.978 | 2 | 0.699 | 46 | |
| | | CentroidMF | 10 | 0.562 | 2 | 10 | 3 | ✓ |

6.4.4 Results

In this section we present the results of the conducted experiments to evaluate the proposed matrix factorization models. First, we analyze several semantic relatedness metrics to compute the inter-domain item similarities. Next, we report the ranking accuracy and diversity of the evaluated recommendation approaches, and study how the size and diversity of the source domain user profile impacts on the target recommendations.

6.4.5 Inter-domain item semantic similarity

The goal of our first experiment is to analyze the performance of several semantic relatedness metrics to compute the inter-domain similarities that we later exploit in our matrix factorization models. We considered the following strategies:

- **TF-IDF.** We use the semantically-enriched item profiles (see 6.4.1 to build TF-IDF vector profiles based on the metadata of each item. The similarity score between a source domain item and a target domain item is computed as the cosine of their corresponding TF-IDF vectors.
- **ESA.** The Explicit Semantic Analysis technique proposed by [153]. Instead of using the semantic metadata, we map each item to its corresponding Wikipedia article. Then, based on the text of the article, ESA extracts a set of other related Wikipedia articles, which represent semantic concepts, and builds a TF-IDF profile from the extracted concepts. Finally, the similarity score between two items is computed as the cosine of their corresponding concept-based vectors.
- **M&W.** The approach proposed by [262] computes the semantic relatedness between two items using the overlap of their sets of inlinks and outlinks in the Wikipedia hyperlink graph.
- **Katz.** Based on Katz’s centrality measure, the relatedness between two items is computed as the accumulated probability of the top shortest paths between their corresponding entities in the semantic network [185].

Table 6.6: MRR of the evaluated semantic relatedness metrics.

| Source | Target | TF-IDF | ESA | M&W | Katz |
|--------|--------|--------------|-------|--------------|-------|
| Books | Movies | 0.058 | 0.030 | 0.123 | 0.092 |
| | Music | 0.028 | 0.015 | 0.042 | 0.022 |
| Movies | Books | 0.054 | 0.011 | 0.031 | 0.013 |
| | Music | 0.030 | 0.011 | 0.028 | 0.009 |
| Music | Books | 0.010 | 0.006 | 0.052 | 0.020 |
| | Movies | 0.013 | 0.018 | 0.088 | 0.006 |

We evaluated the previous semantic relatedness metrics indirectly by comparing their performance in the item recommendation task. For such purpose, we chose a content-based recommendation model with no parameters, so that we can fairly measure the effect of each similarity on the item ranking quality. According to this simple model, the relevance score of an item is computed as the accumulated similarity with the items in the user’s profile:

$$s(u, i) = \sum_{j \in I(u)} s_{ij} \quad (6.18)$$

where s_{ij} is computed any of the methods described above.

The results of our experiment are shown in 6.6. For easier comparison according the methodology from 6.4.2, we averaged the MRR scores for all the cold start sizes in each source-target domain combination. We conclude from the table that M&W is the best performing metric, beating all the other approaches except when considering the movie domain as source, in which case it is still competitive. Hence, in the following experiments we evaluate our proposed matrix factorization models using M&W as the backing semantic similarity. Finally, we note that the low values for MRR are due to the simple recommendation algorithm chosen for this experiment.

6.4.6 Item ranking accuracy

In our second experiment we analyze the accuracy of the item rankings generated by the evaluated recommendation approaches. We aim to understand if cross-domain

variants are in general more effective than single-domain ones, and whether the proposed matrix factorization models are able to outperform the other methods in cold start settings.

6.7 shows the ranking accuracy for book recommendations in terms of MRR. We report the average results for cold start user profiles from sizes 6–10, as we observed that in those cases the trends are stable and, in general, single-domain baselines start to be effective. We remark that, according to the evaluation methodology described in 6.4.2, the number of test users remains constant regardless of the profile size, which we control by iteratively downsampling the training portion of their profile (see 6.1).

We notice from the table that, in general, approaches exploiting cross-domain movies or music preferences provide better recommendations than their single-domain counterparts. In case auxiliary movie preferences are available, we observe that the proposed NeighborMF and CentroidMF models achieve the best performance when only 1–3 book likes are observed. Moreover, in that case, our cross-domain matrix factorization models perform much better than the single-domain baselines. However, once 4 likes are available, CD-INN and single-domain HeteRec are more effective approaches. When the auxiliary preferences consist of music likes, we see that CD-INN is the overall best method, although it is only useful for profiles of size 1. For larger profiles, it is better to use single-domain baselines than any cross-domain method that uses music preferences. In summary, we conclude that music preferences are not useful for book recommendations, whereas movie likes could be used to improve the performance, specially with NeighborMF and CentroidMF for 1–3 book likes. We observe the bad performance of SPRank in cold start situations compared to the other baselines.

In 6.8 we show the results for movie recommendations. We observe that most of the cross-domain approaches are able to provide recommendations better than the most popular items for completely new movie users, and that CD-HeteRec is clearly the best performing approach. If the auxiliary cross-domain data consists of book preferences, we notice that the proposed matrix factorization models outperform the best single-domain baselines. However, in this situation CD-INN is even a better method, clearly providing more accurate recommendations than any other approach

Table 6.7: Accuracy (MRR) for cold start users in the target books domain. The three groups of rows correspond to single-domain, cross-domain with movies as source, and cross-domain with music as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of book <i>likes</i> | | | | | | | |
|---------|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.242 | 0.244 | 0.246 | 0.248 | 0.251 | 0.252 | 0.260 | |
| UNN | | 0.222 | 0.265 | 0.286 | 0.289 | 0.290 | 0.322 | |
| INN | | 0.145 | 0.177 | 0.216 | 0.241 | 0.262 | 0.316 | |
| iMF | | 0.171 | 0.194 | 0.235 | 0.255 | 0.271 | 0.301 | |
| BPR | | 0.110 | 0.116 | 0.136 | 0.154 | 0.157 | 0.193 | |
| FISM | | 0.228 | 0.230 | 0.234 | 0.234 | 0.238 | 0.245 | |
| HeteRec | | 0.218 | 0.244 | 0.279 | 0.297 | 0.316 | 0.351 | |
| SPRank | | 0.048 | 0.055 | 0.070 | 0.065 | 0.062 | 0.059 | |
| Movies | CD-UNN | 0.186 | 0.148 | 0.170 | 0.175 | 0.189 | 0.190 | 0.212 |
| | CD-INN | 0.262 | 0.265 | 0.275 | 0.291 | 0.301 | 0.307 | 0.339 |
| | CD-iMF | 0.261 | 0.262 | 0.268 | 0.272 | 0.275 | 0.274 | 0.287 |
| | CD-BPR | 0.217 | 0.200 | 0.218 | 0.237 | 0.235 | 0.238 | 0.251 |
| | CD-FISM | 0.235 | 0.228 | 0.225 | 0.231 | 0.236 | 0.235 | 0.245 |
| | CD-HeteRec | 0.264 | 0.248 | 0.261 | 0.268 | 0.278 | 0.277 | 0.298 |
| | SimMF | 0.253 | 0.268 | 0.274 | 0.284 | 0.289 | 0.290 | 0.296 |
| | NeighborMF | 0.253 | 0.272 | 0.282 | 0.294 | 0.293 | 0.293 | 0.301 |
| | CentroidMF | 0.252 | 0.271 | 0.283 | 0.289 | 0.293 | 0.295 | 0.301 |
| Music | CD-UNN | 0.136 | 0.103 | 0.115 | 0.120 | 0.138 | 0.140 | 0.157 |
| | CD-INN | 0.259 | 0.260 | 0.266 | 0.278 | 0.296 | 0.302 | 0.329 |
| | CD-iMF | 0.259 | 0.261 | 0.262 | 0.264 | 0.266 | 0.270 | 0.282 |
| | CD-BPR | 0.218 | 0.199 | 0.199 | 0.216 | 0.228 | 0.228 | 0.250 |
| | CD-FISM | 0.230 | 0.228 | 0.227 | 0.229 | 0.236 | 0.233 | 0.245 |
| | CD-HeteRec | 0.266 | 0.249 | 0.251 | 0.259 | 0.270 | 0.267 | 0.281 |
| | SimMF | 0.255 | 0.259 | 0.258 | 0.264 | 0.268 | 0.273 | 0.281 |
| | NeighborMF | 0.253 | 0.258 | 0.258 | 0.263 | 0.267 | 0.273 | 0.280 |
| | CentroidMF | 0.255 | 0.259 | 0.260 | 0.264 | 0.267 | 0.273 | 0.281 |

Table 6.8: Accuracy (MRR) for cold start users in the target movies domain. The three groups of rows correspond to single-domain, cross-domain with books as source, and cross-domain with music as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of movie likes | | | | | | | |
|---------|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.285 | 0.287 | 0.289 | 0.292 | 0.294 | 0.297 | 0.305 | |
| UNN | | 0.332 | 0.320 | 0.318 | 0.330 | 0.348 | 0.405 | |
| INN | | 0.233 | 0.300 | 0.336 | 0.359 | 0.377 | 0.413 | |
| iMF | | 0.256 | 0.291 | 0.314 | 0.334 | 0.348 | 0.388 | |
| BPR | | 0.225 | 0.256 | 0.276 | 0.299 | 0.315 | 0.350 | |
| FISM | | 0.257 | 0.265 | 0.263 | 0.266 | 0.267 | 0.270 | |
| HeteRec | | 0.315 | 0.346 | 0.357 | 0.366 | 0.374 | 0.395 | |
| SPRank | | 0.107 | 0.131 | 0.139 | 0.142 | 0.140 | 0.150 | |
| Books | CD-UNN | 0.219 | 0.169 | 0.185 | 0.219 | 0.256 | 0.292 | 0.385 |
| | CD-INN | 0.344 | 0.347 | 0.371 | 0.386 | 0.398 | 0.410 | 0.435 |
| | CD-iMF | 0.267 | 0.298 | 0.325 | 0.347 | 0.365 | 0.377 | 0.413 |
| | CD-BPR | 0.018 | 0.189 | 0.237 | 0.254 | 0.278 | 0.298 | 0.326 |
| | CD-FISM | 0.338 | 0.267 | 0.263 | 0.283 | 0.273 | 0.287 | 0.282 |
| | CD-HeteRec | 0.479 | 0.320 | 0.349 | 0.359 | 0.367 | 0.375 | 0.396 |
| | SimMF | 0.328 | 0.334 | 0.348 | 0.361 | 0.371 | 0.382 | 0.409 |
| | NeighborMF | 0.330 | 0.335 | 0.348 | 0.361 | 0.371 | 0.383 | 0.409 |
| | CentroidMF | 0.329 | 0.332 | 0.346 | 0.359 | 0.371 | 0.378 | 0.408 |
| Music | CD-UNN | 0.387 | 0.282 | 0.305 | 0.320 | 0.334 | 0.348 | 0.383 |
| | CD-INN | 0.342 | 0.347 | 0.353 | 0.359 | 0.365 | 0.371 | 0.390 |
| | CD-iMF | 0.301 | 0.326 | 0.344 | 0.362 | 0.374 | 0.385 | 0.418 |
| | CD-BPR | 0.352 | 0.305 | 0.316 | 0.332 | 0.332 | 0.343 | 0.361 |
| | CD-FISM | 0.105 | 0.089 | 0.093 | 0.089 | 0.091 | 0.091 | 0.093 |
| | CD-HeteRec | 0.367 | 0.336 | 0.344 | 0.350 | 0.355 | 0.360 | 0.374 |
| | SimMF | 0.339 | 0.351 | 0.361 | 0.374 | 0.384 | 0.396 | 0.419 |
| | NeighborMF | 0.353 | 0.364 | 0.374 | 0.385 | 0.394 | 0.404 | 0.427 |
| | CentroidMF | 0.345 | 0.355 | 0.367 | 0.377 | 0.385 | 0.395 | 0.418 |

from profile sizes 1–10. This is due to the high degree of overlap between the users of books and movies domains (79.7%, see 6.4), which allows CD-INN to compute very accurate item similarities based on the patterns of likes. Instead, when the source domain contains music preferences, we see that NeighborMF, CentroidMF, and SimMF, in that order, are consistently the best performing approaches for sizes 1–10. By regularizing item factors independently, NeighborMF is able to transfer source domain knowledge more effectively, which we also note is due to the greater contribution of cross-domain information (larger values of λ_C in 6.5). In summary, both book and music preferences are helpful for cold start movie recommendations, while our models are more effective when exploiting auxiliary music likes. On a side note, we observe the better performance of UNN over POP on the single domain setting when only 1 like is available. Looking at the results, we found that this is caused by the Jaccard-based similarity, which favors neighbors with small profiles that have rated similar items with high probability. A discussion of this phenomenon is outside of the line of research, and we refer the reader to [50] for a detailed explanation. HeteRec, on the other hand, exploits additional information from item metadata to compute more accurate recommendations than POP while SPRank confirms its bad behavior in cold start scenarios.

Finally, the results for music recommendations are shown in 6.9. As previously, CD-HeteRec is a very good performing approach to provide recommendations for completely new users, in both cross-domain configurations. Once 2 music likes are available, CD-INN is clearly the most competitive approach, independently of the used source domain. Again, we argue that this is due to the high number of music users who also have book and movie preferences, which allows CD-INN to compute very accurate rating-based similarities for items (see last column of 6.4). However, when the source domain consists of book preferences, we see that the proposed NeighborMF and CentroidMF models are slightly better than other cross-domain approaches if only 1 music like is provided. Anyway, even better performance can be achieved in this case simply using the single-domain UNN baseline, which does not need any extra information. Hence, single-domain baselines are compelling approaches for cold start music recommendations, and even though the proposed

Table 6.9: Accuracy (MRR) for cold start users in the target music domain. The three groups of rows correspond to single-domain, cross-domain with books as source, and cross-domain with movies as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of music likes | | | | | | | |
|---------|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.335 | 0.337 | 0.340 | 0.342 | 0.345 | 0.347 | 0.354 | |
| UNN | | 0.422 | 0.389 | 0.389 | 0.419 | 0.448 | 0.517 | |
| INN | | 0.320 | 0.391 | 0.426 | 0.455 | 0.474 | 0.517 | |
| iMF | | 0.347 | 0.396 | 0.427 | 0.451 | 0.471 | 0.517 | |
| BPR | | 0.330 | 0.377 | 0.409 | 0.432 | 0.450 | 0.488 | |
| FISM | | 0.096 | 0.100 | 0.100 | 0.100 | 0.101 | 0.100 | |
| HeteRec | | 0.358 | 0.395 | 0.421 | 0.442 | 0.463 | 0.510 | |
| SPRank | | N/A | N/A | N/A | N/A | N/A | N/A | |
| Books | CD-UNN | 0.290 | 0.244 | 0.266 | 0.300 | 0.344 | 0.387 | 0.487 |
| | CD-INN | 0.310 | 0.368 | 0.416 | 0.442 | 0.465 | 0.482 | 0.522 |
| | CD-iMF | 0.200 | 0.330 | 0.391 | 0.423 | 0.451 | 0.471 | 0.518 |
| | CD-BPR | 0.004 | 0.267 | 0.323 | 0.362 | 0.380 | 0.404 | 0.433 |
| | CD-FISM | 0.153 | 0.124 | 0.105 | 0.126 | 0.116 | 0.118 | 0.113 |
| | CD-HeteRec | 0.514 | 0.367 | 0.407 | 0.432 | 0.453 | 0.474 | 0.516 |
| | SimMF | 0.310 | 0.368 | 0.401 | 0.424 | 0.446 | 0.461 | 0.496 |
| | NeighborMF | 0.328 | 0.372 | 0.402 | 0.425 | 0.445 | 0.461 | 0.496 |
| | CentroidMF | 0.325 | 0.370 | 0.402 | 0.425 | 0.444 | 0.461 | 0.496 |
| Movies | CD-UNN | 0.435 | 0.274 | 0.306 | 0.336 | 0.369 | 0.400 | 0.484 |
| | CD-INN | 0.412 | 0.431 | 0.451 | 0.467 | 0.478 | 0.490 | 0.522 |
| | CD-iMF | 0.293 | 0.356 | 0.398 | 0.428 | 0.454 | 0.474 | 0.516 |
| | CD-BPR | 0.431 | 0.313 | 0.351 | 0.391 | 0.402 | 0.413 | 0.448 |
| | CD-FISM | 0.093 | 0.061 | 0.069 | 0.067 | 0.070 | 0.071 | 0.064 |
| | CD-HeteRec | 0.515 | 0.406 | 0.426 | 0.442 | 0.451 | 0.464 | 0.495 |
| | SimMF | 0.361 | 0.393 | 0.420 | 0.438 | 0.455 | 0.467 | 0.500 |
| | NeighborMF | 0.353 | 0.385 | 0.409 | 0.429 | 0.445 | 0.458 | 0.494 |
| | CentroidMF | 0.354 | 0.386 | 0.413 | 0.431 | 0.447 | 0.460 | 0.495 |

models are able to improve the quality of the item rankings by exploiting cross-domain item metadata, CD-INN, which is purely based on patterns of likes, is the best performing approach.

6.4.7 Recommendation diversity

In this subsection we analyze the diversity of the recommendation lists generated by the methods, as an alternative dimension of ranking quality.

6.10 shows the diversity of book recommendations in terms of the Binomial Diversity metric at cutoff 10 (BinomDiv@10). We observe that, in general, cross-domain approaches provide more diverse recommendations than their single-domain counterparts. However, we note several differences with respect to the accuracy results reported in 6.7. First, CD-UNN is consistently the superior algorithm in terms of diversity, whereas its accuracy results were the poorest among single- and cross-domain approaches. Second, when the source domain consists of movie likes, our proposed models achieve slightly worse diversity than other cross-domain approaches, specially for book profile sizes between 1–3 likes. This is in contrast with the results obtained in 6.7, where our methods performed best precisely in that range. We conclude that there is a clear trade-off between recommendation accuracy and diversity, and that the metric of interest depends on the particular application domain. We argue, however, that in cold start situations providing relevant suggestions may be more useful than recommending diverse, but not relevant items, if the ultimate goal of a system is to keep new users engaged.

The diversity results for movie recommendations are summarized in 6.11. We see that CD-FISM, CD-BPR, and CD-UNN provide the most diverse yet not relevant recommendations. Comparing the sources of auxiliary user preferences, we note that the diversity of the cross-domain baselines is roughly the same as their single-domain versions (comparing, e.g., HeteRec and CD-HeteRec) when considering book likes. In contrast, if the source domain contains music likes their diversity is significantly hurt. By comparing these results with 6.8 we observe once again the accuracy-diversity trade-off. Most methods' MRR greatly benefits from additional music likes at the expense of worse diversity. The exception is CD-FISM,

Table 6.10: Diversity (BinomDiv@10) for cold start users in the books domain. The three groups of rows correspond to single-domain, cross-domain with movies as source, and cross-domain with music as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of book <i>likes</i> | | | | | | | |
|---------|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.739 | 0.674 | 0.690 | 0.702 | 0.703 | 0.710 | 0.736 | |
| UNN | | 0.733 | 0.706 | 0.716 | 0.709 | 0.729 | 0.715 | |
| INN | | 0.655 | 0.674 | 0.654 | 0.665 | 0.672 | 0.669 | |
| iMF | | 0.583 | 0.606 | 0.630 | 0.645 | 0.657 | 0.664 | |
| BPR | | 0.696 | 0.700 | 0.715 | 0.690 | 0.698 | 0.696 | |
| FISM | | 0.513 | 0.692 | 0.708 | 0.686 | 0.706 | 0.719 | |
| HeteRec | | 0.609 | 0.623 | 0.653 | 0.672 | 0.680 | 0.693 | |
| SPRank | | 0.121 | 0.129 | 0.145 | 0.157 | 0.157 | 0.150 | |
| Movies | CD-UNN | 0.792 | 0.833 | 0.816 | 0.791 | 0.778 | 0.784 | 0.746 |
| | CD-INN | 0.740 | 0.676 | 0.683 | 0.684 | 0.680 | 0.692 | 0.695 |
| | CD-iMF | 0.724 | 0.660 | 0.674 | 0.689 | 0.686 | 0.686 | 0.702 |
| | CD-BPR | 0.514 | 0.458 | 0.484 | 0.425 | 0.454 | 0.465 | 0.471 |
| | CD-FISM | 0.453 | 0.464 | 0.479 | 0.488 | 0.495 | 0.499 | 0.533 |
| | CD-HeteRec | 0.747 | 0.673 | 0.672 | 0.680 | 0.690 | 0.704 | 0.709 |
| | SimMF | 0.702 | 0.649 | 0.671 | 0.676 | 0.682 | 0.690 | 0.706 |
| | NeighborMF | 0.690 | 0.652 | 0.660 | 0.671 | 0.680 | 0.682 | 0.702 |
| | CentroidMF | 0.699 | 0.647 | 0.659 | 0.668 | 0.684 | 0.686 | 0.702 |
| Music | CD-UNN | 0.744 | 0.811 | 0.797 | 0.771 | 0.746 | 0.734 | 0.731 |
| | CD-INN | 0.746 | 0.676 | 0.683 | 0.684 | 0.674 | 0.689 | 0.691 |
| | CD-iMF | 0.720 | 0.657 | 0.664 | 0.674 | 0.690 | 0.692 | 0.696 |
| | CD-BPR | 0.303 | 0.333 | 0.374 | 0.384 | 0.358 | 0.374 | 0.380 |
| | CD-FISM | 0.345 | 0.362 | 0.382 | 0.405 | 0.406 | 0.424 | 0.466 |
| | CD-HeteRec | 0.744 | 0.668 | 0.655 | 0.665 | 0.676 | 0.687 | 0.693 |
| | SimMF | 0.724 | 0.656 | 0.675 | 0.684 | 0.692 | 0.692 | 0.708 |
| | NeighborMF | 0.721 | 0.657 | 0.674 | 0.684 | 0.690 | 0.693 | 0.709 |
| | CentroidMF | 0.721 | 0.655 | 0.673 | 0.681 | 0.692 | 0.690 | 0.705 |

Table 6.11: Diversity (BinomDiv@10) for cold start users in the movies domain. The three groups of rows correspond to single-domain, cross-domain with books as source, and cross-domain with music as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of movie likes | | | | | | | |
|---------|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.401 | 0.304 | 0.336 | 0.354 | 0.368 | 0.378 | 0.399 | |
| UNN | | 0.360 | 0.385 | 0.404 | 0.392 | 0.396 | 0.394 | |
| INN | | 0.289 | 0.308 | 0.315 | 0.321 | 0.323 | 0.332 | |
| iMF | | 0.299 | 0.320 | 0.335 | 0.344 | 0.347 | 0.362 | |
| BPR | | 0.590 | 0.608 | 0.628 | 0.644 | 0.650 | 0.653 | |
| FISM | | 0.561 | 0.614 | 0.594 | 0.689 | 0.636 | 0.644 | |
| HeteRec | | 0.311 | 0.328 | 0.334 | 0.337 | 0.341 | 0.348 | |
| SPRank | | 0.218 | 0.242 | 0.260 | 0.262 | 0.254 | 0.269 | |
| Books | CD-UNN | 0.467 | 0.509 | 0.479 | 0.446 | 0.425 | 0.414 | 0.397 |
| | CD-INN | 0.327 | 0.291 | 0.314 | 0.323 | 0.329 | 0.331 | 0.339 |
| | CD-iMF | 0.341 | 0.294 | 0.317 | 0.327 | 0.333 | 0.338 | 0.350 |
| | CD-BPR | 0.646 | 0.677 | 0.645 | 0.668 | 0.624 | 0.643 | 0.666 |
| | CD-FISM | 0.548 | 0.549 | 0.671 | 0.574 | 0.609 | 0.625 | 0.664 |
| | CD-HeteRec | 0.316 | 0.310 | 0.328 | 0.335 | 0.337 | 0.341 | 0.348 |
| | SimMF | 0.308 | 0.265 | 0.297 | 0.307 | 0.320 | 0.325 | 0.339 |
| | NeighborMF | 0.315 | 0.266 | 0.298 | 0.306 | 0.321 | 0.325 | 0.338 |
| | CentroidMF | 0.313 | 0.273 | 0.302 | 0.315 | 0.326 | 0.334 | 0.348 |
| Music | CD-UNN | 0.368 | 0.404 | 0.386 | 0.376 | 0.373 | 0.372 | 0.376 |
| | CD-INN | 0.309 | 0.240 | 0.268 | 0.283 | 0.297 | 0.304 | 0.321 |
| | CD-iMF | 0.270 | 0.231 | 0.270 | 0.289 | 0.302 | 0.315 | 0.332 |
| | CD-BPR | 0.372 | 0.439 | 0.411 | 0.438 | 0.446 | 0.445 | 0.476 |
| | CD-FISM | 0.653 | 0.720 | 0.705 | 0.499 | 0.645 | 0.732 | 0.688 |
| | CD-HeteRec | 0.333 | 0.271 | 0.298 | 0.314 | 0.324 | 0.333 | 0.349 |
| | SimMF | 0.311 | 0.254 | 0.288 | 0.303 | 0.317 | 0.324 | 0.340 |
| | NeighborMF | 0.311 | 0.259 | 0.290 | 0.308 | 0.320 | 0.329 | 0.344 |
| | CentroidMF | 0.302 | 0.246 | 0.279 | 0.297 | 0.310 | 0.319 | 0.338 |

which follows the opposite trend: source music likes lead to significantly worse accuracy but improved diversity. We leave for future work an analysis of CD-FISM to understand which of its characteristics causes this behavior. Finally, we remark the good performance of the NeighborMF method when source music likes are exploited, as it is able to provide a good trade-off of decent diversity and the most accurate recommendations (see 6.8).

Last, we report the diversity results for music recommendations in 6.12. Once again, CD-FISM, which achieved the poorest accuracy in 6.9, provides the most diverse recommendations for all music profile sizes in the 1–10 range. However, for completely new users, we highlight the very good performance of CD-HeteRec, which not only is able to generate diverse recommendations, but also achieved the best accuracy results in terms of MRR. The remaining cross-domain approaches are in general worse than single-domain UNN, independently of the exploited source domain. It is also worth noting the contrasting results for CD-INN. While it provides the best performance in terms of accuracy (see 6.9), its diversity is the worst for books and only average for movies.

In summary, we observe a clear trade-off between accurate and diverse recommendations. In general, when approaches perform well in terms of MRR they tend to suffer in terms of diversity, and vice versa.

6.5 Conclusions and future work

Collaborative filtering approaches have become the most investigated and popular solutions to the cross-domain recommendation problem, as they only mine patterns of user-item preferences (i.e., ratings), and do not require any information about the content of the items to bridge the domains of interest. Some other approaches, however, have shown that content-based relations (e.g., based on social tags) can be exploited to bridge the domains more effectively. In this context, recent initiatives such as the Linked Open Data project provide large interconnected repositories of structured knowledge than can be exploited to relate multiple types of data. Such heterogeneous networks allow establishing content-based links between different

Table 6.12: Diversity (BinomDiv@10) for cold start users in the music domain. The three groups of rows correspond to single-domain, cross-domain with books as source, and cross-domain with movies as source, respectively. Best values for each single- and cross-domain configuration are shown in bold.

| Method | Number of music likes | | | | | | | |
|---------|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6–10 | |
| POP | 0.324 | 0.228 | 0.262 | 0.282 | 0.295 | 0.305 | 0.326 | |
| UNN | | 0.296 | 0.332 | 0.348 | 0.347 | 0.330 | 0.306 | |
| INN | | 0.200 | 0.213 | 0.219 | 0.223 | 0.229 | 0.236 | |
| iMF | | 0.196 | 0.217 | 0.232 | 0.241 | 0.249 | 0.259 | |
| BPR | | 0.539 | 0.577 | 0.589 | 0.590 | 0.594 | 0.619 | |
| FISM | | 0.683 | 0.766 | 0.709 | 0.731 | 0.737 | 0.676 | |
| HeteRec | | 0.227 | 0.264 | 0.280 | 0.288 | 0.296 | 0.304 | |
| SPRank | | N/A | N/A | N/A | N/A | N/A | N/A | |
| Books | CD-UNN | 0.325 | 0.429 | 0.414 | 0.393 | 0.366 | 0.346 | 0.314 |
| | CD-INN | 0.269 | 0.215 | 0.227 | 0.232 | 0.235 | 0.240 | 0.244 |
| | CD-iMF | 0.270 | 0.214 | 0.233 | 0.240 | 0.249 | 0.252 | 0.258 |
| | CD-BPR | 0.570 | 0.585 | 0.578 | 0.602 | 0.592 | 0.603 | 0.609 |
| | CD-FISM | 0.607 | 0.597 | 0.677 | 0.648 | 0.774 | 0.726 | 0.674 |
| | CD-HeteRec | 0.295 | 0.233 | 0.271 | 0.286 | 0.294 | 0.302 | 0.309 |
| | SimMF | 0.274 | 0.220 | 0.240 | 0.249 | 0.257 | 0.264 | 0.275 |
| | NeighborMF | 0.254 | 0.220 | 0.241 | 0.251 | 0.259 | 0.265 | 0.275 |
| | CentroidMF | 0.253 | 0.218 | 0.238 | 0.249 | 0.257 | 0.263 | 0.273 |
| Movies | CD-UNN | 0.296 | 0.411 | 0.380 | 0.358 | 0.347 | 0.329 | 0.312 |
| | CD-INN | 0.277 | 0.231 | 0.255 | 0.264 | 0.270 | 0.272 | 0.275 |
| | CD-iMF | 0.248 | 0.229 | 0.254 | 0.264 | 0.271 | 0.272 | 0.277 |
| | CD-BPR | 0.476 | 0.515 | 0.526 | 0.504 | 0.526 | 0.529 | 0.545 |
| | CD-FISM | 0.601 | 0.664 | 0.541 | 0.757 | 0.578 | 0.771 | 0.669 |
| | CD-HeteRec | 0.372 | 0.271 | 0.314 | 0.331 | 0.342 | 0.349 | 0.360 |
| | SimMF | 0.225 | 0.207 | 0.239 | 0.250 | 0.259 | 0.264 | 0.278 |
| | NeighborMF | 0.252 | 0.226 | 0.251 | 0.265 | 0.269 | 0.274 | 0.283 |
| | CentroidMF | 0.264 | 0.233 | 0.257 | 0.270 | 0.274 | 0.279 | 0.286 |

types of items, and thus providing a new mechanism to bridge domains for cross-domain recommendation.

In this study, we have exploited Linked Open Data to extract metadata about items in three recommendation domains. Using this additional information, we were able to find relations between items in different domains, and ultimately compute inter-domain item similarities. This could be a limit of the presented approaches whenever the underlying LOD knowledge graph does not expose semantic links between items in different domains, e.g., when the source and target domains do not share information, i.e., there is no direct or indirect link between items in different domains or it is not possible to link an item in the catalog to the corresponding entity in the knowledge graph. In fact, in these cases, it is not possible to compute pairwise semantic similarity values between items belonging to different domains.

We then proposed three novel matrix factorization models for cross-domain recommendation that exploit the computed similarities to link knowledge across domains. Experiments in cold start scenarios showed that depending on the involved source and target domains, cross-domain recommendations exploiting item metadata can be more accurate for users with few preferences in the target domain. However, the improved accuracy comes at the cost of less diversity among the recommendations, and approaches thriving in diversity tend to be less accurate. We argue, nonetheless, that in cold start the priority of a system may be keeping the user engaged by delivering relevant recommendations rather than diverse, non relevant ones.

Regarding the categorization presented in 6.2.1, the models proposed in this study belong to the category of **knowledge linkage** cross-domain recommendation approaches. We applied our approaches to the **linked-domain exploitation** task with the goal of **addressing the user cold start** problem. In addition to the results reported here, we conjecture that item metadata may be prove more useful in cross-domain scenarios with low user overlap. In these cases, approaches purely based on collaborative filtering are likely to struggle to compute accurate item-item similarities. Moreover, in our work we relied on advanced Bayesian Optimization

techniques to find the optimal hyperparameters of the models, and in particular the values of the cross-domain regularization λ_C and the item neighborhood size n parameters. It would be interesting, however, to analyze the performance of the models in terms of these parameters to better understand the importance of auxiliary information. We did not report these results due to the high number of possible combinations of different parameter values, source-target domain configurations, cold start profile sizes, and cross-validation folds, which may make it very difficult to extract conclusions that consistently hold through all the possible scenarios.

Chapter 7

Interpretability of Factorization Machines

7.1 Introduction

Research on transparency, and interpretability of predictive models is gaining momentum since it has been recognized as a key element in the next generation of recommendation algorithms. Providing explanations may increase the user awareness in the decision-making process leading to fast (efficiency), conscious and right (effectiveness) decisions. When equipped with interpretability of recommendation results, a system ceases to be just a black-box (transparency) [349, 368, 410] and users are more willing to extensively exploit the predictions [367, 176]. Transparency increases their trust [139] (also exploiting specific semantic structures [132]), and satisfaction in using the system. For a recommender system, the user's trust is becoming more and more important since it leads to better performance [235]. Most of the proposed approaches to interpret recommendation results seek to show how they are related to users' preferences. In a nutshell, we may say that interpreta-

tions for recommendation results can be item-based, user-based or feature-based. Item-based interpretations make use of the shared set of items among users [7]; User-based explanations rely on sets of most similar users, like in [176]; Feature-based explanations exploit features of recommended items as director, genre, and cast [367]. Among interpretable models, we may distinguish between those based on Content-based (CB) approaches and those based on Collaborative filtering (CF) ones. CB algorithms provide recommendations by exploiting the available content and matching it with a user profile [300, 106]. The use of content features makes the model interpretable even though attention has to be paid since a CB approach “*lacks serendipity and requires extensive manual efforts to match the user interests to content profiles*” [418]. It is worth noticing that these features can be dramatically different depending on the considered scenario: a movie recommendation could be based on the director, actors, the producer, the genre whereas a book recommendation may be explained by the author, the book formats or the saga. Sometimes, this prevents the straight adoption of a model independently of the addressed knowledge domain.

When content is missing, the recommender may rely only on the relationships between users and the rates they provide to items in a collaborative fashion. Based on the different collaborative approaches, these relationships may focus on items or users. Interpretation of CF results will then inevitably reflect the approach adopted by the algorithm. For instance, an item-based and a user-based recommendation could be explained, respectively, as “other users who have experienced A have experienced B” or “similar users have experienced B”. Concerning collaborative approaches to recommendation, when we use a kNN approach, either item-based or user-based, we explicitly refer to users and items in the system. Unfortunately, things change when we adopt more powerful and accurate Deep Learning [92], or model-based algorithms and techniques for the computation of a recommendation list. Such approaches project items and users in a new vector space of latent features [224] thus making the final result not directly interpretable. Indeed, it is possible to compute items and users similarities via latent factor exploitation, but we lose entirely any reference to the original user-item interaction and then to an explicit

justification for the computed recommendations.

In the last years, many approaches have been proposed that take advantage of side information to enhance the performance of latent factor models. Side information can refer to items as well as users [390] and can be either structured [358] or semi-structured [419, 40, 97]. Interestingly, in [418] the authors argue about a new generation of knowledge-aware recommendation engines able to exploit information encoded in knowledge graph (KG) to produce meaningful recommendations: *“For example, with knowledge graph about movies, actors, and directors, the system can explain to the user a movie is recommended because he has watched many movies starred by an actor”*.

In this work, we propose a knowledge-aware Hybrid Factorization Machine (kaHFM) to train interpretable models in recommendation scenarios. kaHFM relies on Factorization Machines [310] and it extends them in different key aspects making use of the semantic information encoded in a knowledge graph.

Interpretability without accuracy of results is not sufficient when designing a recommendation engine since it is unlikely that an inaccurate model will be adopted. It would be highly beneficial for users to develop recommender systems that are accurate, and, at the same time, that are built upon an interpretable technique, so that the learned model can be, in case, exploited to generate explanations. In this direction, we show how kaHFM may exploit data coming from knowledge graphs as side information to build a recommender system whose final results are accurate and, at the same time, semantically interpretable. With kaHFM we build a model in which the meaning of each latent factor is bound to an explicit content-based feature extracted from a knowledge graph. Doing this, after the model has been trained, we still have an explicit reference to the original semantics of the features describing the items, thus making possible the interpretation of the final results. Interestingly, we will see that the explicit mapping of latent features to content-based ones makes possible to exploit characteristics of these latter to implement a more effective initialization technique.

We evaluated kaHFM on six different publicly available datasets by getting

content-based explicit features from data encoded in the DBpedia¹ knowledge graph. We analyzed the performance of the approach in terms of accuracy, diversity, and novelty of results by exploiting categorical, ontological and factual features (see Section 7.3.1). For each of them, we used public mappings to DBpedia. Finally, we tested the robustness of kaHFM with respect to its interpretability showing that it ranks meaningful features higher and is able to regenerate them in case they are removed from the original dataset.

With kaHFM we address the following research questions:

- RQ1** Can we develop a model-based recommendation engine whose results are very accurate and, at the same time, can be interpreted with respect to an explicitly stated semantics coming from a knowledge graph?
- RQ2** Can we evaluate that the original semantics of items features is preserved after the model has been trained?
- RQ3** How to measure with an offline evaluation that the proposed model is really able to identify meaningful features by exploiting their explicit semantics?

This investigation can be then summarized as:

- presentation of kaHFM: a framework that exploits data available in knowledge graphs to build semantically interpretable models for recommendation tasks;
- an experimental evaluation based on six different datasets to assess the performance of kaHFM in terms of accuracy, diversity, and novelty of computed results;
- introduction of two metrics, Semantic Accuracy ($SA@K$) and Robustness ($n\text{-Rob}@K$), to measure the interpretability of a knowledge-aware recommendation engine.

The remainder of the chapter is structured as follows: in the next section, we introduce the background technologies behind kaHFM and then we detail the overall approach in Section 7.3. The following section is devoted to the introduction of

¹<http://dbpedia.org>

the two metrics we propose to assess the quality of k_a HFM results in terms of interpretability. In Section 7.4 we describe the experimental setting while in Section 7.2 we report on related work. Conclusions and future lines of research close the chapter.

7.2 Related Work

In recent years, several explainable recommendation models that exploit matrix factorization have been proposed. It is well-known that one of the main issues of matrix factorization methods is that they are not easily explainable (since latent factors meaning is basically unknown). One of the first attempts to overcome this problem was proposed in [416, 419]. In this work, the authors propose Explicit Factor Model (EFM). Products' features and users' opinions are extracted with phrase-level sentiment analysis from users' reviews to feed a matrix factorization framework. Each latent factor is thus mapped with a particular explicit feature. After that, a few improvements to EFM have been proposed to deal with temporal dynamics [420] and to use tensor factorization [97]. In particular, in the latter the aim is to predict both user preferences on features (extracted from textual reviews) and items. This is achieved by exploiting the Bayesian Personalized Ranking (BPR) criterion [312]. Eventually, these preferences are combined to produce recommendation lists. We considered this work interesting because they also adopt a pair-wise learning to rank algorithm, but this is really different from ours since we exploited BPR to explicitly train the feature vectors to rank items. Further advances in MF-based explainable recommendation models have been proposed with Explainable Matrix Factorization (EMF) [6] in which the generated explanations are based on a neighborhood model. We differ from EMF as they do not take advantage of any external data source, other than they use a completely different model. Moreover, they introduced an additional regularizer into the factorization model to constrain users' vectors training. Similarly, in [7] an explainable Restricted Boltzmann Machine model has been proposed. It learns a network model (with an additional visible layer) that takes into account a degree of explainability. To measure it, they defined an ad-hoc *Explain-*

ability Score. Finally, an interesting work took advantage of sentiment analysis of users' reviews. Authors incorporate the sentiments and ratings into a matrix factorization model. The overall approach, named Sentiment Utility Logistic Model (SULM) [40], generates a novel kind of explanations composed by both items and features. In [309] recommendations are computed by generating and ranking personalized explanations in the form of explanation chains. OCuLaR [384] provides interpretable recommendations from positive examples based on the detection of co-clusters between users (clients) and items (products). The recommendation comes with the corresponding user-item co-clusters, which provide much more detailed information than usual collaborative-based explanations. In [183] authors propose a Multi Level Attraction Model (MLAM) in which they build two attraction models, for cast and story. Moreover, the story model is built upon two attraction models: for story level and for sentence level. The interpretability of the model is then provided in terms of attractiveness of Sentence level, Word level, and Cast member. In [301] the authors train a matrix factorisation model to complete the $U \times I$ matrix. They then use the complete (approximated) ratings matrix to compute a set of association rules that explain the obtained recommendations. In [120] the authors prove that, given the conversion probabilities for all actions of customer features, it is possible to transform the original historical data to a new space in order to compute a set of interpretable recommendation rules.

The design of a new explainable model is useless if it does not match with any existing explanations generation technique. For this reason we deepened the different ways of generating feature-based explanations in order to provide a flexible but accurate model. Probably the most used style of explanations of this kind are the content-based ones [360]. In their most simple form, authors consider similarities between items by taking into account both item properties and user ratings. Among the works that exploit content information to produce explainable recommendations, Tagsplanations [383] is worth to mention. It is fed by community tags and it exploits a relevance measure to weight tags w.r.t. items and user preferences. Furthermore, also demographic-based recommendations explanations have been inspected [422], in order to recommend items for particular types (age, loca-

tion, gender) of users. The core of our model is a general Factorization Machines (FM) model [310]. Nowadays FMs are the most widely used factorization models because they offer a number of advantages w.r.t. other latent factors models such as SVD++ [221], PITF [315], FPMC [313]. First of all, FMs are designed for a generic prediction task while the others can be exploited only for specific tasks. Moreover it is a linear model and parameters can be estimated accurately even in high data sparsity scenarios. Nevertheless, several improvements have been proposed for FMs. For instance Neural Factorization Machines [174] have been developed to fix the inability of classical FMs to capture non linear structure of real-world data. This goal is achieved by exploiting the non linearity of neural networks. Furthermore, Attentional Factorization Machines [401] have been proposed that use an attention network to learn the importance of feature interactions. Finally, FMs have been specialized to better work as Context-Aware recommender systems [314]. Usually only top recommended items are provided to the user as suggestions as it is not feasible that a user analyze hundreds of recommended items. For this reason, ranking has become a much more important task than rating prediction [256]. This came with the development of a new class of learning algorithms (in which sorting correctly the recommendation list becomes the key task). These Learning to Rank [88] algorithms can be further categorized in Point-wise [225], Pair-wise [312, 247] and List-wise [344, 343]. In particular, Pair-wise approaches are usually considered as a good trade-off between ordering performances and computational complexity. Among this class of algorithms, Bayesian Personalized Ranking (BPR) [312] is one of the most widely adopted. It is based on a simple stochastic gradient descent algorithm to learn the relative order between positive items (items that a user has experienced in his past history) and negative items (items never rated by the user). BPR can be easily applied to Matrix Factorization and Factorization Machines (as in our work and in [41]).

7.3 Approach

For the sake of completeness, in this section we briefly recap the main technologies we adopted to develop kaHF_M . We introduce Vector Space Models for recommender systems and then we give a quick overview on knowledge graphs and their Linked (Open) Data implementation.

Content-based recommender systems rely on the assumption that it is possible to predict the future behavior of users based on their personalized profile. User profiles can be built by exploiting the characteristics of the items they liked in the past or some other available side information. Several approaches have been proposed, that take advantage of side information in different ways: some of them consider tags [383], demographic data [422] or they extract information from collective knowledge bases [278]. Many of the most popular and adopted CB approaches make use of a Vector Space Model (VSM). In VSM users and items are represented by means of Boolean or weighted vectors in the same space. Their respective positions and the distance, or better the proximity, between them, provides a measure of how these two entities are related or similar. The space containing users and items vectors is composed by dimensions that may refer to characteristics of the users (e.g., the propensity to watch movies in the morning) or of the items as well as to their combination.

The choice of item features may substantially differ depending on their availability and application scenario: crowd-sourced tags, categorical, ontological, or textual knowledge are just some of the most exploited ones.

To sum up, in a CB approach we need (i) to get reliable items descriptions, (ii) a way to measure the strength of each feature for each item, (iii) to represent users and finally (iv) to measure similarities. Regarding the first point, nowadays we can easily get descriptions related to an item from the Web. In particular, thanks to the Linked Open Data initiative a lot of semantically structured knowledge is publicly available in the form of Linked Data datasets.

7.3.1 Knowledge Graphs and Linked Data

In 2012, Google announced its Knowledge Graph² (KG) as a new tool to improve the identification and retrieval of entities in return to a search query. Most of the knowledge encoded in Google Knowledge Graph actually came from Freebase which was a crowdsourced effort to create a base of facts in all possible knowledge domains. Alongside with the development of the above mentioned initiatives, following the original idea of a Semantic Web [60], new technologies have been developed and released with the aim of embedding structured knowledge with unambiguous semantics into Web pages in order to allow software agents to consume and elaborate information in an automated way. The original idea has been modified over the years thus making possible the creation of a full stack of semantic technologies and, more remarkably, gave birth to the Linking Open Data initiative³ where a community of researchers and practitioners devoted an enormous effort to build publicly available knowledge bases of machine-understandable data. A knowledge base exploiting Semantic Web technologies is then represented through a graph (knowledge graph) in which entities are linked to each other by binary relations. In order to manage this knowledge model, several technologies and languages have been developed and among them the basic one is Resource Description Framework (RDF). It provides a simple graph-based data model to encode knowledge in a structured way by means of a triple $\langle \sigma, \rho, \omega \rangle$. In a knowledge graph, each triple represents the connection $\sigma \xrightarrow{\rho} \omega$ between two nodes, named *subject* (σ) and *object* (ω), through the *relation (predicate)* ρ . In an RDF knowledge graph, we usually find different types of encoded information.

- **Factual.** This refers to statements such as *The Matrix was directed by the Wachowskis* or *Melbourne is located in Australia* that describe attributes of an entity;
- **Categorical.** It is mainly used to state something about the subject of an

²<https://googleblog.blogspot.it/2012/05/introducing-knowledge-graph-things-not.html>

³<http://linkeddata.org>

entity. In this direction, the categories of Wikipedia pages are an excellent example. Categories can be used to cluster entities and are often organized hierarchically thus making possible to define them in a more generic or specific way;

- **Ontological.** This is a more restrictive and formal way to classify entities via a hierarchical structure of classes. Differently from categories, sub-classes and super-classes are connected through IS-A (transitive) relations.

If we take a look at the following RDF triples⁴ from the DBpedia knowledge graph

```
dbr:The_Matrix dbo:director dbr:The_Wachowski_Brothers.  
dbr:The_Matrix dct:subject dbc:Dystopian_films.  
dbr:The_Matrix rdf:type dbo:Film.
```

we may see that each of them represents one of the above mentioned types of data. In the first triple we state a fact about the movie *The Matrix* (represented by the corresponding URI `dbr:The_Matrix`) saying that it has been directed (`dbo:director`) by the *Wachowski Brothers* (`dbr:The_Wachowski_Brothers`). The second triple encodes categorical information through the predicate `dct:subject` about the same movie. In particular, here we say that it belongs to the category of dystopian films (`dbc:Dystopian_films`). Finally, with the last triple, we classify *The Matrix* as a Film (`dbo:Film`) thanks to the predicate `rdf:type`.

7.3.2 Formal Model

Factorization models have been proven to be among of the best performing approaches as collaborative filtering methods to build a recommender system [311]. This is due to the high prediction accuracy and the subtle modeling of user-item interactions which let these models operate efficiently even in very sparse settings (compared to other classical collaborative predictive models). Among all the different factorization models, factorization machines propose a unified general model

⁴For the sake of conciseness we will use the CURIE syntax in which URIs are abbreviated using their namespaces. In this study we refer to namespaces as available at <http://prefix.cc>.

to represent most of them. Here we report the definition and results related to a factorization model of order 2 for a recommendation problem involving only implicit ratings. Nevertheless, everything can be easily extended to a more expressive representation by taking into account, e.g., demographic and social information [18], multi-criteria [15], and even relations between contexts [424]. For each user $u \in U$ and each item $i \in I$ we build a binary vector $\mathbf{x}^{ui} \in \mathbb{R}^{1 \times n}$, with $n = |U| + |I|$, representing the interaction between u and i in the original user-item rating matrix. In this modeling, \mathbf{x}^{ui} contains only two 1 values corresponding to u and i while all the other values are set to 0 (see Fig. 7.1). We then denote with $\mathbf{X} \in \mathbb{R}^{n \times m}$ the matrix containing as rows all possible \mathbf{x}^{ui} we can build starting from the original user-item rating matrix as shown in Fig. 7.1.

| | | | | | | | | | | | |
|-------|-------|-------|-------|-------|-----|-------|-------|-------|-------|-------|-----|
| x^1 | 1 | 0 | 0 | 0 | ... | 1 | 0 | 0 | 0 | 0 | ... |
| x^2 | 1 | 0 | 0 | 0 | ... | 0 | 1 | 0 | 0 | 0 | ... |
| x^3 | 1 | 0 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| x^4 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| x^5 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 0 | 1 | 0 | ... |
| x^6 | 0 | 0 | 1 | 0 | ... | 1 | 0 | 0 | 0 | 0 | ... |
| x^7 | 0 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| | U_1 | U_2 | U_3 | U_4 | ... | I_1 | I_2 | I_3 | I_4 | I_5 | ... |
| | User | | | | | Item | | | | | |

Figure 7.1: A visual representation of \mathbf{X} for sparse real valued vectors \mathbf{x}^{ui} .

The factorization machine (FM) for each vector \mathbf{x} can be defined as:

$$\hat{y}(\mathbf{x}^{ui}) = w_0 + \sum_{j=1}^n w_j \cdot x_j + \sum_{j=1}^n \sum_{j'=j+1}^n x_j \cdot x_{j'} \cdot \sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)} \quad (7.1)$$

where the parameters to be learned are: w_0 representing the global bias; w_j giving the importance to every single x_j ; the pair $v_{(j,f)}$ and $v_{(j',f)}$ in $\sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ measuring the strength of the interaction between each pair of variables x_j and $x_{j'}$. The number of latent factors is represented by k . This value is usually selected at design time when implementing the FM.

In order to make the recommendation results computed by `kaHF`M semantically interpretable, we want to inject the knowledge encoded within a knowledge-graph

in a Factorization Machine. Given a set of features retrieved from a KG [277] we first bind them to the latent factors and then, since we address a Top-N recommendation problem, we train the model by using a Bayesian Personalized Ranking (BPR) criterion that takes into account entities within the original knowledge graph.

In [279], the authors originally proposed to encode a Linked Data knowledge graph in a vector space model to develop a CB recommender system. Given a set of items $I = \{i_1, i_2, \dots, i_N\}$ in a catalog and their associated triples $\langle i, \rho, \omega \rangle$ in a knowledge graph \mathcal{KG} we may build the set of all possible features as $F = \{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG} \text{ with } i \in I\}$. Each item can be then represented as a vector of weights $\mathbf{i} = [v_{(i,1)}, \dots, v_{(i,\langle \rho, \omega \rangle)}, \dots, v_{(i,|F|)}]$ where $v_{(i,\langle \rho, \omega \rangle)}$ is computed as the normalized TF-IDF value for $\langle \rho, \omega \rangle$:

$$v_{(i,\langle \rho, \omega \rangle)} = \frac{|\{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG}\}|}{\underbrace{\sqrt{\sum_{\langle \rho, \omega \rangle \in F} |\{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG}\}|^2}}_{TF^{\mathcal{KG}}}} \cdot \underbrace{\log \frac{|I|}{|\{j \mid \langle j, \rho, \omega \rangle \in \mathcal{KG} \text{ and } j \in I\}|}}_{IDF^{\mathcal{KG}}} \quad (7.2)$$

Analogously, when we have a set U of users, we may represent them using the features describing the items they enjoyed in the past. In the following, when no confusion arises, we use f to denote a feature $\langle \rho, \omega \rangle \in F$. Given a user u , if we denote with I^u the set of the items enjoyed by u and we have $\mathbf{u} = [v_{(u,1)}, \dots, v_{(u,f)}, \dots, v_{(u,|F|)}]$ with

$$v_{(u,f)} = \frac{\sum_{i \in I^u} v_{(i,f)}}{|\{i \mid i \in I^u \text{ and } v_{(i,f)} \neq 0\}|}$$

Given the vectors \mathbf{u}_j , with $j \in [1 \dots |U|]$, and $\mathbf{i}_{j'}$, with $j' \in [1 \dots |I|]$, we build the matrix $\mathbf{V} \in \mathbb{R}^{n \times |F|}$ (see Fig. 7.2) where the first $|U|$ rows have a one to one mapping with \mathbf{u}_j while the last ones correspond to $\mathbf{i}_{j'}$. If we go back to Equation (7.1) we may see that, for each \mathbf{x} , the term $\sum_{j=1}^n \sum_{j'=j+1}^n x_j \cdot x_{j'} \cdot \sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ is not zero only once, i.e., when both x_j and $x_{j'}$ are equal to 1. In the matrix depicted in Fig. 7.1, this happens when there is an interaction between a user and an item. Moreover, the summation $\sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ represents the dot product between two vectors \mathbf{v}_j

| | | - dbc:Space_adventure_films | - dbc:Films_set_in_the_future | - dbc:American_science_fiction_action_films | - dbc:1980s_science_fiction_films | - dbc:Paramount_Pictures_films | - dbc:Midlife_crisis_films | - dbc:American_sequel_films |
|----------------|------|-----------------------------|-------------------------------|---|-----------------------------------|--------------------------------|----------------------------|-----------------------------|
| \mathbf{v}_1 | 0 | 0.88 | 0.81 | 0.7 | 0 | 0.60 | 0.53 | ... |
| \mathbf{v}_2 | 1.3 | 1.12 | 0.91 | 0.84 | 0.65 | 0.59 | 0.58 | ... |
| \mathbf{v}_3 | 0.5 | 0 | 0.71 | 0 | 0.28 | 0.35 | 0 | ... |
| \mathbf{v}_4 | 0 | 0 | 0.31 | 0 | 0 | 0 | 0.6 | ... |
| \mathbf{v}_5 | 0 | 0 | 0 | 0 | 0.18 | 0 | 0 | ... |
| \mathbf{v}_6 | 0 | 0.12 | 0.22 | 0 | 0 | 0 | 0 | ... |
| \mathbf{v}_7 | 1.23 | 1.03 | 0.89 | 0.85 | 0.56 | 0.3 | 0.61 | ... |

Figure 7.2: Example of real valued feature vectors for different items v_j .

and $\mathbf{v}_{j'}$ with a size equal to k . Hence, \mathbf{v}_j represents a latent representation of a user, $\mathbf{v}_{j'}$ that of an item within the same latent space, and their interaction is evaluated through their dot product.

In order to inject the knowledge coming from \mathcal{HG} into kaHFM, first of all, we keep Equation (7.1) and we set $k = |F|$. In other words, we impose a number of latent factors equal to the number of features describing all the items in our catalog. We want to stress here that our aim is not representing each feature through a latent vector, but to associate each factor to an explicit feature, obtaining latent vectors that are composed by explicit semantic features. Hence, we initialize the parameters \mathbf{v}_j and $\mathbf{v}_{j'}$ with their corresponding rows from \mathbf{V} which in turn represent respectively \mathbf{u}_j and $\mathbf{i}_{j'}$. In this way, we try to identify each latent factor with a corresponding explicit feature. The intuition is that after the training phase, the resulting matrix $\hat{\mathbf{V}}$ still refers to the original features but contains better values for $v_{(j,f)}$ and $v_{(j',f)}$ that also take into account the latent interactions between users, items and features. It is noteworthy that after the training phase \mathbf{u}_j and $\mathbf{i}_{j'}$ (corresponding to $v_{(j,f)}$ and $v_{(j',f)}$ in \mathbf{V}) contain non-zero values also for features that are not originally in the description of the user u or of the item i .

In Table 7.1 and Table 7.2 we show examples for values after the training (in the

column ka_{HFM}) together with the original TF-IDF ones computed for two movies from the Yahoo! Movies⁵ dataset.

| ka_{HFM} | TF-IDF | Predicate | Object |
|------------|--------|-------------|---|
| 1.3669 | 0.2584 | dct:subject | dbc:Space_adventure_films |
| 1.1252 | 0.2730 | dct:subject | dbc:Films_set_in_the_future |
| 0.9133 | 0.2355 | dct:subject | dbc:American_science_fiction_action_films |
| 0.8485 | 0.3190 | dct:subject | dbc:1980s_science_fiction_films |
| 0.6529 | 0.1549 | dct:subject | dbc:Paramount_Pictures_films |
| 0.5989 | 0.3468 | dct:subject | dbc:Midlife_crisis_films |
| 0.5940 | 0.1797 | dct:subject | dbc:American_sequel_films |
| 0.5862 | 0.2661 | dct:subject | dbc:Film_scores_by_James_Horner |
| 0.5634 | 0.2502 | dct:subject | dbc:Films_shot_in_San_Francisco |
| 0.5583 | 0.1999 | dct:subject | dbc:1980s_action_thriller_films |

Table 7.1: Top-10 features computed by ka_{HFM} for the movie "Star Trek II – The Wrath of Khan".

| ka_{HFM} | TF-IDF | Predicate | Object |
|------------|--------|-------------|---|
| 1.2434 | 0.2858 | dct:subject | dbc:Space_adventure_films |
| 1.0355 | 0.3020 | dct:subject | dbc:Films_set_in_the_future |
| 0.8956 | 0.2605 | dct:subject | dbc:American_science_fiction_action_films |
| 0.8951 | 0.3451 | dct:subject | dbc:Android_(robot)_films |
| 0.7338 | 0.3105 | dct:subject | dbc:Time_travel_films |
| 0.6665 | 0.2701 | dct:subject | dbc:Film_scores_by_Jerry_Goldsmith |
| 0.6581 | 0.2205 | dct:subject | dbc:1990s_action_films |
| 0.6561 | 0.2279 | dct:subject | dbc:1990s_science_fiction_films |
| 0.6118 | 0.1988 | dct:subject | dbc:American_sequel_films |
| 0.5649 | 0.1713 | dct:subject | dbc:Paramount_Pictures_films |

Table 7.2: Top-10 features computed by ka_{HFM} for the movie "Star Trek – First Contact".

7.3.3 Optimization

Factorization machines can be easily trained to reduce the prediction error via gradient descent methods, alternating least-squares (ALS) and MCMC. Since we formulated our problem as a *top-N* recommendation task, ka_{HFM} can be trained using a learning to rank approach like Bayesian Personalized Ranking Criterion (BPR)

⁵http://research.yahoo.com/Academic_Relations

[312]. The BPR criterion is optimized using a stochastic gradient descent algorithm on a set D_S of triples (u, i, j) , with $i \in I^u$ and $j \notin I^u$, selected through a random sampling from a uniform distribution. The BPR optimization criterion can thus be formulated as:

$$\begin{aligned} \text{BPR-OPT} &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{x}_{uij}) - \lambda_{\Theta} \|\Theta\|^2 \\ &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{y}(\mathbf{x}^{\mathbf{ui}}) - \hat{y}(\mathbf{x}^{\mathbf{uj}})) - \lambda_{\Theta} \|\Theta\|^2 \end{aligned} \quad (7.3)$$

In this formulation, $\sigma(\cdot)$ is a sigmoid function, and the update step is defined as:

$$\Theta \leftarrow \Theta + \alpha \left(\frac{e^{-\hat{x}_{uij}}}{1 + e^{-\hat{x}_{uij}}} \cdot \frac{\partial}{\partial \Theta} \hat{x}_{uij} + \lambda \Theta \right) \quad (7.4)$$

where α is the chosen learning rate. Since we are in an implicit feedback setting, we may assume that there is only an instance for the pair user-item. Hence, in our model we can derive \hat{x}_{uij} as:

$$\begin{aligned} \hat{x}_{uij} &= \hat{y}(\mathbf{x}^{\mathbf{ui}}) - \hat{y}(\mathbf{x}^{\mathbf{uj}}) = w_i x_i - w_j x_j + \\ &+ \sum_{f \in F} x_u x_i v_{(u,f)} v_{(i,f)} - x_u x_j v_{(u,f)} v_{(j,f)} \end{aligned} \quad (7.5)$$

For kaHFM, this computation can be performed in an efficient way computing the partial derivatives (to update the factorized parameters) for the only active entities involved in the transactions, w_i , w_j , v_u , v_i , and v_j :

$$\frac{\partial}{\partial \Theta} \hat{x}_{uij} = \begin{cases} 1, & \text{if } \theta = w_i, \\ -1, & \text{if } \theta = w_j, \\ v_{(u,f)}, & \text{if } \theta = v_{(i,f)}, \\ -v_{(u,f)}, & \text{if } \theta = v_{(j,f)}, \\ (v_{(i,f)} - v_{(j,f)}), & \text{if } \theta = v_{(u,f)}, \\ 0, & \text{otherwise} \end{cases} \quad (7.6)$$

Using Equation (7.6) in Equation (7.4) the model parameters can be iteratively updated to maximize the BPR-OPT criterion. The algorithm updates sequentially each sampled triple and continues the training until it reaches the provided number of iterations.

7.3.4 Personalized Recommendation

Once the training phase returns the optimal model parameters, the item recommendation step can take place. We extract the items vectors \mathbf{v}_j from the matrix \mathbf{V} , with the associated optimal values and we use them to implement an Item-kNN recommendation approach. We measure similarities between each pair of items i and j by evaluating the cosine similarity of their corresponding vectors in \mathbf{V} :

$$cs(i, j) = \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \cdot \|\mathbf{v}_j\|}$$

Given a set of neighbors for the item i , denoted as N^i , such that $i \notin I^u$ and a user u we may predict the score assigned by u to i as

$$score(u, i) = \frac{\sum_{j \in N^i \cap I^u} cs(i, j)}{\sum_{j \in N^i} cs(i, j)} \quad (7.7)$$

The proposed approach let us keep the explicit meaning of the “latent” factors computed via a factorization machine thus making possible an interpretation of the recommended results. In the following we propose an automated offline procedure able to assess that the semantics of the features represented in \mathbf{V} is preserved after the training phase. We refer to the values computed by the proposed procedure as *Semantic Accuracy*. A different but related aspect is that of evaluating if kaHFM really assigns a higher value to meaningful features. We refer to this behavior as *Robustness*. Interestingly, both *Semantic Accuracy* and *Robustness* can be evaluated in an offline setting.

7.3.5 Semantic Accuracy

The main idea behind Semantic Accuracy is to evaluate, given an item i , how well kaHFM is able to return its original features available in the returned top-K list \mathbf{v}_i . In other words, given the set of features of i represented by $F^i = \{f_1^i, \dots, f_m^i, \dots, f_M^i\}$, with $F^i \subseteq F$, we check if the values in \mathbf{v}_i corresponding to $f_{m,i} \in F^i$ are higher than those corresponding to $f \notin F^i$. For the set of M features initially describing i we see how many of them appear in the set $\text{top}(\mathbf{v}_i, M)$ representing the top- M features in \mathbf{v}_i . We then normalize this number by the size of F^i and average on all the items within the catalog I .

$$\text{Semantic Accuracy (SA@M)} = \frac{\sum_{i \in I} \frac{|\text{top}(\mathbf{v}_i, M) \cap F^i|}{|F^i|}}{|I|}$$

In many practical scenarios we may have $|F| \gg M$. Hence, we might also be interested in measuring the accuracy for different sizes of the top list. Since items could be described with a different number of features, the size of the top list could be a function of the original size of the item description. Thus, we measured SA@nM with $n \in \{1, 2, 3, 4, 5, \dots\}$ and evaluate the number of features in F^i available in the top- $n \cdot M$ elements of \mathbf{v}_i .

$$\text{SA@nM} = \frac{\sum_{i \in I} \frac{|\text{top}(\mathbf{v}_i, n \cdot M) \cap F^i|}{|F^i|}}{|I|}$$

7.3.6 Robustness

Although SA@nM may result very useful to understand if kaHFM assigns weights according to the original description of item i , we still do not know if a high value in \mathbf{v}_i really means that the corresponding feature is important to define i . In other words, are we sure that kaHFM promotes meaningful features for i ?

In order to provide a way to measure such “meaningfulness” for a given feature, we suppose, for a moment, that a particular feature $\langle \rho, \omega \rangle$ is useful to describe an item i but the corresponding triple $\langle i, \rho, \omega \rangle$ is not represented in the knowledge graph. In case kaHFM was robust in generating weights for unknown features, it

should discover the importance of that feature and modify its value to make it enter the Top- K features in \mathbf{v}_i .

Starting from this observation, the idea to measure robustness is then to “forget” a triple involving i and check if kaHFM can generate it.

In order to implement such process we may proceed by following these steps:

- we train kaHFM thus obtaining optimal values v_i for all the features in F^i ;
- the feature $f_{MAX}^i \in F^i$ with the highest value in v_i is identified;
- we train the model again initializing $f_{MAX}^i = 0$ and we compute v'_i .

After the above steps, if $f_{MAX}^i \in \text{top}(v'_i, M)$ then we can say that kaHFM shows a high robustness in identifying important features.

Given a catalog I , we may then define the *Robustness for 1 removed feature @M* ($1\text{-Rob}@M$) as the number of items for which $f_{MAX}^i \in \text{top}(v'_i, M)$ divided by the size of I .

$$1\text{-Rob}@M = \frac{\sum_{i \in I} |\{i \mid f_{MAX}^i \in \text{top}(v'_i, M)\}|}{|I|}$$

Similarly to $SA@nM$, we may define $1\text{-Rob}@nM$.

7.4 Experimental Evaluation

7.4.1 Experimental Setup

Datasets. In order to provide an answer to RQ1, we evaluated the performance of our method on six datasets belonging to three different domains (Music, Books and Movies). The `Last.fm` dataset [84] corresponds to user-artist plays on Last.fm online music system released during HETRec 2011⁶ Workshop. It contains social networking, tagging, and music artists listening information from a set of 2K users. `LibraryThing` represents books’ ratings collected in the LibraryThing website⁷

⁶<http://ir.ii.uam.es/hetrec2011/>

⁷<https://www.librarything.com/>

community. It contains social networking, tagging and rating information on a [1..10] scale. Yahoo!Movies (Yahoo! Webscope dataset ydata-ymovies-user-movie-ratings-content-v1_0)⁸ contains movies ratings generated on Yahoo! Movies up to November 2003. It provides content, demographic and ratings information on a [1..5] scale, and mappings to MovieLens and EachMovie datasets. Facebook Movies, Facebook Music and Facebook Books datasets have been released for the Linked Open Data challenge co-located with ESWC 2015⁹, and they refer to movies, music and books domains respectively. Only implicit feedback is available for these datasets, but for each item a link to DBpedia is provided. In order to map items in Last.fm and LibraryThing to DBpedia resources, we exploited a freely available mapping¹⁰. For the remaining one (Yahoo!Movies), we extracted all the updated items-features mappings (Yahoo!Movies, LibraryThing, Last.fm, Facebook Movies, Facebook Music and Facebook Books) which are publicly available¹¹.

Datasets statistics are shown in Table 7.3.

| Dataset | #Users | #Items | #Transactions | #Features | Sparsity |
|-----------------|--------|--------|---------------|-----------|----------|
| Yahoo! Movies | 4000 | 2,626 | 69,846 | 988,734 | 99.34% |
| LibraryThing | 7223 | 11,695 | 410,210 | 183,182 | 99.51% |
| Last FM | 1375 | 8,289 | 60,701 | 434,817 | 99.47% |
| Facebook Music | 52068 | 5,749 | 1,374,994 | 345,452 | 99.54% |
| Facebook Movies | 32143 | 3,901 | 689,561 | 180,573 | 99.45% |
| Facebook Books | 1398 | 2,933 | 18,978 | 111,401 | 99.53% |

Table 7.3: Datasets statistics.

Evaluation Protocol and Experimental Setting. "All Unrated Items" [356] protocol has been chosen to compare different algorithms. In All Unrated Items, for each user, all the items that have not yet been rated by the user all over the catalog are considered. We split the dataset using Hold-Out 80-20 retaining for every user the 80% of their ratings in the training set and the remaining 20% in the test set.

⁸http://research.yahoo.com/Academic_Relations

⁹<https://2015.eswc-conferences.org/program/semwebeval.html>

¹⁰<https://github.com/sisinflab/LODrecsys-datasets>

¹¹<https://github.com/sisinflab/LinkedDatasets/>

Moreover, a temporal split has been performed [163] whenever timestamps associated to every transaction is available. We tested our approach against the most related content-based and collaborative algorithms in terms of Accuracy, Diversity and Novelty [23]. We compared k_{aHFM} ¹² w.r.t. a canonical 2 degree Factorization Machine (users and items are intended as features of the original formulation) by optimizing the recommendation list ranking via BPR (BPR-FM). In order to preserve the fairness of the comparison, we used the same parameters adopted for k_{aHFM} and the same number of hidden factors (see the "Selected" column in Table 7.5). Moreover, since we use items similarity in the last step of our approach (see Equation (7.7)), we compared k_{aHFM} against an *Attribute Based Item-kNN* (ABItem-kNN) algorithm, where each item is represented as a vector of weights, computed through a TF-IDF model. In this model, the attributes are computed via Equation (7.2). For the sake of completeness we also compared k_{aHFM} against a pure Item-kNN, that is an item-based implementation of the k-nearest neighbors algorithm. It finds the k-nearest item neighbors based on Cosine Similarity. Items in the neighborhood are then used to predict a score for each user-item pair. Regarding BPR parameters, *learning rate*, *bias regularization*, *user regularization*, *positive item regularization*, and *negative item regularization* have been set respectively to 0.05, 0, 0.0025, 0.0025 and 0.00025 while a sampler "without replacement" has been adopted in order to sample the triples as suggested by authors[312]. For the sake of reproducibility, the BPR parameters are chosen from `mymedialite` implementation. Moreover, the same parameters are used in all the algorithms that make use of BPR to avoid affecting the results, and to guarantee a fair comparison. We also compared k_{aHFM} against the corresponding User-based nearest neighbor scheme, and Most-Popular, a simple baseline that shows high performance in specific scenarios [108]. In our context, we considered mandatory to also compare against a pure knowledge-graph content-based baseline based on Vector Space Model (*VSM*) [279].

¹²Implementation available at <https://github.com/sisinflab/HybridFactorizationMachines/>

7.4.2 Features extraction

The feature extraction is one of the most sensitive steps in our approach. A wrong feature selection may result in noisy data, or in the lack of some important features. This preprocessing was basically divided in three steps: (i) ”**Extraction**”, in which we retrieve data from the DBpedia knowledge graph, (ii) ”**Selection**” where only features involved in the specific experiment are selected, and (iii) ”**Filtering**” in which uninformative features are removed [277].

Extraction. Thanks to the publicly available mappings, all the items from the datasets represented in Table 7.3 come with a DBpedia link. Exploiting this reference, we retrieved all the $\langle \rho, \omega \rangle$ pairs. Some noisy features (based on the following predicates) have been excluded already in this early extraction. In particular we removed: `owl:sameAs`, `dbo:thumbnail`, `prov:wasDerivedFrom`, `foaf:depiction`, `foaf:isPrimaryTopicOf`. Behind this choice, the main reason is that they give us only information about the nature of the entity in the specific knowledge base (e.g., the links between DBpedia and Wikipedia pages) or are linked to multimedia data or even external datasets. After this cleaning step, all the features have been indexed, saved separately and associated to the *id* of the item.

Selection. We performed our experiments with three different settings to analyze the impact of the different kind of features in the recommendation accuracy and diversity. The features have been chosen as they are present in all the different domains and because of their factual, categorical or ontological meaning:

- **Categorical Setting (CS):** We selected only the features containing the property `dcterms:subject`.
- **Ontological Setting (OS):** In this case the only features we considered are: `rdf:type`, `dbo:genre`, `dcterms:subject`.
- **Factual Setting (FS):** We considered all the features but those involving the properties selected in OS.

| Dataset | Threshold |
|-----------------|-----------|
| Yahoo! Movies | 99.62 |
| LibraryThing | 99.91 |
| Last FM | 99.88 |
| Facebook Music | 99.83 |
| Facebook Movies | 99.74 |
| Facebook Books | 99.66 |

Table 7.4: Threshold values adopted to filter out irrelevant properties for each dataset.

Filtering. This last step corresponds to the removal of irrelevant features, that bring little value to the recommendation task, but, at the same time, pose scalability issues. The pre-processing phase has been done following [277], and [299] with a unique threshold. The corresponding thresholds (tm [277], and p [299] for missing values) for each dataset are represented in Table 7.4.

We discarded features for which we had more than tm (or, equivalently p) missing values. For a fair comparison the features used for the baseline Attribute-based Item kNN (ABItem-kNN) are the same used in our approach and the number of latent factors for FMs has been set equal to the number of features involved in the specific setting. The characteristics of each datasets (varying the setting) in terms of considered features are reported in Table 7.5.

| Datasets | Categorical Setting | | Ontological Setting | | Factual Setting | |
|------------------------|---------------------|----------|---------------------|----------|-----------------|----------|
| | Total | Selected | Total | Selected | Total | Selected |
| Yahoo!Movies | 26155 | 747 | 38699 | 1240 | 950035 | 3186 |
| LibraryThing | 9443 | 1169 | 14585 | 1934 | 168597 | 5826 |
| Last.fm | 16422 | 1315 | 30734 | 3032 | 404083 | 9413 |
| Facebook Music | 15016 | 1057 | 27988 | 2531 | 317464 | 7881 |
| Facebook Movies | 8843 | 1103 | 13828 | 1848 | 166745 | 5427 |
| Facebook Books | 6231 | 263 | 9881 | 592 | 101520 | 1315 |

Table 7.5: Considered features in the different settings

7.4.3 Accuracy, Diversity and Novelty with kaHFM

In order to evaluate our approach, we measured accuracy, novelty and aggregate diversity metrics. Accuracy metrics were measured through Precision@N (*Prec@N*) and Normalized Discounted Cumulative Gain (*nDCG@N*). This latter measures the usefulness of an item based on its position in the recommendation list. Hence, it has been computed only for datasets that have explicit ratings (i.e., `LibraryThing` and `Yahoo!Movies`). Since the recommended results may vary in length depending on the user, cumulative gain at each position is normalized across users. *EPC* (Expected Popularity Complement) is used to measure novelty, or more precisely the ability of the algorithm to select items that belong to the long tail. Finally, diversity has been measured through *Catalog Coverage* (aggregate diversity in *top-N* list), *Gini Index* and *Shannon entropy*. In particular, *Catalog Coverage* denotes the ability of a system in selecting as many elements as possible from the whole catalog while *Gini index* (Gini) and *Shannon entropy* are used to measure the distributional inequality with different approaches. Both accuracy and novelty metrics have been computed by averaging their values per-user. To compute those metrics we used the implementation provided by RankSys¹³ framework. The evaluation has been performed considering Top-10 ([88, 108, 344]) recommendations for all the datasets. When a rating score was available, a *Threshold-based relevant items* condition [83] was adopted in order to take into account only relevant items. In particular, a relevance threshold of 4/5 and 8/10 has been set for `Yahoo!Movies` and `LibraryThing` respectively.

Tables 7.6, and 7.7 show the results of our experiments regarding accuracy and diversity.

In all the tables we highlight in **bold** the best result while we underline the second one. Statistically significant results are denoted with a * mark. We used a Student's paired t-test with a 0.05 level.

`LibraryThing` experiments show that our approach outperforms the competing algorithm for all the considered metrics. It is worth to notice that the *nDCG@N*

¹³<http://ranksys.org/>

| Categorical Setting (CS) | Facebook Movies | | | | | Facebook Music | | | | | Facebook Books | | | | |
|--------------------------|-----------------|---------------|-------------|---------------|----------------|----------------|---------------|-------------|---------------|----------------|----------------|---------------|-------------|---------------|----------------|
| | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0173* | 0.0196* | 3566 | 0.2217 | 9.5537 | 0.0200* | 0.0195* | 4961 | 0.1800 | 10.2039 | 0.0060* | 0.0056* | 1343 | 0.1498 | 9.0882 |
| BPR-FM | 0.0158* | 0.0149* | 179 | 0.0043 | 4.5639 | 0.0138* | 0.0107* | 544 | 0.0061 | 5.4518 | 0.0036* | 0.0034* | 78 | 0.0062 | 4.6815 |
| MostPopular | 0.0118* | 0.0099* | 27 | 0.0029 | 3.8543 | 0.0146* | 0.0089* | 30 | 0.0020 | 3.8628 | 0.0032* | 0.0030* | 17 | 0.0034 | 3.6193 |
| ItemKnn | 0.0262* | 0.0270* | 2554 | 0.0963 | 8.7266 | 0.0279* | 0.0269* | 3916 | 0.0960 | 9.2819 | 0.0041* | 0.0042* | 1437 | 0.1759 | 9.2526 |
| UserKnn | 0.0168* | 0.0157* | 466 | 0.0104 | 5.7206 | 0.0130* | 0.0109* | 961 | 0.0157 | 6.7938 | 0.0101* | 0.0095* | 269 | 0.0150 | 5.8499 |
| VSM | 0.0185* | 0.0205* | 3326 | 0.1769 | 9.5856 | 0.0289* | 0.0325* | 4582 | 0.1395 | 9.6626 | 0.0104* | 0.0112* | 1833 | 0.2631 | 9.8735 |
| kaHFM | 0.0296 | 0.0324 | 3560 | 0.2493 | 10.1243 | 0.0338 | 0.0353 | 4373 | 0.1166 | 9.5861 | 0.0129 | 0.0136 | 1905 | 0.3128 | 10.1855 |
| Ontological Setting (OS) | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0172 | 0.0188 | 3646 | 0.2589 | 9.9645 | 0.0295 | 0.0314 | 5152 | 0.2021 | 10.2691 | 0.0109* | 0.0113* | 2004 | 0.2954 | 10.0712 |
| BPR-FM | 0.0155* | 0.0128* | 144 | 0.0040 | 4.4288 | 0.0083 | 0.0065 | 207 | 0.0036 | 4.8614 | 0.0037* | 0.0034* | 85 | 0.0062 | 4.6887 |
| MostPopular | 0.0118* | 0.0099* | 27 | 0.0029 | 3.8543 | 0.0146* | 0.0089* | 30 | 0.0020 | 3.8628 | 0.0032* | 0.0030* | 17 | 0.0034 | 3.6193 |
| ItemKnn | 0.0263* | 0.0271* | 2557 | 0.0963 | 8.7270 | 0.0280 | 0.0270 | 3919 | 0.0960 | 9.2813 | 0.0041* | 0.0042* | 1437 | 0.1759 | 9.2526 |
| UserKnn | 0.0168* | 0.0157* | 466 | 0.0104 | 5.7207 | 0.0130 | 0.0109 | 961 | 0.0157 | 6.7938 | 0.0101* | 0.0095* | 269 | 0.0150 | 5.8493 |
| VSM | 0.0181* | 0.0198* | 3255 | 0.1765 | 9.5424 | 0.0257 | 0.0274 | 4231 | 0.1180 | 9.3889 | 0.0072* | 0.0085* | 1582 | 0.1691 | 9.1551 |
| kaHFM | 0.0273 | 0.0300 | 3631 | 0.2500 | 10.1188 | 0.0276 | 0.0288 | 5391 | 0.2537 | 10.5843 | 0.0148 | 0.0162 | 2099 | 0.3183 | 10.1984 |
| Factual Setting (FS) | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0234 | 0.0275 | 3605 | 0.2003 | 9.5874 | 0.0246 | 0.0258 | 5232 | 0.2336 | 10.4700 | 0.0147* | 0.0163* | 2098 | 0.3105 | 10.1480 |
| BPR-FM | 0.0157 | 0.0133 | 110 | 0.0039 | 4.3774 | 0.0119 | 0.0091 | 357 | 0.0049 | 5.2379 | 0.0041* | 0.0037* | 90 | 0.0067 | 4.7874 |
| MostPopular | 0.0123 | 0.0102 | 26 | 0.0029 | 3.8531 | 0.0114 | 0.0070 | 31 | 0.0020 | 3.8674 | 0.0033* | 0.0031* | 17 | 0.0034 | 3.6195 |
| ItemKnn | 0.0273 | 0.0283 | 2617 | 0.0983 | 8.7616 | 0.0289 | 0.0278 | 3993 | 0.0982 | 9.3085 | 0.0041* | 0.0040* | 1656 | 0.2117 | 9.4983 |
| UserKnn | 0.0176 | 0.0165 | 470 | 0.0106 | 5.7393 | 0.0136 | 0.0114 | 987 | 0.0159 | 6.8209 | 0.0109* | 0.0101* | 271 | 0.0157 | 5.9059 |
| VSM | 0.0219 | 0.0245 | 2812 | 0.0909 | 8.5810 | 0.0348 | 0.0366 | 3846 | 0.0925 | 9.1179 | 0.0126* | 0.0140* | 1862 | 0.2641 | 9.9322 |
| kaHFM | 0.0240 | 0.0268 | 3619 | 0.2434 | 10.1562 | 0.0313 | 0.0336 | 5350 | 0.2491 | 10.4870 | 0.0179 | 0.0189 | 2211 | 0.3523 | 10.3441 |

Table 7.6: Accuracy, Diversity and Novelty results for Facebook Movies, Facebook Music and Facebook Books

| Categorical Set. | LibraryThing | | | | | Yahoo!Movies | | | | | Last.fm | | | | | | |
|------------------|---------------|---------------|---------------|-------------|---------------|----------------|---------------|---------------|---------------|-------------|---------------|----------------|---------------|---------------|-------------|---------------|----------------|
| | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0408* | 0.0460* | 0.0424* | 6335 | 0.1607 | 11.1013 | 0.0421* | 0.1174* | 0.0528* | 2447 | 0.3640 | 10.0302 | 0.0223 | 0.0210 | 3642 | 0.1823 | 10.6194 |
| BPR-FM | 0.0151* | 0.0162* | 0.0138* | 527 | 0.0029 | 5.2783 | 0.0189* | 0.0344* | 0.0184* | 123 | 0.0056 | 4.3270 | 0.0314 | 0.0327 | 158 | 0.0022 | 4.6219 |
| MostPopular | 0.0056* | 0.0058* | 0.0051* | 34 | 0.0009 | 3.8302 | 0.0154* | 0.0271* | 0.0149* | 48 | 0.0043 | 3.9046 | 0.0252 | 0.0233 | 35 | 0.0012 | 3.7101 |
| ItemKnn | 0.0425* | 0.0590* | 0.0477* | 5939 | 0.1396 | 10.8123 | 0.0203* | 0.0427* | 0.0193* | 1442 | 0.1486 | 8.9548 | 0.0449 | 0.0405 | 2318 | 0.0947 | 9.8270 |
| UserKnn | 0.0213* | 0.0346* | 0.0226* | 1330 | 0.0103 | 6.7102 | 0.0231* | 0.0474* | 0.0232* | 729 | 0.0336 | 6.6712 | 0.0394 | 0.0443 | 1127 | 0.0248 | 7.6857 |
| VSM | 0.0367* | 0.0472* | 0.0393* | 7431 | 0.2106 | 11.4202 | 0.0385* | 0.1129* | 0.0496* | 2320 | 0.2893 | 9.7604 | 0.0025 | 0.0023 | 4582 | 0.0967 | 9.6626 |
| kaHFM | 0.0639 | 0.0913 | 0.0726 | 9367 | 0.3139 | 12.1071 | 0.0524 | 0.1399 | 0.0613 | 2433 | 0.3406 | 9.9831 | 0.0354 | 0.0364 | 4732 | 0.2976 | 11.5368 |
| Ontological Set. | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0446* | 0.0539* | 0.0485* | 7669 | 0.2139 | 11.4349 | 0.0427* | 0.1223* | 0.0545* | 2461 | 0.3634 | 10.0776 | 0.0297 | 0.0308 | 4278 | 0.2507 | 11.1528 |
| BPR-FM | 0.0121* | 0.0126* | 0.0108* | 396 | 0.0023 | 5.0055 | 0.0199* | 0.0356* | 0.0195* | 111 | 0.0053 | 4.2386 | 0.0287 | 0.0304 | 155 | 0.0022 | 4.5842 |
| MostPopular | 0.0056* | 0.0058* | 0.0051* | 34 | 0.0009 | 3.8302 | 0.0154* | 0.0271* | 0.0149* | 48 | 0.0043 | 3.9046 | 0.0252 | 0.0233 | 35 | 0.0012 | 3.7101 |
| ItemKnn | 0.0425* | 0.0591* | 0.0477* | 5939 | 0.1396 | 10.8125 | 0.0203* | 0.0427* | 0.0193* | 1442 | 0.1486 | 8.9548 | 0.0449 | 0.0405 | 2318 | 0.0947 | 9.8270 |
| UserKnn | 0.0213* | 0.0346* | 0.0226* | 1330 | 0.0103 | 6.7102 | 0.0232* | 0.0474* | 0.0232* | 729 | 0.0336 | 6.6711 | 0.0526 | 0.0577 | 823 | 0.0163 | 7.2512 |
| VSM | 0.0367* | 0.0472* | 0.0393* | 7431 | 0.2106 | 11.4202 | 0.0349* | 0.1083* | 0.0450* | 2216 | 0.2706 | 9.7345 | 0.0025 | 0.0026 | 4231 | 0.0819 | 9.3889 |
| kaHFM | 0.0635 | 0.0912 | 0.0728 | 9083 | 0.3081 | 12.0563 | 0.0521 | 0.1380 | 0.0608 | 2444 | 0.3442 | 10.0086 | 0.0371 | 0.0381 | 4853 | 0.3200 | 11.6318 |
| Factual Set. | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | nDCG@10 | EPC | AD@10 | Gini | SE | P@10 | EPC | AD@10 | Gini | SE |
| ABItem-kNN | 0.0488* | 0.0596* | 0.0513* | 7419 | 0.1968 | 11.2464 | 0.0619 | 0.1764 | 0.0777 | 2433 | 0.3177 | 9.7847 | 0.0376 | 0.0362 | 4179 | 0.2379 | 11.0850 |
| BPR-FM | 0.0087* | 0.0087* | 0.0080* | 184 | 0.0015 | 4.5400 | 0.0177 | 0.0305 | 0.0171 | 116 | 0.0054 | 4.2903 | 0.0233 | 0.0253 | 1815 | 0.0347 | 6.4129 |
| MostPopular | 0.0056* | 0.0058* | 0.0051* | 34 | 0.0009 | 3.8301 | 0.0154 | 0.0271 | 0.0149 | 48 | 0.0043 | 3.9046 | 0.0252 | 0.0233 | 35 | 0.0012 | 3.7103 |
| ItemKnn | 0.0436* | 0.0615* | 0.0491* | 6162 | 0.1471 | 10.9008 | 0.0203 | 0.0427 | 0.0193 | 1442 | 0.1486 | 8.9548 | 0.0426 | 0.0372 | 2351 | 0.0955 | 9.8092 |
| UserKnn | 0.0217* | 0.0349* | 0.0228* | 1399 | 0.0108 | 6.7459 | 0.0232 | 0.0474 | 0.0232 | 729 | 0.0336 | 6.6711 | 0.0403 | 0.0451 | 1163 | 0.0256 | 7.7004 |
| VSM | 0.0456* | 0.0575* | 0.0496* | 6860 | 0.1928 | 11.3774 | 0.0627 | 0.1725 | 0.0752 | 2203 | 0.2251 | 9.1183 | 0.0430 | 0.0419 | 3325 | 0.1590 | 10.4706 |
| kaHFM | 0.0627 | 0.0906 | 0.0713 | 9089 | 0.3134 | 12.0663 | 0.0564 | 0.1434 | 0.0639 | 2394 | 0.3511 | 10.1138 | 0.0339 | 0.0352 | 4788 | 0.3139 | 11.6257 |

Table 7.7: Accuracy, Diversity and Novelty results for LibraryThing, Yahoo!Movies and Last.fm

is almost doubled when using `kaHFM`. The worst results are achieved by Most-Popular, followed by BPR-FM and this happens in almost all the experiments. This is probably due to the high number of latent factors that does not make the method perform efficiently. `Yahoo!Movies` experiments show that in Categorical and Ontological settings our method is the most accurate, while the diversity performance of `kaHFM` w.r.t. `ABItem-kNN` are quite similar. In the `Yahoo!Movies` mapping a strong popularity bias is present and it is interesting to notice that this affects only the Factual setting leading our approach to be less precise than `ABItem-kNN` while our method proposes a bit more personalized recommendations as we can see through Gini index and Shannon entropy values. In `Last.fm`, Categorical and Ontological settings show that our method outperforms the others whereas in Factual setting the results are almost identical. In terms of catalog coverage and distributional inequality our approach achieves good results. In `Facebook Movies` we see very a good improvement in terms of accuracy as it almost doubles up the `ABItem-kNN` algorithm values. Diversity results show no relevant differences between `kaHFM` and `ABItem-kNN`. In `Facebook Music` the accuracy improvements are clear in Categorical and Factual settings while the Ontological setting seems to be the least descriptive setting because accuracy results come reduced w.r.t. the other settings, and they are quite similar to `ABItem-kNN`. Finally, in `Facebook Books`, `kaHFM` shows the best results for all the considered metrics. Let us discuss the baselines more related to our approach. We compared `kaHFM` against `ABItem-kNN` to check if the collaborative trained features may lead to better similarity values. This hypothesis seems to be confirmed since in former experiments `kaHFM` beats `ABItem-kNN` 16 times over 18. This suggests that collaborative trained features achieve better accuracy results. Moreover, we want to check if a knowledge-graph-based initialization of latent factors may improve the performance of Factorization Machines. `kaHFM` beats BPR-FM 18 times over 18, and in our opinion, this happens since the random initialization takes a while to drive the Factorization machine to reach good performance. Finally, we want to check if collaborative trained features lead to better accuracy results than a purely informativeness-based Vector Space Model even though it is in its knowledge-

graph-aware version. This seems to be confirmed in our experiments, since `kaHFM` beats `VSM` 15 times over 18.

In order to strengthen the results we got, we computed recommendations with 0,1,5,10,15,30 iterations. For the sake of brevity we report here only the plots related to Categorical setting shown in Figure 7.3. Results of the full experiments are available online¹⁴.

It is worth to notice that in every case we considered, we show the best performance in one of this iterations. Moreover, the positive influence of the initialization of the feature vectors is particularly evident in all the datasets, with performances being very similar to the ones depicted in [312]. Given the obtained results we may say that the answer to RQ1 is positive when adopting `kaHFM`.

7.4.4 Semantic Accuracy

The previous experiments showed the effectiveness of our approach in terms of accuracy, diversity and novelty. In practical terms, we proved that: (i) content initialization generally lead to better performance with our method, (ii) the obtained items vectors are fine-tuned better than the original ones for a *top-N* item recommendation task, (iii) results may depend on the features we extract from the Knowledge Graph. However, we still do not know if the original semantics of the features is preserved in the new space computed after the training of `kaHFM` (as we want to assess by posing RQ2). In Section 7.3.5 we introduced `Semantics Accuracy` ($SA@nM$) as a metric to automatically check if the importance computed by `kaHFM` and associated to each feature reflects the actual meaning of that feature.

Thus, we measured $SA@nM$ with $n \in \{1,2,3,4,5\}$ and $M = 10$, and evaluated the number of ground features available in the *top-nM* elements of \mathbf{v}_i for each of the six datasets.

Table 7.8 shows the results for all the different datasets computed in the Categorical setting. In general, the results we obtain are noteworthy. We now examine the worst one to better understand the actual meaning of the values we get. In Ya-

¹⁴<https://github.com/sisinflab/papers-results/tree/master/kahfm-results/>

| Semantics Accuracy | @M | @2M | @3M | @4M | @5M | F.A. |
|------------------------|-------|-------|-------|-------|-------|--------|
| Yahoo!Movies | 0.847 | 0.863 | 0.865 | 0.868 | 0.873 | 12.143 |
| LibraryThing | 0.960 | 0.996 | 0.998 | 0.999 | 0.999 | 3.820 |
| Last.fm | 0.960 | 0.987 | 0.991 | 0.994 | 0.995 | 6.615 |
| Facebook Music | 0.892 | 0.948 | 0.962 | 0.970 | 0.974 | 7.113 |
| Facebook Movies | 0.864 | 0.883 | 0.889 | 0.894 | 0.899 | 12.856 |
| Facebook Books | 0.995 | 1 | 1 | 1 | 1 | 3.133 |

Table 7.8: Semantics Accuracy results for different values of M. F.A. denotes the Feature Average number per item.

`hoo!Movies` Categorical setting, 747 different features compose each item vector (see Table 7.5). After the training phase, on average, more than 10 (equal to 0.847×12.143) over 12 features (last column in Table 7.8) are equal to the original features list. This means that `kaHFM` was able to compute almost the same features starting from hundreds of them. Even then, given the obtained results we may provide a positive answer to RQ2.

7.4.5 Generative Robustness

The previous experiment showed that the features computed by `kaHFM` keep their original semantics if already present in the item description. In section 7.3.6, we introduced a procedure to measure the capability of `kaHFM` to compute meaningful features. Here, we computed $1 - \text{Rob}@nM$ for the six adopted datasets. Results are represented in Table 7.9.

| 1-Robustness | @M | @2M | @3M | @4M | @5M | F.A. |
|------------------------|-------|-------|-------|-------|-------|--------|
| Yahoo!Movies | 0.487 | 0.645 | 0.713 | 0.756 | 0.793 | 12.143 |
| LibraryThing | 0.275 | 0.481 | 0.554 | 0.597 | 0.632 | 3.820 |
| Last.fm | 0.125 | 0.281 | 0.346 | 0.394 | 0.430 | 6.615 |
| Facebook Music | 0.714 | 0.893 | 0.935 | 0.955 | 0.966 | 7.113 |
| Facebook Movies | 0.821 | 0.945 | 0.970 | 0.980 | 0.984 | 12.856 |
| Facebook Books | 0.315 | 0.516 | 0.605 | 0.682 | 0.745 | 3.133 |

Table 7.9: 1-Robustness for different values of M. Column F.A. denotes the Feature Average number per item.

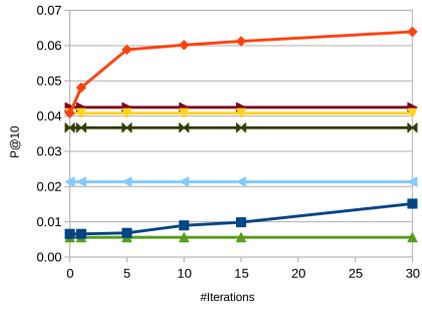
Even here, we focus on the CS setting. For a better understanding of the ob-

tained results, we start by focusing on `Yahoo!Movies` which apparently has a bad behavior. As we said before, Table 7.8 shows that `kaHFM` was able to guess 10 on 12 different features for `Yahoo!Movies`. In this experiment, we remove a feature thus making `kaHFM` to guess an average of 9 over 12 features. What we measure now is if `kaHFM` is able to guess the removed feature in the remaining 3 “slots”. Results in Table 7.9 show that our method is able to put the removed feature in one of the three slots the 48.7% of the times over 747 overall features. This example should help the reader to appreciate even more `Facebook Music` and `Facebook Movies` results. For the remaining datasets the situation is even much harder because there were no free slots (see Table 7.8). Thus, our method has only one missing slot to fill with the right feature. Let us take `Facebook Books` as an example: there are 263 different features in the item vector (see Table 7.5) and a very low average number of features per item (3.133). `kaHFM` is able to fill the missing slot with the right feature 31% of the times. The obtained results shows that `kaHFM` is able to propose meaningful features as we asked with RQ3.

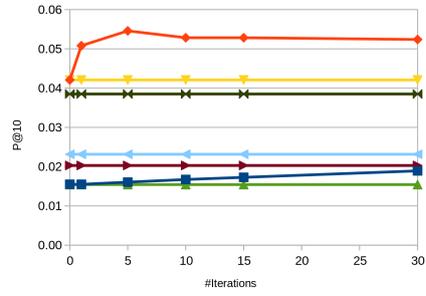
7.5 Conclusion

In this work, we have proposed an explainable method, `kaHFM`, in which we bind the meaning of latent factors for a Factorization machine to data coming from a knowledge graph. We evaluated `kaHFM` on six different publicly available datasets and compared it against state-of-the-art algorithms showing that our approach outperforms the other approaches with respect to accuracy, diversity, and novelty on different sets of semantics-aware features. In particular, we considered Ontological, Categorical and Factual information coming from a freely available knowledge graph. We have shown that the generated recommendation lists are more precise and personalized, and they select more items from the long tail. Summing up, performed experiments show that: (RQ1) the learned model shows very good performance in terms of accuracy, novelty and diversity and, at the same time, is effectively explainable; (RQ2) the computed features are semantically meaningful; (RQ3) the model is robust regarding computed features.

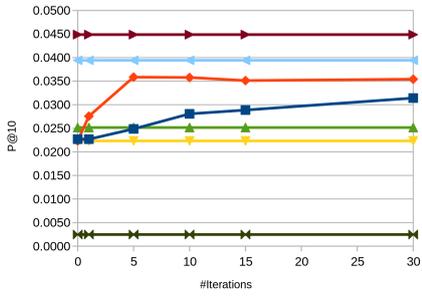
In the future we want to test the k_{aHFM} performance in different scenarios, other than recommender systems. Moreover, the model can be improved in many different ways. First of all, a stopping condition based on a validation set could be introduced to avoid wasteful training steps. Different relevance metrics could be beneficial in different scenarios, as the method itself is agnostic to the specific adopted measure. This work focused on the items' vector; however, an interesting key point would be analyzing the learned users' vectors to extract more accurate profiles. Furthermore, it would be useful to exploit k_{aHFM} in order to provide suggestions to knowledge graphs maintainers while adding relevant missing features to the knowledge base. In this direction, we would like to evaluate our approach in knowledge graph completion tasks.



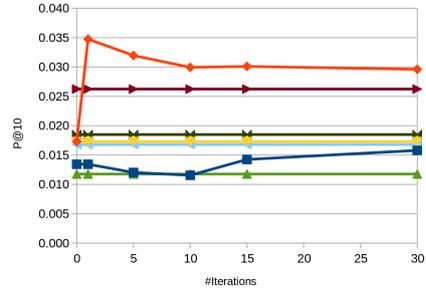
(a) Precision on LibraryThing



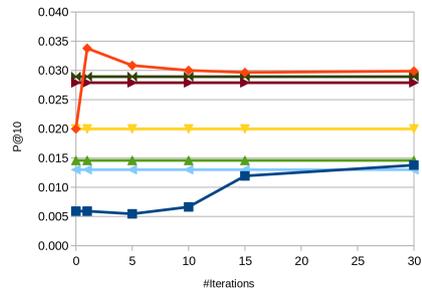
(b) Precision on Yahoo!Movies



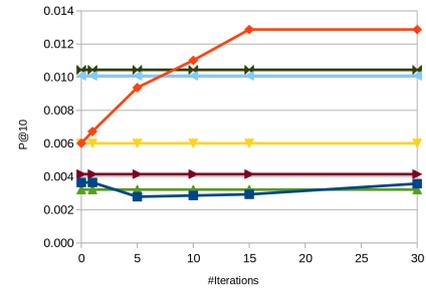
(c) Precision on Last.fm



(d) Precision on Facebook Movies



(e) Precision on Facebook Music



(f) Precision on Facebook Books



(g) legend

Figure 7.3: Precision@10 varying # iterations 0, 1, 5, 10, 15, 30

Chapter 8

Reasoning about Preferences

8.1 Introduction

In this line of research, the focus is on model-based preference reasoning, which relies on specific assumptions about the structure of the preference relation [152]. The simplest assumption that can be made is that the target ranking of a set of resources, described in terms of multiple attributes, can be represented as a lexicographical order [124]. Lexicographical preference models define orders of importance on the attributes that describe the objects of a domain. As an example, consider the choice of a movie. Typically, the most important attribute that one considers is the genre of the movie (e.g., drama, superhero, etc.). Then, among the movies of the preferred genre, the choice can rely on the movie's actor (e.g., Tom Hanks, Christian Bale, etc.). The assumption of a lexicographical order restricts significantly the hypothesis space, but induces a bias rarely justified in practical applications. In fact, preferences on individual attributes are generally not independent of each other. With reference to the movie domain, Tom Hanks may be preferred to Christian Bale, if the movie genre is drama, while Christian Bale may be preferred in case of a su-

perhero film. *CP-nets* [71] offer a language to express preferences on the values of single attributes, and, at the same time, allow to model dependencies of this type. A CP-net is a qualitative graphical representation that reflects conditional dependence and independence of preferences under a *ceteris paribus* (all else being equal) interpretation. It is a compact representation of a complex preference relation (partial order), where each node refers to a single attribute and is associated with a function that assigns a preference relation on the values of that attribute to each combination of the values of the parent attributes. More precisely, CP-nets require that the user specifies (i) for any attribute A of interest, which other attributes can impact her preferences for values of A (the parents of A), and (ii) for each instantiation of the parent attributes, the preference ordering over values of A . Points (i) and (ii) could make CP-nets a rather rigid formalism, compared to the expressive needs of a user. Furthermore, some statements that are very natural for the user to assess cannot be represented within a CP-net. *Conditional preference theories* (or *CP-theories*) [395] are a more general and flexible formalism for qualitative preferences that allows to go beyond the expressiveness limitations of CP-nets.

In our previous work [320], we focused on the well-known CP-net graphical language and have addressed the problem of preference representation and reasoning with Linked Data from different perspectives. We have proposed a vocabulary to represent statements formulated according to the *ceteris paribus* semantics and have shown how to encode a CP-net by means of this vocabulary. Inspired by [162], we have also explained how to embed such a compact preference model into a SPARQL 1.1 query in order to access semantic data in a personalized way. The current investigation extends the leading motivation and the approach of a previous work [320], but embraces the more general and flexible formalism of CP-theories. We point out that the approach proposed here only deals with context-uniform conditional (cuc) acyclic CP-theories [395], which are a special type of CP-theories exposing nice polynomial computational properties while comparing outcomes.

This study heavily extends the approach presented in [320] in many points. First of all, here we deal with the much more expressive CP-theories instead of CP-nets. A new extended vocabulary as well as a completely new algorithm to encode

CP-theories in SPARQL is also proposed. Moreover, an implementation of the overall framework is presented together with experimental evaluations targeted at assessing the users' experience in representing their preferences as CP-theories and the performance of the tool. The main contributions of this work can be summarized as follows:

- presentation of RDF vocabularies to represent qualitative preference statements over Linked Data, built on top of the vocabulary proposed in [320], but adjusted for more general preference statements;
- an encoding into RDF triples of the qualitative preferential information represented by a CP-theory and the exploitation of theoretical results of [394] to compute a partial order over items for acyclic cases;
- a procedure to translate conditional preference statements into a SPARQL 1.1 query able to retrieve a **ranked list** of resources whose order reflects the user preferences;
- an application framework that meets a user's needs while representing her preferences as a CP-theory encoded in RDF and that eventually allows a SPARQL-enabled software agent to retrieve a ranked list of resources according to the user's tastes;
- An experimental evaluation to verify the effectiveness of the proposed approach.

The rest of the chapter is structured as follows. Section 8.3 presents the motivating scenario that fostered the overall approach. The semantics of CP-theories and a recap on some relevant results and theorems has been provided in background, Section 2.2.1. In Section 8.4, we propose an RDF vocabulary to represent a CP-theory with the preferential statements of a user, and then we show how to embed the RDF version of the CP-theory into a SPARQL query able to retrieve a ranked list of results ordered according to user's preferences. Section 8.7 describes the tool that supports the user both in the formulation of her preferences under the CP-theories semantics

and in retrieving the resources of interest ordered according to her preferences. The results on a user study are reported in Section 8.8, where we also measure the performance of the implemented system on a synthetic dataset. Section 8.2 provides an overview of related work about preference reasoning and enabling query languages with preferences. Conclusions close the chapter.

8.2 Related work

The ability to infer, model, and reason with user preferences has been recognized as a prominent research direction in many fields, especially artificial intelligence (AI) [131, 303]. Preferences are generally classified as quantitative, if they make use of a scoring function to assess an order over the available resources, resulting in a total order, or qualitative, if they are treated independently, resulting in incomparable resources and a partial preference order. Much work has focused on qualitative approaches, since these are closer to how people express their preferences; among the earliest logic-based approaches is von Wright's [385]. Following the overview over qualitative multi-attribute preference reasoning approaches provided by [284], in AI, there are, in particular, (i) methods adopting graphical structures to represent and reason about preferences, e.g., CP-nets [71] and TCP-nets [72]; (ii) methods that extend constraints satisfaction problems and incorporate soft constraints, as in the approximation of CP-nets with soft constraints described in [128, 127]; and (iii) methods that use specific logic-based languages to represent qualitative preferences and derive utility functions, exploiting, e.g., machine learning techniques, such as support vector machines [126, 125].

Conceptually close in spirit to our investigation is in particular [276], where ontological knowledge expressed via existential rules in Datalog[±] is combined with CP-theories to represent qualitative conditional preferences along with domain knowledge, and to perform preference-based answering of conjunctive queries. Another related work [102] combines Datalog with CP-theories, but only considers atomic queries. Our work, in contrast, focuses on SPARQL queries in a more restricted ontological context and conditional preferences specified via *cuc*-acyclic CP-theories.

There is also a large body of work on handling preferences in logic programming, e.g., asprin [75], which is a framework for handling preferences among the stable models of a logic program. Similarly, the qualitative choice logic [74] is a propositional logic for representing a preference relation among models, which allows to specify alternative, ranked options for problem solutions. The above two and similar works on handling preferences in logic programming are fundamentally different from our approach, as they are about preferences for ordering models of a logic program, rather than preferences for ordering the answers to a query subject to all models of a knowledge base.

Databases are another research area where preferences have been investigated. In a relational database management system, for example, the *top-k* (or *ranking*) queries represent a quantitative approach, since they return the k best matches according to a numerical score. In [234], a formalism supporting ranking queries for a relational database is presented. With reference to the qualitative approach instead, *skyline* queries [67] extend the notion of best matching to contexts, where multiple independent scores have to be taken into account. The result of a *skyline* query is a set of objects that are no worse than any other across all dimensions of a set of independent Boolean or numerical preferences [67]. Within the database community, both Chomicki [98, 99] and independently Kießling and colleagues [213, 214] formalized the first examples of *preference-based querying* languages, that is, extensions of SQL that support the specification of quantitative and qualitative queries.

The notion of preference is of primary importance also in the Linked Open Data context. The provision of means to enable users to look for data sources (e.g., SPARQL endpoints) and data content that is tailored to their individual preferences is one of the target of the original project by Tim Berners-Lee et al. Even the motivating example proposed in the introductory article about the Semantic Web [60] can be interpreted as a preference-based search, as extensively discussed in [347]. Based on this insight, in [347], the authors add preference-based querying capabilities to the most known Semantic Web query language, SPARQL. However, when the paper was published, it was not possible to specify multiple (independent) preference dimensions in SPARQL, and consequently the authors had to introduce

the `PREFERRING` solution modifier. For example, the following query provides a preference-enabled SPARQL query for a user who is searching for an appointment, preferring excellent therapist, appointments out of the rush hour and later appointments over earlier ones, if both are equal with respect to the rush hour constraint.

```

1 SELECT ?appointment WHERE {
2   ?therapist :rated ?rating;
3   :offers ?appointment.
4   ?appointment :starts ?start;
5   :ends ?end.
6   PREFERRING (?rating = excellent AND
7   ?end < 1600 || ?start > 1800
8   CASCADE HIGHEST(?start))
9 }

```

At line 6, the `PREFERRING` clause behaves as a solution modifier, and the `AND` keyword separates independent preference dimensions. The `CASCADE` keyword at line 8 allows to give higher priority to the left-hand preference over the right-hand one. In their paper, the authors state that within the same SPARQL query, the use of a `LIMIT k` statement in combination with `PREFERRING` ones could inform the query evaluator to go deeper in the retrieval of skyline solutions, thus allowing the system to return a set of results ordered by user preferences.

A mapping operation between an OWL ontology, called *OWLPref*, and the SPARQL Preference syntax of [347] has been proposed by [231]. *OWLPref* allows for representing in a declarative, domain-independent and machine-interpretable way several kinds of preferences, namely, *SimplePreference*, *CompositePreference* (which makes compositions of the former), and *ConditionalPreference* (which models preferences that vary according to the context, thanks to a property *OnCondition*). However, considering the unavailability of conditional preferences in the SPARQL Preference syntax of [347], the use of *ConditionalPreference* in *OWLPref* seems of marginal utility.

The *PrefSPARQL* syntax of [162] keeps the goal of identifying the Pareto-optimal set, but introduces preferences at the level of filters. It still uses the `AND` to separate independent dimensions and to build what the authors call *MultidimensionalPref*. Each “dimension” is either a *conditional* preference (`IF-THEN-ELSE`) or

an *atomic* preference, which in turn can be a simple expression or can involve more complex constructs [162]. Besides the support for conditional preferences and the substitution of `CASCADE` with `PRIOR TO`, the main innovative point of [162] with respect to [347] is perhaps that the proposed preference-enabled query can be completely rewritten using SPARQL 1.1 characteristics. In particular, [162] uses the `FILTER NOT EXISTS`. The translation of the previous query according to the *PrefSPARQL* query rewriting is given below.

```

1 SELECT ?appointmentA WHERE {
2 ?terapistA :rated ?ratingA;
3 :offers ?appointmentA.
4 ?appointmentA :starts ?startA;
5 :ends ?endA.
6 BIND ((?ratingA = :excellent) AS ?Pref1A)
7 BIND ((?endA < 16 || ?startA > 18:00) AS ?Pref2A)
8 BIND ((?startA) AS ?Pref3A)
9 FILTER NOT EXISTS {
10 ?terapistB :rated ?ratingB;
11 :offers ?appointmentB.
12 ?appointmentB :starts ?startB;
13 :ends ?endB.
14 BIND ((?ratingB = :excellent) AS ?Pref1B)
15 BIND ((?endB < 1600 || ?startB > 1800) AS ?Pref2B)
16 BIND ((?startB) AS ?Pref3B)
17 FILTER (
18 ((?Pref1B > ?Pref1A) &&
19 !((?Pref2B < ?Pref2A) ||
20 (?Pref3B < ?Pref3A && ?Pref2B = ?Pref2A)))
21 ||
22 (!(?Pref1B < ?Pref1A) &&
23 ((?Pref2B > ?Pref2A) ||
24 (?Pref3B > ?Pref3A && ?Pref2B = ?Pref2A))))}
25 }

```

The query looks for appointments `?appointmentA` satisfying a certain pattern expressed in lines 2–5. The research is carried out checking that there is no `?appointmentB` that verifies the same pattern (lines 10–13) and dominates `?appointmentA` in any preference dimension. The example refers to only two independent preference dimensions and the situations when `?appointmentB` dominates `?appointmentA` are represented in the branches of `||` symbol at line 21, that is, lines 18–20 and lines 22–24. For the sake of completeness, `?appoint-`

ment_B would dominate ?appointment_A if it was better in one dimension (line 18 or line 23–24) and no worse in the other one (line 19–20 or line 22). The `PRIOR TO` preference relation is encoded in lines 19–20 and 23–24 through the `||` operator.

Although *PrefSPARQL* allows the user to encode conditional preferences in a SPARQL 1.1 query, it differs from the approach that we presented here in at least three main aspects: (i) in *PrefSPARQL*, the focus is on computing the most preferred solution (undominated outcome), given a set of conditional preferences, while we provide a list of results ordered by user preferences; (ii) they deal with conditional preferences in the form $u_\varphi : x_\varphi > x'_\varphi$, while our approach is able to manage CP-statements in the form $u_\varphi : x_\varphi > x'_\varphi [W_\varphi]$, which result to be much more expressive for finite and discrete domains even in their ccc-acyclic version that we consider here; (iii) we provide an ontological vocabulary and a procedure to automatically encode preferences in a SPARQL query. However, thanks to its more agile structure, differently from our approach, *PrefSPARQL* allows the user to express preferences on variables with continuous domains as well as the usage of comparison operators.

In [371], the authors present SPREFQL, an extension of the SPARQL language that allows for appending a “PREFER” clause, which expresses soft preferences over the query results obtained by the main body of the query. The main ideas behind the approach are to associate relations of tuples with preference formulas, and to select the relations’ most preferred tuples via a so-called winnow operator. Consequently, the approach does not allow for expressing conditional ceteris paribus preferences as in CP-theories.

As a general remark, previous SPARQL-related works on preference reasoning have been mainly devoted to preference representation and the retrieval of undominated outcomes. In principle, one may encode each of them in a procedural approach able to compute the first level of undominated solutions, then the second one and so on. At each iteration step, the procedure should be able to filter out the results coming from the “higher levels” of results computed in the previous steps.

Less closely related are approaches to preference-based query answering over graph databases. In particular, [149] presents regular languages for graph queries, where answers are partially ordered via a partial order on the strings of the lan-

guages. In the same vein, [159] introduces preferential regular path queries for enhanced querying of semi-structured databases. Query symbols are annotated with preference weights for scaling up or down the importance of matching a symbol against a database edge label. The paper studies (progressive) query answering, (certain) query answering in LAV data-integration systems, and query containment and equivalence. A similar approach in [146] introduces a graph query language that enables to declaratively express preferences. None of the above approaches to preference-based query answering over graph databases (which are intuitively based on (potentially recursive) pattern-recognition-style regular expressions), however, allows for expressing conditional *ceteris paribus* preferences as in CP-theories. Less closely related are also information retrieval systems based on manipulating fuzzy truth values (which may also be interpreted as quantitative preferences), such as the fuzzy multimedia retrieval system in [138].

8.3 Motivating scenario

The leading scenario behind the framework that we propose here is that of a distributed system where a user may pose a query to a SPARQL endpoint and have the returned results ordered with respect to a set of personal preferences on a specific knowledge domain. For instance, a user might be willing to get a list of books to read by querying DBpedia and then have it ranked according to a set of preferences hosted on their own Web page. A possible implementation of the approach that we propose is depicted in Fig. 8.1,¹ where the main building blocks required to implement the whole framework are shown:

- a reference model to encode and reason with preferences (CP-theory in our case);
- an ontological vocabulary to represent preferences by adopting Web languages;

¹For ease of presentation, we use DBpedia as the main dataset to query. The approach can be adapted to any Linked Data dataset.

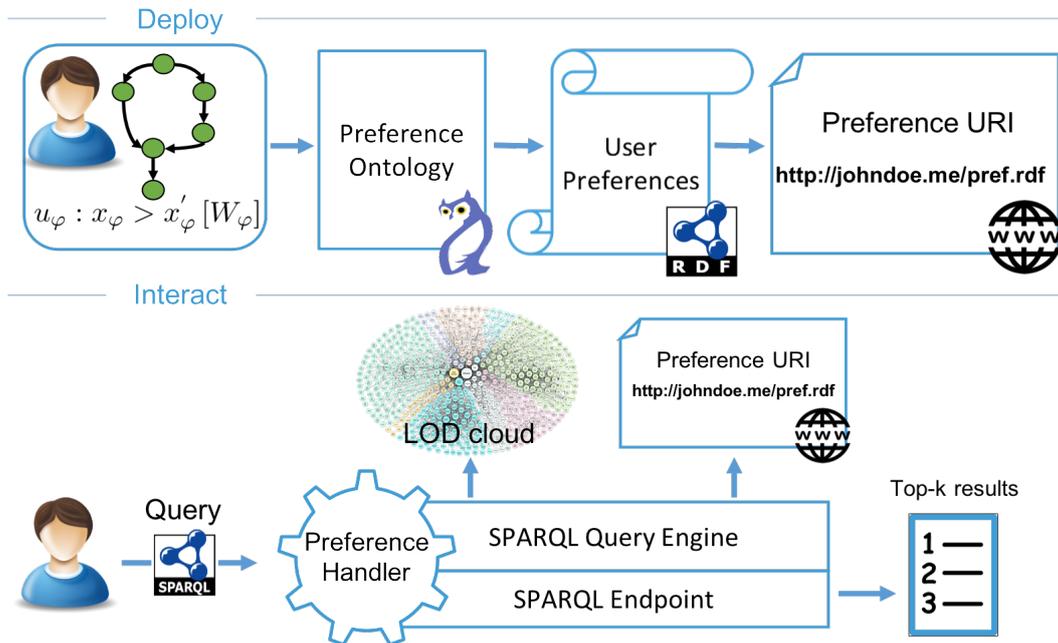


Figure 8.1: A graphical representation of the proposed approach. The Deploy and Interact steps only rely on the adoption of standard technologies and languages.

- a tool able to handle and manage preferences as well as to encode them in a set of SPARQL queries.

In order to implement the whole framework, we propose the following deployment and interaction steps.

Deploy.

- *Model user preferences.* We adopt CP-theories [394] as reference model to represent user preferences. As we will see in Section 2.2.1, they are a formalism to represent and reason with sets of qualitative preferences.²
- *Encode preferences in RDF by means of a preference ontology.* We developed an OWL ontology to represent CP-theory statements encoding user preferences (see Section 8.4).

²The interest here is not whether such preferences are automatically learned from data or manually modeled and set by a human agent.

- *Publish user preferences publicly on the Web.* The aim is to foster a principled adoption of user preferences by systems interested in providing a personalized access to data available in the Linked Data cloud. User preferences can be encoded and published in an RDF file.

Interact.

- *Use SPARQL to get user preferences.* A preference handler loads³ the preference model of the user encoded in RDF by means of a general purpose SPARQL query engine.
- *The preference handler formulates SPARQL 1.1 queries able to retrieve and order resources by taking into account user preferences.* Standard SPARQL queries are formulated in order to rank the result set with reference to the preferences expressed by the user.

A strong requirement that we had in mind while developing our solution was to use only standard Web technologies and languages to implement the overall framework. Indeed, we selected RDF to model user preferences and the power of SPARQL 1.1 to perform preference-based reasoning in order to rank the results of a query. In fact, as we will see in Section 8.4.2, advanced features of the last version of the SPARQL query language can be employed to perform preference reasoning over a model based on CP-theories. In Section 8.3, all the details needed to implement the Deploy and Interact steps in a pure Linked Data setting will be provided.

8.4 CP-theories and Linked Open Data

As we stated in Section 8.1, the target of this work is twofold. On the one hand, we want to supply the user with a vocabulary to represent qualitative statements formulated in terms of *ceteris paribus* semantics. On the other hand, we aim to provide an encoding of user preferences that can be used in a *top-k* query answering scenario. In this section, we start by proposing a first ontology that allows a system to

³See *SPARQL 1.1 Update* specification at <https://www.w3.org/TR/sparql11-update/#load>.

represent preferential statements according to CP-theories in a very straightforward way and an extended version to manage the directed graph $J_o(\Gamma)$ on V introduced in Section 2.2.1 and exploited in Theorem 2. Note that RDF triples encoding the directed graph can be automatically derived from the original preferential statements. In Section 8.7, we provide a description of an implemented tool to infer the RDF version of the directed graph, starting from a set of conditional preferences. We deal with the complementary target in Section 8.4.2, where we show how to employ a user profile represented as an instantiation of the extended ontology to encode the corresponding preferences in a standard SPARQL query able to retrieve and rank resources in a personalized way.

Figure 8.2 shows the ontology that we modeled to express user profiles in terms of CP-theory statements⁴. The main idea behind the modeling of the ontology is that we may express preferences on properties of items that the user is looking for, e.g., `dbo:literaryGenre`, `dbo:country`, `dbo:subsequentWork`, or potentially `dbo:filmVersion`. Hereafter, the ontology in Figure 8.2 will be referred to as the *lite* ontology. The aim of the *lite* ontology is that of creating an ontological vocabulary providing all the elements to syntactically represent conditional preference statements in a theory Γ . By means of this ontology, it is possible to encode whatever preference φ in its general form $u_\varphi : x_\varphi > x'_\varphi [W_\varphi]$.

The ontology is composed by four main classes and nine properties. The class `Value` represents possible values of a variable. If we look at the book *GoodKnyght!* in Example 5, we see that the “actual values” for which the user expresses a preference are composed by both a *property* (e.g., `dbo:country`) and its related *object* (e.g., `db:UnitedKingdom`). This is the reason why the class `Value` is the domain of the two properties `attribute` and `value`. The former mapping the property, the latter mapping the object. `Condition` is used to express the conditional part u_φ of a preference statement $u_\varphi : x_\varphi > x'_\varphi [W_\varphi]$, which is also the condition for the relative importance [72] of the variable X_φ over variables in W_φ in case $W_\varphi \neq \emptyset$. It is the domain of the property `contains`, whose range

⁴The corresponding OWL file is available at http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_light.owl

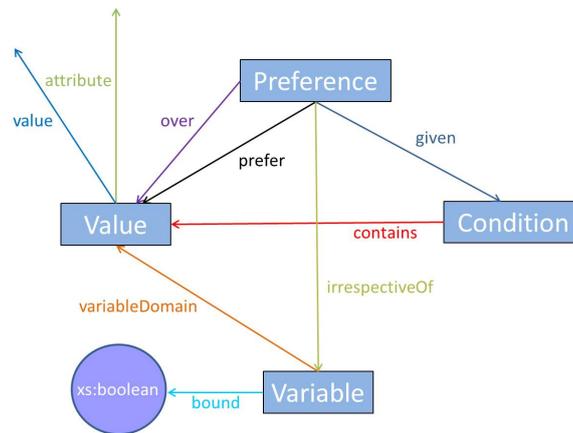


Figure 8.2: A graphical representation of the *lite* version of the ontology proposed to represent conditional statements.

is Value. The class Preference represents the whole conditional statement φ . The properties having Preference as domain reflect the structure of the preferential statement “given a Condition, I prefer a Value over another Value, optionally irrespectiveOf some other variables”. The class Variable is used to model the variables of a CP-theory, and it is domain of variableDomain, whose range is Value. Finally, we have the bound property needed to explicitly state if the value associated to an attribute is an actual value (as for `dbo:country` and `dbo:literaryGenre`) or if it represents the situation that we (do not) have a triple involving the attribute, as for `dbo:subsequentWork` or `dbo:filmVersion`. Here, we adopted the modeling choice of representing directly the conditions generated by the combination for the values of variables in U_φ instead of relating the variables themselves. As an example, in case we have “*Given an English book that is part of a saga, I prefer...*”, with reference to the previous example, the corresponding encoding will be⁵

```
cpl:combined-cond a cpl:Condition ;
cpl:contains cpl:c1 ;
cpl:contains cpl:s1 .
```

⁵Here, we use the prefix `cpl` to denote `<http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_light.owl\#>`

Here, `cpl:combined-cond` represents u_φ as a whole, while `cpl:c1` and `cpl:s1` represent the values x_1 and x_2 composing $u_\varphi = x_1x_2$.

We will see how this modeling choice will be useful when embedding the CP-theory into a SPARQL query.

Notice that such classes and predicates are sufficient for the user to express the CP-theory with her preferential statements, and to facilitate the user experience even further, we built a user-friendly tool, where preferences related to a specific domain (e.g., books or movies) can be added. The aforementioned tool will be extensively described in Section 8.7.

Example 5 (Books cont'd)

With respect to Giorgio's preference "*given an English book, he prefers those being part of a saga*", if we look in DBpedia, we may find, for instance, the book *GoodKnyght!*. Indeed, we have:

```
@prefix db: <http://dbpedia.org/resource/>
@prefix dbo: <http://dbpedia.org/ontology/>

db:GoodKnyght! a dbo:Book ;
dbo:country db:United_Kingdom ;
dbo:subsequentWork db:Whizzard! .
```

From the previous RDF statements, we see that Giorgio's preference refers to values of objects in a triple with reference to a specific predicate. Indeed, given a set of resources of type `dbo:Book` such that the value for `dbo:country` is `db:United_Kingdom`, he prefers those with an associated triple whose predicate is `dbo:subsequentWork`. In order to be fully compliant with the Linked Data technological stack, we need a vocabulary/ontology that allows users to represent their preferences on different attributes of resources that they might be interested in. Hence, with reference to the ontology in Fig. 8.2, we have the following RDF triples modeling the preference introduced at the beginning of this example.

```
@prefix cpl: <http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_light.owl#>
@prefix db: <http://dbpedia.org/resource/>
```

@prefix dbo: <http://dbpedia.org/ontology/>

Variables

cpl:country a cpl:Variable;
cpl:bound true;
cpl:variableDomain cpl:c1,cpl:c2.

cpl:subsequentWork a cpl:Variable;
cpl:bound false;
cpl:variableDomain cpl:s1, cpl:s2.

Values allowed for each variable

cpl:c1 a cpl:Value;
cpl:attribute dbo:country;
cpl:value db:United_Kingdom.

cpl:c2 a cpl:Value;
cpl:attribute dbo:country;
cpl:value db:France.

cpl:s1 a cpl:Value;
cpl:attribute dbo:subsequentWork;
cpl:value cpl:subsequentWorkYes.

cpl:s2 a cpl:Value;
cpl:attribute dbo:subsequentWork;
cpl:value cpl:subsequentWorkNo.

Condition

cpl:cond a cpl:Condition;
cpl:contains cpl:c1.

Preference

cpl:pref a cpl:Preference;
cpl:given cpl:cond;
cpl:prefer cpl:s1;
cpl:over cpl:s2.

■

Once we have defined and modeled Γ in RDF, in order to compare two outcomes α and β as in Theorem 2, we may build the RDF version of the directed graph \triangleright_α on V (see Section 2.2.1).

To this aim, the *lite* ontology is extended as shown in Figure 8.3⁶ to what we call the *full* ontology. Hence, starting from a set of preferences represented via the *lite* ontology, we derive its *full* version such that, once instantiated with an outcome α , it represents \triangleright_α . The derivation step is performed by adding new statements via the following relations:

moreImportantThan has `Variable` both as domain and as range, and it is used to model the transitive closure of dependency and (unconditional) relative importance information, i.e., for the transitive closure of edges $U_\varphi \rightarrow \{X_\varphi\} \cup W_\varphi$, and, if $U_\varphi = \emptyset$, for $\{X_\varphi\} \rightarrow W_\varphi$, for any φ .

conditionallyMoreImportantThan takes into account conditional relative importance information. Its range is an instance of the new class `InstanceOfRelativeImportance`, which is used to represent a pair (u_φ, Z) , $Z \in W_\varphi$, for a statement φ with $U_\varphi \neq \emptyset$.

In fact, the class `InstanceOfRelativeImportance` is the domain of the property `hasCondition`, whose range is the `Condition` instance representing u_φ , and of the `hasLessImportantVariable` property, whose range is the `Variable` instance for Z . In the following, we say “ X is `conditionallyMoreImportantThan` Y under the `Condition` C ” when X is linked via `conditionallyMoreImportantThan` to an instance of the class `InstanceOfRelative-Importance`, which, in turn, is linked by `has-Condition` to a `Condition` C and by `has-LessImportant-Variable` to a `Variable` Y .

⁶The corresponding OWL file is available at http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_full.owl

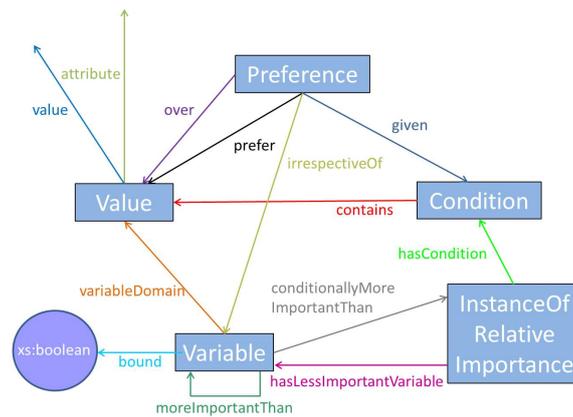


Figure 8.3: A graphical representation of the *full* ontology proposed to represent conditional statements.

Both properties must act as transitive relations. Therefore, we add more statements to the original preferences by means of the following rules involving the transitive closure of property `moreImportantThan` and then, for pairs of variables (Y,Z) not linked by it, the transitive closure of property `conditionallyMoreImportantThan`.

- if variable Y is `moreImportantThan` a variable X and at the same time variable X is `conditionallyMoreImportantThan` a variable Z , under a condition C not involving values of Y , then we add the RDF statements representing that “ Y is `conditionallyMoreImportantThan` Z under the condition C ”;
- if Y is `conditionallyMoreImportantThan` a variable X under a condition C , X is at the same time `moreImportantThan` a variable Z , and the condition C does not involve any value of Z , then the additional fact we add is “ Y is `conditionallyMoreImportantThan` Z under the condition C ”;
- if Y is `conditionallyMoreImportantThan` a variable X under a condition C , and X is `conditionallyMoreImportantThan` a variable Z under a condition C' , then the additional fact added to the knowledge base is “ Y is `conditionallyMoreImportantThan` Z under the condition

C'' , where $C'' = C \wedge C'$ is the condition joining C and C' , but only if $C \wedge C'$ does not contains values of Y, Z or two different values of any other variable.

According to the properties just introduced, given a pair of outcomes (α, β) , if there exist no variable $X' \in \Delta(\alpha, \beta)$ that is moreImportantThan X and no variable $X'' \in \Delta(\alpha, \beta)$ that is conditionallyMore-ImportantThan (u, X) , with α extending u , then a variable X is in $\Theta'(\alpha, \beta)$.

Example 6 (Books cont'd)

The *full* encoding corresponding to $\Gamma_{C-LG-SW-F}$ for Giorgio's preferences in Table 2.1, is represented in Listing 8.1⁷. ■

```

@prefix cpt: <http://sisinflab.poliba.it/semanticweb/
lod/ontologies/cpt_full.owl#>
@prefix db: <http://dbpedia.org/resource/>
@prefix dbo: <http://dbpedia.org/ontology/>

cpt:country1 a cpt:Value;
cpt:attribute dbo:country;
cpt:value db:United_Kingdom.
cpt:country2 a cpt:Value;
cpt:attribute dbo:country;
cpt:value db:France.
cpt:genre1 a cpt:Value;
cpt:attribute dbo:literaryGenre ;
cpt:value db:Crime_fiction.
cpt:genre2 a cpt:Value;
cpt:attribute dbo:literaryGenre ;
cpt:value db:Autobiographical_novel.
cpt:sw1 a cpt:Value;
cpt:attribute dbo:subsequentWork;
cpt:value cpt:subsequentWorkYes.
cpt:sw2 a cpt:Value;
cpt:attribute dbo:subsequentWork;
cpt:value cpt:subsequentWorkNo.
cpt:film1 a cpt:Value;
cpt:attribute dbo:filmVersion;
cpt:value cpt:filmVersionYes.
cpt:film2 a cpt:Value;
cpt:attribute dbo:filmVersion;
cpt:value cpt:filmVersionNo.

cpt:conditionC1 a cpt:Condition;
cpt:contains cpt:country1.
cpt:conditionC2 a cpt:Condition;
cpt:contains cpt:country2.
cpt:conditionSW1 a cpt:Condition;
cpt:contains cpt:sw1.
cpt:conditionSW2 a cpt:Condition;

```

⁷For conciseness, the prefix `cpt` is always assumed in the following for `<http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_full.owl\#>`.

```

cpt:contains cpt:sw2.

cpt:country a cpt:Variable;
cpt:bound true;
cpt:variableDomain cpt:country1,cpt:country2;
cpt:moreImportantThan cpt:literaryGenre;
cpt:moreImportantThan cpt:subsequentWork;
cpt:moreImportantThan cpt:filmVersion.
cpt:literaryGenre a cpt:Variable;
cpt:bound true;
cpt:variableDomain cpt:genre1,cpt:genre2;
cpt:conditionallyMoreImportantThan
cpt:instanceOfRelativeImportance1;
cpt:conditionallyMoreImportantThan
cpt:instanceOfRelativeImportance3.
cpt:subsequentWork a cpt:Variable;
cpt:bound false;
cpt:variableDomain cpt:sw1, cpt:sw2;
cpt:moreImportantThan cpt:filmVersion;
cpt:conditionallyMoreImportantThan
cpt:instanceOfRelativeImportance2.
cpt:filmVersion a cpt:Variable;
cpt:bound false;
cpt:variableDomain cpt:film1, cpt:film2.

cpt:instanceOfRelativeImportance1
a cpt:instanceOfRelativeImportance;
cpt:hasCondition cpt:conditionC1;
cpt:hasLessImportantVariable cpt:subsequentWork.
cpt:instanceOfRelativeImportance2
a cpt:instanceOfRelativeImportance;
cpt:hasCondition cpt:conditionC2;
cpt:hasLessImportantVariable cpt:literaryGenre.
cpt:instanceOfRelativeImportance3
a cpt:instanceOfRelativeImportance;
cpt:hasCondition cpt:conditionC1;
cpt:hasLessImportantVariable cpt:filmVersion.

cpt:preference1 a cpt:Preference;
cpt:prefer cpt:country2;
cpt:over cpt:country1.
cpt:preference2 a cpt:Preference;
cpt:given cpt:conditionC1;
cpt:prefer cpt:genre2;
cpt:over cpt:genre1;
cpt:irrespectiveOf cpt:subsequentWork.
cpt:preference3 a cpt:Preference;
cpt:given cpt:conditionC2;
cpt:prefer cpt:sw2;
cpt:over cpt:sw1;
cpt:irrespectiveOf cpt:literaryGenre.
cpt:preference4 a cpt:Preference;
cpt:given cpt:conditionC1;
cpt:prefer cpt:sw1;
cpt:over cpt:sw2.
cpt:preference5 a cpt:Preference;
cpt:given cpt:conditionC2;
cpt:prefer cpt:genre1;
cpt:over cpt:genre2.
cpt:preference6 a cpt:Preference;
cpt:given cpt:conditionSW1;
cpt:prefer cpt:film1;
cpt:over cpt:film2.
cpt:preference7 a cpt:Preference;

```

```
cpt:given cpt:conditionSW2;  
cpt:prefer cpt:film2;  
cpt:over cpt:film1.
```

Listing 8.1: The RDF version of the CP-theory $\Gamma_{\text{C-LG-SW-F}}$ in Table 2.1, according to the *full* ontology.

In Section 8.4.2, we will see how to pose a SPARQL query against the *full* version of Γ in order to compute \gg_{Γ} , thus retrieving a ranked list of semantic resources ordered according to a user’s preferences.

8.4.1 The special case of CP-nets

Before moving to the description of how to encode a CP-theory in a SPARQL query for personalized results ranking, we point out that the ontological model just described subsumes the vocabulary introduced in [320] to represent CP-nets models. There, we proposed how to model the information encoded in a CP-net through an ontology and how to formulate the query able to order outcomes accordingly in a consistent way. On the other hand, we know that CP-nets are a special and simple case of CP-theories and, therefore, we want to show how the new ontology in Figure 8.3 can deal with a CP-net.

The theoretical result exploited in [320] (namely Corollary 4 of [71]) orders an outcome o over another one o' , consistently with a CP-net \mathcal{N} , if there exists a variable X such that o and o' assign the same values to all ancestors of X in \mathcal{N} and given the assignment provided by o (and o') to the parents of X , i.e., $Pa(X)$, o assigns a more preferred value to X than that assigned by o' (according to the conditional preference table of X). The sufficient condition of Corollary 4 of [71] may be reformulated asking for a variable X such that there does not exist any variable `moreImportantThan` X different in o and o' , and, given the assignment provided by o (and o') to $Pa(X)$, o assigns a more preferred value to X than that assigned by o' . The predicate `moreImportantThan`, in fact, covers dependency information and, when applied to a CP-net, allows to define a set of variables coincident with the ancestor set. Moreover, for a CP-net, the predicate `given` for any instance of `Preference` ordering the values of a variable may be used to define

the parent set of that variable. The ontological model proposed in [320] for CP-nets is hence subsumed by the *full* ontology of Figure 8.3. On the other hand, if one wants to represent a CP-net according to the *full* ontology, one has to consider that the dependency information is the only kind of information required for CP-nets, because there is no relative importance encoded. This implies that the RDF version of the CP-net, in terms of the *full* ontology, would not contain any predicate `conditionallyMoreImportantThan` and `irrespectiveOf` or any instance of the class `InstanceOfRelativeImportance`.

8.4.2 Ordering SPARQL results via CP-theories

In the following, we assume that users are looking for the best k items satisfying some requirements and that the choice for the best ones is led by their preferences, formulated according to a CP-theory Γ , on a set of variables $V = \{X_1, \dots, X_n\}$. Hence, we aim at solving a top- k query answering problem, where the ordering criterion is encoded in Γ . In the presented approach, we concentrate on acyclic CP-theories, to preserve the nice computational properties introduced in Section 2.2.1 and exploit the algorithmic approach suggested by Theorem 2.

The ultimate goal of our proposal can be summarized by the following query formulated in a meta-language on top of SPARQL:

```
SELECT ?item
WHERE {
    ?item satisfies user requirements
}
ORDER BY  $\Gamma$ 
LIMIT  $k$ 
```

Here, *user requirements* are represented by a SPARQL graph pattern where at least one triple has `?item` as subject. In the following, we use the notation $\mathcal{R}(\text{?item})$ to denote the user requirements associated with the variable `?item`.

Example 7 (Books cont'd)

“Giorgio really wants to relax, and so he is looking only for books with more than 300 pages”. In this case, Giorgio’s requirements $\mathcal{R}(?item)$ are represented by:

```
?item a dbo:Book.
?item dbo:numberOfPages ?page.
FILTER(?page>300).
```

■

The computation of an answer to the previous query, goes through the exploitation of the *full* version of Γ . The overall approach is composed of two main steps.

Step 1. Here, we compute a representation of α and β , starting from each $\varphi \in \Gamma$ and, by exploiting Theorem 2, an ordering based on \gg_{Γ} for all possible outcomes is eventually returned. The representation of α and β as URIs goes through a string concatenation (we use the `GROUP_CONCAT` aggregate function of SPARQL).

Step 2. This step deals with ranking the items matching $\mathcal{R}(?item)$ according to \gg_{Γ} as computed in the previous step.

Both steps are detailed in the following.

Ordering the Outcomes (Step 1). From Theorem 2 in Section 2.2.1, we know how to build a strict partial order on a set of outcomes \mathcal{O} extending $>_{\Gamma}$ by comparing outcomes via \gg_{Γ} . By means of the same meta-language that we used before, the ordering of outcomes can be done via the following query. There, we see the outcomes are ordered according to a counter representing the number of outcomes that they are able to \gg_{Γ} -dominate.

```
SELECT ?outcome-Dominating
      (COUNT(?outcome-dominated) AS ?counter)
WHERE {
  FILTER { ?outcome-Dominating  $\gg_{\Gamma}$  ?outcome-dominated }
}
```

```

GROUP BY ?outcome-Dominating
ORDER BY DESC(?counter)

```

In order to compute values for the two variables `?outcome-Dominating` and `?counter`, the previous query should act on one pair (α, β) per time by checking if $\alpha \gg_{\Gamma} \beta$.

The preference-based reasoning is performed exclusively by means of the SPARQL 1.1 query *OrderingQuery* whose generation is detailed in Algorithm 2 of Appendix 8.5.

8.5 Query formulation algorithm for CP-theories

Algorithm 2 has the user's preferences graph \mathcal{G}_{user} as input, that is, the RDF version of the user's CP-theory Γ in terms of the *full* ontology and returns the SPARQL query able to order outcomes according to \gg_{Γ} , a strict partial order extending $>_{\Gamma}$ (see Theorem 2). Line 3 computes, for each outcome o , the values (of variables in V) that it is composed of, through the string *Outcome_D_values* built with the for cycle that ends on line 2. The string *Outcome_D_values* contains just the names of variables in V with a suffix *D*. Line 3 also computes the number of outcomes o' that o dominates according to \gg_{Γ} , referred to as `?counter`. The counting is made possible by the combination of the `COUNT` in line 3 and the `GROUP BY` in line 18. The `?counter` variable is then used by the `ORDER BY` in line 19 to rank the result set. It is worth to notice that line 3 asks for `URI (?outcome_D)` and not just for `?outcome_D`, since this entity will be employed as the subject of triples at the beginning of **Step 2** described in Section 8.4.2. Given a pair of outcomes (o, o') , lines 4 to 16 are used to identify the set of variables $\Theta'(o, o')$ and among them the set of variables $\{X \in \Theta'(o, o'): o(X) \succ_o^X o'(X)\}$. The first nested subquery (lines 5 to 13) considers one pair of outcomes at a time, `?outcome_D` and `?outcome_d`, where `D` and `d` stand respectively for *Dominating* and *dominated*. The for loop of lines 4 until 9 allows to consider first `?outcome_D` and then `?outcome_d`. For each of them, the nested loop of lines 5 to 13 introduces the values corresponding to

the variables in V . For each variable $X_i \in V$, Algorithm 2 looks for values $?X_i-y$ filtering only elements of the set $\{value(x_{i1}), value(x_{i2})\}$ in the binary case, or $\{value(x_{i1}), \dots, value(x_{in})\}$ elsewhere, through the VALUES assignment at line 6. The algorithm requires that a list W of variables is defined from the set of variables V of Γ : the variables that refer to each instance of Condition must appear in the same order in W , optionally allowing some recurrence. At lines 7-9, the outcome is explicitly built by concatenating, according to the order imposed by W , the values extracted for various $?X_i-y$ together with $attribute(X_i)$, for all members of W . Line 10 is added to verify that the pair of outcomes to compare is made of distinct elements.

At line 11, the patterns and the FILTER are used to identify the variables $?V$ with different values $?value1$ and $?value2$ in the outcomes $?outcome_D$ and $?outcome_d$, namely the variables in the set $\Delta(?outcome_D, ?outcome_d)$. As imposed by the couple of FILTER NOT EXISTS of lines 12 and 13, these variables $?V$ are such that:

- there does not exist any variable $?V2$ in the set $\Delta(?outcome_D, ?outcome_d)$ that is `more-ImportantThan` $?V$;
- there does not exist any variable $?V3$ in the set $\Delta(?outcome_D, ?outcome_d)$ that is `conditionallyMoreImportantThan` $?V$ under a Condition extended by $?outcome_D$.

In particular, the nested subquery appearing within the FILTER NOT EXISTS of line 13 extracts only the instances of Condition extended by $?outcome_D$, building the representative strings of the conditional values, concatenating them in $?Concatenated$ and checking the inclusion of the string $?Concatenated$ in $?outcome_D$ through FILTER and CONTAINS. The pair of FILTER NOT EXISTS of lines 12 and 13 allows therefore to identify the set of variables $\Theta'(?outcome_D, ?outcome_d)$.

The UNION of Query 1 and Query 2 is added at line 14. It returns a set of quadruples of the general form $\langle ?V, ?ConcatenatedParent, ?Prefer, ?Over \rangle$, able to order an outcome over another one, locally with respect to $?V$. For

each variable $?V$ within the set $\Theta'(?outcome_D, ?outcome_d)$, the IF of line 15 verifies if one of the quadruples on variable $?V$ can be used to order $?outcome_D$ over $?outcome_d$ locally with respect to $?V$. In particular, for quadruples with a missing value for $?ConcatenatedParent$, it is sufficient to verify if $?outcome_D$ contains the better value of a preference, i.e., $?Prefer$, and $?outcome_d$ contains the relative worse value $?Over$. Instead, for quadruples with a bound value for $?ConcatenatedParent$, it must happen that $?outcome_D$ contains the value of $?ConcatenatedParent$, as well as $?Prefer$ and $?outcome_d$ contain the value of $?Over$. If one of the $||$ (or) conditions happens, the BIND instantiate the value of $?counterBind$ to 1, otherwise to 0. The $?counterBind$ value is summed up across all instantiation of $?V$ in $\Theta'(?outcome_D, ?outcome_d)$, and it resolves into $?counterV$ at line 4. The same line also computes the cardinality of $\Theta'(?outcome_D, ?outcome_d)$, namely, the value $?counterVundominated$. The FILTER at line 17 verifies if the pair of values $?counterVundominated$ and $?counterV$ coincides, that is, if $?outcome_D(X) \succ_{?outcome_D}^X ?outcome_d(X)$ for all variables X in $\Theta'(?outcome_D, ?outcome_d)$. In conclusion, if the FILTER of line 17 returns true then $?outcome_D$ dominates $?outcome_d$ with respect to \gg_Γ , and its $?counter$ value is incremented of a unit. Only distinct dominated outcomes are counted, through the solution modifier DISTINCT at line 3.

If we consider the query in Appendix 8.6 resulting from the running example on Giorgio's preferences, we see that the variables involved in CP-statements as well as the corresponding values are encoded in the initial part of the query (line 9-50) and in the GROUP BY statement (line 130). The remaining of the query is quite standard and does not depend on the underlying CP-theory Γ . As for the initial part of the query, it contains a number of $2 \cdot |V|$ VALUES statements, where we assign all the allowed values to variables in V and a BIND statement used to compose the strings representing all possible outcomes. We emphasize that the whole query is automatically generated by Algorithm 2, starting from the *full* version of Γ , and then the process is completely transparent to the user.

8.6 Ordering Query for the Book Example

```
1 prefix cpt:<http://sisinflab.poliba.it/semanticweb/lod/ontologies/cpt_full.owl#>
2 prefix dbpedia-owl:<http://dbpedia.org/ontology/>
3 prefix dbpedia:<http://dbpedia.org/resource/>
4 prefix g:<http://sisinflab.poliba.it/semanticweb/graphs/>
5 SELECT (URI(?outcome_D) AS ?URIOutcome) ?genre_D ?country_D ?subwork_D
6 ?filmVersion_D (COUNT(DISTINCT ?outcome_d) AS ?counter)
7 WHERE
8 {
9 { SELECT DISTINCT ?outcome_D ?outcome_d ?genre_D ?country_D ?subwork_D
10 ?filmVersion_D (COUNT(DISTINCT ?V) AS ?counterVundominated)
11 (SUM((?counterBind)) AS ?counterV)
12 WHERE
13 {
14 {SELECT DISTINCT ?outcome_D ?outcome_d ?V ?genre_D ?country_D ?subwork_D
15 ?filmVersion_D
16 WHERE
17 {
18 VALUES (?genre_D) {
19 (dbpedia:Crime_fiction) (dbpedia:Autobiographical_novel)
20 }
21 VALUES (?country_D) {
22 (dbpedia:France) (dbpedia:United_Kingdom)
23 }
24 VALUES (?subwork_D) {
25 (cpt:subsequentWorkYes) (cpt:subsequentWorkNo)
26 }
27 VALUES (?filmVersion_D) {
28 (cpt:filmVersionYes) (cpt:filmVersionNo)
29 }
30 BIND (CONCAT(STR(dbpedia-owl:country),STR(?country_D),
31 STR(dbpedia-owl:literaryGenre),STR(?genre_D),
32 STR(dbpedia-owl:subsequentWork),STR(?subwork_D),
33 STR(dbpedia-owl:filmVersion),STR(?filmVersion_D)) AS ?outcome_D).
34 VALUES (?genre_d) {
35 (dbpedia:Crime_fiction) (dbpedia:Autobiographical_novel)
36 }
37 VALUES (?country_d) {
38 (dbpedia:France) (dbpedia:United_Kingdom)
39 }
40 VALUES (?subwork_d) {
41 (cpt:subsequentWorkYes) (cpt:subsequentWorkNo)
42 }
43 VALUES (?filmVersion_d) {
44 (cpt:filmVersionYes) (cpt:filmVersionNo)
45 }
```

```

46 BIND (CONCAT (STR(dbpedia-owl:country), STR(?country_d),
47 STR(dbpedia-owl:literaryGenre), STR(?genre_d),
48 STR(dbpedia-owl:subsequentWork), STR(?subwork_d),
49 STR(dbpedia-owl:filmVersion), STR(?filmVersion_d)) AS ?outcome_d).
50 FILTER(?outcome_D!=?outcome_d).
51
52 ?V a cpt:Variable.
53 ?V cpt:variableDomain ?variable1. ?variable1 cpt:value ?value1.
54 ?V cpt:variableDomain ?variable2. ?variable2 cpt:value ?value2.
55 FILTER (!(?value1=?value2) && CONTAINS(?outcome_D, STR(?value1))
56 && CONTAINS(?outcome_d, STR(?value2))).
57 FILTER NOT EXISTS{
58 ?V2 cpt:moreImportantThan ?V.
59 ?V2 cpt:variableDomain ?vd1. ?vd1 cpt:value ?v1.
60 ?V2 cpt:variableDomain ?vd2. ?vd2 cpt:value ?v2.
61 FILTER (!(?v1=?v2) && CONTAINS(?outcome_D, STR(?v1)) &&
62 CONTAINS(?outcome_d, STR(?v2))).
63 }
64 FILTER NOT EXISTS{
65 ?V3 cpt:conditionallyMoreImportantThan ?instanceOfRelativeImportance.
66 ?instanceOfRelativeImportance cpt:hasCondition ?C.
67 ?instanceOfRelativeImportance cpt:hasLessImportantVariable ?V.
68 ?V3 cpt:variableDomain ?vd13. ?vd13 cpt:value ?v13.
69 ?V3 cpt:variableDomain ?vd23. ?vd23 cpt:value ?v23.
70 { SELECT DISTINCT ?C
71 (GROUP_CONCAT (CONCAT (STR(?attr), STR(?value)); separator =""))
72 AS ?Concatenated) WHERE{
73 ?C cpt:contains ?c.
74 ?c cpt:attribute ?attr; cpt:value ?value.
75 }
76 GROUP BY ?C
77 }
78 FILTER (CONTAINS (?outcome_D, ?Concatenated)).
79 FILTER (!(?v13=?v23) && CONTAINS (?outcome_D, STR(?v13)) &&
80 CONTAINS (?outcome_d, STR(?v23))).
81 }
82 }
83 }
84 { SELECT DISTINCT ?V ?ConcatenatedParent ?Prefer ?Over {
85 { SELECT ?V
86 (CONCAT (STR(?attrPrefer), STR(?valuePrefer)) AS ?Prefer)
87 (CONCAT (STR(?attrPrefer), STR(?valueOver)) AS ?Over) WHERE
88 {?preference cpt:prefer ?p;
89 cpt:over ?o.
90 FILTER NOT EXISTS {?preference cpt:given ?condition.}
91 ?V cpt:variableDomain ?p.
92 ?p cpt:attribute ?attrPrefer;

```

```

93  cpt:value ?valuePrefer.
94  ?o cpt:value ?valueOver.
95  }
96  }
97  UNION
98  { SELECT DISTINCT ?V ?ConcatenatedParent ?Prefer ?Over WHERE
99  {
100 SELECT DISTINCT ?condition ?V
101 (CONCAT(STR(?attrPrefer),STR(?valuePrefer)) AS ?Prefer)
102 (CONCAT(STR(?attrPrefer),STR(?valueOver)) AS ?Over)
103 (GROUP_CONCAT(CONCAT(STR(?attr),STR(?value));separator = "")
104 AS ?ConcatenatedParent) WHERE
105 { ?preference cpt:given ?condition.
106 ?preference cpt:prefer ?p;
107 cpt:over ?o.
108 ?V cpt:variableDomain ?p.
109 ?p cpt:attribute ?attrPrefer;
110 cpt:value ?valuePrefer.
111 ?o cpt:value ?valueOver.
112 ?condition cpt:contains ?c.
113 ?c cpt:attribute ?attr;
114 cpt:value ?value.
115 }
116 GROUP BY ?condition ?V ?attrPrefer ?valuePrefer ?valueOver
117 }
118 }
119 }
120 }
121 BIND(IF( ( (!BOUND(?ConcatenatedParent) &&
122 CONTAINS(?outcome_D,?Prefer)&&
123 CONTAINS(?outcome_d,?Over))
124 ||
125 (BOUND(?ConcatenatedParent) && ?ConcatenatedParent!="" &&
126 CONTAINS(?outcome_D,?ConcatenatedParent) &&
127 CONTAINS(?outcome_D,?Prefer) &&
128 CONTAINS(?outcome_d,?Over)) ) ,1,0) AS ?counterBind )
129 }
130 GROUP BY ?outcome_D ?genre_D ?country_D ?subwork_D ?filmVersion_D ?outcome_d
131 }
132 FILTER(?counterV=?counterVundominated)
133 }
134 GROUP BY ?outcome_D ?genre_D ?country_D ?subwork_D ?filmVersion_D
135 ORDER BY DESC (?counter)

```

The algorithm takes as input the *full* version of Γ and computes a query able to return a list of outcomes ordered according to `?counter`. In particular, the query

returns for each outcome, a numerical score representing its position in the ranking imposed by \gg_{Γ} .

For a better clarification, the reasoning procedure under the comparison between the pair (α, β) is summarized in the following:

1. The query computes the set $\Theta'(\alpha, \beta)$ by considering the variables X in the set $\Delta(\alpha, \beta)$ for which there do not exist: (i) a variable $X' \in \Delta(\alpha, \beta)$ linked to X by property `cpt:moreImportant-Than` and (ii) a variable $X'' \in \Delta(\alpha, \beta)$ which is `cpt:conditionallyMoreImportant-Than` than X under a condition extended by α .
2. It then counts the number of variables X from the set $\Theta'(\alpha, \beta)$ that let to state $\alpha(X) \succ_{\alpha}^X \beta(X)$ and compares it to the cardinality of $\Theta'(\alpha, \beta)$;
3. If the numerical values coincide, which means that for each variable X in $\Theta'(\alpha, \beta)$, $\alpha(X) \succ_{\alpha}^X \beta(X)$ holds, then the query concludes that $\alpha \gg_{\Gamma} \beta$.

In order to get all the information needed to check $\alpha(X) \succ_{\alpha}^X \beta(X)$ from the *full* version of Γ , *OrderingQuery* embeds Query 1 and Query 2 reported in the following. They return a set of quadruples $\langle ?V, ?ConcatenatedParent, ?Prefer, ?Over \rangle$, with `?ConcatenatedParent` optionally not instantiated, able to locally order α over β with respect to variable `?V`.

Given a variable $X_i \in V$ with $dom(X_i) = \{x_{i1}, x_{i2}, \dots, x_{in}\}$, we use the following notation relative to the corresponding instances `cpt:xi1`, `cpt:xi2,...`, `cpt:xin` of the class `cpt:Value`:

- $value(x_{ij})$ is the object of the triple `cpt:xij cpt:value` object;
- $attribute(x_{ij})$ denotes the object of the triple `cpt:xij cpt:attribute` object;
- we call *representative string* of x_{ij} the concatenation of the two strings represented by $attribute(x_{ij})$ and $value(x_{ij})$ respectively. The combination of $attribute(x_{ij})$ and $value(x_{ij})$ is used to represent x_{ij} , as they uniquely identify a value in the domain of a variable. Indeed, in case we used only $value(x_{ij})$, ambiguous situations could arise when it is used in combination with different attributes.

Finally, for an instance `cpt:c` of the class `cp:Condition`, we call *conditional values* of `cpt:c` all the objects of the triples `cpt:c cpt:contains object`.

Query 1

```

1 SELECT ?V
2 (concat(str(?attrPrefer),str(?valuePrefer)) as ?Prefer)
3 (concat(str(?attrPrefer),str(?valueOver)) as ?Over)
4 WHERE
5 {
6   ?preference cpt:prefer ?p;
7   cpt:over ?o.
8   FILTER NOT EXISTS {?preference cpt:given ?condition.}
9   ?V cpt:variableDomain ?p.
10  ?p cpt:attribute ?attrPrefer;
11  cpt:value ?valuePrefer.
12  ?o cpt:value ?valueOver.
13 }

```

Query 1 processes elements φ of Γ with $u_\varphi = \top$. Within the query, they are represented by the variable `?preference`. The selection is made possible by the `FILTER NOT EXISTS` on the pattern `{?preference cpt:given ?condition.}`

(line 8). Considering that the objects of properties `cpt:prefer` and `cpt:over` must be distinct values of the same variable, the query firstly extracts the variable that the preference acts on, i.e., `?V` (line 1 and line 9). Then, it computes the *representative strings*, `?Prefer` and `?Over` (lines 2–3) for the objects `?p` and `?o` of the two triples involving `cpt:prefer` and `cpt:over` (lines 10–12).

Query 2

```

1 SELECT DISTINCT ?V ?ConcatenatedParent ?Prefer ?Over WHERE{
2 SELECT ?condition ?V
3 (GROUP_CONCAT(concat(str(?attr),str(?value)); separator="")
4 as ?ConcatenatedParent)
5 (concat(str(?attrPrefer),str(?valuePrefer)) as ?Prefer)
6 (concat(str(?attrPrefer),str(?valueOver)) as ?Over)
7 WHERE
8 {
9   ?preference cpt:given ?condition;
10  cpt:prefer ?p;
11  cpt:over ?o.
12  ?V cpt:variableDomain ?p.
13  ?p cpt:attribute ?attrPrefer;
14  cpt:value ?valuePrefer.
15  ?o cpt:value ?valueOver.
16  ?condition cpt:contains ?c.

```

```

17 ?c cpt:attribute ?attr;
18 cpt:value ?value.
19 }
20 GROUP BY ?condition ?V ?attrPrefer ?valuePrefer ?valueOver
21 }

```

Differently from the previous query, Query 2 is used to process statements φ belonging to Γ with $u_\varphi \neq \top$. The selection is made via the pattern $\{?preference\text{cpt:given } ?condition.\}$ (line 9). Let us consider first the nested subquery in lines 2–20. For each instance of class `Preference`, such query extracts the variable `?V` that the preference is about (lines 2 and 12) and considers the `cpt:given condition ?condition` (line 9), extracting its corresponding *conditional values* (line 16). The *representative strings* of such *conditional values* are then computed (lines 17–18) and concatenated at lines 3–4 in `?ConcatenatedParent`, grouping by condition. The variables `?Prefer` and `?Over` are defined similarly to Query 1. The external query is just used to restrict the result set to variables `?V`, `?ConcatenatedParent`, `?Prefer`, `?Over`.

Ordering the Items (Step 2). Given the information on outcomes returned by the *OrderingQuery* at the previous step, on both values of variables and position in the ranking, an external RDF dataset, e.g., `DBpedia`, may be queried, asking for items satisfying the hard constraints ($\mathcal{R}(?item)$) imposed by the user and such that, when limiting the attention on variables in V , they match the description of an outcome. Items are then ordered according to the ranking over corresponding outcomes.

We are well aware that the one we propose is just a possible rewriting of a CP-theory in a SPARQL query and other encodings are possible, ever more efficient from a computational perspective. Moreover, we may see that the performance of the overall approach decreases when the size of variable domains grows and, in its current version, the approach is not able to handle continuous domains as for distance and time. Nevertheless, we believe that the proposed approach is a good starting point to reason with preferences in a pure Linked Data environment, as it is a straight implementation of theoretical results coming from previous works [395].

8.6.1 Instantiation of the framework

The procedure to retrieve items ordered according to user's preferences is made up of four phases:

- the *loading* of user's preferences;
- an *insert* to add information about outcomes;
- the execution of a *federated query*;
- the (optional) *dropping* of user's preferences.

First of all, the user's preferences file representing the *full* version of Γ is loaded in the SPARQL server through a LOAD operation and becomes the user's graph of preferences \mathcal{G}_{user} .

Example 8 (Book cont'd)

If the path to the RDF file encoding Giorgio's preferences (see Listing 8.1) is generally denoted as *path_to_ttl_file*, the load operation is executed as follows:

```
prefix g:<http://sisinflab.poliba.it/semanticweb/graphs/>

LOAD path.to.ttl.file INTO GRAPH g:Giorgio.preferences
```



The *OrderingQuery* able to order the outcomes according to \gg_{Γ} is then executed. The information returned by the *OrderingQuery* is used to integrate the graph of user preferences \mathcal{G}_{user} with additional triples on outcomes. Specifically, we add information about the score of an outcome and its description. Hence, the following triples are defined for each outcome:

- a triple satisfying the pattern

```
?URIOutcome cpt:hasScore ?score
```

- a set of triples instantiating the pattern

```
?URIOutcome cpt:hasValueForX
?ValueForX
```

for every variable X of Γ .

Such information are added to \mathcal{G}_{user} through an INSERT query.

Example 9 (Book cont'd)

For the CP-theory $\Gamma_{C-LG-SW-F}$ with Giorgio's preferences, the INSERT operation would behave as follows:

```
prefix cpt:<http://sisinflab.poliba.it/semanticweb/
lod/ontologies/cpt-full.owl#>
prefix dbo:<http://dbpedia.org/ontology/>
prefix db:<http://dbpedia.org/resource/>
prefix g:<http://sisinflab.poliba.it/semanticweb/graphs/>

INSERT { GRAPH g:Giorgio_preferences
{?URIOutcome cpt:hasScore ?counter .
?URIOutcome cpt:hasValueForCountry ?country_D;
cpt:hasValueForLiteraryGenre ?genre_D;
cpt:hasValueForSubsequentWork ?subwork_D;
cpt:hasValueForFilmVersion ?filmVersion_D.
}
}
where { GRAPH g:Giorgio_preferences {
```

OrderingQuery

```
}
}
```

Here, *OrderingQuery* denotes the ordering query over outcomes returned by Algorithm 2 in 8.5. The output of Algorithm 2 applied to Listing 8.1 is available in 8.6. ■

The next step is the execution of a federated query, composed by two subqueries. The first subquery retrieves the items satisfying the requirements imposed by the user, i.e., $\mathcal{R}(?item)$, and for each item it looks for the values of variables in Γ . The retrieval of values grounds on the VALUES construct for variables that are `cpt:bound true` and on the combination of BIND, IF and EXISTS otherwise.

The second subquery on the user's preferences graph \mathcal{G}_{user} , retrieving for each outcome its score and the variables values. A matching between items and outcomes is hence performed through such values and the items are finally ordered according to the position in the ranking of the relative outcome.

The main reason behind the federation of two (or more) endpoints is that: while the graph containing the RDF version of the user's preferences is encoded in the corresponding document (available at *Preference URI* in Fig. 8.1), all the information about the items that we want to retrieve and rank is encoded in a separate dataset, e.g., DBpedia. The main assumption here is that the user's preferences are expressed with respect to a reference dataset/vocabulary, which can be queried via a SPARQL endpoint.

Example 10 (Book cont'd)

Suppose that Giorgio is interested in the top-5 list of books matching his hard constraints (see Example 7), ordered according to his preferences encoded in the CP-theory $\Gamma_{C-LG-SW-F}$ of Table 2.1. The federated query to carry out the searching task would be as follows:

```

prefix cpt:<http://sisinflab.poliba.it/semanticweb/
lod/ontologies/cpt.full.owl#>
prefix dbo:<http://dbpedia.org/ontology/>
prefix db:<http://dbpedia.org/resource/>
prefix g:<http://sisinflab.poliba.it/semanticweb/graphs/>

SELECT ?item_D ?score WHERE {
{SERVICE <http://dbpedia.org/sparql> {
SELECT DISTINCT ?item_D ?genre_D ?country_D
?subwork_D?filmVersion_D WHERE{
?item_D a dbo:Book;
dbo:numberOfPages ?page_D.
FILTER(?page_D>300).
?item_D dbo:literaryGenre ?genre_D;
dbo:country ?country_D.
VALUES (?genre_D) {
(db:Crime_fiction)
(db:Autobiographical_novel)
}
VALUES (?country_D) {
(db:France)
(db:United_Kingdom)
}
}
BIND(IF(EXISTS{?item_D dbo:subsequentWork ?object},
cpt:subsequentWorkYes, cpt:subsequentWorkNo
AS ?subwork_D).
BIND(IF(EXISTS{?item_D dbo:filmVersion ?object},
cpt:filmVersionYes,cpt:filmVersionNo)
AS ?filmVersion_D).

```

```

}
}
}
{graph g:Giorgio_preferences {
SELECT ?score ?genre_D ?country_D ?subwork_D
?filmVersion_D WHERE{
?s cpt:hasScore ?score;
cpt:hasValueForCountry ?country_D;
cpt:hasValueForLiteraryGenre ?genre_D;
cpt:hasValueForSubsequentWork ?subwork_D;
cpt:hasValueForFilmVersion ?filmVersion_D.
}
}
}
}
ORDER BY DESC(?score)
LIMIT 5

```



Finally, the graph with the user's preferences \mathcal{G}_{user} can be optionally eliminated with a DROP operation.

Please note that all SPARQL queries are executed in a *simple entailment between RDF graphs* on the *full* version of Γ as well as on the external dataset for the federated query. Hence, all the queries are executed under the *RDF Entailment Regime* of SPARQL.

Example 11 (Book cont'd)

The graph related to Giorgio's preferences may be dropped as follows:

```

prefix g:<http://sisinflab.poliba.it/semanticweb/graphs/>

DROP GRAPH g:Giorgio_preferences

```



8.7 Application

We now describe a tool⁸ implementing the framework described in previous sections and aimed at supporting the end-user in retrieving a list of semantic resources

⁸The tool is available at <http://cptheorysparql.cloudapp.net:10002/>

Figure 8.4: Preferences Insertion.

ordered according to her preferences formulated under the CP-theory formalism. The tool just asks for preferential statements formulated under the CP-theory formalism, i.e., “given u_φ , x_φ is strictly preferred to x'_φ , all else being equal, but irrespective of the values of variables in W_φ ”. With reference to the ontologies introduced in Section 8.4, this means that for this preliminary step of preference definition, the interested user only has to deal with classes and properties of the *lite* ontology of Figure 8.2. In particular, after the selection of the domain of interest, the user inserts her preferences as depicted in Figure 8.4. The interface manages both instances of variables `cpt:bound true` and `cpt:bound false`, as introduced in Section 8.4. In the former case, the variable which the preference is “about” and the couple of values separated by the word “over” must be introduced; in the latter case, the user has to specify whether the presence or the absence of a variable is preferred. Optionally, she can insert a “Condition” under which the above order holds and make explicit, in the “Irrespective” section, the set of variables for the (conditional) relative importance. She can insert as much preferences as she wants with the “add Another Preference” button or complete the insertion procedure with the “Insert Preferences” button.

When this second button is pressed, the tool takes care that the user has defined the transitive closure of those preferential statements related to multiple values of a variable. More specifically, if $\varphi_1 = u_{\varphi_1} : x > \hat{x} [W_{\varphi_1}]$ and $\varphi_2 = u_{\varphi_2} : \hat{x} > \bar{x} [W_{\varphi_2}]$ have been inserted, with $x, \hat{x}, \bar{x} \in \text{dom}(X)$, and u_{φ_1} and u_{φ_2} do not contain two different values of the same variable, then the rule $\varphi_3 = u_{\varphi_3} : x > \bar{x} [W_{\varphi_3}]$ is added, if missing

(where u_{φ_3} is the condition joining u_{φ_1} and u_{φ_2} , and W_{φ_3} is the intersection of sets W_{φ_1} and W_{φ_2}). As an example in the movie domain, one may state that (φ_1) the actor Hugh Grant is preferred over Colin Firth for comedy films irrespective of the country of production and that (φ_2) Colin Firth is preferred over Joaquin Phoenix for movies directed by Woody Allen. In this case, the additional fact to add would be that (φ_3) Hugh Grant is preferred over Joaquin Phoenix for comedy movies directed by Woody Allen *ceteris paribus* and with no relative importance specification, since the intersection of sets W_{φ_1} and W_{φ_2} is empty. The tool then exploits the *lite* version of Γ to generate its *full* version by managing the transitive closure of `moreImportantThan` and `conditionallyMoreImportantThan`, as described in Section 8.4⁹.

The tool also helps the user to understand if her CP-theory is *cuc*-acyclic or not. It returns an error message to the user if the CP-theory is not locally consistent or the directed graph $J_{u_\varphi}(\Gamma)$ on V , for any u_φ introduced by the user in her preferential statements, is cyclic. Otherwise, it returns the encoding that can be used to query the DBpedia dataset.

As an alternative to the manual insertion of preferences, the user can decide to upload a file of preferences written according to the *lite* ontology, using the specific top right button. The transitivity and *cuc*-acyclicity checking and the introduction of the additional class and properties of the *full* ontology is performed in this case as well.

At this point, a file with the *full* version of Γ is available and can be visualized by pushing the button in Figure 8.5 (a) or exploited directly to formulate a query against DBpedia through the button in Figure 8.5 (b).

More specifically, when the button of Figure 8.5 (b) is pressed, an interface as the one depicted in Figure 8.6 appears. The interface mimics the query presented at the beginning of Section 8.4.2. There, the users may insert their own requirements $\mathcal{R}(\text{?item})$ and specify the number k of results to use in the `LIMIT` modifier, which

⁹The rules that allow the system to manage the definition and the transitive closure of both properties `moreImportantThan` and `conditionallyMoreImportantThan` have been implemented in Prolog and are available at <http://sisinflab.poliba.it/semanticweb/lod/ontologies/rules.pl>.



Figure 8.5: The buttons to display the Full RDF File (a) or to formulate a preference-based query (b).

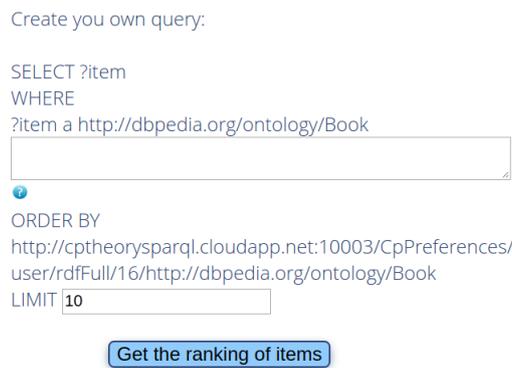


Figure 8.6: The interface to build the query.

by default is set to 10. The URL displayed after the `ORDER BY` clause represents the location of the RDF file containing the *full* version of Γ .

The query returns the top- k list of items belonging to the domain of interest (e.g., books as in Figure 8.6), satisfying $\mathcal{R}(?item)$ and ordered according to user's preferences.

Example 12 (Book cont'd)

Introducing Giorgio's preferences, contained in the CP-theory of Table 2.1, in the proposed tool would produce the top-5 list of results shown in Table 8.1. The results refer to the release of DBpedia 2015-04¹⁰ (also known as: 2015 A).

By looking in DBpedia, one can observe an exact matching with the expected order shown in Example 4, according to the following triples:

```
@prefix db: <http://dbpedia.org/resource/>
@prefix dbo: <http://dbpedia.org/ontology/>

db:An_Uncertain_Place dbo:country db:France ;
dbo:literayGenre db:Crime_fiction .
db:Requiem_for_a_Fish dbo:country db:France ;
```

¹⁰<http://wiki.dbpedia.org/dbpedia-data-set-2015-04>

```

dbo:literayGenre db:Crime_fiction .
db:Blood_Red_Rivers dbo:country db:France ;
dbo:literayGenre db:Crime_fiction .
db:Tropic_of_Capricorn_(novel) dbo:country db:France ;
dbo:literayGenre
db:Autobiographical_novel .
db:Have_Mercy_on_Us_All dbo:country db:France ;
dbo:literayGenre db:Crime_fiction ;
dbo:subsequentWork
db:Wash_This_Blood_Clean_from_My_Hand.

```

■

8.8 Experiments

In order to assess the effectiveness of the presented approach and the implemented tool, we set up two different experiments.

The first experiment consisted of 20 real users using the tool to express their preferences. After the test, users were asked to fill up a questionnaire (reported in Table 8.4). The dataset adopted consisted of a subset of DBpedia 2015-04 related to the four popular domains of: Movies, Food, Music and Books. The statistics of the dataset used for experiments are detailed in Table 8.2.

It is worth noticing that CP-statements can also be automatically extracted from users data [244, 249, 226, 243]. Nevertheless, we set up the previous experiment to have a hint on the average number of CP-statements φ needed to model a user profile as well as on what is, from a user perspective, the most tricky version of φ to represent among:

- $\top : x_\varphi > \hat{x}_\varphi[\emptyset]$,
- $u_\varphi : x_\varphi > \hat{x}_\varphi[\emptyset]$,
- $u_\varphi : x_\varphi > \hat{x}_\varphi[W_\varphi]$.

The second experiment consisted of simulating 168 users using the platform and expressing an overall number of 6720 preferences and 6720 queries to retrieve the resources ranked by taking into their preferences. The aim of this experiment was that of evaluating the response time of the overall system in retrieving a list of resources based on a set of user preferences.

8.8.1 Test on Real Users

In order to test the capability of a user to exploit the platform and even to test if human users unaware of CP-theories were able to express their preferences, we selected 20 users that did not know anything about CP-theories and, after a 5 minutes tutorial, we asked them to express their preferences by using our tool. We asked them to insert as many preferences as they wanted for each domain on the platform, and we then asked them to fill up a post-experience questionnaire in order to acquire some feedback about the experience. The motivation of this experiment is twofold: the first information that we wanted to collect was the number of preferences that a user is prone to explicitly express. The result for this evaluation is shown in Table 8.3. The users provided an overall number of 322 preferences. The average numbers per user are quite similar among the different domains (between 4 and 6) with a little higher propensity to express preferences over books w.r.t. songs. The similar average values, and the similar standard deviations, suggest that there exists a commonality in the number of expressed preferences over a specific domain.

The second relevant information that we wanted to collect is how much the CP-theories expressiveness may fit a “natural” way of expressing preferences by a human being. To this aim, we submitted a small questionnaire with 10 questions whose relative answers in aggregate form are shown in Table 8.4 and Fig. 8.7. All the questions but Q.2 needed to express a value in a 5-star rating scale, with 1 being the worst answer and 5 the best one.

Users felt that representing preferences was not a trivial operation (3.2 corresponds to the lowest value of the overall questionnaire), but this perceived difficulty is clearly dependent on the type of preference (it is worth to notice that for every specific kind of preference, the score is higher than the overall score).

Thanks to the survey, we can list in an increasing order of difficulty the different kinds of preferences:

- “*About a property, I prefer a Value over another Value*” corresponding to $\top : x_\varphi > \hat{x}_\varphi[\emptyset]$.
- “*Given a condition, I prefer a Value over another Value*” corresponding to

$$u_\varphi : x_\varphi > \hat{x}_\varphi[\emptyset].$$

- “About a property, I prefer a Value over another Value irrespectively to a property” corresponding to $u_\varphi : x_\varphi > \hat{x}_\varphi[W_\varphi]$.

Moreover, if we look at the pie chart in Fig. 8.7, it emerges that the most difficult part of the process was to detect the properties (variables V) on which the preferences should be expressed. Another information that we wanted to collect was if the possible difficulty in expressing preferences is stable or it progressively vanishes as the number of expressed preferences increases. Questions 7, 8 and 9 show that the first preference was quite hard to express, but, as the experience goes on, it becomes much easier reaching an average value of 4.1.

The last relevant information that we wanted to collect is how much the expressiveness of CP-theories can correspond to a perceived “natural” way to express preferences. Even here, the result is interesting, because CP-theories are perceived as a quite good way of expressing preferences with a high value of 3.8.

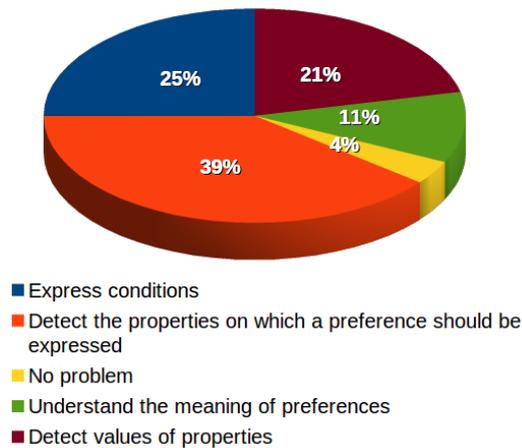


Figure 8.7: Pie chart depicting the results of the second question of the survey.

8.8.2 Test on Simulated Users

In order to closely simulate the behavior of a real user, we designed a tool able to perform the classical operations of expressing a preference and asking the system

for an ordered list of relevant resources. For each domain of interest, the simulated users randomly extract (with a uniform distribution) a property that they might be interested in, and then randomly select the other components of the preference (e.g., in case of a simple preference, $\top : x > \hat{x}[\emptyset]$, they select either the more liked resource x and the less liked one \hat{x}). The composed preference is then sent to the server to be processed and stored. The system checks if the preference produced a cycle, eventually warning the user (in case of a cycle, a new preference is produced). Once the preference is correctly inserted, the simulated user performs a query to the system to retrieve an ordered list of the 100 most relevant resources. The system continues, inserting a new preference for the same domain, and asking the system for a new list. The process ends when 10 preferences are inserted for each domain and the 10 related queries accomplished. Based on the previous experiment, we considered 10 as a representative number of preferences per user. Fig. 8.8 shows the average execution time for an increasing number of preferences related to the simulated users.¹¹ The SPARQL engine adopted for the experimental evaluation is Jena Fuseki v. 2.3.1 running on a Linux server (kernel v. 4.4.0-28-generic) with an Intel Xeon @ 2.30GHz CPU and 8 GB RAM, while the local version of DBpedia had been loaded in a Virtuoso Server (v. 07.20.3212), running on a Linux server (kernel v. 4.2.0-23-generic) with an Intel Xeon @ 2.40 GHz CPU and 56 GB RAM.

The results show that queries based on a number of preferences lower than six take approximately less than one second to return results to the user. This is even more interesting if we consider results of the previous experiments, where we saw that users tend to express an average number of preferences between 4 and 6.

8.9 Conclusion and future work

In this study, we have investigated how user preferences can be taken into account while querying Linked Open Data datasets. Having realized that the Pareto-

¹¹For those interested in a more fine-grained view of the data, a report of the execution times is publicly available at <https://github.com/sisinflab/CP-theories-SPARQL/blob/master/evaluation/evaluationResults.tsv>

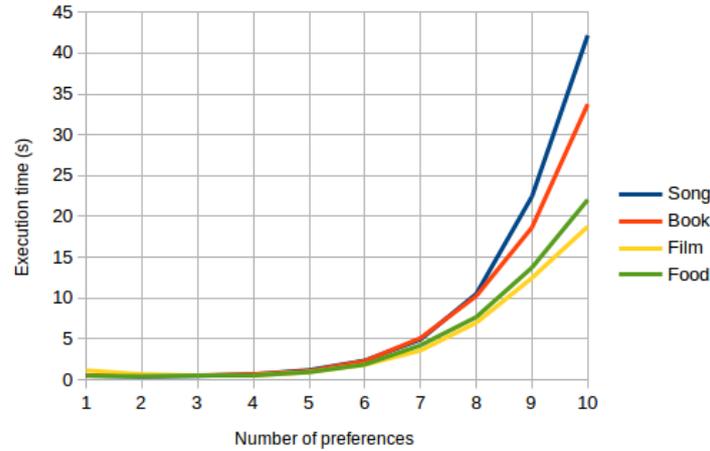


Figure 8.8: Average execution time for increasing number of preferences (1 to 10) in the four domains of: Song, Book, Film and Food.

optimal set identification is not enough, we moved beyond it, proposing an approach to retrieve a ranked list of semantic resources, ordered according to a user’s soft constraints. We focused on qualitative preferences, which are closer to how a user makes decisions, especially in a multi-attribute context, and integrated the partial order implied by a qualitative approach into a top- k scenario, that is, returning to the user, who is formulating qualitative preferential statements, a ranked list of resources, optionally limited in its size, ordered according to her preferences. Among qualitative approaches to preference reasoning, we relied on CP-theories, a general and well-known formalism based on the *ceteris paribus* semantics. We proposed an ontological vocabulary to model CP-theories by means of `RDF` statements under the *ceteris paribus* semantics. Then, we presented an algorithm able to build a standard `SPARQL` 1.1 query encoding the CP-theory and able to retrieve a ranked set of resources satisfying the corresponding preferential constraints. To our knowledge, this is the first attempt to encode the semantics of a CP-theory into a `SPARQL` query and, along with [320], the first approach that lets `SPARQL` to retrieve a ranked list of resources ordered according to a user’s preferences.

We intend the proposed approach as a starting point for many future directions to reason with preferences in a pure Linked Data setting. More efficient encodings

for the proposed queries could be investigated, able to mitigate some performance problems, related, for example, to the increase in the size of variables domains. Moreover, in its current version, the algorithm first computes the complete partial order of the outcomes, and then it matches them with items satisfying the users' hard requirements. As the computation of the partial order is computationally expensive, an improvement could surely be to compute and order only those outcomes matching the users' requirements, thus reducing the number of comparisons needed to return query results to the user. As a future direction of our research, we are also working on approaches proposed for automated CP-nets and CP-theories elicitation [121, 161]. Another interesting topic for future research is to explore how our present work can be extended by integrating approaches to preference-based query answering over graph databases, such as the ones in [149, 159, 146], as well as how to deal with variables having continuous domains.

```

1: procedure GENERATEORDERINGQUERYFORCP-THEORIES( $\mathcal{G}_{user}$ )
2:   forall  $X_i \in V$  do
|   Outcome_D_values += ? $X_i$ _D
|   end
3:   OrderingQuery = SELECT (URI(?outcome_D) AS ?URIOutcome)
|   Outcome_D_values (COUNT(DISTINCT ?outcome_d) AS ?counter)
|   WHERE{ ;
4:   OrderingQuery += {SELECT DISTINCT ?outcome_D ?outcome_d
|   Outcome_D_values (count(DISTINCT ?V) AS ?counterVundominated)
|   (sum((?counterBind)) AS ?counterV) WHERE{ ;
5:   OrderingQuery += {SELECT DISTINCT ?outcome_D ?outcome_d ?V
|   Outcome_D_values WHERE{ ;
|   for  $y \in \{D, d\}$  do
|   forall  $X_i \in V$  do
|   §: OrderingQuery += VALUES (? $X_i$ -y) { (value( $x_{i1}$ )) (value( $x_{i2}$ )) } ;
|   end
7:   OrderingQuery += BIND (CONCAT (STR( ; for  $i = 1, \dots, |W| - 1$  do
|   §: OrderingQuery += attribute( $X_i$ ), STR (? $X_i$ -y), ;
|   end
9:   OrderingQuery += STR (attribute( $X_{|W|}$ ), STR (? $X_{|W|}$ -y)) AS
|   ?outcome_y). ;
|   end
10:  OrderingQuery += FILTER (?outcome_D != ?outcome_d). ;
11:  OrderingQuery += ?V a cpt:Variable.
|   ?V cpt:variableDomain ?variable1.
|   ?variable1 cpt:value ?value1.
|   ?V cpt:variableDomain ?variable2.
|   ?variable2 cpt:value ?value2.
|   FILTER (!(?value1=?value2)) &&
|   contains(?outcome_D, str(?value1)) &&
|   contains(?outcome_d, str(?value2)). ;
12:  OrderingQuery += FILTER NOT EXISTS{
|   ?V2.
|   ?V2 cpt:variableDomain ?vd1.
|   ?vd1 cpt:value ?v1.
|   ?V2 cpt:variableDomain ?vd2.
|   ?vd2 cpt:value ?v2.
|   FILTER((!(?v1=?v2)) && contains(?outcome_D, str(?v1)) &&
|   contains(?outcome_d, str(?v2))). } ;

```

Algorithm 2: Algorithm

```

13:   OrderingQuery += FILTER NOT EXISTS{ ?V3
cpt:conditionallyMoreImportantThan
?instanceOfRelativeImportance.
    ?instanceOfRelativeImportance cpt:hasCondition ?C.
    ?instanceOfRelativeImportance cpt:hasLessImportantVariable
?V.
    ?V3 cpt:variableDomain ?vd13.  ?vd13 cpt:value ?v13.
    ?V3 cpt:variableDomain ?vd23.  ?vd23 cpt:value ?v23.
    {Select distinct ?C
(GROUP_CONCAT(CONCAT(str(?attr),str(?value)); separator ="")) as
?Concatenated)
    where{ ?C cpt:contains ?c.
    ?c cpt:attribute ?attr; cpt:value ?value.  }
    GROUP BY ?C }
FILTER(CONTAINS(?outcome_D,?Concatenated)).
FILTER(!(?v13=?v23)&& contains(?outcome_D,str(?v13)) &&
contains(?outcome_d,str(?v23))).  } } ;
14:   OrderingQuery += {SELECT DISTINCT ?V ?ConcatenatedParent
?Prefer ?Over WHERE{ {Query 1} UNION {Query 2} } } ;
15:   OrderingQuery += BIND( IF( (( !BOUND(?ConcatenatedParent) &&
contains(?outcome_D,?Prefer)&& contains(?outcome_d,?Over))
|| ( BOUND(?ConcatenatedParent) && ?ConcatenatedParent!="")
&& contains(?outcome_D,?ConcatenatedParent) &&
contains(?outcome_D,?Prefer) && contains(?outcome_d,?Over)))
,1,0) as ?counterBind );
16:   OrderingQuery += } GROUP BY ?outcome_D Outcome_D_values
?outcome_d } ;
17:   OrderingQuery += FILTER(?counterV=?counterVundominated)} ;
18:   OrderingQuery += GROUP BY ?outcome_D Outcome_D_values ;
19:   OrderingQuery += ORDER BY DESC (?counter) ;
20:   return OrderingQuery

```

| ?item_D | ?score |
|--------------------------------|--------|
| db:An_Uncertain_Place | 15 |
| db:Requiem_for_a_Fish | 15 |
| db:Blood_Red_Rivers | 15 |
| db:Tropic_of_Capricorn_(novel) | 13 |
| db:Have_Mercy_on_Us_All | 9 |

Table 8.1: The top-5 list of items retrieved for preferences in the CP-theory of Table 2.1.

| Classes | Instances | Properties | 1 hop resources |
|--------------|---------------|------------|-----------------|
| dbo:Book | 31172 | 36 | 12964 |
| dbo:Film | 90063 | 31 | 82922 |
| dbo:Food | 6003 | 21 | 2367 |
| dbo:Song | 7195 | 27 | 2220 |
| Total | 134433 | 115 | 100473 |

Table 8.2: Dataset Statistics

| | Total | Min | Max | Mean | Std Dev |
|-----------------|-------|-----|-----|--------|---------|
| dbo:Book | 109 | 1 | 11 | 5.7368 | 2.1562 |
| dbo:Film | 78 | 1 | 17 | 4.5882 | 3.8900 |
| dbo:Food | 75 | 1 | 12 | 4.1667 | 2.6844 |
| dbo:Song | 60 | 1 | 12 | 3.5294 | 2.6485 |

Table 8.3: User Experiments Statistics

| Q.N. | Questions | Min | Max | Mean | St Dev |
|-------------|--|--------------|------------|-------------|---------------|
| <i>Q.1</i> | <i>How easy has been to represent your preferences?</i> | 1 | 5 | 3.1667 | 1.1100 |
| <i>Q.2</i> | <i>Which among these did you consider the hardest?</i> | See Fig. 8.7 | | | |
| <i>Q.3</i> | <i>How easy is to represent a preference like "About a property I prefer a Value over another Value"?</i> | 1 | 5 | 4.1667 | 1.3048 |
| <i>Q.4</i> | <i>How easy is to represent a preference like "I prefer a resource that has a certain property"?</i> | 2 | 5 | 3.5556 | 1.0966 |
| <i>Q.5</i> | <i>How easy is to represent a preference like "Given a condition I prefer a Value over another Value"?</i> | 2 | 5 | 3.8889 | 0.9852 |
| <i>Q.6</i> | <i>How easy is to represent a preference like "About a property I prefer a Value over another Value irrespectively to a property"?</i> | 2 | 5 | 3.6111 | 1.2005 |
| <i>Q.7</i> | <i>How easy was to represent the first preference?</i> | 1 | 5 | 3.2222 | 1.3492 |
| <i>Q.8</i> | <i>After the first preference how easy was to represent the next two ones?</i> | 2 | 5 | 3.9444 | 0.9852 |
| <i>Q.9</i> | <i>After the first three preferences how easy was to represent the next ones?</i> | 1 | 5 | 4.1111 | 1.2783 |
| <i>Q.10</i> | <i>How much this way of expressing preferences is similar to your own?</i> | 1 | 5 | 3.7778 | 1.1448 |

Table 8.4: User Survey Statistics

Chapter 9

Hybrid Relevance: How to enhance traditional relevance weighting schemes

9.1 Introduction

In the last years, many recommendation approaches have been proposed that take advantage of side information to enhance the performance of latent factor models. Side information can refer to items as well as users [390] and can either be structured [358] or semi-structured [419, 40, 97]. Interestingly, in [418] the authors argue about a new generation of knowledge-aware recommendation engines able to exploit information encoded in *knowledge graphs (KG)* to produce meaningful recommendations: “*For example, with knowledge graph about movies, actors, and directors, the system can explain to the user a movie is recommended because he has watched many movies starred by an actor*”. Hence, the use of side information could not only be used to improve the recommendations but also the whole user

experience, by providing explanations or some kind of reasoning from the system.

Nonetheless, in recommendation scenarios where user interactions with the system is of key importance, we cannot deal only on content/context-aware data, but we do need to take into account also collaborative one. The point here is: how to enhance, in a principled way, knowledge about the relevance of each attribute by also encoding collaborative information?

In this research line, we tackle this problem by drawing conceptual and methodological techniques from *Information Retrieval (IR)* and *Recommender Systems (RS)* in our collaborative-aware relevance framework $CRe1-FM$. We start from a generic relevance measure in *IR* (such as *TF-IDF* [35] or *BM25* [201, 202, 319]) to formally represent the importance or informativeness of the attributes extracted from user and item descriptions [35]. Then, we make use of *Bayesian Personalized Ranking (BPR)* optimization procedure [312] to combine the relevance measurements obtained before with the collaborative patterns extracted from the user-item interactions. This augmented representation of items (or users) attributes can be eventually used to compute better item-item (or user-user) similarity values exploited in recommendation engines where the main goal is to present the best possible (most interesting) items to the user as a *top-N* ranking (as opposed to other works where the task is to predict the rating of a user towards an item). We believe that a better understanding of this hybrid representation of attributes relevance will not only improve the quality of recommendations returned to the user, but it will also open up many other possibilities, such as generating better explanations for users or building better suited item and user (latent) representations, by integrating more complex relevance measures or even other optimization techniques tailored to particular recommendation tasks.

More specifically, we address the following main research question: **(RQ1)** In personalized scenarios, whenever we have content-based and collaborative information, how should we exploit them together in a principled way to measure the relevance of an attribute? If a positive answer is found, we then analyze the following: **(RQ2)** May collaborative information improve the quality of attributes' relevance in terms of accuracy, diversity, and novelty of results in a recommendation setting?

The main contributions thus include:

- A principled approach to integrate content and collaborative information by exploiting the relevance of an attribute.
- Extensive experiments on three real-world datasets using state-of-the-art recommendation algorithms.
- Positive results in terms of accuracy, diversity, and novelty for our approach based on hybrid relevance in a majority of the tested scenarios, evidencing its generalization capabilities and potential to capture user and item preferences in different domains.

The remainder of this chapter is structured as follows: in the next section we introduce the proposed framework CRe1-FM for collaborative-aware computation of attributes relevance. Then, in Section 9.3 we describe the experimental setting to prove the effectiveness of CRe1-FM and discuss the obtained results. We continue with a description of the main related works in Section 9.4. Conclusion and future work close the chapter.

9.2 CRe1-FM : Collaborative-aware relevance for recommendation

We present here CRe1-FM , a collaborative-aware relevance framework particularly well-suited for recommendation scenarios. It mainly bases on *Factorization Machines (FM)* and exploits content-based relevance measures that can be plugged-in the overall framework.

In the following, we will consider a typical recommendation setting where we have two separate sets of users $u \in U$ and items $i \in I$ such that $U \cap I = \emptyset$ and the corresponding user-item interaction matrix containing all the implicit ratings (i.e., interactions) users have given to items. In case we have access to some form of items description, such as keywords, tags, attributes, etc., as well as users' one, for recommendation purposes, we usually exploit them to compute similarity values

through a relevance measure ρ . For instance, in case items are endowed with a textual representation we may use $\rho = TF-IDF$ and associate a relevance value ω to each keyword in the text. Each item i has then a vector-based representation

$$i^{TF-IDF} = (\omega_0^{TF-IDF}, \dots, \omega_{|A|}^{TF-IDF})$$

with A (for attributes) being the set of all possible keywords. In an analogous way, for the computation of u^{TF-IDF} we may consider all the items enjoyed by u and, even here, compute the ω values associated to each keyword.

This can be further generalized. Indeed, given a relevance measure ρ and a set of attributes A , we can always represent an item as

$$i^\rho = (\omega_0^\rho, \dots, \omega_{|A|}^\rho)$$

The same can be done with u^ρ . In case we want to compute an attribute-based similarity value between items or users, both u^ρ and i^ρ are then the perfect candidates.

In recommendation scenarios, the adoption of factorization models has shown its effectiveness since its initial introduction [310]. In fact, thanks to their subtle modeling of user-item interactions, such models are very precise and effective even in very sparse settings. *FM* have been proposed as a unifying framework to represent all the different factorization models in a general and theoretically sound framework. Without loss of generality, in the following we describe our proposal of *CREL-FM* by looking at *FM* of order 2 for a recommendation problem involving only implicit ratings. The model can be easily extended to a more expressive representation by taking into account, e.g., demographic and social information [18], multi-criteria [15], and even relations between contexts [424].

For each pair of user u and item i we build the binary vector $\mathbf{x}^{ui} \in \{0, 1\}^{1 \times (|U|+|I|)}$ representing the interaction between u and i in the original user-item rating matrix. In this modeling, \mathbf{x}^{ui} contains only two 1 values corresponding to u and i while all the other values are set to 0. Based on all possible vectors \mathbf{x}^{ui} we then build the matrix $\mathbf{X} \in \{0, 1\}^{(|U| \times |I|) \times (|U|+|I|)}$ containing as rows all possible \mathbf{x}^{ui} we can build starting from the original U - I matrix as shown in Fig. 9.1.

| | | | | | | | | | | | |
|-------|-------|-------|-------|-------|-----|-------|-------|-------|-------|-------|-----|
| x^1 | 1 | 0 | 0 | 0 | ... | 1 | 0 | 0 | 0 | 0 | ... |
| x^2 | 1 | 0 | 0 | 0 | ... | 0 | 1 | 0 | 0 | 0 | ... |
| x^3 | 1 | 0 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| x^4 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| x^5 | 0 | 1 | 0 | 0 | ... | 0 | 0 | 0 | 1 | 0 | ... |
| x^6 | 0 | 0 | 1 | 0 | ... | 1 | 0 | 0 | 0 | 0 | ... |
| x^7 | 0 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 0 | 0 | ... |
| | u_1 | u_2 | u_3 | u_4 | ... | l_1 | l_2 | l_3 | l_4 | l_5 | ... |
| | User | | | | | Item | | | | | |

Figure 9.1: A visual representation of the interaction matrix \mathbf{X} .

$$\hat{y}(\mathbf{x}^{\mathbf{ui}}) = w_0 + \sum_{j=1}^{|U|+|I|} w_j \cdot x_j + \sum_{j=1}^{|U|+|I|} \sum_{p=j+1}^{|U|+|I|} x_j \cdot x_p \cdot \sum_{f=1}^k v_{(j,f)} \cdot v_{(p,f)} \quad (9.1)$$

The *FM* score $\hat{y}(\mathbf{x}^{\mathbf{ui}})$ for each vector $\mathbf{x}^{\mathbf{ui}}$ is defined as in Equation (9.1), where the parameters to be learned are, respectively: w_0 representing the global bias; w_j giving the importance to every single x_j ; and the pair $v_{(j,f)}$ and $v_{(p,f)}$ in $\sum_{f=1}^k v_{(j,f)} \cdot v_{(p,f)}$ measuring the strength of the interaction between each pair of variables: x_j and x_p . The summation $\sum_{f=1}^k v_{(j,f)} \cdot v_{(p,f)}$ represents the dot product between two vectors: \mathbf{v}_j and \mathbf{v}_p with a size equal to k . Hence, \mathbf{v}_j represents a latent representation of a user, \mathbf{v}_p that of an item within the same latent space, and their interaction is evaluated through their dot product¹. The number of latent factors is represented by the hyper-parameter k whose value is usually selected at design time. In order to compactly represent all the vectors \mathbf{v}_j and \mathbf{v}_p we introduce the matrix $V \in \mathbb{R}^{(|U|+|I|) \times k}$ having \mathbf{v}_j in the first $|U|$ rows and \mathbf{v}_p in the last $|I|$ ones. For training purposes, V moves from an initial value V_0 to the final \hat{V} where the values in V_0 can be initialized following different strategies [180]. After the training, the first $|U|$ rows of \hat{V} represent a latent representation of each user u with respect to the k latent features while the last $|I|$ ones are the latent representation of all the items i in the catalog. Such vectors can be eventually used to compute the similarity between items or users which is usually exploited in recommendation algorithms. We may say that each value in \hat{V} represents a collaborative-based relevance value of the corresponding latent attributes of u (first $|U|$ rows) or i (last $|I|$ rows). If we look at the mathematical

¹In our case, we have $v_{(j,f)} \cdot v_{(p,f)} \neq 0$ only when $x_j = 1$ and $x_p = 1$ at the same time. Hence, the final value represents the strength in the interaction between the corresponding u and i .

formulation of FM , we see it is able to encode only collaborative information thus being completely agnostic to the nature of catalog items the user interacts with. As a consequence, the same holds for the final similarity values between users or items.

The main idea behind $CRel-FM$ is to combine the attribute-based relevance representation of u^ρ and i^ρ with the collaborative one computed by FM in a principled and effective way. As a first step, given a relevance measure ρ and a set of attributes A , for each $i \in I$ and $u \in U$ we compute the corresponding vectors i^ρ and u^ρ . Then, we build the matrix $W \in \mathbb{R}^{(|U|+|I|) \times |A|}$ where the first $|U|$ rows correspond to u^ρ vectors and the last $|I|$ rows to i^ρ vectors. Eventually, we set $k = |W|$ in Equation (9.1) and initialize the corresponding matrix V with $V_0 = W$. In other words, we impose the number of latent factors equal to the number of all the attributes A and we then inject the relevance values ω for both u and i within the factorization model. As a consequence, at the end of model training we have the matrix $\hat{V} \in \mathbb{R}^{(|U|+|I|) \times |A|}$ containing a representation of the original content-based values of u^ρ and i^ρ enhanced with collaborative information coming from the FM modeling.

Our intuition is that these collaborative-aware relevance values are by far more representative of users and items with respect to the pure collaborative latent representation in the original formulation of FM or the pure content-based representation contained in u^ρ and i^ρ . As an example, in Table 9.1 we refer to some attributes extracted from the knowledge graph `DBpedia`² and show an example of values obtained after the training (in the column $CRel-FM$) together with the original $TF-IDF$ ones [27, ?] computed for a movie from the `Yahoo!Movies`³ dataset. The attributes considered here are the categories associated to a movie coming from the corresponding `Wikipedia` ones. It is interesting to compare the attributes ranking coming from the values of $TF-IDF$ and the one coming from $CRel-FM$ and to see how this latter looks more meaningful than former one.

²<http://dbpedia.org>

³http://research.yahoo.com/Academic_Relations

| CRel-FM | TF-IDF | Attribute | CRel-FM | TF-IDF | Attribute |
|---------|--------|---------------------------------------|---------|--------|---------------------------------------|
| 1.3669 | 0.2584 | Space_adventure_films | 1.2434 | 0.2858 | Space_adventure_films |
| 1.1252 | 0.2730 | Films_set_in_the_future | 1.0355 | 0.3020 | Films_set_in_the_future |
| 0.9133 | 0.2355 | American_science_fiction_action_films | 0.8956 | 0.2605 | American_science_fiction_action_films |
| 0.8485 | 0.3190 | 1980s_science_fiction_films | 0.8951 | 0.3451 | Android_robot_films |
| 0.6529 | 0.1549 | Paramount_Pictures_films | 0.7338 | 0.3105 | Time_travel_films |
| 0.5989 | 0.3468 | Midlife_crisis_films | 0.6665 | 0.2701 | Film_scores_by_Jerry_Goldsmith |
| 0.5940 | 0.1797 | American_sequel_films | 0.6581 | 0.2205 | 1990s_action_films |
| 0.5862 | 0.2661 | Film_scores_by_James_Horner | 0.6561 | 0.2279 | 1990s_science_fiction_films |
| 0.5634 | 0.2502 | Films_shot_in_San_Francisco | 0.6118 | 0.1988 | American_sequel_films |
| 0.5583 | 0.1999 | 1980s_action_thriller_films | 0.5649 | 0.1713 | Paramount_Pictures_films |

Table 9.1: Top-10 features computed by CRel-FM for the movies "Star Trek II - The Wrath of Khan" and "Star Trek - First Contact".

9.3 Experimental Evaluation

In this section, we aim at assessing if there is empirical evidence on the usefulness of adopting a hybrid relevance measure to feed recommender systems.

Datasets. To provide an answer to research questions posed in Section 9.1, we have evaluated the performance of our method on three well-known datasets for recommender systems belonging to different domains. The `Last.fm` dataset [84] corresponds to user-artist plays on `Last.fm` online music system released during *HETRec 2011⁴ Workshop*. It contains social networking, tagging, and music artists listening information from a set of *2K* users. `LibraryThing` represents books' ratings collected in the `LibraryThing` website⁵ community. It contains social networking, tagging and rating information on a [1..10] scale. `Yahoo!Movies` (Yahoo! Webscope dataset `ydata-ymovies-user-movie-ratings-content-v1_0`)⁶ contains movies ratings generated on `Yahoo! Movies` up to November 2003. It provides content, demographic and ratings information on a [1..5] scale, and mappings to `MovieLens` and `EachMovie` datasets. **Feature Selection.** In order to get attributes related to items in the datasets we exploited the freely available mapping⁷ that links each item to an entity in the `DBpedia` knowledge graph. Our

⁴<http://ir.ii.uam.es/hetrec2011/>

⁵<https://www.librarything.com/>

⁶http://research.yahoo.com/Academic_Relations

⁷<https://github.com/sisinflab/LinkedDatasets>

assumption is that the usage of such well-curated features does not introduce any informative bias (both positive or negative) in the values computed by CReL-FM and the eventual recommendations. Following [277], and [299] we filtered out some irrelevant features with a unique threshold for missing values (corresponding to tm [277], and p [299]). Datasets statistics are shown in Table 9.2.

| Dataset | Threshold | #Users | #Items | #Transactions | #Features | Sparsity |
|---------------|-----------|--------|--------|---------------|-----------|----------|
| Last FM | 99.86 | 1,375 | 7,312 | 46,982 | 1,315 | 99.53% |
| LibraryThing | 99.91 | 7,221 | 10,605 | 313,069 | 1,169 | 99.59% |
| Yahoo! Movies | 99.60 | 4,000 | 2,528 | 55,711 | 747 | 99.45% |

Table 9.2: Datasets statistics.

Experimental Setting. “All Unrated Items” [356, 49] protocol has been adopted to compare different algorithms. We have split the dataset using Hold-Out 80-20 retaining for every user the 80% of their ratings in the training set and the remaining 20% in the test set. Moreover, a temporal split has been performed [163] whenever timestamps associated to every transaction is available.

As expressed in our research questions, we want to check if the adoption of our hybrid relevance measure CReL-FM is beneficial or not in a recommendation scenario. For this reason each algorithm has been fed using: **i)** the original information contained in R ; **ii)** the relevance values as computed adopting a relevance measure; **iii)** the new relevance information as defined in Section 9.2. As for point **ii)** we have used $TF-IDF$ since we wanted to start with the simplest relevance measure possible.

Algorithms. We evaluate three different families of algorithms that could be fed with relevance information: two *Neighborhood-based* algorithms (ItemKNN [333, 334] and UserKNN [73]), *Factorization Machines* as a representative of latent factors models (BPR-FM [312, 310]), and *Vector Space Model* as a representative of content-based recommender systems (VSM [279]). Additionally, we have compared against two non-personalized baselines, i.e., Random and MostPopular . As for this latter, it is acknowledged that popularity ranking typically shows very good performance because of statistical biases in the data [49] and it is an important baseline to compare against [108].

Table 9.3: Results for Yahoo!Movies.

| Recommender | Source | P | R | nDCG | EPC | Gini | SE | IC | UC |
|-------------|-----------|----------------|----------------|----------------|----------------|----------------|---------------|----------------|-------|
| Random | R | 0.001 | 0.004 | 0.003 | 0.001 | † 0.826 | † 11.3 | † 2,528 | 4,000 |
| MostPopular | R | 0.015 | 0.037 | 0.027 | 0.015 | 0.004 | 3.9 | 48 | 4,000 |
| ItemkNN | R | 0.040 | 0.169 | 0.109 | 0.042 | 0.172 | 9.2 | 1,808 | 3,998 |
| | V_0 | 0.043 | 0.151 | 0.119 | 0.054 | 0.353 | 10.0 | 2,456 | 4,000 |
| | \hat{V} | † 0.055 | 0.197 | † 0.151 | † 0.067 | 0.300 | 9.7 | 2,405 | 4,000 |
| UserkNN | R | 0.031 | 0.134 | 0.086 | 0.033 | 0.041 | 6.8 | 736 | 3,998 |
| | V_0 | 0.030 | 0.140 | 0.092 | 0.033 | 0.055 | 7.4 | 809 | 4,000 |
| | \hat{V} | 0.037 | 0.149 | 0.098 | 0.041 | 0.083 | 7.9 | 1,062 | 4,000 |
| FM | R | 0.027 | 0.076 | 0.050 | 0.026 | 0.008 | 4.7 | 276 | 4,000 |
| | V_0 | 0.039 | 0.140 | 0.113 | 0.050 | 0.289 | 9.8 | 2,320 | 4,000 |
| | \hat{V} | 0.041 | 0.142 | 0.088 | 0.040 | 0.021 | 5.6 | 594 | 4,000 |
| VSM | V_0 | 0.039 | 0.140 | 0.113 | 0.050 | 0.289 | 9.8 | 2,320 | 4,000 |
| | \hat{V} | 0.052 | † 0.198 | 0.149 | 0.063 | 0.297 | 9.8 | 2,357 | 4,000 |

For all the considered recommendation engines we have performed a grid search to tune the hyperparameters. We considered the range of values as suggested by the original authors or by varying the parameters values around the ones showed in the original papers as the best performing ones.

Metrics. In order to evaluate the algorithms, we have measured accuracy through *Precision@N* ($P@N$), *Recall@N* ($R@N$) and *normalized Discounted Cumulative Gain* ($nDCG@N$). *EPC* (*Expected Popularity Complement*) [91] is used to measure novelty, or more precisely the ability of a system to recommend relevant long-tail items. Finally, diversity has been measured through *Item Coverage* (aggregate diversity in top- N list, $IC@N$), *Gini Index* ($Gini@N$) and *Shannon entropy* ($SE@N$). To measure the ability of producing recommendation lists for each user, *User Coverage* (UC) is also computed. The evaluation has been performed considering Top-10 [108] recommendations for all the datasets. A *Threshold-based relevant items* condition [83] of 4/5 has been set for Yahoo!Movies and 8/10 for LibraryThing, and Last.fm respectively in order to take into account only relevant items.

Table 9.4: Results for Last.fm.

| Recommender | Source | P | R | nDCG | EPC | Gini | SE | IC | UC |
|-------------|-----------|----------------|----------------|----------------|----------------|----------------|---------------|----------------|-------|
| Random | R | 0.001 | 0.002 | 0.001 | 0.001 | † 0.547 | † 12.4 | † 6,395 | 1,375 |
| MostPopular | R | 0.025 | 0.064 | 0.035 | 0.023 | 0.001 | 3.7 | 35 | 1,375 |
| ItemkNN | R | 0.031 | 0.069 | 0.042 | 0.032 | 0.154 | 10.5 | 3,235 | 1,375 |
| | V_0 | 0.022 | 0.039 | 0.029 | 0.022 | 0.265 | 11.3 | 4,487 | 1,373 |
| | \hat{V} | 0.037 | 0.070 | 0.051 | 0.041 | 0.216 | 11.0 | 4,044 | 1,370 |
| UserkNN | R | 0.052 | 0.122 | 0.084 | 0.057 | 0.016 | 7.2 | 831 | 1,375 |
| | V_0 | 0.054 | 0.130 | 0.086 | 0.059 | 0.012 | 7.0 | 420 | 1,375 |
| | \hat{V} | 0.060 | 0.136 | 0.093 | 0.066 | 0.015 | 7.4 | 517 | 1,375 |
| FM | R | 0.042 | 0.091 | 0.065 | 0.046 | 0.004 | 5.0 | 313 | 1,375 |
| | V_0 | 0.031 | 0.055 | 0.038 | 0.032 | 0.188 | 10.8 | 3,660 | 1,375 |
| | \hat{V} | † 0.064 | † 0.138 | † 0.099 | † 0.070 | 0.020 | 7.2 | 897 | 1,375 |
| VSM | V_0 | 0.031 | 0.055 | 0.038 | 0.032 | 0.188 | 10.8 | 3,660 | 1,375 |
| | \hat{V} | 0.050 | 0.087 | 0.062 | 0.053 | 0.187 | 10.7 | 3,691 | 1,375 |

9.3.1 Results

Results in Tables 9.3-9.5 show the performance of the different recommender systems fed with different sources. We have marked with **bold** and † the best and the second-best algorithm. The tables correspond respectively to the experiments conducted on Yahoo!Movies, Last.fm, and LibraryThing datasets. To answer RQ2, we analyzed the behavior of the recommender systems on a per dataset basis, focusing on accuracy, diversity, and novelty. While this section is devoted to highlighting the more interesting results, Section 9.3.2 is devoted to drawing some considerations based on the findings depicted here.

Yahoo!Movies. As expected, the *Random* recommender shows the highest dataset values in terms of *Gini index*, *Shannon entropy*, and *Item Coverage*. It is worth noticing that *Item-kNN* almost reaches the same performance. *Item-kNN* also shows the best results, at a dataset level, considering *Precision* and *nDCG*. In this case, the behavior of the recommender varying its knowledge source is very clear. Feeding it with V_0 , and then \hat{V} , we note progressive increases of *Precision*, *Recall*, *nDCG*, and *EPC*. Interestingly, from V_0 to \hat{V} the diversity decreases, whereas the novelty keeps increasing. This is a signal that when we use *R* there is still room for improving all

Table 9.5: Results for LibraryThing.

| Recommender | Source | P | R | nDCG | EPC | Gini | SE | IC | UC |
|-------------|-----------|----------------|----------------|----------------|----------------|----------------|---------------|-----------------|-------|
| Random | R | 0.001 | 0.001 | 0.001 | 0.001 | † 0.774 | † 13.4 | † 11,560 | 7,223 |
| MostPopular | R | 0.006 | 0.006 | 0.006 | 0.005 | 0.001 | 3.8 | 34 | 7,223 |
| ItemkNN | R | † 0.080 | † 0.146 | † 0.124 | † 0.093 | 0.147 | 10.9 | 6,378 | 7,221 |
| | V_0 | 0.047 | 0.077 | 0.059 | 0.051 | 0.217 | 11.4 | 8,009 | 7,221 |
| | \hat{V} | 0.076 | 0.134 | 0.112 | 0.089 | 0.245 | 11.7 | 8,258 | 7,221 |
| UserkNN | R | 0.033 | 0.070 | 0.059 | 0.038 | 0.022 | 7.8 | 1,994 | 7,221 |
| | V_0 | 0.035 | 0.076 | 0.061 | 0.039 | 0.040 | 8.8 | 2,489 | 7,221 |
| | \hat{V} | 0.035 | 0.071 | 0.056 | 0.038 | 0.044 | 9.0 | 2,603 | 7,221 |
| FM | R | 0.024 | 0.038 | 0.031 | 0.023 | 0.007 | 6.2 | 1,122 | 7,221 |
| | V_0 | 0.037 | 0.065 | 0.047 | 0.039 | 0.211 | 11.4 | 7,437 | 7,221 |
| | \hat{V} | 0.044 | 0.077 | 0.065 | 0.047 | 0.033 | 8.0 | 2,687 | 7,221 |
| VSM | V_0 | 0.037 | 0.065 | 0.047 | 0.039 | 0.211 | 11.4 | 7,434 | 7,221 |
| | \hat{V} | 0.052 | 0.086 | 0.072 | 0.060 | 0.265 | 11.9 | 8,012 | 7,221 |

the metrics. Instead, when we use \hat{V} , we increase accuracy exploiting more items that come from the long tail. Regarding accuracy, the behavior of *User-kNN* is similar, with the same progressive increases. On the other side, the different mechanics of the user-based schema are reflected on an increase in diversity when using \hat{V} . *FM* shows different behavior. *Precision* and *Recall* progressively increase changing the knowledge source. *nDCG*, novelty, and diversity show an increase using V_0 , but they decrease when \hat{V} is exploited. A possible reason for that is the low number of mapped items. These items are hence described by a low number of features. After the training, these features are over-trained and they affect the performance. This behavior is not present in *k-NNs* algorithms and *VSM*. The *k-NNs* approaches exploit the Neighborhood to alleviate this effect. On the other hand, *VSM* alleviates it by estimating scores based on vector similarity. Since vector are composed of all the features, this effect is absent.

Last.fm. It is a completely different dataset, that contains a lesser number of users, and a higher number of items. Moreover, the average number of features that describe an item is 6.6, against the 12.1 of *Yahoo!Movies*. These characteristics affect the behavior of the recommenders. *Random* recommender still shows the

highest values for diversity. Differently from the previous experiment, only *Item-kNN* and *VSM* show a good accuracy-diversity trade-off. In *User-kNN* we may notice that are still present the progressive increases for all the accuracy and novelty metrics, but with a vary small margin. *Item-kNN* and *FM* show an interesting commonality. For both recommenders, the exploitation of V_0 produces the negative effect of decreasing accuracy. This may be due to two reasons. First, we have noticed the limited number of features that describe items. Second, maybe those features are not sized for a recommendation task. Nevertheless, after the training, both recommenders show important improvements. Finally, varying the knowledge source in *VSM* leads to remarkable improvements.

LibraryThing. This dataset shows a different behavior. In `LibraryThing`, we curiously note that the best recommender is *Item-kNN* fed with R matrix. This is definitely due to the small number of features per item: 3.8. However, let us focus on what happens when feeding it with \hat{V} . We may notice that accuracy decreases a bit, but the diversity values become much more interesting. Regarding accuracy, *FM* and *VSM* show the same improvements depicted before, with some differences concerning diversity. In *VSM* the item coverage keeps increasing, while *FM* sacrifices diversity in favor of accuracy. The only recommender which shows a confusing behavior is *User-kNN*, in which we may notice an increase of *Precision* and diversity but a decrease of *Recall* and *nDCG*.

9.3.2 Discussion

Once we have detailed the results of the experiments, we can focus on the general findings. If we analyze the behavior of the recommender over the different datasets, we may notice that the exploitation of *TF-IDF* relevance is useful for `Yahoo!Movies` and `LibraryThing`. In those datasets, the extra-knowledge leads to better-tuned recommenders. However, this effect is absent in `Last.fm`, where the *TF-IDF* is not able to overpass the R -based models. Since it happens only in this dataset, in our opinion the reasons for the behavior have to be found in the characteristics of the dataset and a poor description of items. On the other side,

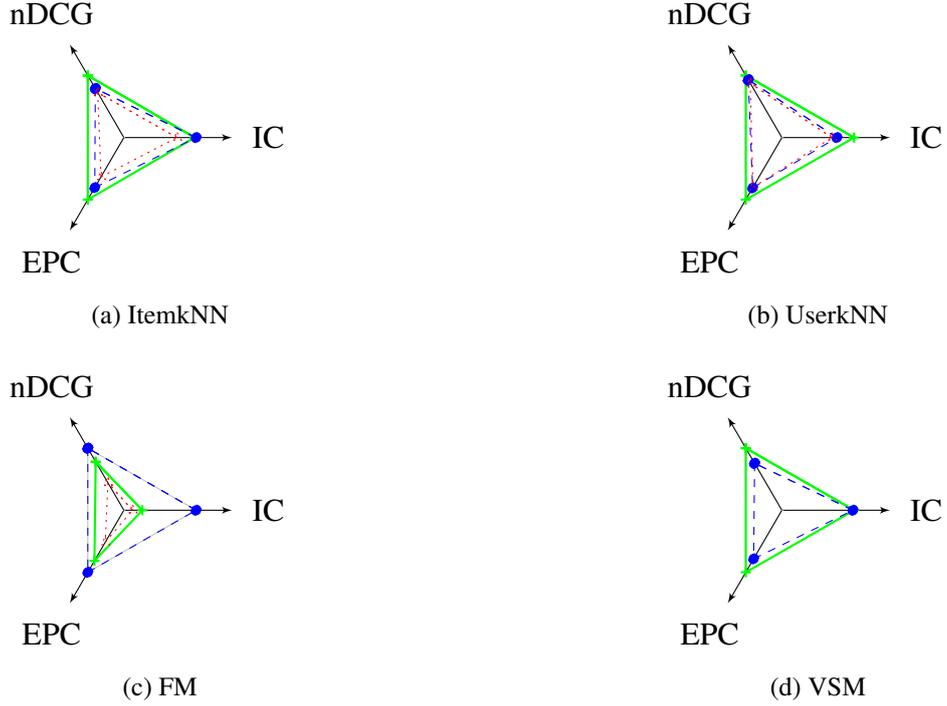


Figure 9.2: Graphical representation of Accuracy, Diversity, and Novelty results for Yahoo!Movies. The different knowledge sources, R , V , and \hat{V} are represented respectively through dotted red line, dashed blue line, and solid green line.

the proposed hybrid measure for relevance in almost all cases leads to improvements in accuracy, novelty, and diversity. The training of the features effectively modify the original relevance weights also taking into account a signal of the popularity of the features. Moreover, another trend deserves attention. *Item-kNN* and *FM* show higher values for *Item Coverage* with V_0 . This value typically decreases with \hat{V} , while accuracy increases. Irrespectively to possible accuracy increases, we observe that *TF-IDF* introduces a noise we consider responsible for those values. After the training phase, we observe a reduction of the recommended items because the noise is partially removed. Some other considerations concern the comparison with [30]. The authors have also considered Yahoo!Movies dataset in their analysis, thus we consider mandatory to compare our findings against it. In particular, among the settings they consider, the most similar is Categorical Setting. Their *kaHFM* corresponds in our experiments to *Item-kNN* with \hat{V} , while *Attribute-based*

Item-kNN corresponds to *Item-kNN* with V_0 . In both experiments, we note that *Item-kNN* with \hat{V} is the best performing one, while *Item-kNN*, *Item-kNN* with V_0 and *VSM* show good performance. Since their work is limited to *Precision* and *nDCG*, we can not compare the other dimensions of the analysis. Moreover, since we have considered many different applications of hybrid relevance, further comparison would be difficult. However, the results for `Yahoo!Movies` can be analyzed considering a trade-off perspective. In Fig. 9.2, we may see the four considered models represented on three axes. The idea here is to observe the different trade-offs in terms of accuracy (*nDCG*), diversity (*Item Coverage*), and novelty (*EPC*). The models fed with \hat{V} are denoted by a solid green line. The V_0 models are represented with a dashed blue line and R models with a dotted red line. It is clear that for *Item-kNN*, *User-kNN*, and *VSM* there is no trade-off. The \hat{V} model overpasses the other along each dimension. Only *FM* shows different behavior, with \hat{V} shape which is containing R shape, but it is contained by the V_0 model. However, here it is clear that the \hat{V} accuracy values are similar to V_0 , while the *Item Coverage* is particularly low.

9.4 Related Work

In recent years, some works have proposed to measure relevance exploiting jointly the content information and the collaborative information. In [237], the authors face some known issues of Collaborative-Filtering approaches like *sparsity* and *Cold-Start*. They propose a Bayesian generative model called *Collaborative Variational AutoEncoder (CVAE)*. *CVAE* learns implicit relationships between items and users from both content and rating information. In their experiments, they show that *CVAE* outperforms the *Collaborative Topic Regression model (CTR)*, the *Collaborative Deep Learning model (CDL)*, and the *Deep Music* neural model. They point out that the generative nature of the method eases to capture the key points of items and users. In a subsequent work [171], the authors propose another Deep Neural Network architecture to consider jointly side information and rating information. In detail, they propose *CAVAE*, a *Collaborative Additional Variational AutoEncoder*, in which content (tag information) feeds an additional variational

autoencoder. They show that their Bayesian probabilistic model outperforms the *CTR* model, the *CDL* model, and *CVAE*. In *entity2rec* [289], the authors consider a knowledge graph that models *user – item* and *item – item* relations. They exploit it to learn property-specific latent representations. From these representations, the authors compute *user – item* relatedness features to feed two Learning to Rank algorithms (*Adarank* and *LambdaMart*). In [130] the authors combine a *bag-of-words* (*BOW*) representation with a *Convolutional Neural Network* (*CNN*) to create an overall model to retrieve semantically similar questions. Since each original question is a transaction, the collaborative information is naturally injected during the training phase.

An interesting work is [70], in which the authors summarize the contributions of the *Social Information Retrieval* community. Their work is particularly interesting for us because they explicitly exploit collaborative information to improve Information Retrieval techniques. The main difference is that they exploit the *user – user* connections in a social graph which are a completely different signal from *user – item* transactions. In a different scenario, the social tagging task, a specific research line focuses on exploiting personal tagging preferences. Even then, relevance information is affected by user behavior. Relevant work is [335] where the authors compute local interaction networks among users exploiting comments. Instead, in [236], the tagging preferences are represented through tag statistics. More recently, in [241], the authors learn a user-specific embedding space. Collaborative information is also exploited in [329], where the authors propose a framework that takes into account user affinities. Interestingly, in [308], the goal is discovering latent associations between users, images, and tags. Another work that focuses on modifying tag relevance using user information is [238], in which two simple methods are proposed to update relevance: *Borda Count* and *UniformTagger*. The latter, proposed by the authors, exploits a neighbor voting algorithm to combine multiple tag relevance estimations.

Another relevant research line is Relevance Feedback Environment, in which feedback information is exploited to correct relevance values. This line of investigation is not closely related to ours. Nevertheless, we have decided to mention

it because, even here, the content information is combined with user transactions. Relevance feedback is a well-established technique for improving Interactive Information Retrieval. We can find a specific reference to feature relevance since 1994, in [77], where authors have proposed a modified Rocchio’s relevance feedback approach to improve accuracy. This idea is then reflected in [397], in which the authors propose an *idealized relevance feedback (IRF)* framework. To further improve the previous work, they perform a term normalization and selection based on the ranking position.

CRe1-FM is also related to the exploitation of explicit features in Matrix Factorization. One of the first models of this kind is *Explicit Factor Model (EFM)* [419]. In this model, products’ features and users’ opinions are extracted with phrase-level sentiment analysis from users’ reviews to feed a matrix factorization framework. Since then, a few improvements to *EFM* have been proposed to deal with temporal dynamics [420] and to use tensor factorization [97]. In particular, in the latter the aim is to predict both user preferences on features (extracted from textual reviews) and items. This is achieved by exploiting the *Bayesian Personalized Ranking (BPR)* criterion [312].

In CRe1-FM, we exploit the *Factorization Machines (FM)* data model [310]. FMs are the most widely used factorization models because they offer a number of advantages compared against other latent factors models such as *SVD++* [221], *PITF* [315], *FPMC* [313]. The main advantage is that *FMs* are designed for a generic prediction task. On the other side, other mentioned techniques are task-specific. Moreover, *FM* is a linear model where parameters can be estimated accurately even in a high data sparsity scenario. Nevertheless, many researchers have proposed variants of *FMs*. In [174], the authors propose *Neural Factorization Machines* to increase the expressiveness of *FM*. In particular, they enable *FMs* to capture the non-linear structure of real-world data exploiting the non-linearity of neural networks. Furthermore, *Attentional Factorization Machines* [401] use an attention network to learn the relevance of feature interactions. Finally, in [314], the authors propose a specialized Context-Aware variant of *FMs*.

In a *top-N* recommendation task, recommenders provide users a shortlist of rec-

ommendations. For this reason, correctly sorting recommendations has gradually replaced the rating prediction task [256]. These *Learning to Rank* [88] algorithms can be further categorized in *Point-wise* [225], *Pair-wise* [312, 247] and *List-wise* [344, 343]. In particular, *Pair-wise* approaches are usually considered as a good trade-off between ordering performance and computational complexity. Among this class of algorithms, *Bayesian Personalized Ranking (BPR)* [312] is one of the most widely adopted. It is based on a simple stochastic gradient descent algorithm to learn the relative order between positive items (items that a user has experienced in his past history) and negative items (items never rated by the user). *BPR* can be easily applied to *Matrix Factorization* and *Factorization Machines* (as in CReL-FM and in [41]).

9.5 Conclusions and Future Work

In this chapter, we propose and test CReL-FM , a principled method to compute the relevance of item attributes. Its crucial insight is that of exploiting collaborative information in a recommendation scenario to fine-tune the relevance of users/items attributes. The overall approach goes through the exploitation of *Factorization Machines* to map latent factors to actual attributes and the training of the model adopting a *Learning to Rank* optimization criterion to inject collaborative information. Through an extensive experimental evaluation, we have found that the proposed notion of hybrid relevance provides significant improvements to Recommender Systems even when we adopt the very simple *TF-IDF* relevance measure combined with collaborative information via CReL-FM . Furthermore, the exploitation of different signals has led to more accurate and diverse recommendations proving that this could be a solution to the classical accuracy-diversity trade-off [377]. In this work, we have shown the effectiveness of the approach by relying only on one traditional relevance measure (*TF-IDF*). As for our future work, we plan to investigate the behavior of CReL-FM with *Entropy-based* [389], *BM25* [201, 202], and *BM25F* [201, 202, 319] weighting schemes. Another investigation direction would be that of replacing the *BPR* optimization with other *pair-wise* and *list-wise Learning to*

Rank optimization algorithms like *AdaRank* [402] and *LambdaMart* [400]. Moreover, we plan to evaluate the method in a more general Information Retrieval setting and to enable explanation services.

Part III

Modeling and evaluating without an explicit knowledge representation

Chapter 10

Introduction

This chapter is devoted to all the approaches that make use of non-completely structured knowledge in recommender systems. Here we study some dimensions like the Popularity of items and Time that are present in real case scenarios that are encoded differently. Popularity is usually a derived quantity exploited by some non-personalized recommenders. On the other side, Time is a piece of optional information that can be exploited for improving recommendations. Since when they are available they are represented in the same way, we decided to face these dimensions in this research. Our study has led to face other classic problems, like the definition of similarity, the tuning of parameters, and the consequent evaluation. However, the evaluation itself is nowadays under observation, because it should monitor "unfair" recommendations. For this reason, a study on fairness evaluation closes the chapter.

In modern recommender systems, diversity has been widely acknowledged as an important factor to improve user experience and, more recently, intent-aware approaches to diversification have been proposed to provide the user with a list of recommendations covering different aspects of her behavior. In this line of research,

we propose and analyze the performances of two diversification methods taking into account temporal aspects of the user profile: in the first one we adopt a temporal decay function to emphasize the importance of more recent items in the user profile while in the second one we perform an evaluation based on the identification and analysis of temporal sessions. The two proposed methods have been implemented as temporal variants of the well-known xQuAD framework. In both cases, experimental results on Netflix 100M show an improvement in terms of accuracy-diversity balance .

Items popularity is a strong signal in recommendation algorithms. It strongly affects collaborative filtering approaches and it has been proven to be a very good baseline in terms of results accuracy. Even though we miss an actual personalization, global popularity can be effectively used to recommend items to users. In this second line of research, we introduce the idea of **time-aware personalized popularity** in recommender systems by considering both items' popularity among neighbors and how it changes over time. An experimental evaluation shows a highly competitive behavior of the proposed approach, compared to state of the art model-based collaborative approaches, in terms of results accuracy.

In recommendation scenarios, similarity measures play a fundamental role in memory-based nearest neighbor approaches. They recommend items to a user based on the similarity of either items or users in a neighborhood. In this line of research, **we argue that similarity** between users or items, although it keeps leading importance in computing recommendations, **should be paired with a value of dissimilarity** (computed not just like the complement of the similarity one). We formally modeled and injected this notion in some of the most used similarity measures and evaluated our approach in a recommendation scenario showing its effectiveness for accuracy and diversity results on three different datasets.

Hyperparameters tuning is a crucial task to make a model perform at its best. However, despite the well-established methodologies, some aspects of the tuning remain unexplored. As an example, it may affect not just accuracy but also Novelty as well as it may depend on the adopted dataset. Moreover, sometimes it could be sufficient to concentrate on a single parameter only (or a few of them) instead

of their overall set. In this line of research, we report on our **investigation on hyperparameters tuning** by performing an extensive 10-Folds Cross-Validation on MovieLens and Amazon Movies for three well-known baselines: User-kNN, Item-kNN, BPR-MF. We adopted a grid search strategy considering approximately 15 values for each parameter, and we then evaluated each combination of parameters in terms of accuracy and novelty. **We investigated the discriminative power of $nDCG$, Precision, Recall, MRR, EFD, EPC**, and, finally, we analyzed the role of parameters on model evaluation for Cross-Validation.

One common characteristic of research works focused on fairness evaluation (in machine learning) is that they call for some form of parity (equality) either in treatment – meaning they ignore the information about users’ memberships in protected classes during training – or in impact – by enforcing proportional beneficial outcomes to users in different protected classes. In the recommender systems community, fairness has been studied for both users’ and items’ memberships in protected classes defined by some sensitive attributes (e.g., gender or race for users, revenue in a multi-stakeholder setting for items).

The fairness has been commonly interpreted as some form of *equality*. This implies to measure if the system is meeting the information needs of all its users *in an equal sense*. In this line of research, **we propose a probabilistic framework based on Generalized Cross-Entropy (GCE) to measure the fairness** of a given recommendation model. The framework comes with a suite of advantages: first, it allows the system designer to define and measure fairness for both users and items and can be applied to any classification task; second, it can incorporate various notions of fairness as it does not rely on specific and pre-defined probability distributions but they can be defined at design time.

In its design, it uses a gain factor, which can be flexibly defined to contemplate different accuracy-related metrics to measure fairness upon such as decision-support metrics (e.g., precision, recall) or rank-based measures (e.g., NDCG, MAP). Results on three real-world datasets show the nuances capture by our metric regarding fairness on different user and item attributes, where nearest-neighbor recommenders tend to obtain good results under equality constraints.

Chapter 11

A re-ranking algorithm exploiting Time and Diversity

11.1 Introduction

Recommender systems are designed to meet the users' needs suggesting relevant items in a personalized fashion. As recommendations are usually presented in form of list or group, the user experience strongly depends on the overall quality of such recommendations and, the diversity among them has been identified as one of the most important quality factor [256, 91]. Generally, accuracy and diversity are considered as contrasting properties, due to the demonstrated trade-off between them in offline evaluation [91]. In spite of that, a recent user's behaviour study proved that diversity in recommendations has an important positive impact on user satisfaction [134]. Moreover, in [283] the authors show that by taking into account the users propensity towards diversity, it is possible to foster the recommendation diversity without affecting accuracy or even slightly improve it. The diversity issue has been originally addressed in the Information Retrieval field. As user queries are often am-

ambiguous and their intent is not clear, proposing a set of answers covering different intents may increase the probability that users find at least one relevant document [330, 94]. The concept of intent-aware diversification has then been applied to the Recommender Systems field [379] and extensively studied thus producing new algorithms and evaluation metrics [378, 376, 391]. Here, user intents as defined in Information Retrieval have been mapped to user interests with reference to item characteristics.

Very often, in the design of the model behind a recommendation engine, the user profile is considered as a static snapshot without taking into proper account its temporal dimension. Actually, the importance of analyzing temporal aspect for user modeling has been proved to affect the final recommendation results [218, 222, 200].

Based on the above observations, we investigate the effect on the trade-off between accuracy and diversity of a recommendation list when dealing with temporal aspects of the user profile. The intuition behind our idea is that temporal dynamics might allow to better understand the user interests with respect to the items characteristics and then provide a more accurate intent-aware diversification. Therefore, this work presents two intent-based modeling methods that exploit the time dimension in a user profile. The first one analyzes the frequency of interaction between the users and the items features using a temporal decay function in order to emphasize *persistence* and *recency* of an intent. The other method is based of a new session analysis technique of user ratings for intent modeling. Considering that a session is usually defined as a set of consecutive ratings with a very small gap of time among them (e.g., less than one hour in music [218]), we provided a wide definition of session tailored for movie ratings. In particular, such method is designed to highlight *importance*, *persistence* and *recency* of an intent among user sessions. We experimentally evaluated such methods with the large scale movie dataset published in the context of Netflix Prize Context [52]. The experimental results demonstrated that the analysis of temporal aspects in the user profile leads to better accuracy-diversity balance and intent-aware diversity compared to the original xQuAD. As an additional benefit, the aggregate diversity results improved too thus demonstrating

to produce more personalized recommendations.

To the best of our knowledge, this is the first attempt in the investigation of how temporal dynamics affect diversity in recommendations.

11.2 Related Work

A list of recommendations can be diversified in an implicit or explicit manner [42]. While the implicit diversification is used to increase the average distance between pairs of items in the recommendation list, the explicit method diversifies the recommendations trying to cover the user interests represented via categories or other descriptive information of the items. Therefore, the explicit diversification is known as Intent-Aware since it considers the likeness of user intents in information retrieval and user interests in recommender systems [379]. Explicit Query Aspect Diversification (xQuAD) is one of the most well-known intent-aware framework originally proposed for query results diversification [330] and then adapted to the recommendation [378] field. In a nutshell, xQuAD aims at maximizing the coverage of the inferred interests while minimizing their redundancy in a recommendation list.

The importance of taking into account the temporal dynamics in recommender systems has been recently pointed out in different works for diverse recommendation domains. A method to model user sessions in music domain was proposed in [218]. It considers as session each set of consecutive ratings without an extended time gap between them. Considering that there are various psychological phenomena that lead to a set of ratings to be grouped into a single session, such method captures these effects by means of user session biases. [222] presented a collaborative filtering algorithm able to model time drifting of user preferences and the results on the Netflix dataset indicated the importance of unveiling temporal effects in order to produce more accurate recommendations. A more recent method proposed to take advantage of temporal information in user behavior is called Time-based Markov Embedding [200], used to find the best next-song recommendation via Latent Markov Embedding.

In this work we aim at exploring the exploitation of temporal dynamics in user

intents to provide a better intent-aware diversification.

11.3 Intent-aware diversification for recommendations

Typically, a recommender system produces a list of personalized recommendations for each user. According to [14], a re-ranking of such list can be applied to improve its diversity, without modifying the recommendation process. However, finding the most diverse results is a NP-hard problem and hence several heuristics have been proposed [228]. Most previous diversification approaches are based on a greedy selection strategy [90, 330, 33]. Such strategy selects the next most relevant item only if that item is diverse with respect to the items already selected [228]. Algorithm

Data: The original list \mathbf{R} , $N \leq n$

Result: The re-ranked list \mathbf{S}

$\mathbf{S} = \langle \rangle;$

while $|\mathbf{S}| < N$ **do**

$i^* = \operatorname{argmax}_{i \in \mathbf{R} \setminus \mathbf{S}} f_{obj}(i, \mathbf{S}, u);$

$\mathbf{S} = \mathbf{S} \circ i^*;$

$\mathbf{R} = \mathbf{R} \setminus \{i^*\}$

end

return \mathbf{S} .

Algorithm 3: The greedy strategy.

3 describes the working scheme of a greedy selection method. For the purpose of this work, we consider x_{QuAD} , one of the most well-known intent-aware greedy heuristics. It maximizes the coverage of the inferred interests while minimizing their redundancy. x_{QuAD} was proposed for search diversification in information retrieval by Santos et al. [330], as a probabilistic framework to explicitly model an ambiguous query as a set of sub-queries that are supposed to cover the potential aspects of the initial query. More recently, it has been adapted for recommendation diversification by Vargas and Castells [378], replacing query and relative aspects

with user and items features, respectively. The expression of the $\times\text{QuAD}$ objective function is

$$f_{obj}(i, \mathbf{S}, u) = \lambda \cdot r^*(u, i) + (1 - \lambda) \cdot \text{div}(i, \bar{\mathbf{S}}, u) \quad (11.1)$$

with $\bar{\mathbf{S}}$ representing the set of the items belonging to \mathbf{R} not already in \mathbf{S} and $\text{div}(i, \bar{\mathbf{S}}, u)$ defined as

$$\text{div}(i, \bar{\mathbf{S}}, u) = \sum_f p(i|f, u) \cdot p(f|u) \cdot \prod_{j \in \bar{\mathbf{S}}} (1 - p(j|f, u)) \quad (11.2)$$

In Equation (11.2) $p(i|f, u)$ represents the likelihood of item i being chosen given the feature f and is computed as a binary function that returns 1 if the item contains f , 0 otherwise; $p(f|u)$ represents the interest of user u in the feature f and is usually computed as the relative frequency of the feature f on the items rated by user u . In other words, $\times\text{QuAD}$ fosters the idea of promoting items that are simultaneously highly related to at least one of the features of interest for the user and slightly related to the features of the items already recommended. In particular, this work focuses on the intent modeling in the $\times\text{QuAD}$ framework, namely the aforementioned $p(f|u)$ component in Equation (11.2).

11.4 Intent modeling with Temporal Dynamics

In this section, we propose two methods to exploit temporal analysis for intent modeling in diversification that we call **session-based** and **time-based intent modeling**. Both relies on the intuition that user intent can change during the interaction with the system and the simple computation of features frequency in the user profile may not represent the current user interests.

11.4.1 Time-Based Intent Modeling

In order to valorize *persistence* and *recency* of an intent, we propose to analyze the frequency of interaction between the user u and the feature f and to weight each interaction by a temporal decay function. More formally, the following formula computes the interest of the user u with respect to the feature f :

$$p(f|u) = \frac{\sum_{i \in R(u)} cov(f, i) disc(u, i)}{\sum_{i \in R(u)} disc(u, i)} \quad (11.3)$$

where $R(u)$ indicates the set of rating provided by the user u ; $cov(f, i)$ is a binary function returning 1 if the item i is associated with the feature f , 0 otherwise; $disc(i, u)$ is a temporal decay function returning lower values for older ratings, and higher values for the most recent ones. Inspired by [222], as decay function we adopted the following exponential function

$$disc(u, i) = e^{-\beta \cdot |t_{u, last} - t_{u, i}|} \quad (11.4)$$

where $t_{u, last}$ indicates the timestamp of the last rating of the user u and $t_{u, i}$ the timestamp when user u rated i ; $\beta > 0$ controls the decay rate.

In our experimental setting we adopted the Netflix dataset which contains ratings from October 1998 until October 2005. This information affects the choose of β because a too big value of beta could make the initial ratings not valuable at all. We decided to consider all the ratings valuable with a very smooth curve and a minimum value of 10^{-6} . Hence β value was set to $1/200$ obtaining $2,831 * 10^{-6}$ for the farrest rating over time.

11.4.2 Session-Based Intent Modeling

User sessions definition. Session analysis is quite common in music domain, since users are used to listen to many songs in sequence. There a session is represented by a set of consecutive ratings with a small gap of time between them [218]. Conversely, sessions are not easy to find in movie domain, since users typically watch a small number of movies in brief timeslots and the temporal gap among visions or ratings could be large (sometimes several days or even months). In our setting, in order to identify user sessions we propose an EM clustering used to train two univariate Gaussian Mixture Models (with equivariance and variable variance). The number of clusters has been evaluated based on the Bayesian Information Criterion considering the fitted models with a number of clusters from 1 up to 300. In order to remove outliers from each session s , for each computed cluster we do not consider

ratings falling outside the interval $[\mu_s - \sigma_s, \mu_s + \sigma_s]$, with μ_s and σ_s being respectively the mean and the standard deviation of ratings distribution for the session s .

Intent modeling. Once user sessions are determined, they can be used to analyze the user activities taking into account the temporal dynamics. In this work we present an approach to model the users intents over time, by considering three key properties: *importance*, *persistence* and *recency* of an intent among the user sessions. The first property indicates the importance of an intent in each session computed as the percentage of items covering that intent. The second property considers how many sessions the intent is important for, therefore it sums the importance of the intent for each sessions. Finally, the third property focuses on the intent freshness, penalizing old sessions with a temporal decay function.

More formally, the following formula computes the interest of the user u with respect to the feature f :

$$p(f|u) = \frac{\sum_{s \in S(u)} \frac{\sum_{i \in I(s)} cov(f,i)}{|F(s)|} disc(s,u)}{\sum_{s \in S(u)} disc(s,u)} \quad (11.5)$$

where $S(u)$ indicates sessions computed for the user u ; $I(s)$ is the set of items in s ; $cov(f,i)$ is a binary function returning 1 if the item i is associated with the feature f , 0 otherwise; $F(s)$ represents the set of features associated with all the items in s ; $disc(s)$ is the temporal decay function adapted to handle the sessions instead of the items, considering a session as an item in Equation (11.4) where the session date is that of the last rated item in such session. As for the previous case, β value was set to $1/200$.

11.5 Experimental setting

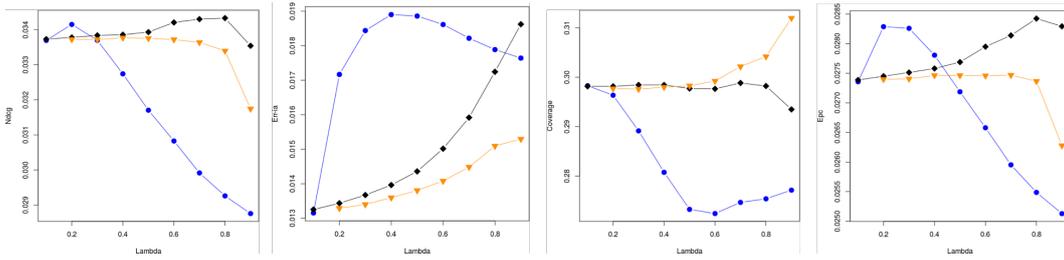
Dataset. In order to verify our research questions and evaluate our proposal, we used the popular movie datasets derived from the Netflix Prize Context [52]. Netflix dataset contains over 100 million ratings provided by $\sim 480,000$ users on $\sim 17,000$ movies. Such ratings were collected between 1998 and 2005 and associated with

the relative date. However, such dataset contains noise added on purpose for reasons of privacy, as explained in the Netflix Prize Rules¹: "some of the rating data for some customers in the training and qualifying sets have been deliberately perturbed in one or more of the following ways: deleting ratings; inserting alternative ratings and dates; and modifying rating dates". Indeed, we found that some users rated an exaggerated number of movies in some days: 30% of all the users have rated at least 61 movies in the most *prolific* day. Therefore, in order to train the session-based user model on a clean subset of the dataset, we selected a sample of users removing the outliers by means of the following steps: (i) we discarded the users with less than 20 ratings as at this stage of our study we are not interested in evaluating the cold users behavior; (ii) we ordered the users in decreasing order of the maximum number of daily interactions and discarded the top 30%; (iii) we ordered the users in decreasing order of the average number of daily interactions and discarded the top 30%. The dataset used for training the models contains 233,452 users, 18,104,476 rating and 17,763 movies. We built training and test sets by employing a 80%-20% holdout temporal split for each user.

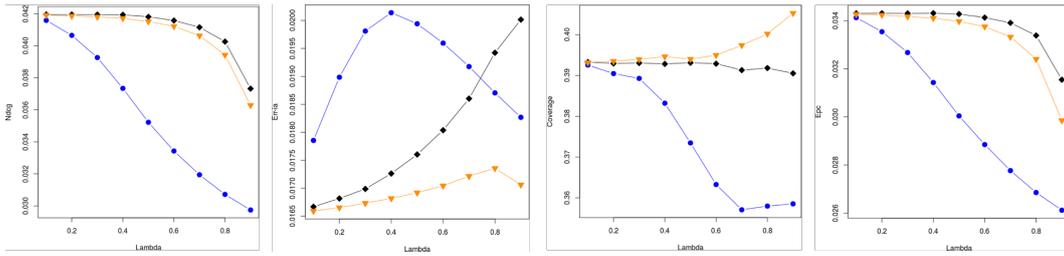
Recommendation algorithms. We evaluated our approaches w.r.t. the x_{QuAD} baseline via a re-ranking of the BPRMF [312] and BPRSLIM [273] algorithms resulting recommendations. We trained both models using the MyMediaLite² implementation upon the dataset described in Section 11.5 to produce for each user a Top-300 recommendations list (\mathbf{R} in Algorithm 3, used to compute $r^*(u, i)$ in Equation (11.1)), and then we re-ranked those lists using x_{QuAD} . The resulting recommendations lists are the baselines we compare against. x_{QuAD} uses side information to lead to diversified recommendations. In this work the diversification is based on the movie genre feature. In Netflix dataset this information is not explicitly provided. In order to extract the genre information, we mapped each movie with the corresponding Freebase resource by means of its title and year of release and we then selected the corresponding genres. Overall, the number of distinct genres extracted is 266.

¹<http://netflixprize.com/rules.html>

²<http://www.mymedialite.net>



(a) Results for BPRMF



(b) Results for BPRSLIM

Figure 11.1: Curves obtained by varying λ from 0.1 to 0.9 (step 0.1), using BPRMF and BPRSLIM as recommendation algorithm. From left to right we plot the values of nDCG, ERR-IA, Catalog Coverage and EPC. The blue line represents values for the base xQuAD evaluation while the black and yellow lines represent respectively values for the time-based and session-based version of xQuAD

In our evaluation, the time-based and session-based intent modelings proposed in Section 11.4 are used as alternatives to the pure frequency based intent modeling in the original xQuAD. These two variations of xQuAD, are denoted as: TB_xQuAD, SB_xQuAD, where *TB* stands for time-based and *SB* for session-based. The resulting evaluated algorithms are then:

- BPRMF + xQuAD (baseline)
- BPRMF + Time-based xQuAD variant (TB_xQuAD)
- BPRMF + Session-based xQuAD variant (SB_xQuAD)
- BPRSLIM + xQuAD (baseline)
- BPRSLIM + Time-based xQuAD variant (TB_xQuAD)

- BPRSLIM + Session-based \times QuAD variant (SB- \times QuAD)

Metrics. In order to evaluate *accuracy*, we measured nDCG. As for *individual diversity*, namely the degree of dissimilarity among all items in the list provided to a user, was measured by ERR-IA as it has been shown [377] to be the metric targeted by \times QuAD, while for *aggregate diversity* we computed the Catalog Coverage (percentage of items recommended at least to one user). An evaluation on the *novelty* of computed results has been done through EPC (Expected Popularity Complement) [377]. For all the aforementioned metrics we used the implementation provided by RankSys framework³ on the Top-10 recommendation list.

Results Discussion. Charts in Figure 11.1a and Figure 11.1b show the curves for nDCG, ERR-IA, Catalog Coverage and EPC, for both BPRMF and BPRSLIM variants. Very interestingly both the time- and session-based version of \times QuAD improves results in terms of accuracy, aggregate diversity as well as of novelty independently of the recommendation algorithm adopted. Generally, the time-based variant of \times QuAD performs better than the session-based one but for Catalog Coverage where we have better results for the session-based implementation of \times QuAD. It is worth noticing that the base version of \times QuAD outperforms its time-based variants up to a certain value of λ for both BPRMF and BPRSLIM. For the former this value lies between 0.8 and 0.9 while for the latter between 0.7 and 0.8. Hence, in case we are interested in higher values of diversity, time may play an important role. This observation is also strengthened by the higher values we obtain in terms of precision, catalog coverage and novelty. In Table 11.1 we see that with $\lambda = 0.8$ we obtain the best result in terms of trade-off among the various metrics we measured in our experiments.

11.6 Conclusions and future work

This line of research has been devoted to investigate the role of temporal information while modeling a user profile in computing diversified recommendations. We propose two different time-dependent user modelings which take into account also

³<https://github.com/RankSys/RankSys>

| Algorithm | Ndcg@10 | ERR IA@10 | Coverage@10 | EPC@10 |
|------------------|----------------|----------------|----------------|----------------|
| BPRMF+XQUAD | 0.029264 | 0.01789 | 0.27540 | 0.02549 |
| BPRMF.SB_XQuAD | 0.03340 | 0.01510 | 0.30417 | 0.02737 |
| BPRMF.TB_XQuAD | 0.03433 | 0.01724 | 0.29820 | 0.02843 |
| BPRSLIM+XQUAD | 0.03072 | 0.01870 | 0.35799 | 0.02686 |
| BPRSLIM.SB_XQuAD | 0.03943 | 0.01736 | 0.40021 | 0.03240 |
| BPRSLIM.TB_XQuAD | 0.04026 | 0.01942 | 0.39183 | 0.03339 |

Table 11.1: Comparative results in terms of accuracy, individual diversity and aggregate diversity with $\lambda = 0.8$

the user rating history. One of the two proposed methods bases on a new session analysis technique by considering those periods where the user interacted in a more constant way with the system. Experimental results demonstrated that considering temporal dynamics leads to better accuracy-diversity balance and better intent-aware diversification. The results we obtained in this preliminary investigation are quite promising and we are currently extending our experimental evaluation to different datasets in diverse domains as well as to other metrics in order to measure the quality of the recommendations not just in terms of accuracy when time is taken into account.

Chapter 12

A study over the dimensions of Time, and Popularity

12.1 Introduction

Collaborative-Filtering (CF) [310] algorithms, more than others, have gained a key-role among recommendation approaches and have been effectively implemented in commercial systems to help users in dealing with the information overload problem. Some of them also use additional information (hybrid approaches) to build a more precise user profile in order to serve a much more personalized list of items [27, 140].

However, it is well known [193] that all the algorithms based on a CF approach are affected by the so called “popularity bias” meaning that popular items tend to be recommended more frequently than those in the long tail. Initially considered as a shortcoming of collaborative filtering algorithms and then not useful to produce good recommendations [191], in some works items popularity has been intentionally penalized [285]. Very interestingly, a recommendation algorithm purely based

on most popular items, has been proven to be a strong baseline [108] although it does not exploit any actual personalization. More recently, popularity has also been considered as a natural aspect of recommendation that, by measuring the user tendency to diversification, can be exploited to balance the recommender optimization goals [204]. The study of popularity in user tendencies is not completely new in the recommender systems field. Some interesting works explored these criteria for re-ranking purposes [204, 285], and multiple goals optimization [191].

In the approach we present here, we introduce a more fine-grained personalized version of popularity by assuming that it is conditioned by the items that a user u already experienced *in the past*. To this extent, we look at a specific class of neighbors, that we name *Precursors*, defined as the users who already rated the same items of u in the past. This led us to the introduction of a time-aware analysis while computing a recommendation list for u . As time is considered a contextual feature, most of the works dealing with temporal aspects are considered as a subclass of Context-Aware RS (CARS) [18]: Time-Aware RS (TARS) [430, 16, 222]. In TARS, the freshness of different ratings is often considered as a discriminative factor between candidate items. Usually, a time window [229] is adopted to filter out all the ratings that stand before (and/or after) a certain time relative to the user or the item. Recently, an interesting work that makes use of time windows has been proposed in [51] where the authors focus on the last common interaction between the target user and her neighbors to populate the candidate items list. In [38] social information and time are integrated dealing with the interests of the users as a series of temporal matrices. Probabilistic matrix factorization technique are adopted to learn latent factors. Regarding sequences and recommendation it is worth to mention [399], in which the authors combine an LSTM network with a low-rank matrix factorization algorithm to produce recommendation lists. A pioneer work was proposed more than a decade ago in [122] which used an exponential decay function $e^{-\lambda t}$ to penalize old ratings. An exponential decay function [222] was then used to integrate time in a latent factors model. In the last years, several Item-kNN [246, 122] with a temporal decay function have been deployed. Another interesting work was proposed in [2] where three different kinds of time decay were adopted:

exploiting concave, convex and linear functions.

This research line has led to the development of `TimePop`, an algorithm that combines the notion of personalized popularity conditioned to the behavior of users' neighbors while taking into account the temporal dimension. It is worth noticing that `TimePop` works with implicit feedback to compute recommendations. Differently from some of the approaches previously described, in `TimePop` we avoid both the use of a time window and the selection of a fixed number of candidate items. Indeed, while on the one hand, a time window may severely restrict the selection of candidates, on the other hand, a fixed number of candidate items may heavily affect the algorithm results.

12.2 Time-aware Local Popularity

The leading intuition behind `TimePop` is that the popularity of an item has not to be considered as a global property but it can be personalized if we consider the popularity in a neighborhood of users. We started from this observation to formulate a form of personalized popularity, and then we added the temporal dimension to strengthen this idea.

In `TimePop`, given a user u the first step is then the identification of user's neighbors who rated the same items as u but before u . We name these users *Precursors*. In our intuition, Precursors represent a community of users u relies on to choose the items to enjoy. In a neighborhood of u , the same item is enjoyed by users in different time frames. This leads us to the second ingredient behind `TimePop`: personalized popularity is a function of time. The more the ratings about an item are recent, the more its popularity is relevant for the specific user. Hence, in order to exploit the temporal aspect of these ratings, the contributions of *Precursors* can be weighted depending on their freshness.

We now introduce some basic notation that will be used in the following. We use $u \in U$ and $i \in I$ to denote users and items respectively. Since we are not just interested in the items a user rated but also at when the rating happened, we have that for a user u the corresponding user profile is $P_u = \{(i_1, t_{ui_1}), \dots, (i_n, t_{ui_n})\}$ with

$P_u \subseteq I \times \mathfrak{R}$, being t_{ui} a timestamp representing when u rated i .

Definition 1 (Candidate Precursor and Precursor)

Given $(i, t_{ui}) \in P_u$ and $(i, t_{u'i}) \in P_{u'}$, we say that u' is a **Candidate Precursor** of u if $t_{u'i} < t_{ui}$. We use the set $\hat{\mathcal{P}}^u$ to denote the set of Candidate Precursors of u . Given two users u' and u such that u' is a Candidate Precursor of u and a value $\tau_u \in \mathfrak{R}$ we say that u' is a **Precursor** of u if the following condition holds.

$$|\{i \mid (i, t_{ui}) \in P_u \wedge (i, t_{u'i}) \in P_{u'} \wedge t_{u'i} < t_{ui}\}| \geq \tau_u$$

We use \mathcal{P}^u to denote the set of **Precursors** of u .

A user u' is a *Candidate Precursor* of u if u' rated at least one common item i before u . Although this definition catches the intuition behind the idea of *Precursors*, it is a bit weak as it also considers users u' who have only a few or even just one item in common with u and rated them before she did. Hence, we introduced a threshold taking somehow into account the number of common items in order to enforce the notion of *Precursors*. The threshold parameter τ_u in Definition 1 can also be computed automatically as:

$$\tau_u = \frac{\sum_{u' \in \hat{\mathcal{P}}^u} |\{i \mid (i, t_{ui}) \in P_u \wedge (i, t_{u'i}) \in P_{u'} \wedge t_{u'i} < t_{ui}\}|}{|\hat{\mathcal{P}}^u|} \quad (12.1)$$

To give an intuition on the computation of Precursors and of τ_u let us describe the simple example shown in Figure 12.1.

Here, for the sake of simplicity, we suppose that there are only four users and six items and u is the user we want to provide recommendations to. Items that users share with u are highlighted in blue and items with a dashed red square are the ones that have been rated before u . We see that $\hat{\mathcal{P}}^u = \{u_2, u_4\}$. Indeed, although u_3 rated some of the items also rated by u they have been rated after. By Equation (12.1) we have $\tau_u = \frac{3}{2} = 1.5$. Then, only u_2 results to be in \mathcal{P}^u because she has $2 > 1.5$ shared items rated before those of u . As for u_3 , it is more likely that u is a Precursor of u_3 and not vice versa.

Temporal decay. As the definition of Precursor goes through a temporal analysis of user behaviors, we may look at the timestamp of the last rating provided by

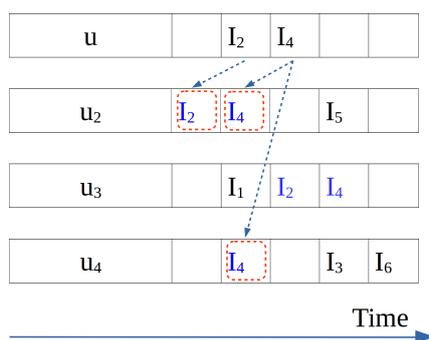


Figure 12.1: Example of Precursors computation.

a Precursor in order to identify how active she is in the system. Intuitively, the contribution to popularity for users who have not contributed recently with a rating is lower than “active” users. On the other side, given an item in the profile of a Precursor we are interested in the freshness of its rating. As a matter of fact, old ratings should affect the popularity of an item less than newer ratings. Summing up, we may classify the two temporal dimensions as **old/recent user** and **old/recent item**. In order to quantify these two dimensions for Precursors we introduce the following timestamps: \mathbf{t}_0 this is the reference timestamp. It represents the “now” in our system; $\mathbf{t}_{u'i}$ is the time when u' rated i ; $\mathbf{t}_{u'l}$ represents the timestamp associated to the last item l rated by the user u' . Different temporal variables are typically used [122, 222], and they mainly focus on **old/recent items**. ΔT may refer to the timestamp of the items with reference to the last rating of u' [122] with $\Delta T = \mathbf{t}_{u'l} - \mathbf{t}_{u'i}$ or to the reference timestamp [222] with $\Delta T = \mathbf{t}_0 - \mathbf{t}_{u'i}$. As we stated before, our approach captures the temporal behavior of both **old/recent users** and **old/recent items** at the same time. We may analyze the desired ideal behavior of ΔT depending on the three timestamps previously introduced as represented in Table 12.1.

Let us focus on each case. In the upper-left case we want ΔT to be as small as possible because both u' and the rating for i are “recent” and then highly representative for a popularity dimension. In the upper-right case, the rating is recent but the user is old. The last item has been rated very close to i but a large value of ΔT should remain because the age of u' penalizes the contribution. The lower-left case denotes a user that is active on the system but rated i a long time ago. In this case

| | recent user ($t_0 \approx t_{u'1}$) | old user ($t_0 \gg t_{u'1}$) |
|---|---|--|
| recent item ($t_{u'1} \approx t_{u'i}$) | ≈ 0 | $t_0 - t_{u'1}$ |
| old item ($t_{u'1} \gg t_{u'i}$) | $t_{u'1} - t_{u'i}$ | $t_0 - t_{u'1}$ |

Table 12.1: *Ideal values of ΔT w.r.t. the Precursor characteristics*

the contribution of this item is almost equal to the age of its rating. The lower-right case is related to a scenario in which both the rating and u' are old. In this scenario, the differences between the reference timestamp minus the last interaction and the reference timestamp minus the rating of i are comparable: $(t_0 - t_{u'1}) \approx (t_0 - t_{u'i})$. In this case, we wish the contribution of ΔT to consider the elapsed time from the last interaction (or the rating) until the reference timestamp. All the above observations lead us to define $\Delta T = |t_0 - 2t_{u'1} + t_{u'i}|$. In order to avoid different decay coefficients, in our experimental evaluation, all ΔT s are transformed in days (from milliseconds) as a common practice.

The Recommendation Algorithm. We modeled our algorithm `TimePop` to solve a *top-N* recommendation problem. Given a user u , `TimePop` computes the recommendation list by executing the following steps:

1. Compute \mathcal{P}^u ;
2. For each item i such that there exists $u' \in \mathcal{P}^u$ with $(i, t_{u'i}) \in P_{u'}$ compute a score for i by summing the number of times it appears in $P_{u'}$ multiplied by the corresponding decay function;
3. Sort the list in decreasing order with respect to the score of each i .

For sake of completeness, in case there were no precursors for a certain user, a recommendation list based on global popularity is returned to u . Moreover, if `TimePop` is able to compute only m scores, with $m < N$, the remaining items are returned based on their value of global popularity.

12.3 Experimental Evaluation

In order to evaluate `TimePop` we tested our approach considering datasets related to different domains. Two of them related to the movie domain —the well-known Movielens1M dataset and Amazon¹ Movies — and a dataset referring to toys and games —Amazon Toys and Games, with 2M ratings and a sparsity of 99.99949%. “All Unrated Items” [356] protocol has been chosen to compare different algorithms where, for each user, all the items that have not yet been rated by the user all over the platform are considered. In order to evaluate time-aware recommender systems in an offline experimental setting, a typical k-folds or hold-out splitting would be ineffective and unrealistic. To be as close as possible to an online real scenario we used the fixed-timestamp splitting method [83, 163], also used in [51] but with a *dataset centered* base set. The basic idea is choosing a single timestamp that represents the moment in which test users are on the platform waiting for recommendations. Their past corresponds to the training set, and the performance is evaluated with data coming from their future. In this work, we select the splitting timestamp that maximizes the number of users involved in the evaluation by setting two constraints: the training set must keep at least 15 ratings, and the test set must contain at least 5 ratings. Training set and test set for the three datasets are publicly available² along with the splitting code for research purposes. In order to evaluate the algorithms we measured *normalized Discount Cumulative Gain@N* ($nDCG@N$) using *Time-independent rating order condition* [83]. The metric was computed per user and then the overall mean was returned using the RankSys framework and adopting *Threshold-based relevant items* condition [83]. The threshold used to consider a test item as relevant has been set to the value of 4 w.r.t. a 1-5 scale for all the three datasets.

Baselines. We evaluated our approach w.r.t CF and time-aware techniques. **Most-Popular** was included as `TimePop` is a time-aware variant of “Most Popular”. From model-based collaborative filtering approaches we selected some of the best

¹<http://jmcauley.ucsd.edu/data/amazon/>

²<https://github.com/sisinflab/DatasetsSplits>

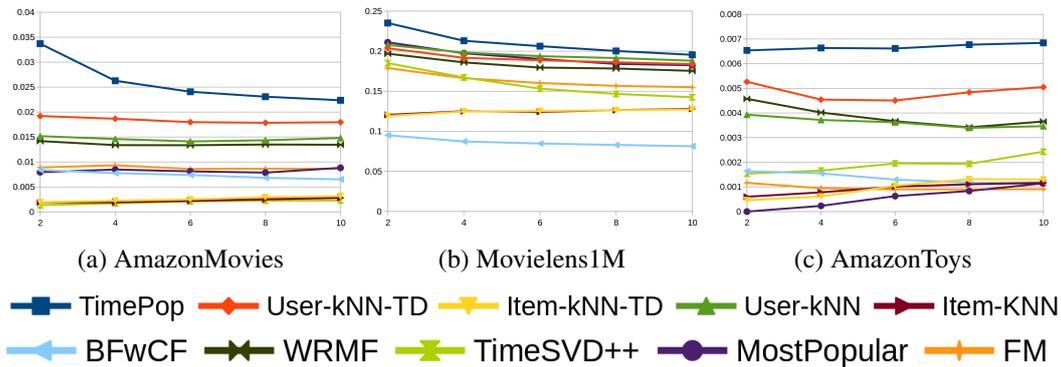


Figure 12.2: nDCG @N varying N in 2..10

performing matrix factorization algorithms **WRMF** trained with a regularization parameter set to 0.015, α set to 1 and 15 iterations, and **FM**³[310], computed with an ad-hoc implementation of a 2 degree factorization machine considering users and items as features, trained using Bayesian Personalized Ranking Criterion[312]. Moreover, we compared our approach against the most popular memory-based kNN algorithms, **Item-kNN**³ and **User-kNN**³ [333], together with their time-aware variants (**Item-kNN-TD**³, **User-kNN-TD**³)[122]. We included **TimeSVD++**³ [222] in our comparison even though this latter has been explicitly designed for the rating prediction task. All model-based algorithms were trained using 10, 50, 100, and 200 factors; only best models are reported in the evaluation: for Movielens1M WRMF 10, FM 10; for Amazon Movies WRMF 100, FM 200; for Amazon Toys and Games WRMF 100, FM 50. Finally **BFwCF** [51] is an algorithm that takes into account interaction sequences between users and it uses the last common interaction to populate the candidate items list. In this evaluation we included the BFwCF variant that takes advantage of similarity weights per user and two time windows, left-sided and right-sided (Backward-Forward). BFwCF was trained using parameters from [51]: 100 neighbors, *indexBackWards* and *indexForwards* set to 5, normalization and combination realized respectively via *DummyNormalizer* and *SumCombiner*. Recommendations were computed with the implementation publicly provided by authors. In order to guarantee a fair evaluation, for all the time-based variants the β

³<https://github.com/sisinflab/recommenders>

coefficient was set to $\frac{1}{200}$ [222]. TimeSVD++ was trained using parameters used in [222].

Results Discussion. Results of experimental evaluation are shown in Figure 12.2 which illustrate nDCG (12.2a, 12.2b, 12.2c) curves for increasing number of top ranked items returned to the user. Significance tests have been performed for accuracy metrics using Student’s t-test and p-values and they result consistently lower than 0.05. By looking at Figure 12.2a we see that TimePop outperforms comparing algorithms in terms of accuracy on AmazonMovies dataset. We also see that algorithms exploiting a Temporal decay function perform well w.r.t. their time-unaware variants (User-kNN and Item-kNN) while matrix factorization algorithms (WRMF, TimeSVD++ and FM) perform quite bad. The low performance of MF algorithms is very likely due to the temporal splitting that makes them unable to exploit collaborative information. We may assume that the good performance of TimePop w.r.t. kNN algorithms are due to the adopted threshold, that emphasizes the popular items, and hence increases accuracy metrics values. Results for Amazon Toys and Games dataset are analogous to those computed for Amazon Movies. Results for Movielens additionally show that the high number of very popular items make neighborhood-based approaches perform similarly.

12.4 Conclusion

In this investigation we have presented TimePop, a framework that exploits local popularity of items combined with temporal information to compute top-N recommendations. The approach relies on the computation of a set of time-aware neighbors named Precursors that are considered the referring population for a user we want to serve recommendations. We compared TimePop against state-of-art algorithms showing its effectiveness in terms of accuracy despite its lower computational cost in computing personalized recommendations.

Chapter 13

Beyond similarities: Dissimilarity and Asymmetric Similarities

13.1 Introduction

Neighborhood-based approaches have been the first family of algorithms developed for collaborative filtering recommender systems. They identify similar users or items, and they provide users with a list of items they could be interested in by exploiting the degree of similarity. Though there have been around for many years, it has been shown that neighborhood-based approaches may perform better than latent model-based methods to solve the *top-N* recommendation problem [19, 119, 205, 273]. In the *top-N* recommendation task, the focus is on providing an accurate ranked list rather than minimizing the rating prediction error. Among neighborhood-based methods the best-known are user-kNN, item-kNN and Sparse Linear Methods (SLIM) [273]. User-based and item-based schemes have proven to be effective in different settings although they use the same logic behind the scenes. In details item-kNN and SLIM (which uses an item-based scheme) have

shown to outperform user-kNN to solve the *top-N* recommendation problem [100] and several algorithms have been proposed in the literature to enhance neighbors models like GLSLIM [100] and Weighted kNN-GRU4REC [195], taking advantage of personalized models and recurrent neural networks. Moreover, we also have approaches focusing on the injection of time in neighborhood models [51] and modeling similarities by directly optimizing the pair-wise preferences error [312]. All in all, under the hood, what makes neighborhood-based methods work is a similarity measure.

Several similarity measures have been proposed and used extensively such as Jaccard [129][307] and Tanimoto [113] coefficients, Cosine Vector similarity [36, 63, 13], Pearson Correlation [177], Constrained Pearson correlation [342], Adjusted Cosine similarity [334], Mean Squared Difference similarity [342], Spearman Rank Correlation [211], Frequency-Weighted Pearson Correlation [73], Target item weighted Pearson Correlation [37]. In the vast majority of cases, all these similarity measures are based on two assumptions:

1. the correlation between i and j is the same correlation between j and i (symmetry of similarity);
2. the correlation between two entities only captures how much they are similar to each other without taking into account their degree of dissimilarity.

To the best of our knowledge, a few works have been proposed in the past years related to asymmetric similarity, and they are mainly designed for a user-based scheme. Dissimilarity was first suggested in 1999 [381, 380] when Varian described the value of introducing diversity into search results. An Asymmetric User Similarity has been proposed in [261] where the authors underline that a similarity measure should distinguish between a user with a rich profile and a cold user. Thus, given two users u and v , they slightly modify the Jaccard index in order to consider exclusively the number of ratings of the current user (instead of the overall number of both users). In HYBRTyco [210] the similarity proposed by Millan [261] is combined with the Sørensen index [354]. HYBRTyco [210] is a hybrid recommender system which combines matrix factorization with an asymmetric similarity model

to realize a typicality-based collaborative filtering recommender system. The same approach is exploited in another asymmetric user similarity model [306] to feed a user-user similarity matrix that is then completed using a matrix factorization algorithm. Additionally, both of them provide an extension for the latter similarity measure based on explicit numerical feedbacks (ratings). Despite previous works are focused on the user-based scheme, we already underlined that item-kNN shows excellent performance in *top-N* recommendation task. Moreover, when the number of users exceeds the number of items, as in most of the cases, item-based recommendation approaches require much less memory and time to compute the similarity weights than user-based ones, making them more scalable. Due to these reasons, both approaches have been considered in this work.

In this study, we investigate the effect on recommendation accuracy when we go beyond the above two assumptions and define (and include) the concepts of dissimilarity and asymmetry in similarity measures. In our proposal, we start from a probabilistic interpretation of similarity to define symmetric and asymmetric dissimilarities. The dissimilarity measures are then combined with traditional similarity values using additive and multiplicative strategies. The experimental evaluation shows that our approach outperforms the non-dissimilarity-aware counterparts improving the accuracy of results or diversity or both.

The rest of the chapter is organized as follows: Section 13.2 presents the motivation behind our work and the proposed approach. Section 13.3 presents the evaluation protocol, metrics, datasets and performance of the method. Finally, in Section 13.4 concluding considerations are provided.

13.2 Dissimilarity in Recommendation

13.2.1 Motivation

The main idea behind our proposal is that symmetric similarity may not be sufficient to capture subtle interactions between items. We assume that representing the similarity through traditional measures can lead to imperfect results as important

information might not be properly considered. Let us consider some examples in an item-kNN scenario. Suppose we are dealing with a dataset containing rating data from the book domain on the following books:

| Title | Short name | Author | # Votes |
|----------------------|------------|-----------------|---------|
| A Game of Thrones | GoT | G. R. R. Martin | 100 |
| A Dance with Dragons | DwD | G. R. R. Martin | 10 |
| Shroud of Eternity | SoE | T. Goodkind | 120 |

By looking at the previous data, we see that both *A Game of Thrones* and *A Dance with Dragons* belong to the same saga *A song of ice and fire* and they are, respectively, the first and the fifth volume. *Shroud of Eternity* is the second volume of Nicci Chronicles' saga. We may assume that all the users who rated DwD also rated GoT, i.e., $U_{DwD} \subseteq U_{GoT}$. Analogously, since the topic of the book is mostly the same, we may assume that a number of readers of SoE also voted GoT, $U_{SoE} \cap U_{GoT} \neq \emptyset$. Suppose now that we have $U_{SoE} \cap U_{GoT} = 20$ and $U_{DwD} \cap U_{GoT} = 10$. If we compute the Jaccard similarity between the pairs SoE, GoT and DwD, GoT we have

$$\begin{aligned}
 JS(DwD, GoT) &= \frac{|U_{DwD} \cap U_{GoT}|}{|U_{DwD} \cup U_{GoT}|} = 0.1 \\
 JS(SoE, GoT) &= \frac{|U_{SoE} \cap U_{GoT}|}{|U_{SoE} \cup U_{GoT}|} = 0.1
 \end{aligned}$$

In our opinion, much relevant information has been lost in this simple example. The scenario is shown graphically in Figure 13.1.

It is clear that U_{DwD} is a proper subset of U_{GoT} and, on the contrary, there are many users in U_{SoE} that have not experienced GoT. This information is mostly lost in the computation of the similarity values even though a piece of this information is retained in the denominator of the Jaccard coefficient through the overall value of $U_{GoT} \cup U_{DwD}$ and $U_{GoT} \cup U_{SoE}$.

In order to clarify the reason why we do not consider this remaining information sufficient, let us consider a more formal description of the scenario. Many similarity measures (like Jaccard in the previous example) mainly rely on a value that denotes the similarity between two items normalized by their overall weight. This can be

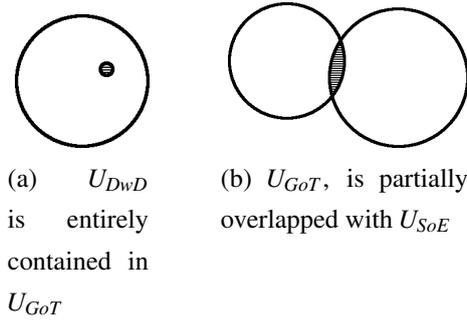


Figure 13.1: Representation of the motivation example

represented as a probability. Let $\mathcal{S}(u) = \{\langle i, r_{ui} \rangle \mid u \text{ rated } i \text{ with } r_{ui}\}$ be the user profile containing the pairs item-rating and $\mathcal{U}(i) = \{\langle u, r_{ui} \rangle \mid \langle i, r_{ui} \rangle \in \mathcal{S}(u)\}$ be the set of users that experienced i . If we consider Jaccard similarity we see it represents the following probability:

$$p^{JS}(i, j) = p(\langle u, r_{ui} \rangle \in \mathcal{U}(i) \wedge \langle u, r_{uj} \rangle \in \mathcal{U}(j)) \quad (13.1)$$

which can be read as the probability that u experienced both i and j . Under a probabilistic lens, we may define dissimilarity measures with two different probabilities:

- the probability that a generic user u experienced the item i but never experienced the item j
- the probability that a generic user u experienced the item j but never experienced the item i

Since we are introducing an asymmetric behavior in computing the similarity between i and i as well as between j and i , we see that the two probabilities have a different role while computing a similarity measure.

In a classical memory-based Item-kNN, $sim(i, j)$ is used to compare i to different j to find out which js are the most similar to i . In practical terms, we are interested in how much j is similar to i . If we focus on Figure 13.1a, we realize that DwD is more similar to GoT than the opposite situation. The reason behind this behavior is not directly related to the size of the involved sets but it depends on the probability that a user who experienced GoT did not experience DwD.

Though some interesting asymmetric similarities have been proposed in the last years [306, 210], to our knowledge, no one focused on this probability that represents a negative asymmetric dissimilarity.

13.2.2 Metrics

In this work, we propose a general asymmetric similarity model in which items i, j similarities are computed by taking into account the probability that users that experienced j never experienced i . The idea behind our work is preserving the core meaning of a specific similarity, applying a corrective factor encoding the dissimilarity we mentioned before.

We introduced this correction into two binary symmetric similarities: Jaccard and Sørensen index, and in the asymmetric variant to the Jaccard coefficient proposed in [261]. We tested this correction both as an additive and as a multiplicative factor.

For the sake of completeness, we reintroduce *Jaccard coefficient similarity (JS)* that, for a memory-based item-kNN model can be expressed as:

$$JS(i, j) = \frac{|\mathcal{U}(i) \cap \mathcal{U}(j)|}{|\mathcal{U}(i) \cup \mathcal{U}(j)|} \quad (13.2)$$

The probability that users who experienced j never experienced i can be modeled as the complementary probability of Equation (13.1) w.r.t. $\mathcal{U}(i)$ over $|\mathcal{U}(i) \cup \mathcal{U}(j)|$. This probability, we name *Jaccard Asymmetric Dissimilarity (JAD)* can be formulated as follows:

$$JAD(i, j) = \frac{|\mathcal{U}(j)| - |\mathcal{U}(i) \cap \mathcal{U}(j)|}{|\mathcal{U}(i) \cup \mathcal{U}(j)|} \quad (13.3)$$

Again, Equation (13.3) can be seen in terms of probability as

$$p^{JAD}(i, j) = p(\langle u, r_{uj} \rangle \in \mathcal{U}(j)) - p^{JS}(i, j)$$

We now propose to modify the original similarity by injecting the former negative correction weighted with a parameter λ that can be easily customized. The overall similarity, *Additive Adjusted Jaccard (AAJ)*, is then formulated as:

$$AAJ(i, j) = JS(i, j) - \lambda \cdot JAD(i, j) \quad (13.4)$$

where λ is a parameter that could depend on many factors such as the number of users and items in the dataset and the intensity of interactions between them.

Recalling that the overall formula should represent a degree of similarity between the two different items we defined the Multiplicative Adjusted Jaccard as the product of Jaccard similarity with the inverse of Jaccard Asymmetric Dissimilarity (IJAD). In other words, we use $\frac{1}{p^{JAD}(i,j)}$ as a corrective factor for $p^{JS}(i,j)$. We then define the *Multiplicative Adjusted Jaccard (MAJ)* as:

$$MAJ(i,j) = \frac{JS(i,j)}{JAD(i,j)} = JS(i,j) \cdot IJAD(i,j) \quad (13.5)$$

In the multiplicative variant, in order to avoid division by zero the minimum value for JAD is set to $\frac{1}{|\mathcal{U}(i) \cup \mathcal{U}(j)|}$. As mentioned in Section 13.2.1 a symmetric variant of Jaccard coefficient (named *Jaccard Symmetric Dissimilarity (JSD)*) can be used composing Equation (13.3) with the probability that a user u experienced the item i but never experienced the item j :

$$\begin{aligned} JSD(i,j) &= \frac{(|\mathcal{U}(j)| - |\mathcal{U}(i) \cap \mathcal{U}(j)|) + (|\mathcal{U}(i)| - |\mathcal{U}(i) \cap \mathcal{U}(j)|)}{|\mathcal{U}(i) \cup \mathcal{U}(j)|} \\ &= \frac{|\mathcal{U}(i)| + |\mathcal{U}(j)| - 2 \cdot |\mathcal{U}(i) \cap \mathcal{U}(j)|}{|\mathcal{U}(i) \cup \mathcal{U}(j)|} \end{aligned}$$

leading to the corresponding probability

$$p^{JSD}(i,j) = p^{JAD}(i,j) + p^{JAD}(j,i)$$

Thus the *Symmetric Additive Adjusted Jaccard (S-AAJ)* and *Symmetric Multiplicative Adjusted Jaccard (S-MAJ)* can be defined as follows:

$$\begin{aligned} S-AAJ(i,j) &= JS(i,j) - \lambda \cdot JSD(i,j) \\ S-MAJ(i,j) &= \frac{JS(i,j)}{JSD(i,j)} \end{aligned}$$

In order to test our idea, we applied all the variants previously introduced for Jaccard similarity to two popular similarity measures: *Asymmetric Jaccard Similarity (AJS)* and *Sørensen coefficient (SOR)*. All the derived variants are represented, respectively, in Table 13.1 and Table 13.2.

All the above metrics have been introduced having in mind an item-kNN approach but, without loss of generality, they can be applied to user-kNN model as well.

Table 13.1: Asymmetric Jaccard considered variants.

| Short name | Extended | Formula |
|-------------|--------------------------------------|---|
| AJS(i,j) | Asymm. Jaccard Similarity | $\frac{ \mathcal{U}(i) \cap \mathcal{U}(j) }{ \mathcal{U}(i) }$ |
| AJD(i,j) | Asymm. Jaccard Dissimilarity | $\frac{ \mathcal{U}(j) - \mathcal{U}(i) \cap \mathcal{U}(j) }{ \mathcal{U}(i) }$ |
| AAAJ(i,j) | Additive Adjusted Asymm. Jaccard | $AJS(i, j) - \lambda \cdot AJD(i, j)$ |
| MAAJ(i,j) | Multiplicative Adjus. Asymm. Jaccard | $AJS(i, j) \cdot IJAD(i, j)$ |
| S-AAAJ(i,j) | Symmetric AAAJ | $AJS(i, j) - \lambda \cdot \frac{ \mathcal{U}(i) + \mathcal{U}(j) - 2 \cdot \mathcal{U}(i) \cap \mathcal{U}(j) }{ \mathcal{U}(i) }$ |
| S-MAAJ(i,j) | Symmetric MAAJ | $AJS(i, j) \cdot IJSD(i, j)$ |

13.3 Experimental Evaluation

The experimental evaluation has been carried out on three publicly available datasets and with different values of k and λ .

Datasets. We evaluated the effectiveness of our approach on the three datasets shown in Table 13.3 belonging to different domains (Music, Books, and Movies). The `Last.fm` dataset [84] corresponds to transactions with Last.fm online music system released in HetRec 2011¹. It contains social networking, tagging, and music artist listening information from a set of 2K users. `LibraryThing` represents books ratings collected in the LibraryThing community website. It contains social networking, tagging, and rating information on a [1..10] scale. `Yahoo!Movies` (Yahoo! Webscope dataset `ydata-ymovies-user-movie-ratings-content-v1_0`)² contains movies ratings generated by Yahoo! Movies up to November 2003. It provides content, demographic, rating information, and mappings to `MovieLens` and `EachMovie` datasets.

¹<http://ir.ii.uam.es/hetrec2011/>

²http://research.yahoo.com/Academic_Relations

Table 13.2: Sørensen similarity considered variants.

| Short name | Extended | Formula |
|------------|--|--|
| SOR(i,j) | Sørensen Similarity | $\frac{ U(i) \cap U(j) }{ U(i) + U(j) }$ |
| ASD(i,j) | Asymm. Sørensen Dissimilarity | $\frac{ U(j) - U(i) \cap U(j) }{ U(i) + U(j) }$ |
| AAS(i,j) | Additive Adjusted Assym Sørensen | $SOR(i, j) - \lambda ASD(i, j)$ |
| MAS(i,j) | Multiplicative Adjusted Assym Sørensen | $SOR(i, j) \cdot IJAD(i, j)$ |
| S-AAS(i,j) | Symmetric AAS | $SOR(i, j) - \lambda \frac{ U(i) + U(j) - 2 \cdot U(i) \cap U(j) }{ U(i) + U(j) }$ |
| S-MAS(i,j) | Symmetric MAS | $SOR(i, j) \cdot IJSD(i, j)$ |

Table 13.3: Datasets statistics.

| Dataset | #Users | #Items | #Transactions | Sparsity |
|---------------|--------|--------|---------------|----------|
| Yahoo! Movies | 7642 | 11,916 | 221,367 | 99.76% |
| LibraryThing | 7279 | 37,232 | 2,056,487 | 99.24% |
| Last FM | 1850 | 11,247 | 59,071 | 99.72% |

Columns corresponding to #Users, #Items and #Transactions show the number of users, number of items and number of transactions, respectively, in each dataset. The last column shows the sparsity of the dataset.

Evaluation Protocol and Experimental Setting with Parameters tuning. The evaluation protocol we adopted in our experiments is *all unrated items* [356]. With this protocol, the recommendation list is computed from a candidate list given by the cartesian product between users and items minus the items each user experimented in the training set. We performed a temporal 64-16-20 hold-out split (when temporal information is available) retaining the last 20% of ratings as test set and 16% as validation set. We measured the performance by computing Precision@N ($Prec@N$) for *top-N* recommendation lists as accuracy metric. Precision has been computed on a per-user basis, and the returned results have been averaged. As

Precision needs relevant items to be computed, we set the relevance threshold to 8 over 10 for `LibraryThing` and `Yahoo!Movies`, and to 0 for `Last.fm` since in this latter no ratings are provided but the number of user-item transactions. We measured Diversity through *catalog coverage* (aggregate diversity in *top-N* list). The *catalog coverage*, also called diversity-in-top-N ($D@N$) [14], is measured by computing the overall number of different items recommended within the complete recommendation list. It represents the propensity of a system to recommend always the same items.

Baselines. We compared our approaches in both User-kNN and Item-kNN settings. The former finds the k -nearest user neighbors based on a similarity function and then exploits them to predict a score for each user-item pair. The latter is the item-based version of the k -nearest neighbors algorithm that uses the k -nearest items to compute the predictions. For both schemes we used the validation set to find the optimal hyperparameters. However, we are not interested in the algorithm itself but on the similarity measures that are used to compute neighbors and predictions. As baseline to compare with, we used both symmetric and asymmetric measures, namely, *Jaccard (JS)* and *Sorensen (SOR)* (for symmetric measures) and *asymmetric Jaccard (AJS)* and *asymmetric Jaccard weighted with the Sorensen Index (ASOR)* [210] (for asymmetric measures). For all the similarities that make use of λ we evaluated them varying λ in $\{0.2, 0.4, 0.6, 0.8\}$ whereas we considered a number of *Neighbors* varying in $\{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$. We ran the algorithms with all possible combinations, and we selected the best performing ones with respect to Precision@N. The best parameters for user-based and item-based schemes are represented, respectively in Table 13.4 and 13.5.

Table 13.4: Best parameters for User-kNN scheme

| Precision - P@10 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|------------------|----|-----------|-----|-----------|-----|-----------|-------|-----------|-------|-----------|------|-----|-----------|-----|-----------|-------|-----------|-------|-----------|-----|-----------|---|-----------|-----|-----------|---|-----------|-----|------|----|
| | JS | | AAJ | | MAJ | | S-AAJ | | S-MAJ | | ASOR | SOR | | AAS | MAS | S-AAS | | S-MAS | | AJS | AAAJ | | MAAJ | | S-AAAJ | | S-MAAJ | | P@10 | |
| Datasets | k | λ | k | λ | k | λ | k | λ | k | λ | k | k | λ | k | λ | k | λ | k | λ | k | λ | k | λ | k | λ | k | λ | k | P@10 | |
| LibraryThing | 50 | 0.2 | 100 | | 50 | 0.2 | 30 | | 100 | | 90 | | 50 | 0.2 | 50 | | 50 | 0.2 | 30 | | 100 | | 20 | 0.2 | 100 | | 40 | 0.2 | 50 | 20 |
| Yahoo | 90 | 0.4 | 100 | | 40 | 0.2 | 100 | | 90 | | 100 | | 90 | 0.2 | 40 | | 40 | 0.2 | 100 | | 90 | | 10 | 0.6 | 80 | | 30 | 0.2 | 100 | 70 |
| Last FM | 90 | 0.2 | 100 | | 100 | 0.2 | 50 | | 90 | | 90 | | 100 | 0.2 | 100 | | 100 | 0.2 | 50 | | 90 | | 30 | 0.2 | 100 | | 100 | 0.2 | 50 | 20 |

Performance of the proposed methods. Results in Table 13.6 show the per-

Table 13.5: Best parameters for Item-kNN scheme

| Precision - P@10 | | | | | | | | | | | | | | | | | | | | | | |
|------------------|----|-----------|-----|-------|-----------|------|-----|-----|-----------|-------|-------|-----|------|-----------|--------|-----------|-----|------|----|-----|----|----|
| | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ | | | | | | |
| Datasets | k | λ | k | k | λ | k | k | k | λ | k | k | k | k | λ | k | λ | k | P@10 | | | | |
| LibraryThing | 10 | 0.2 | 20 | 20 | 0.4 | 10 | 20 | 70 | 10 | 0.2 | 20 | 20 | 0.4 | 10 | 20 | 60 | 0.2 | 10 | 10 | 0.2 | 10 | 60 |
| Yahoo | 10 | 0.2 | 20 | 10 | 0.2 | 10 | 10 | 10 | 10 | 0.2 | 10 | 10 | 0.2 | 10 | 10 | 20 | 0.2 | 20 | 10 | 0.2 | 10 | 10 |
| Last FM | 10 | 0.2 | 20 | 30 | 0.2 | 20 | 10 | 10 | 10 | 0.2 | 30 | 30 | 0.2 | 20 | 10 | 10 | 0.2 | 40 | 20 | 0.2 | 70 | 10 |

Table 13.6: Comparison in terms of Precision and Aggregate Diversity for User-kNN scheme with best parameters

| Precision - P@10 | | | | | | | | | | | | | | | | |
|----------------------------|---------|-------------|----------------|-------------|----------------|----------------|----------------|---------|----------------|-------------|---------|---------|-------------|----------------|-------------|----------------|
| Datasets | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ |
| LibraryThing | 0.03025 | 0.04005 | 0.04675 | 0.00416 | 0.03102 | 0.03565 | 0.03031 | 0.04687 | 0.04687 | 0.00416 | 0.03088 | 0.03768 | 0.03871 | 0.04670 | 0.00429 | 0.03796 |
| Yahoo | 0.04165 | 0.05196 | 0.06284 | 0.03892 | 0.04149 | 0.04261 | 0.04151 | 0.06354 | 0.06354 | 0.03917 | 0.04170 | 0.05465 | 0.04234 | 0.06176 | 0.03705 | 0.05207 |
| Last FM | 0.02704 | 0.02163 | 0.02532 | 0.00601 | 0.02773 | 0.02927 | 0.02747 | 0.02532 | 0.02532 | 0.00609 | 0.02695 | 0.02901 | 0.02120 | 0.03039 | 0.00592 | 0.03107 |
| Aggregate Diversity - D@10 | | | | | | | | | | | | | | | | |
| Datasets | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ |
| LibraryThing | 2735 | 7330 | 3093 | 3341 | 1953 | 2687 | 2660 | 2999 | 2999 | 3337 | 1879 | 2519 | 7443 | 3037 | 3337 | 2834 |
| Yahoo | 698 | 2087 | 1126 | 2181 | 796 | 882 | 662 | 1074 | 1074 | 2179 | 732 | 974 | 1903 | 1147 | 2177 | 580 |
| Last FM | 1136 | 1665 | 1485 | 1156 | 1203 | 1283 | 1097 | 1446 | 1446 | 1154 | 1145 | 764 | 1748 | 889 | 1163 | 1025 |

Table 13.7: Comparison in terms of Precision and Aggregate Diversity for Item-kNN scheme with best parameters

| Precision - P@10 | | | | | | | | | | | | | | | | |
|----------------------------|----------------|----------------|--------------|---------|----------------|---------|----------------|---------|--------------|----------------|----------------|--------------|----------------|---------|---------|---------|
| Datasets | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ |
| LibraryThing | 0.08180 | 0.07449 | 0.03189 | 0.06482 | 0.08048 | 0.07835 | 0.07949 | 0.02462 | 0.02462 | 0.06111 | 0.07894 | 0.04556 | 0.08748 | 0.02256 | 0.06293 | 0.06097 |
| Yahoo | 0.05105 | 0.04976 | 0.00272 | 0.04966 | 0.05187 | 0.05013 | 0.05043 | 0.00196 | 0.00196 | 0.04915 | 0.05141 | 0.01598 | 0.05031 | 0.00175 | 0.04842 | 0.02437 |
| Last FM | 0.02146 | 0.02549 | 0.00326 | 0.02489 | 0.02120 | 0.02052 | 0.02069 | 0.00275 | 0.00275 | 0.02403 | 0.02077 | 0.00876 | 0.02910 | 0.00240 | 0.02506 | 0.01373 |
| Aggregate Diversity - D@10 | | | | | | | | | | | | | | | | |
| Datasets | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ |
| LibraryThing | 11945 | 12338 | 22097 | 11654 | 11551 | 10521 | 11774 | 21365 | 21365 | 11447 | 11450 | 18604 | 12899 | 13042 | 13604 | 16018 |
| Yahoo | 3262 | 5159 | 4644 | 3466 | 3091 | 3539 | 3340 | 4742 | 4742 | 4351 | 3725 | 6864 | 3946 | 3645 | 4514 | 5565 |
| Last FM | 2867 | 4394 | 3884 | 2920 | 2763 | 3389 | 2917 | 3648 | 3654 | 3446 | 3341 | 5150 | 3680 | 3078 | 3429 | 4447 |

formance of all algorithms with a user-based scheme. Concerning accuracy, it is clear that both asymmetric and symmetric multiplicative variants are the best-performing ones. MAJ and S-MAJ achieve good performance, outperforming JS. The same trend is shown between MAAJ, S-MAAJ, and AJS. The same behavior can be observed in the Sorensen algorithms block, in which MAS outperforms SOR in LibraryThing and Yahoo!Movies datasets with the only exception of Last.fm datasets. Concerning Diversity the proposed variants constantly outperform the base variants. It is worth to note that AAJ and AAAJ algorithms that can recommend much more items with a little loss of Precision. Table 13.7 shows

the results for the item-based scheme. Performance is much different here, and we can note that the multiplicative asymmetric variants have a worse behavior. However, the Last.fm results show that the additive variants outperform the base ones. We can observe the same behavior also in the whole asymmetric Jaccard similarity block for all the three datasets. The Jaccard similarity and the Sørensen block, for LibraryThing and Yahoo!Movies show no clear champion concerning the performance of JS/SOR and S-MAJ/S-MAS that result very similar. This behavior may be due to the tuning results that are very close to each other, and this probably prevented us from selecting the best parameters.

Experimental Setting with a fixed number of neighbors. In tables 13.4 and 13.5 we see that best values of λ and k are different depending on the adopted approach. Hence, we also tested the different algorithms with a fixed number of neighbors. In other words, we checked: How would the different algorithms perform if, i.e., we fix the number of neighbors? Once again we employed the *all unrated items* evaluation protocol to evaluate the methods. We performed a temporal 80-20 hold-out split retaining the last 20% of ratings as test set using temporal information when available.

Baselines. Even in this experiment, we compared our approaches with both User-kNN and Item-kNN settings considering, for all the algorithms, the number of neighbors fixed and set to $k = 80$. For all the similarities that make use of λ we evaluated them varying λ in 0.2,0.4,0.6,0.8. In Tables 13.8 and 13.9 we show the best results we obtained³. The best values regarding Precision and Aggregate Diversity are highlighted in bold. We computed significance tests for precision results, and we found they are statistically significant at the 0.05 level w.r.t. their respective baselines.

Performance of the proposed methods. Results in Table 13.8 show that our approach always outperforms baseline variants in the User-kNN scheme. In details **additive** asymmetric similarity and **multiplicative** asymmetric similarity significantly perform better than JS, SOR and ASOR for all three dataset. Among

³The complete results are publicly available at <https://github.com/sisinflab/The-importance-of-being-dissimilar-in-Recommendation>.

Table 13.8: Comparison in terms of Precision and Aggregate Diversity for User-kNN scheme

| Precision - P@10 | | | | | | | | | | | | | | | | | | | | | | |
|------------------|--------|-----------|--------|---------------|-----------|--------|--------|--------|-----------|-------|--------|---------------|-----------|---------------|--------|-----------|------|--------|---------------|-----|--------|--------|
| | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ | | | | | | |
| Datasets | P@10 | λ | P@10 | P@10 | λ | P@10 | P@10 | P@10 | λ | P@10 | P@10 | P@10 | λ | P@10 | P@10 | λ | P@10 | | | | | |
| LibraryThing | 0.0363 | 0.2 | 0.0582 | 0.0627 | 0.4 | 0.0364 | 0.0010 | 0.0394 | 0.0363 | 0.2 | 0.0586 | 0.0629 | 0.2 | 0.0405 | 0.0364 | 0.0364 | 0.2 | 0.0558 | 0.0603 | 0.2 | 0.0057 | 0.0375 |
| Yahoo | 0.0437 | 0.4 | 0.0561 | 0.0676 | 0.2 | 0.0403 | 0.0433 | 0.0442 | 0.0438 | 0.4 | 0.0573 | 0.0687 | 0.2 | 0.0568 | 0.0435 | 0.0607 | 0.6 | 0.0529 | 0.0664 | 0.2 | 0.0378 | 0.0529 |
| Last FM | 0.0164 | 0.4 | 0.0242 | 0.0257 | 0.2 | 0.0037 | 0.0160 | 0.0215 | 0.0166 | 0.4 | 0.0241 | 0.0253 | 0.2 | 0.0275 | 0.0164 | 0.0248 | 0.2 | 0.0228 | 0.0323 | 0.2 | 0.0037 | 0.0242 |

| Aggregate Diversity - D@10 | | | | | | | | | | | | | | | | | | | | | | |
|----------------------------|------|-----------|-------------|-------------|-----------|-------------|------|------|-----------|-------|-------------|-------------|-----------|------|--------|--------|-----|-------------|------|-----|-------------|------|
| | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ | | | | | | |
| Datasets | D@10 | λ | D@10 | D@10 | λ | D@10 | D@10 | D@10 | λ | D@10 | D@10 | D@10 | λ | D@10 | D@10 | D@10 | | | | | | |
| LibraryThing | 2136 | 0.2 | 7367 | 2406 | 0.4 | 2246 | 873 | 2819 | 2083 | 0.2 | 7330 | 2292 | 0.2 | 5528 | 2171 | 1299 | 0.2 | 7501 | 2060 | 0.2 | 3331 | 1486 |
| Yahoo | 734 | 0.4 | 2070 | 835 | 0.2 | 2187 | 825 | 936 | 695 | 0.4 | 2048 | 786 | 0.2 | 1114 | 784 | 451 | 0.6 | 1773 | 717 | 0.2 | 2175 | 587 |
| Last FM | 1449 | 0.4 | 1345 | 1654 | 0.2 | 1203 | 1460 | 1433 | 1324 | 0.4 | 1324 | 1625 | 0.2 | 1410 | 1420 | 774 | 0.2 | 1665 | 1077 | 0.2 | 1207 | 941 |

Table 13.9: Comparison in terms of Precision and Aggregate Diversity for Item-kNN scheme

| Precision - P@10 | | | | | | | | | | | | | | | | | | | | | | |
|------------------|--------|-----------|---------------|--------|-----------|---------------|--------|---------------|-----------|-------|---------------|--------|-----------|---------------|--------|-----------|------|---------------|--------|-----|---------------|--------|
| | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ | | | | | | |
| Datasets | P@10 | λ | P@10 | P@10 | λ | P@10 | P@10 | P@10 | λ | P@10 | P@10 | P@10 | λ | P@10 | P@10 | λ | P@10 | | | | | |
| LibraryThing | 0.0869 | 0.2 | 0.0949 | 0.0451 | 0.4 | 0.0822 | 0.0944 | 0.1019 | 0.0815 | 0.2 | 0.0914 | 0.0374 | 0.4 | 0.0826 | 0.0901 | 0.0598 | 0.2 | 0.1011 | 0.0303 | 0.4 | 0.0830 | 0.0792 |
| Yahoo | 0.0331 | 0.2 | 0.0510 | 0.0016 | 0.2 | 0.0535 | 0.0373 | 0.0447 | 0.0297 | 0.2 | 0.0504 | 0.0014 | 0.2 | 0.0531 | 0.0352 | 0.0046 | 0.4 | 0.0509 | 0.0012 | 0.2 | 0.0527 | 0.0083 |
| Last FM | 0.0141 | 0.2 | 0.0248 | 0.0036 | 0.2 | 0.0230 | 0.0158 | 0.0127 | 0.0124 | 0.4 | 0.0195 | 0.0031 | 0.2 | 0.0230 | 0.0155 | 0.0036 | 0.2 | 0.0237 | 0.0032 | 0.2 | 0.0229 | 0.0068 |

| Aggregate Diversity - D@10 | | | | | | | | | | | | | | | | | | | | | | |
|----------------------------|------|-----------|-------------|--------------|-----------|-------------|-------|-------------|-----------|-------|-------------|--------------|-----------|-------|--------|--------------|-----|-------------|-------|-----|-------|-------|
| | JS | AAJ | MAJ | S-AAJ | S-MAJ | ASOR | SOR | AAS | MAS | S-AAS | S-MAS | AJS | AAAJ | MAAJ | S-AAAJ | S-MAAJ | | | | | | |
| Datasets | D@10 | λ | D@10 | D@10 | λ | D@10 | D@10 | D@10 | λ | D@10 | D@10 | D@10 | λ | D@10 | D@10 | D@10 | | | | | | |
| LibraryThing | 9745 | 0.2 | 12004 | 17510 | 0.4 | 12304 | 10249 | 10399 | 9556 | 0.2 | 11727 | 16557 | 0.4 | 13306 | 9998 | 17361 | 0.2 | 11221 | 11132 | 0.4 | 12315 | 14819 |
| Yahoo | 3103 | 0.2 | 3718 | 3463 | 0.2 | 4315 | 2541 | 3346 | 3184 | 0.2 | 3935 | 3397 | 0.2 | 3541 | 2475 | 2948 | 0.4 | 4092 | 2334 | 0.2 | 3756 | 2295 |
| Last FM | 3508 | 0.2 | 3538 | 3486 | 0.2 | 3462 | 2911 | 3779 | 3265 | 0.4 | 305 | 3175 | 0.2 | 3509 | 2887 | 2992 | 0.2 | 3935 | 2644 | 0.2 | 3692 | 3457 |

these two variants of similarity, the multiplicative variant is the best-performing one. Quite interestingly, modifying AJS, which is asymmetric in its inner nature with our asymmetric dissimilarity factor leads to an improvement irrespective of the considered dataset. It is worth noticing that, other than the accuracy improvements, aggregate diversity also increases due to the dissimilarity injection. In details, the asymmetric additive variant achieves the best results and triples catalog coverage values for LibraryThing and Yahoo!Movies.

Table 13.9 shows Precision and Catalog Coverage results for an item-based scheme. Obtained results are quite interesting for many reasons. First of all, it is clear that the same similarities can lead to very different results depending on the adopted scheme. In particular, asymmetric Jaccard (AJS) performs very badly for the item-kNN algorithm. Under the dissimilarities perspective, we have the same behavior, and the multiplicative approach performs badly. Quite surprisingly, the

additive version can always outperform the base variants. This suggests that adopting an additive strategy for item-kNN may lead to better results. This may be due to the wide number of items pairs without any common user. Focusing on additive symmetric and asymmetric similarities we can note that aggregate diversity results reflect the same improvements observed in accuracy values. The only case that appears to behave differently is AAAJ that registered a catalog coverage lower than AJS. This happens as we considered the best performing λ for precision. In case of $\lambda \in \{0.4, 0.6\}$ we obtain aggregate diversity values of 16,205 and 18,071, respectively, with precision results constantly higher than AJS (0.09374 and 0.08820). We may observe another interesting pattern on the `Yahoo!Movies` row: the symmetric version outperforms the asymmetric one. This could be due to some dataset's characteristic. By looking at the data in Table 13.3 we see that the ratio of the number of items to the number of users is much higher in `LibraryThing` and `Last.fm` (≈ 5 and ≈ 6) with respect to `Yahoo!Movies` (≈ 1.5). This suggests that the more the ratio is, the more is convenient to adopt an asymmetric scheme. However, this consideration needs to be further investigated.

13.4 Conclusion and Future Work

In this work, we propose a method to improve the performance of neighborhood-based models, by capturing subtle interactions between users and items, which cannot be appreciated using a traditional similarity measure. We defined a dissimilarity measure, that can be used combined with traditional user-based and item-based schemes. The proposed approach takes into account the single asymmetric components, leading to an improvement in both precision and aggregate diversity results. We performed a comparative experimental evaluation using three well-known datasets, varying the tuning parameter λ and k . Experiments show that our approach outperforms competing algorithms, denoting the usefulness of incorporating symmetric and asymmetric dissimilarity in neighborhood-based models. We are currently working on an extension of our idea that also takes into account user ratings and not just set-based measures. As a further extension, we are also interested in

making the approach even more personalized by weighting dissimilarity with user-centered values of λ .

Chapter 14

An investigation over hyperparameters tuning: A Discriminative Power perspective

14.1 Introduction

Recommenders Systems now play an important role in the lives of users. These systems avoid massive data overwhelm users and help them in finding a path to relevant information [317]. To enhance the expressiveness of the models or to improve the learning phase, these models can be equipped with a special class of parameters, named hyperparameters. Since many recommender systems come with one or more hyperparameters, the goodness of the system depends on how these parameters are selected. Although several strategies are available [61, 93, 350, 252, 188], the choice of the metric to evaluate them is not manifest. Are there particular measures that are well-suited for hyperparameters tuning? Are the changes in the metric's values significant or they are fundamentally originated by chance?

Up to the Netflix prize [52], the research community widely considered the recommendation problem as a rating prediction task [425, 363]. Consequently, the optimization goal was the minimization of the prediction error [178, 339]. However, in real recommender systems applications, only a small subset of relevant items are provided to users [178]. Indeed, several studies acknowledged that rating prediction optimization was not able to produce the optimal top-N recommendation lists [256]. Recommendation problem was hence re-defined as a top-N recommendation task [108], in which the optimization goal shifted to items ranking.

From this new perspective, many Information Retrieval metrics came to play to evaluate recommender systems. After decades, accuracy is still broadly considered as the key element in evaluation. Nonetheless, new dimensions as novelty and diversity of recommendation [91, 377, 186] became progressively important. In compliance with the purpose of the system, accuracy, novelty, and diversity metrics are used both to evaluate the recommender and tune the hyperparameters. Although recommender systems are evaluated using an online or offline setting, hyperparameters are usually tuned in an offline setting [339]. In this setting, to evaluate the competing models (or hyperparameters candidate values), past users interactions are split adopting distinct strategies like Hold-Out [245, 107, 364] and k-Folds Cross-Validation (CV) [111, 219, 189]. In the former, the training set is split into two further sets: training, and validation set. In the latter, data is divided into k sets, retaining in rotation one of them as the validation set and the remaining ones as the training set.

The choice of hyperparameters values to test has also been deeply investigated. Among all the most adopted techniques are Random [53, 305, 54], Bayesian optimization [353, 76], and Grid Search [150, 55, 187]. Nevertheless, even though a recommender system's hyperparameter tuning is wisely designed to achieve more robust results, some aspects need further investigation. For instance, the behaviour of different metrics when varying the folds is still an almost unexplored field. As an example, if the metric is not able to capture significant differences when different values of parameters are set, that metric is not the ideal one to tune hyperparameters. A recent study [375] depicted a new interesting methodology to establish whether

a metric is discriminative, robust and the authors also performed a metric-to-metric comparison. Even though the authors designed it to measure the robustness of metrics to changes in the cut-off (previously explored in Information Retrieval in [251]), the overall procedure, along with the Discriminative Power and Robustness definition inspired us to propose a new procedure to evaluate metrics and study the metric-hyperparameter combined behaviour.

Our experiments on two well-known datasets show that *Precision* and *normalised Discounted Cumulative Gain* are the most discriminative accuracy metrics to use in model selection. We additionally show that also considering the standard deviation variations over the folds these metrics still behave better than *Recall* and *Mean Reciprocal Rank*. We show that *Expected Free Discovery* and *Expected Popularity Complement* metrics are discriminative and are significantly influenced by changes in hyperparameters values. Finally, if more than one hyperparameters have to be set, it is possible to investigate how changes in the specific parameters affect the considered metric.

In this research line we have focused on:

1. a study on the discriminative power of accuracy and novelty metrics for the k-Folds Cross-Validation hyperparameter tuning for three well-known collaborative models;
2. a procedure for models with more than one parameter to check if variations in one of the parameters were particularly relevant for accuracy of recommendation;
3. a study on the impact of "Number of latent factors", "Number of iterations", and "learning rate" on accuracy of recommendation for BPR-MF.

The remainder of the chapter is structured as follows: Section 14.2 introduces the setting of our experiments describing the adopted methodologies; then we focus on the discriminative power of metrics and their variations across folds; in Section 14.5 we study separately the hyperparameters of BPR-MF. Conclusions are drawn at the end of the chapter.

14.2 Experimental Settings

Discriminative Power (DP) is a metric proposed by [375] to measure the discriminative power of an evaluation metric over a set of competing algorithms. The procedure was originally presented by [326] in 2006. Given a metric, a dataset, and a set of recommender systems, the authors perform a statistical test considering all the possible system pairs. The obtained *p-values* can be sorted by decreasing value and plotted. The resulting curve is the *p-values* curve of the considered metric. Analogously, the corresponding *p-values* curve can be obtained for each considered metric. Since lower values of *p-values* denote statistical differences between system pairs, the metric with the lower area under the curve can be considered as the most discriminative. DP consists of the summation of all the *p-values* for a given metric, and it can be considered as an approximation of the area under the curve for that metric. Interestingly, in [375], the authors extend the idea of competing algorithms to a set composed of instances of the same algorithm considering different cut-off values. However, this idea can be further extended to consider a set of instances of the same algorithm considering different hyper-parameter values. This idea is the starting point of our work.

| Dataset | Users | Items | Ratings | Sparsity |
|---------------|-------|--------|---------|----------|
| MovieLens-1M | 6040 | 3706 | 1000209 | 95.53% |
| Amazon Movies | 16141 | 111537 | 858314 | 99.95% |

Table 14.1: Datasets statistics

Datasets. To conduct our study we exploited two different datasets in the Movies domain: Amazon Movies¹ and MovieLens-1M. Both datasets contain explicit ratings on a 1-5 scale. For Amazon Movies dataset we removed users with less than 20 interactions and items with less than 25 votes. Then we sampled the resulting dataset to generate a subset that preserves the original distribution of data [253, 254]. Experiments were conducted on a dedicated server equipped with an *Intel Xenon* with 32 cores, and 256GB RAM memory. The sampling step was necessary to perform

¹<https://snap.stanford.edu/data/web-Movies.html>

experiments in a reasonable time. The characteristics of datasets are reported in Table 14.1.

Over the years, several splitting methodologies have been proposed [32, 51, 339]. We decided to split our data in training and test set using a temporal Hold-Out splitting [339]. For each user, the first 80% of the interactions are considered as the training set, whereas the remaining 20% is used as the test set. The training set is further divided using a 10-Folds Cross-Validation.

Evaluation protocol. Offline evaluation in recommender systems is a well-studied field. To evaluate the approaches we decided to use "All Unrated Items" protocol [356], in which the set of candidate items for user u is composed of all items i not rated in u 's training set. Many metrics make use of binary relevance. Since we use datasets with explicit ratings, a relevance threshold τ [83] should be set to establish whether the items in each user's test set are relevant or not. We set τ to 4 for both datasets: only items with a rate above τ are considered as relevant during evaluation.

Algorithms. To study the hyperparameters influence on the discriminative power of metrics we decided to consider two distinct families of algorithms: Neighborhood models, and Matrix Factorization models. For the former, we considered both the *User-based* and *Item-based* scheme [333, 31]. For the latter, *BPR-MF* [312] was considered. In BPR-MF the classic MF model is optimized adopting the Bayesian Personalized Ranking Criterion, a well-known pairwise "*learning to rank*" algorithm.

Metrics. We decided to study the discriminative power of some widely used metrics along two dimensions: Accuracy and Novelty. In order to evaluate the accuracy of the algorithms, we measured *normalised Discount Cumulative Gain@N* ($nDCG@N$) [198], *Precision* ($Prec@N$), *Recall* ($Rec@N$), and *Mean Reciprocal Rank* ($MRR@N$). To evaluate Novelty, we decided to measure *Expected Free Discovery* ($EFD@N$) [377], and *Expected Popularity Complement* ($EPC@N$) [91]. These metrics were computed per user to perform the *Student's t* statistical test. The metrics values and the overall mean was computed using the RankSys framework [91].

Grid Search. To study the DP of the metrics for the different algorithms, we conducted a grid search exploration. This procedure is very common for hyperparameters tuning. However, based on how much exhaustive this search is, the operation can be time and space consuming. For this reason, we needed to determine the boundaries of this grid. The number of hyperparameters we decided to explore was 1 for Neighborhood models (the *number of Neighbors*), and 3 for Matrix Factorization (*latent factors, iterations, learning rate*). For Matrix Factorization we assumed user and item regularization to be dependent on learning rate, with a scale factor of respectively $1/20$ and $1/200$. We computed the values of the grid exploiting an exponential function with base 2 [56, 117]. To determine the evolution of latent factors, we used as an exponent for our function values within the range $[3.321, 10.821]$ with a step of 0.5. The same procedure and the same step have been used to define all the hyperparameters values. The difference basically consists of the considered range, chosen coherently with literature. For the number of iterations, we considered a range of $[0, 7]$ as the exponent. Finally, for the learning rate the exponent is in the range $[-2.3219, -16.3219]$. This choice led to learning values in the range $[0.00001220726897, 0.2000038948]$ considering 5 orders of magnitude. Summing up, for Matrix Factorization we had a grid of dimensions $15 \times 15 \times 14$. This grid generates 3150 different configurations of hyperparameters for each fold. Since the exploration on Amazon would have required several months we extracted a sub grid ($5 \times 3 \times 3$) with the best value for each hyper-parameter (considering $nDCG@10$) ± 2 neighbors in the grid. The overall hyperparameters values are depicted in Table 14.2.

14.3 Discriminative Power of metrics on Hyperparameters

The discriminative power (DP) of each metric using hyper parameters depicted in Table 14.2 can be studied exploiting the same strategy proposed in [375]. We compute the *p-value* across all folds in our k-fold cross-validation setting. For each algorithm, we generate all possible combinations of pairs of hyperparameters and

| Latent factors | Iterations | Learning rate | Nearest neighbors |
|----------------|------------|---------------|-------------------|
| 10 | 1 | 0.20000389 | 10 |
| 14 | 2 | 0.10000195 | 14 |
| 20 | 3 | 0.05000097 | 20 |
| 28 | 4 | 0.02500049 | 28 |
| 40 | 6 | 0.01250024 | 40 |
| 57 | 8 | 0.00625012 | 57 |
| 80 | 11 | 0.00312506 | 80 |
| 113 | 16 | 0.00156253 | 113 |
| 160 | 23 | 0.00078127 | 160 |
| 226 | 32 | 0.00039063 | 226 |
| 320 | 45 | 0.00019532 | 320 |
| 452 | 64 | 0.00009766 | 452 |
| 640 | 91 | 0.00004883 | 640 |
| 905 | 128 | 0.00002441 | 905 |
| 1809 | | 0.00001221 | 1809 |

Table 14.2: Hyperparameters grid

we randomly take 25 combinations. Thus, for these pairs, we compute the p -values for each fold. The p -values of the paired statistical tests are sorted by decreasing value and the corresponding values are averaged over the folds. For memory-based algorithms, the pairs correspond to pairs of systems with a different number of nearest neighbors. Even for BPR-MF, these pairs are pairs of systems. However, each

| MovieLens | EFD@N | EPC@N | MRR@N | nDCG@N | Prec@N | Rec@N |
|-----------------|-------|-------|--------------|--------------|--------------|--------------|
| Item-kNN | 1.884 | 1.840 | 1.766 | 1.710 | 1.855 | <u>2.385</u> |
| User-kNN | 1.688 | 1.576 | <u>2.196</u> | 1.665 | 1.869 | 1.993 |
| BPR-MF | 0.875 | 0.776 | 0.803 | 0.594 | 0.585 | <u>1.314</u> |
| Amazon | EFD@N | EPC@N | MRR@N | nDCG@N | Prec@N | Rec@N |
| Item-kNN | 0.188 | 0.188 | 0.188 | 0.213 | 0.161 | <u>0.281</u> |
| User-kNN | 4.738 | 5.717 | <u>6.646</u> | 6.310 | 5.474 | 6.281 |
| BPR-MF | 3.771 | 3.841 | <u>4.518</u> | 3.989 | 3.289 | 4.434 |

Table 14.3: Metrics Discriminative Power.

system is defined by a triple: *Latent Factors*, *Iterations*, *Learning Rate*. The sorted and averaged values can be plotted and they correspond to the averaged p -values curve. For each Accuracy and Novelty metrics, these plots are depicted in Fig-

ure 14.1. For each metric, the corresponding curve gives us an intuition on how that metric is discriminative with respect to the evolution of hyperparameters. Table 14.3 summarises the DP values of the metrics respectively on MovieLens and Amazon. We used **bold** to highlight the best-performing metric and underline to highlight the worst one. On MovieLens the trends of accuracy metrics seem to be clear: $nDCG@N$ always performs better than the others while $MRR@N$ and $Rec@N$ seem to be the worst metrics. It is interesting to notice that Novelty metrics achieve good DP values. This could be a signal that changes in parameters values lead to significant differences in terms of Novelty. On Amazon, the Best Accuracy metric is definitely $Prec@N$, while, even here, for User-kNN and BPR-MF the worst metric is $MRR@N$. In [375], it is suggested that this kind of behaviour can be due to the complexity of the metric. We agree with the authors and we reckon that this could also be due to some characteristics of the dataset, like the average number of rating per user (lower than MovieLens's), the sparsity of the dataset (higher than MovieLens's), and the low average number of ratings per item.

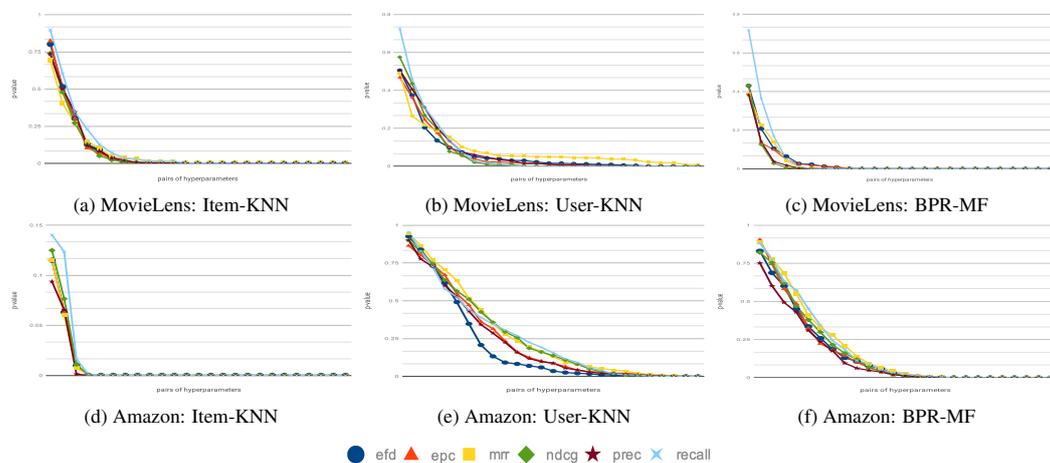


Figure 14.1: Discriminative Power of Accuracy and Novelty metrics

14.4 Metrics Confidence

In the previous section, we compared the Discriminative Power of different metrics and we found that $nDCG@N$ and $Prec@N$ are two good choices to select the best hyperparameters value for Neighborhood-based models and BPR-MF. In details, we consider the best hyperparameters value as the value which ensures the best performance with respect to the most discriminative metric. However, since we computed the averaged p -values curves across 10 Folds, and hence the averaged DP, it is still possible that these metrics could be much less discriminative on some folds. If this is true, the choice of an elected metric to conduct hyperparameters learning could be argued. For this reason, now we study the variations across different folds of the metrics p -values. Given the sorted p -values for each fold, we computed the standard deviation for each ordered pair across folds. These values can be exploited to define two additional p -values curves, which represent reasonable boundaries of p -values for each metric. Moreover, it is possible to compute the corresponding DP values, for each metric \pm the standard deviation. This could give us an intuition of the metric's DP robustness across folds. Table 14.4 shows the results for respectively Amazon and Movielens dataset. On Amazon, we may notice the good performance of $Prec@N$ with Item-kNN model. Although Amazon is a very sparse dataset with a large number of items, $Prec@N$ is able to capture significant differences between similar models with a different number of neighbours. Moreover, if we observe the DP value considering the standard deviation, this extreme value is still very close to the DP of the worst metric. This behaviour is also present on Movielens, for both Item-kNN and BPR-MF. We considered the worst scenario in which we added the standard deviation value to each pair in comparison.

However, it seems clear that the metrics chosen with the previous procedure show good performance across the different folds.

| MovieLens | Best metric | Worst metric | Best avg | Best + Std Dev | Worst avg |
|-----------------|-------------|--------------|----------|----------------|-----------|
| Item-kNN | nDCG@N | Recall@N | 1.710 | 2.940 | 2.385 |
| User-kNN | nDCG@N | MRR@N | 1.665 | 2.958 | 1.704 |
| BPR-MF | Prec@N | Recall@N | 0.585 | 1.126 | 1.314 |

| Amazon | Best metric | Worst metric | Best avg | Best + Std Dev | Worst avg |
|-----------------|-------------|--------------|----------|----------------|-----------|
| Item-kNN | Prec@N | Recall@N | 0.161 | 0.287 | 0.281 |
| User-kNN | Prec@N | MRR@N | 5.474 | 7.410 | 6.646 |
| BPR-MF | Prec@N | MRR@N | 3.289 | 4.689 | 4.518 |

Table 14.4: Metrics Discriminative Power deviation.

| Latent factors | MovieLens | | | | | | | | | | | | | | | | | Amazon | | | | |
|----------------------|-----------|--------------|--------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|--------------|---------|--------------|--------------|--------------|-------|--|--|
| | 10 | 14 | 20 | 28 | 40 | 56 | 80 | 113 | 160 | 226 | 320 | 452 | 640 | 905 | 1809 | 113 | 160 | 226 | 320 | 452 | | |
| DP | 0.703 | 0.664 | 0.602 | 0.512 | 0.592 | 0.707 | 0.575 | 0.535 | 0.635 | 0.639 | 0.572 | 0.499 | 0.303 | 0.565 | 0.290 | 2.775 | 2.681 | 2.348 | 2.127 | 2.542 | | |
| Iterations | 1 | 2 | 3 | 4 | 6 | 8 | 11 | 16 | 23 | 32 | 45 | 64 | 91 | 128 | | 64 | 92 | 128 | | | | |
| DP | 1.633 | 0.796 | 0.707 | 0.939 | 0.978 | 0.781 | 1.037 | 1.338 | 1.045 | 1.210 | 1.222 | 1.196 | 1.008 | 0.939 | | 4.106 | 3.299 | 2.357 | | | | |
| Learning rate | 0.20000 | 0.10000 | 0.05000 | 0.02500 | 0.01250 | 0.00625 | 0.00312 | 0.00156 | 0.00078 | 0.00039 | 0.00019 | 0.00009 | 0.00004 | 0.00002 | 0.00001 | 0.20000 | 0.10000 | 0.05000 | | | | |
| DP | 0.765 | 0.702 | 1.918 | 2.667 | 2.131 | 2.429 | 2.811 | 1.552 | 1.044 | 0.803 | 1.083 | 1.225 | 0.950 | 2.335 | 3.548 | 3.941 | 3.423 | 4.803 | | | | |

Table 14.5: DP w.r.t. hyperparameters evolution

14.5 Dominant Hyperparameter

In the previous section, we mainly focused on the Discriminative Power of the metrics to find out the best metric for hyper-parameter tuning taking into account the recommendation model and the considered dataset. In this section, we fix the metric to study the specific hyperparameters. Differently from nearest neighbors models, in BPR-MF we decided to explore three different hyperparameters: *number of latent factors*, *number of iterations*, *learning rate*. Usually, during the tuning phase, these dimensions are handled in the same way. Indeed, irrespective of the adopted search strategy, all the hyperparameters are equally important and should be explored. However, it is straightforward that one or more hyperparameters changes could influence more the accuracy of the provided recommendation list. This led us to pose two additional research questions:

- Is there one or more hyperparameters that affect more the accuracy of recommendations?
- Could be established a procedure to check if different values for a specific hyperparameter can lead to significant differences?

To answer these research questions, we decided to analyze the three hyperparameters of BPR-MF separately. In details, for a certain parameter, we want to define a procedure to check if different values of that hyperparameter lead to systems which show significant differences in accuracy of recommendation. Let us suppose we fix the metric and the number of latent factors: we still have two other parameters that can vary. Similarly to the procedure defined in the previous sections, we computed all the possible combinations of the remaining hyperparameters. From the set of combinations, we randomly chose 25 pairs of combinations. We recall that a pair of combinations corresponds to a pair of systems that share the number of latent factors, and differ in the number of iterations and learning rate. Now we compute the *p-values* of all these pairs and we order them by decreasing value. These values correspond to a *p-values* curve which is peculiar for the considered metric and number of latent factors. Consequently, for the curve, the corresponding Discriminative Power can be computed. This procedure can be repeated to analyze the discriminative power of various values of latent factors parameter. Thus, the whole procedure can be repeated to analyze the remaining hyperparameters. The results of this study is depicted in Tables 14.5. We used **bold** to highlight the best DP value for each hyperparameter analysis. As suggested in [375], the results between the two tables are not comparable. However, for both datasets, the *number of latent factors* seems to be the dimension on which variations in hyperparameter value lead to significant differences in recommendation accuracy. Moreover, for Movielens dataset, the DP value is much lower than the best values for "*Number of iterations*" and "*Learning rate*" hyper-parameter analysis. This clearly suggests that "*Number of latent factors*" is dominant with respect to the other hyperparameters. "*Learning rate*" dimension shows a different behaviour on the two datasets: on Movielens it shows big variations in terms of DP values, while on Amazon it shows oscillating performance. Finally, the DP values denote that the choice to conduct the study on the sub-grid was a reasonable choice.

14.6 Conclusions and Future Work

In this work, we explored the behavior of accuracy and novelty metrics in response to changes in hyperparameters values (we focused on the k-Folds Cross-Validation hyperparameter tuning). We found that nDCG@N and Precision@N represent a good choice for hyperparameters tuning for Neighborhood-based models and BPR-MF. Novelty metrics also show good DP values suggesting that these metrics are very sensitive to changes in hyperparameters values. We proposed a general procedure for models with more than one parameter to check if variations in one of the parameters were particularly relevant for accuracy of recommendation. We explored separately the BPR-MF hyperparameters and we found that the number of latent factors is dominant with respect to the learning rate and the number of iterations. Following this research direction, we want to explore other well-known algorithms and datasets to check if our findings can be further generalized.

Chapter 15

Generalized Cross-Entropy for Fairness Evaluation

15.1 Introduction

The use of recommender systems (RS) has exploded over the last decade, mostly due to their huge business value. According to the statistics revealed by Netflix, 75% of the downloads and rentals come from their recommendation service. This is a clear mark of the strategic importance of such a service in several companies [3, 196]. The success of RS is commonly measured by how well they are capable of making accurate predictions, i.e., items that users will likely interact with, purchase, or consume. Hence, the main effort of the research community over the last decade has been devoted to improve the *utility of recommendations* often measured in terms of effectiveness as well as to address beyond-accuracy aspects (e.g., novelty or diversity).

Models based on collaborative filtering (CF) (e.g., pureSVD, SVD++, WRMF, MMMF, SLIM, NeuralCF) lie at the core of most real-world recommendation en-

gines due to their state-of-the-art recommendation quality. In addition, a growing number of research works have leveraged different types of contextua information or external knowledge sources, such as knowledge bases/graphs, multimedia, user-generated tags and reviews among others, as additional information beyond the user-item interaction matrix to *further enhance* the final utility/quality of recommendation.

While recommendation models have reached a remarkable level of maturity in terms of effectiveness/performance in many application scenarios, at the same time, concerns have been recently raised on *fairness* of the recommendation models. As a matter of fact, recommendation algorithms, like other machine learning algorithms, are prone to imperfections due to algorithmic biases or biases in data. As stated by Barocas et al. [39] “*data can imperfect the algorithms in ways that allow these algorithms to inherit the prejudices of prior decision makers*”. Since RS assist users in many decision-making and mission-critical tasks such as medical, financial, or job-related ones [382, 355], unfair recommendation could have far-reaching consequences, impacting peoples lives and putting minority groups at a major disadvantage. This is particularly relevant for recommender systems as they are machine learning applications that directly interact with users, which are not necessarily experts in the field, and provide their outcome to them.

In the past, the notion of unfair recommendation was often associated with a non-uniform distribution of the benefits among different groups of users and items. Interestingly, many works in the last years have gone beyond this view and, nowadays, fairness and, analogously, unfairness can be defined adopting more fine-grained and non-uniform perspectives. As a consequence, measuring fairness is becoming more complex especially if one wants to quantify it.

Furthermore, according to [408, 409], we can classify the most popular notions of unfairness used in the literature as *disparate treatment* and *disparate impact*. Their common characteristic is that they both call for some type of parity (equality), either by ignoring user’s membership in protected classes (parity in treatment) or enforcing parity in fraction of users belonging to different protected classes, receiving beneficial outcomes (parity in impact). Under an operational lens, we may

say that parity in treatment refers to the training phase of a model while parity in impact to its usage. Although they look tightly connected, actually, we know that a parity in treatment does not necessarily imply a parity in impact.

From a recommender systems perspective, where users are first class citizens, there are multiple stakeholders and then fairness issues can be studied for more than one group of participants [78]. Previous work on fairness evaluation in RS has mostly interpreted fairness as some form of *equality* across multiple groups (e.g., gender, race). For example, Ekstrand et al. [135] studied whether RS produce *equal utility* for users of different demographic groups. In addition, Yao and Huang [404] studied various types of unfairness that can occur in collaborative filtering models where, to produce fair recommendations, the authors proposed to penalize algorithms producing disparate distributions of prediction error. Nonetheless, although less common, there are a few works where fairness has been defined beyond uniformity [62, 348, 412]. For instance, Biega et al. [62] concentrates on discovering the relation between *relevance* and *attention* in search (information retrieval). During a search session, searchers are subject to a high degree of *positional bias* due to paying much more attention to the top-ranked items than lower-ranked items. As a consequence, despite having a proper ranking based on relevance, lower-ranked items receive disproportionately less attention than they deserve. Their proposed approach promotes the notion that ranked subjects should receive attention that is proportional to their *worthiness* in a given search scenario and achieve fairness of attention by making exposure proportional to relevance. These research works however have focused on fairness from different perspectives and for different purposes.

In the present work, we argue that fairness does not necessarily imply equality between groups, but instead proper distribution of utility (benefits) based on merits and needs. Starting from this idea, we mainly focus on quantifying unfairness in recommendation systems, and we propose a probabilistic framework based on Generalized Cross Entropy (GCE) to measure fairness (or unfairness) of a given recommendation model that can be applicable to diverse recommendation scenarios. This is a general approach that can be easily adapted to any classification task. Our framework allows the designer to define and measure fairness for groups of

users (samples in a generic classification task) and for groups of items (target in a classification task). Moreover, the proposed framework is particularly flexible in the definition of different notions of fairness as it does not rely on specific and pre-defined probability distributions but they can be defined at design time. This lets the designer consider equality- and non-equality-based fairness notions adopting a single and unified perspective. The main characteristics of the proposed framework can be summarized as:

- One important characteristic of good evaluation metrics is their interpretability and explainability power. Generalized cross entropy is designed based on theoretical foundations, which makes it easy to understand and interpret.
- A large portion of previous work defines fairness as some form of **equality** across multiple groups (e.g., gender, race) [135]. However, as pointed out by some researchers [167, 404], fairness is not necessarily equivalent to equality. The proposed framework is sufficiently flexible to allow designers in the definition of fairness based on any arbitrary probabilistic distribution (in which uniform distribution is equivalent to equality in fairness).
- As a general remark, the proposed fairness-evaluation metric comes with a suite of other advantages compared to prior art:
 1. It incorporates a **gain factor** in its design, which can be flexibly defined to contemplate different accuracy-related metrics to measure fairness upon. Examples of such measures are *recommendation count* (focused on global count of recommendations), *decision-support metrics* (e.g., precision, recall) or *rank-based metrics* (e.g., nDCG, MAP). Prior art usually focuses on one of these aspects (see Section 15.2.4 for more details), which makes our approach more encompassing and general (cf. Section 15.3.1).

This choice derives from the assumption that the user satisfaction can be defined in many different ways. Based on the specific scenario, a certain metric could be more useful than others and, as a consequence,

the considered gain factor should differ. Additionally, the generalization of the gain factor allows the designer also adopting ranking-based gains like nDCG. This opens up new interesting perspectives. Let us suppose we are measuring fairness for different groups of items adopting nDCG as a gain factor. If the adopted probability distribution is not equal among groups, the GCE value will be related to the average position of the items of specific groups in the recommendation lists. The GCE will then measure if a RS is promoting relevant items from specific groups to users.

2. Unlike most previous work that solely focused on either user fairness or item fairness, the proposed framework integrates user-related and item-related gain factors. Also, we choose to evaluate fairness considering the various item and user attributes, showing how the different RSs behave in this respect. This brings our work closer to multiple stakeholder settings where benefits of multiple parties involved in the recommendation process should be considered (see Sections 15.2.3 and 15.3.1 for more details).

We have developed our investigation around the following research questions:

RQ1: How to define a fairness evaluation metric that considers different notions of fairness (not only equality)? We propose a probabilistic framework for evaluating RS fairness based on attributes of any nature (e.g., sensitive or insensitive) for both items or users. We show that the proposed framework is flexible enough to measure fairness in RS by considering it as equality or non-equality among groups, as specified by the system designer or any other parties involved in multi-stakeholder settings.

RQ2: How do classical recommendation models behave in terms of such an evaluation metric, especially under non-equality definitions of fairness? Some studies have been developed under different definitions of fairness, however in this study we shall focus on comparing the effect that equality

vs non-equality notions of fairness may have on classical families of recommendation algorithms.

RQ3: Which user and item attributes are more sensitive to different notions of fairness? Which ones impact more on different families of recommendation algorithms? Since fairness can be defined according to different user or item attributes, we aim to study the sensitivity of recommendation algorithms with respect to these parameters under the proposed probabilistic framework.

We answered the research questions performing extensive experiments on three well-known datasets: Amazon Toys & Games, Amazon Video Games, and Amazon Electronics. We tuned seven well-known baseline recommenders (ItemKNN, UserKNN, SVD++, BPR-MF, BPR-Slim, MostPopular, Random) and we evaluated them exploiting the proposed framework to measure fairness. To address the second research question, we considered uniform and non-uniform distributions among groups. This gave us a clear idea about how these classic recommenders behave. The third research question was addressed considering an adequate number of items and users attributes. We considered two attributes for items (Price and Popularity), and three attributes for users (Happiness, Helpfulness, and Interactions). While Popularity, Happiness, and Interactions are derived from the original user-item matrix, Price and Helpfulness are two attributes that are, at the same time, dataset-specific, and sensitive attributes. This research question imposed to re-evaluate all the baseline five times. However, this effort is paid back by results. On the one hand, they show that some recommenders make large use of popularity and they show a non-uniform behavior. On the other hand, some interesting similarities between different attributes emerged, resulting in recommenders that are more or less prone to produce better recommendations for groups of users or items, according to these attributes.

15.2 Background and Prior Art

In this section, we briefly review different notions of fairness and the trade-off between fairness and accuracy-oriented metrics explored in the literature.

15.2.1 Fairness notions

Machine learning (ML) is now involved in life-affecting decision points such as criminal risk prediction, credit risk assessments, housing allocation, loan qualification prediction, or hiring decision making [355, 382]. As ML is increasingly being employed to ease or automate decision making for human, some concerns have been recently raised on *fairness* of such models. Over the last decade, a growing number of research articles in the ML community have focused on defining appropriate notions of fairness and then developing models to ensure fairness in automated decision making (DM). Awareness on fairness and ethics in information retrieval has been raised by Belkin and Robertson already in 1976 [43]. By and large, the current notions of fairness are mainly influenced by the concept of discrimination in social sciences, law and economy [105]. For instance, back in the 90's there was interest to measure the distribution of personal characteristics such as income or wealth for a given population. As a result, the concept of unfairness (or discrimination) referred to disproportionate distribution of these resources.¹

We are concerned with fairness in algorithmic DM. One common characteristic of the fairness notion in the ML literature is that they all call for some form of parity (i.e., equality), either in treatment or in impact or both [408, 409]. Here we review these two popular notions of fairness used in the ML literature:

1. **Treatment parity:** Parity in treatment means the DM system does not use the information about user's membership in protected classes (e.g., gender, race), which are protected by anti-discrimination laws [1, 39]. The use of group-conditional DM systems is often prohibited to avoid this kind of disparity.

¹The terms "poverty", "welfare" or "inequality" were used interchangeably in the economy literature [104, 105] when referring to discrimination or unfairness.

As a result, a DM system whose outcome varies based on a change in the protected feature value is called to suffer from *disparate treatment*.

2. **Impact parity:** Parity in impact means the DM system needs to ensure parity in the fraction of users belonging to different protected groups (e.g., men, women) receiving beneficial outcomes. As a result, a DM system, which grants disproportionately large beneficial outcomes (or positive classification) favoring certain sensitive feature groups is called to suffer from *impact disparity*.

It is possible to extend the concept for one of the parity notions above to accommodate other classification outcomes. For example, a DM system whose error rate changes for different protected classes is called to suffer from disparity mistreatment [408].

Although the above notions proposed in prior studies provide an attractive viewpoint, they often lack flexibility with respect to one or more of the following aspects:

1. They are specifically designed for classification problems and define fairness based on the results of confusion matrix.
2. Fairness is measured with respect to instances of the training data.

15.2.2 Fairness and accuracy trade-off

Recommender systems help users in many decision-making and mission-critical tasks such as entertainment, medical, financial, or job-related applications. One of the key success indicators of RS is linked with the fact that they can alleviate the information overload pressure on information seekers by offering suggestions that match their tastes or preferences. It is common to measure the quality of a personalized recommendation algorithm in terms of *relevance* (e.g., personalized ranking) metrics. In domains such as news, books, movies and music where the individual preference is paramount, providing personalized recommendations can increase users' trust in and engagement with the system. These are important factors to motivate users to stay in and keep receiving recommendations, resulting in loyalty in

the long term and offering benefits to different parties involved in a recommendation setting such as consumers, suppliers, the system designer and other related services. Even in sensitive domains such as job recommendation, where fair opportunities to job seekers is desired, personalization can be relevant, e.g., a job-seeker might be willing to compensate salary with the distance factor or other benefits.

Nonetheless, blindly optimizing for accuracy-oriented metrics (or consumer relevance) may have adverse or unfavorable impacts on the fairness aspect of recommendations [258], e.g., in the employment recommendation context, certain genders or users from certain areas might be more likely to be recommended a job due to their behavioral differences. For example, male users or users from certain regions with high-speed Internet connection may produce more clicks compared to the others. A system optimizing for consumer relevance, might be unfair to less active users such as female or people from areas with less Internet activity thereby placing these groups at an unfair disadvantage. On the other hand, exposing all users equally might have detrimental impact on relevance and eventual consumer satisfaction. This inadvertently leads to a trade-off between relevance/personalization and fairness.

Zafar et al. [407] propose a framework for modeling the trade-off between fairness and accuracy in a classifier that suffers from disparate mistreatment-based fairness. The proposed model treats the under-researched variables (fairness and accuracy) in a joint fashion, e.g., by casting them in a convex-concave optimization constraint. This results in a fair classification in which disparate mistreatment on false positive and false negative rates are eliminated and can be tailored to measure, e.g., false discovery and false omission rates depending on a specific application. The framework provides the advantage to measure unfairness in situations where sensitive attributes of protected classes might not be accessible for reasons such as privacy or disparate treatment laws [39] prohibiting their use. In another study, Grgic-Hlaca et al. [160] propose a fairness-aware decision making system that unlike previous work focuses on fairness of outcomes (of a ML system), and results in a new dimension of fairness named *fairness in decision-making* (aka process fairness). The work introduces different measures to model individual users' moral

sense in deciding whether it is fair to use various input features in the decision making process. The authors show that it is possible to find a reasonable trade-off between process fairness and accuracy of a classifier over the set of features and provide the empirical evidence.

Studying fairness in information retrieval and recommendation is not limited to the aforementioned works. Relevance-fairness trade-off has been also studied from other angles by Abebe et al. [10] that propose an approach based on fair division of resources, by Mehrotra et al. [257] that focus on auditing search engine performance for fairness, and by Biega et al. [62] as well as Singh and Joachims [348] that study fairness in ranking.

Majority of the above works focused on fairness from the perspective of users (or user fairness). On the research works that focus on other fairness recipient we can name of [258], which exclusively focuses on supplier fairness in marketplaces.

15.2.3 From reciprocal recommendation to multiple stakeholders

Reciprocal recommendation views RS as systems fulfilling *dual* goals; the first goal is associated with satisfying *customers' preference* (user-centered utility) and the other goal is quite often associated with the *value of recommendations to the vendors* (vendor-centered utility aka profitability) [20]. Reciprocal recommendation regards the recommendation in most scenarios similar to a transaction and states that in generating recommendation *bilateral considerations* should be made meaning that the recommendation list must be acceptable to both parties involved in a transaction. On the domains, which use reciprocal recommendation we can name of on-line dating, on-line advertising, scientific collaboration and so on [78]. Maintaining a balance between the user and the vendor-centered utilities is the main attention focus of RS acknowledging this viewpoint to recommendation. In [20] the authors propose ValuePick, a framework which integrates the proximity to a target user and the global value of a network to recommend relevant nodes within a network. Several approaches have been proposed to combine the aforementioned util-

ities to either optimize profitability or to generate a win-win situation for providers and consumers [192] – according to which recommended items are ranked, see, e.g., [192, 96, 293]. Various approaches based on heuristic scoring model [96], mathematical optimization model [20, 34, 115], reinforcement learning [340, 216] and more complex models are proposed and used for this purpose. Some approaches have attempted to place into a mathematical optimization problem additional constraints such as consumer budget and other decision factors, for example, customer satisfaction levels [387].

When recommendations must account for the needs of more than just the two transacting parties, we move beyond reciprocal recommendation setting to a *multi-stakeholder recommendation setting*. Systems designed to meet the requirements of multiple stakeholders are referred to as Multi-stakeholder Recommender Systems (MRS). For instance, Etsy [5] is an e-commerce website focused on hand-made products and craft supplies. The recommender system platform in Etsy provide recommendation from small-scale artisans to consumers (shoppers). Hence, the recommender system on such a website needs to deal with the needs of both consumers and sellers [248]. According to [79], one can classify multiple stakeholders involved in a MRS into three main groups: consumers, providers and the platform (the system). The concept of MRS has been reintroduced in the *context of fairness* to highlight the fact that fairness is a *multisided concept* in which the impact of recommendation on multiple groups of individuals must be considered. In this context, [79] proposes to study the fairness issues relative to each one of these groups: (i) consumers (C-fairness), (ii) providers (P-fairness), and (iii) both (CP-fairness) [79]. The authors propose a balanced neighborhood as a mechanism to balance personalization vs. fairness of recommendation outcomes.

Several works have been proposed for evaluating recommendations in MRS. In [9, 82, 423] they suggest a utility-based framework for representing multiple stakeholders. As an example, in [423] the authors suggest a utility-based framework for MRS for personalized learning. Specifically, a recommender system is built for suggesting course projects to students by accounting both the student preferences and the instructors expectations in the model. The model aims to address the chal-

lence of over-expectations (by instructors) and under-expectations (by students) in the utility-based MRS. [359] approaches the MRS issue differently by relying on a constraint-based integer programming (IP) optimization model, where different sets of constraints can be used to characterize the objectives of different stakeholders. Finally, [8] provides a good frame of reference for the origins of multi-stakeholder recommendation, and the landscape of system designs.

15.2.4 Evaluating fairness in recommender systems

Even though research on fairness has been a very active topic in ML in general, and in RS in particular, there are not too many works—to the best of our knowledge—where authors address the goal we aim to achieve here: propose an evaluation metric that is capable of measuring fairness in RS. The closest work is [372], where the authors define a metric to capture the relative change of biases between the recommendations produced by an algorithm and those biases inherently found in the data. For this, the authors need to categorize both users and items, hence, it is not possible to measure only user or item fairness as allowed by our framework. Moreover, the most important disadvantage of the proposed metric is that the authors do not provide a single value for a given recommender, but a table (similar to a confusion matrix or contingency table) with all the possible combinations for pairs of user and item categories.

Nonetheless, even though we have not found other investigations specifically tackling the problem of defining a fairness evaluation metric, works that propose algorithms tailored for fairness need to be evaluated somehow, and these metrics, although usually based on heuristics, can also be considered to evaluate fairness. We start by describing the purely theoretical survey presented in [382], where the authors collect many definitions from the literature about the concept of fairness. The following three could be easily applied in a recommendation context: group fairness (equal probability of being assigned to the positive predicted class), predictive parity (correct positive predictions should be the same for both classes), and overall accuracy equality (groups have equal prediction accuracy). The last two could be computed by measuring the precision or the error in each class and somehow

comparing those values across all the groups. This is exactly the idea behind MAD (absolute difference between the mean ratings of different groups) used in [427]. Here, the authors also use in their experiments the Kolmogorov-Smirnov statistics of two distributions (predicted ratings for groups) as a comparison. The main problem with these two approaches and with some of the definitions in [382] is that they are only valid for 2 groups and are focused on ratings —and, consequently, only valid for the rating prediction task, which has been displaced by the community because it does not correlate with the user satisfaction [163, 256]—, mostly because fairness is addressed as a classification problem in ML. We find the same situation in [404] where the authors define several unfairness quantities (non-parity, value, absolute, underestimation, overestimation, and balance unfairness) that can only be applied to 2 groups of users and based on prediction errors.

Finally, we found other types of metrics not directly based on prediction errors. On the one hand, in [248] the authors define a metric tailored for P-fairness (fairness from the perspective of the providers in a multi-stakeholder setting) based on the provider coverage, that is, the number of providers covered by a recommendation algorithm. On the other hand, in [331] the authors use the Matthews correlation coefficient, since it allows to quantify the performance of an algorithm at a threshold while, at the same time, it penalizes the classifier for classifying all samples as the target class. In the paper, as some of the metrics presented above, the coefficient is defined only for the binary case where the attribute has two possible values, however it is possible to compute a multiclass version. In any case, as proposed by the authors, it can only be applied to user attributes.

Summing up, several metrics have been used to evaluate recommender systems under different notions of fairness. All of them assume equality as fairness and, usually, they are limited to user attributes with only two values (such as gender or other binary attributes). Therefore, we believe the framework we present could open up several possibilities in the field, since it overcomes all the above-mentioned limitations.

15.3 A probabilistic framework to evaluate fairness

We now present a probabilistic framework for evaluating RS fairness based on attributes of any nature (e.g., sensitive or insensitive) for both items or users and show that the proposed framework is flexible enough to measure fairness in RS by considering fairness as equality or non-equality among groups, as specified by the system designer or any other parties involved in multi-stakeholder setting.

15.3.1 Using Generalized Cross Entropy to measure user and item fairness

In this section, we propose a framework based on generalized cross entropy for evaluating fairness in recommender systems. Let U and I denote a set of users and items, respectively and A be a set of sensitive attributes in which fairness is desired. Each attribute can be defined for either users, e.g., gender and race, or items, e.g., item provider (or stakeholder). Given a set M (for models) of recommendation systems, we define the *unfairness measure* as the function

$$\omega : M \times A \rightarrow \mathbb{R}^+$$

The goal is to find a function ω that produces a non-negative real number for a recommender system that represents and measures its (un)fairness. A recommender system $m \in M$ is considered less unfair (i.e., more fair) than $m' \in M$ with respect to the attribute $a \in A$ if and only if $\omega(m, a) < \omega(m', a)$. Previous works have used *inequality* measures to evaluate algorithmic unfairness, however, we argue that fairness does not always imply equality. For instance, let us assume that there are two types of users in the system – regular (free registration) and premium (paid) – and the goal is to compute fairness with respect to the users' subscription type. In this example, it might be more fair to produce better recommendations for paid users, therefore, equality is not always equivalent to fairness – note that, in any case, the goal is to ensure that premium users receive good (or better) recommendations without affecting the experience of regular users.

We define fairness of a recommender system with respect to an attribute a using the *Csiszar generalized measure of divergence* as follows [112]:

$$\omega(m, a) = \int p_m(a) \cdot \varphi \left(\frac{p_f(a)}{p_m(a)} \right) da \quad (15.1)$$

where p_m and p_f respectively denote the probability distribution of the model m 's performance and the fair probability distribution, both with respect to the attribute a [69]. A distinguishing property of this measure is that conceptually there are no differences for the case in which p_m and p_f are discrete densities, in such a case the integral is simply replaced by the sum. Csiszar's family of measures subsumes all of the information-theoretic measures used in practice (see [208, 169]). We restrict our attention to the case when $\varphi(x) = \frac{x^\alpha - x}{\alpha \cdot (\alpha - 1)}$ and $\alpha \neq 0, 1$ for some parameter α ; then, the family of divergences indexed by α boils down to *Generalized Cross Entropy*:

$$f(m, a) = GCE(m, a) = \frac{1}{\alpha \cdot (1 - \alpha)} \left[\int p_f^\alpha(a) \cdot p_m^{(1-\alpha)}(a) da - 1 \right] \quad (15.2)$$

The unfairness measure I is minimized with respect to attribute a when $p_m = p_f$, meaning that the performance of the system is equal to the performance of a fair system. In the next sections, we discuss how to obtain or estimate these two probability distributions.

Note that the defined unfairness measure indexed by α includes the Hellinger distance for $\alpha = 1/2$, the Pearson's χ^2 discrepancy measure for $\alpha = 2$, Neymann's χ^2 measure for $\alpha = -1$, the Kullback-Leibler divergence in the limit as $\alpha \rightarrow 1$, and the Burg CE distance as $\alpha \rightarrow 0$. Figure 15.1 illustrates simulation of how GCE changes across different α values.

If the attribute a is discrete or categorical (as typical attributes, such as gender or race), then the unfairness measure is defined as:

$$GCE(m, a) = \frac{1}{\alpha \cdot (1 - \alpha)} \left[\sum_{a_j} p_f^\alpha(a_j) \cdot p_m^{(1-\alpha)}(a_j) - 1 \right] \quad (15.3)$$

The role of α in the definition of GCE is critical, as we show in Figure 15.1. We observe, for instance, that at extreme values of p_m , GCE obtains larger values for

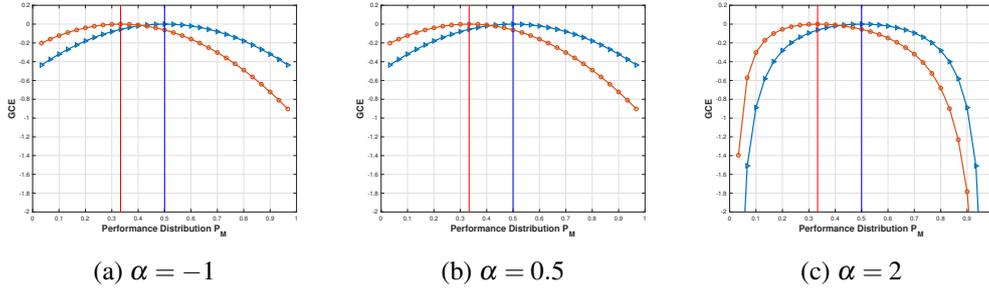


Figure 15.1: Simulations of values obtained using GCE fairness evaluation metric for different fair distribution types p_f and performance distributions p_m and different α values. For example, when x-axis is 0.3 then $p = [0.3, 0.7]$. The **blue** curve represents $p_f = [0.5, 0.5]$ while the **red** represents $p_f = [0.3, 0.7]$. It can be noted when fairness means equality the representative **blue** curve is used which is maximized at 0.5; this is while when fairness means non-equality the representative **red** curve should be used which is maximized at a point non-equal to 0.5 (here 0.3).

lower values of α . Besides, according to [69], Pearson's χ^2 measure (which corresponds to $\alpha = 2$) is more robust to outliers than other typical divergence measures such as Kullback-Leibler divergence; hence, in the rest of the dissertation, unless stated otherwise, we shall use this value for parameter α .

Defining the fair distribution p_f

The definition of a fair distribution p_f is problem-specific and should be determined based on the problem or target scenario in hand. For example, one may want to ensure that premium users, who pay for their subscription, would receive more relevant recommendations, because running complex recommendation algorithms might be costly and not feasible for all users². In this case, p_f should be non-uniform across the user classes (premium versus free users). In other scenarios, a uniform definition of p_f might be desired. Generally, when fairness is equivalent to equality, then p_f should be uniform and in that case, the generalized cross entropy

²These scenarios are becoming more and more realistic especially in edge computing settings where computational resources are often quite limited.

would be the same as generalized entropy (see [355] for more information).

Note that p_f can be seen as a more general utility distribution, and the goal is to observe such distribution in the output of the recommender system. In this chapter, since we focus on recommendation fairness, we refer to p_f as the fair distribution.

Finding the fair distribution p_f is challenging. It is task-specific and a fair distribution in one domain is not necessarily a fair distribution in another. However, generalized cross entropy is a general framework that allows researchers and practitioners in different domains to define the fairness definition based on their needs. We leave discussions on various definition of p_f in different domains for future.

Estimating the model distribution p_m

The model distribution p_m should be estimated based on the output of the recommender system on a test set. In the following, we explain how we can compute this distribution for item attributes. We define the recommendation gain (rg_i) for each item $i \in I$ as follows:

$$rg_i = \sum_{u \in U} \phi(i, Rec_u^K) \cdot g(u, i, r) \quad (15.4)$$

where Rec_u^K is the set of top- K items recommended by the system to the user $u \in U$. $\phi(i, Rec_u^K) = 1$ if item i is present in Rec_u^K ; otherwise $\phi(i, Rec_u^K) = 0$. The function $g(u, i, r)$ is the gain of recommending item i to user u with the rank r . Such gain function can be defined in different ways. In its simplest form, if always $g(u, i, r) = 1$, the recommendation gain in Eq. (15.4) would boil down to recommendation count (i.e., $rg_i = rc_i$).

A binary gain in which $g(u, i, r) = 1$ when item i recommended to user u is relevant and $g(u, i, r) = 0$ otherwise, is another simple form of the gain function based on relevance. The gain function g can be also defined based on ranking information, i.e., recommending relevant items to users in higher ranks is given a higher gain. In such case, we propose to use the discounted cumulative gain (DCG) function that is widely used in the definition of nDCG [?], given by $\frac{2^{\text{rel}(u,i)-1}}{\log_2(r+1)}$ where $\text{rel}(u, i)$ denotes the relevance label for the user-item pair u and i . We can further normalize the above formula based on the ideal DCG for user u to compute the gain

function g .

As we can see in the definition of the gain function for items, it is possible to flexibly specify the constraint under which fairness needs to be satisfied (e.g., based on recommendation count, relevance, ranking, or a combination thereof). As such, our approach extends considerably the previous approaches, e.g., [62, 348, 412] which focused on a single aspect of fairness, e.g., either based on error or ranking.

Then, the model probability distribution p_m^I is computed proportionally to the recommendation gain for the items associated to an item attribute value a_j . Formally, the probability $p_m^I(a_j)$ used in Eq. (15.3) is defined as:

$$p_m^I(a_j) = \frac{\sum_{i \in a_j} rg_i}{Z} \quad (15.5)$$

where Z is a normalization factor set equal to $Z = \sum_i rg_i$ to make sure that $\sum p_m^I(a_j) = 1$.

Under an analogous formulation, we could define a variation of fairness for users $u \in U$ based on Eq. (15.4):

$$rg_u = \sum_{i \in I} \phi(i, Rec_u^K) \cdot g(u, i, r) \quad (15.6)$$

where in this case, the gain function cannot be reduced to 1, otherwise, all users would receive the same recommendation gain rg_u . Then, to compute $p_m^U(a_j)$, we normalize these gains in a similar way as shown in Eq. (15.5).

It should be noted that, to avoid zero probabilities, we smoothed the previous computations by using the Jelinek-Mercer method [413] as follows, where p_m^e corresponds to either p_m^I or p_m^U depending if rg_i or rg_u are used:

$$\begin{aligned} \tilde{p}_m^e(a_j) &= \frac{\sum_{i \in a_j} rge}{Z} \\ \hat{p}_m^e(a_j) &= \lambda \cdot \tilde{p}_m^e(a_j) + (1 - \lambda) \cdot p_C \\ \hat{Z} &= \sum_j \hat{p}_m^e(a_j) \\ p_m^e(a_j) &= \frac{\hat{p}_m^e(a_j)}{\hat{Z}} \end{aligned}$$

Here, smoothing is applied in the second equation, where we use a background probability p_C . In the experiments, we used $\lambda = 0.95$ and $p_C = 0.0001$. Additionally, to obtain more robust values of the probabilities estimated using the recommendation gains, a slightly more complicated version of these formulations could be used where the probabilities consider the average of gains rg_e in a user-basis instead of such gains directly, since this is how typical evaluation metrics are computed in the recommender systems literature. For the sake of space we avoid including such formulation here.

15.3.2 Toy example

For the illustration of the proposed concept, in Table 15.1 we provide a toy example on how our approach for fairness evaluation framework could be applied in a real recommendation setting. A set of six users belonging to two groups (each group is associated with an attribute value a_1 (red) or a_2 (green)) who are interacting with a set of items are shown in Table 15.1. Let us assume the red group represents users with a *regular* (free registration) subscription type on an e-commerce website while the green group represents users with a *premium* (paid) subscription type. A set of recommendations produced by different systems (**Rec0**, **Rec1**, and **Rec2**) are shown in the last columns. The goal is to compute fairness using the proposed fairness evaluation metric based on GCE given by Eqs. (15.3) and (15.6). The results of evaluation using three different evaluation metrics are shown in Table 15.2. The metrics used for the evaluation of fairness and accuracy of the system include: (i) GCE, (ii) Precision, and (iii) Recall, all at cutoff 3. Note that $GCE = 0$ means the system is completely fair, and the closer the value is to zero, the more fair the respective system is.

By looking at the recommendation results from **Rec0**, one can note that *if fairness is defined as equality between two groups*, defined through fair distribution $p_f = [\frac{1}{2}, \frac{1}{2}]$, then **Rec0** is not a completely fair system, since $GCE = -0.09 \neq 0$. In contrast, *if fairness is defined as providing recommendation of higher utility (usefulness) to green users who are users with paid premium membership type*, (e.g., by setting $p_{f_2} = [\frac{1}{3}, \frac{2}{3}]$) then, since GCE is smaller, we can say that recommendations

Table 15.1: A set of 6 users belonging to groups g_1 and g_2 and 10 items along with their true labels marked by ✓ and recommended items by recommenders Rec 0, Rec 1, Rec 2. Rec 0 produces 3 and 6 relevant items for free and premium users (in total) respectively; Rec 1 generates 1 relevant item for each user; Rec 2 produces recommended items that are all relevant for all users.

| | | i_1 | i_2 | i_3 | i_4 | i_5 | i_6 | i_7 | i_8 | i_9 | i_{10} | Rec 0 | Rec 1 | Rec 2 |
|-------|---------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|---------------------|------------------------|------------------------|
| a_1 | user 1 | ✓ | | ✓ | | | | ✓ | | | | $\{i_1, i_6, i_8\}$ | $\{i_1, i_5, i_9\}$ | $\{i_1, i_3, i_7\}$ |
| a_1 | user 2 | | | | | ✓ | | | ✓ | | | $\{i_2, i_5, i_9\}$ | $\{i_2, i_5, i_7\}$ | $\{i_1, i_5, i_8\}$ |
| a_1 | user 3 | | ✓ | | | | | ✓ | | | | $\{i_1, i_6, i_7\}$ | $\{i_2, i_5, i_9\}$ | $\{i_2, i_7, i_9\}$ |
| a_2 | user 4 | | | ✓ | ✓ | | | | | ✓ | | $\{i_3, i_4, i_9\}$ | $\{i_4, i_5, i_6\}$ | $\{i_3, i_4, i_9\}$ |
| a_2 | user 5 | | | | | ✓ | | ✓ | | | ✓ | $\{i_1, i_5, i_7\}$ | $\{i_1, i_2, i_{10}\}$ | $\{i_5, i_7, i_{10}\}$ |
| a_2 | user 6 | ✓ | | ✓ | | | ✓ | | | ✓ | | $\{i_2, i_6, i_9\}$ | $\{i_1, i_5, i_8\}$ | $\{i_3, i_6, i_9\}$ |

produced by **Rec0** are more fair for this type of users and also with respect to the other recommenders. Both of the above conclusions are drawn with respect to attribute “subscription type” (with categories free/paid premium membership). This is an interesting insight which shows the evaluation framework is flexible enough to capture fairness based on the interest of system designer by defining what she considers as fair recommendation through the definition of p_f . While in many application scenarios we may define fairness as equality among different classes (e.g., gender, race), in some scenarios (such as those where the target attribute is not sensitive, e.g., regular v.s. premium users) fairness may not be equivalent to equality.

Furthermore, by comparing the performance results of **Rec1** and **Rec2**, we observe that, even though precision and recall improve for **Rec2** and becomes the most accurate recommendation list, it fails to keep a decent amount of fairness with respect to any parameter settings of GCE, as in all the cases it is outperformed by the other methods. Moreover, GCE only reaches the optimal value for **Rec1** and p_{f_0} , since that recommender produces the same number of relevant items (one) for every user, independently of the user group; in the other cases, since there are more relevant items on green than red users, the results reflect the amount of inherent biases in the data due to the unequal distribution of resources among classes.

Table 15.2: Fairness of different recommenders in the toy example presented in Table 15.1 according to proposed GCE and individual-level accuracy metrics. Note that $p_{f_0} = [\frac{1}{2}, \frac{1}{2}]$, $p_{f_1} = [\frac{2}{3}, \frac{1}{3}]$, and $p_{f_2} = [\frac{1}{3}, \frac{2}{3}]$ characterize the fair distribution as uniform or non-uniform distribution (of resources) among two groups.

| | GCE ($p_f, p_m, \alpha = 2$) | | | P@3 | R@3 |
|--------------|---------------------------------------|-----------|-----------|---------------|--|
| | p_{f_0} | p_{f_1} | p_{f_2} | | |
| Rec 0 | -0.0952 | -0.3201 | -0.0026 | $\frac{1}{2}$ | $\frac{1}{6} \cdot \frac{19}{6} = 0.530$ |
| Rec 1 | 0 | -0.0556 | -0.0556 | $\frac{1}{3}$ | $\frac{1}{6} \cdot \frac{9}{4} = 0.375$ |
| Rec 2 | -0.0079 | -0.1067 | -0.0220 | 1 | $\frac{1}{6} \cdot \frac{23}{4} = 0.958$ |

This evidences that optimizing an algorithm to produce relevant recommendations does not necessarily result in more fair recommendation rather, conversely, a trade-off between the two evaluation properties can be noticed.

15.4 Experimental settings

15.4.1 Datasets

To address the research questions presented before, we use datasets from different domains with more or less sensitive attributes. This allows us to evaluate several notions of fairness under user and item dimensions. More specifically, we have used multiple product categories of the Amazon Review dataset [173, 4]. This dataset is a collection of product reviews aggregated at the category level, which also includes metadata from Amazon; in total it contains 142.8 million reviews spanning from May 1996 to July 2014. Beyond ratings, these datasets include reviews (which consist of ratings, text, timestamp, and votes from other users to determine how helpful a review is), product metadata (descriptions, category information, price, brand, and image features [255]), and links (graphs with information about also viewed/also bought items).

Overall the Amazon Dataset contains 24 different category-level datasets. Based on the number of users, items, and transactions we have selected the fol-

lowing three datasets to conduct our study. The smallest one is `Amazon Video Games`, with more than 1 million ratings, devoted to videogames sold on the Amazon Store. The second dataset is `Amazon Toys & Games`, with more than 2 million transactions of toys and tangible games. The last and largest dataset is `Amazon Electronics`, with almost 8 million overall ratings.

15.4.2 Evaluation Protocol and Temporal Splitting

The experimental evaluation is conducted adopting the so-called “All Items” evaluation protocol [49] in which, for each user, all the items that are not rated yet by the user are considered as candidates when building the recommendation list.

To simulate an online real scenario as realistically as possible, we use the fixed-timestamp splitting method [32, 29], initially suggested in [83, 163]. The core idea is choosing a single timestamp that represents the moment in which test users are on the platform waiting for recommendations. Their past will correspond to the training set, whereas the performance is evaluated exploiting data which occurs after that moment. In this work, we select the splitting timestamp that maximizes the number of users involved in the evaluation by setting two reasonable constraints: the training set of each user should keep at least 15 ratings, while the test set should contain at least 5 ratings. Training set and test set for the three datasets are made publicly available for research purposes, along with the splitting code³.

Finally, the statistics of the training and test datasets used in the experiments are depicted in Table 15.3, where the difference in the number of transactions between the original datasets (see previous section) and the ones used in the experiments is due to the constraints imposed in the splitting process. It is important to note that, in any case, the processed datasets keep very small density values – between 0.054% and 0.48% – as it is standard in the literature.

³<https://github.com/sisinflab/DatasetsSplits/>

Table 15.3: Statistics about the datasets used in the experiments.

| Training Set | | | | | | |
|---------------------|--------|--------|---------------|----------|------------|------------|
| Dataset | #Users | #Items | #Transactions | Sparsity | From | To |
| Amazon Electronics | 5,351 | 56,727 | 164,375 | 99.94584 | 07/14/1999 | 05/14/2013 |
| Amazon Toys & Games | 1,108 | 24,158 | 38,317 | 99.85685 | 07/22/2000 | 08/30/2013 |
| Amazon Videogames | 479 | 8,892 | 20,369 | 99.52177 | 11/18/1999 | 10/28/2011 |
| Test Set | | | | | | |
| Dataset | #Users | #Items | #Transactions | Sparsity | From | To |
| Amazon Electronics | 5,351 | 28,792 | 74,090 | 99.95191 | 05/15/2013 | 07/23/2014 |
| Amazon Toys & Games | 1,108 | 9,192 | 15,169 | 99.85106 | 08/31/2013 | 07/22/2014 |
| Amazon Videogames | 479 | 4,171 | 8,114 | 99.59387 | 10/29/2011 | 07/21/2014 |

15.4.3 Attribute selection and discretization

In this work, we follow an attribute-based analysis of fairness in RS. In particular, we assume that users and items are associated with some attributes. Each attribute partitions the users/items into a number of groups (classes) where users/items in each group have the same attribute value (e.g., male or female for users) or (e.g., low-priced or high-priced items). One of the main objectives in the attribute-based study of fairness is to avoid discrimination against protected groups; as such these attributes are quite often chosen as nontrivial or (in some cases) sensitive attributes. Therefore, in this section we describe which user and item attributes were selected and how they were discretized in a limited number of groups or classes.

We start by selecting some attributes that we feel they are common enough to be found in almost any recommender system, in this way, the presented analysis could be relevant for both researchers that use domains not addressed in this work and industry practitioners with different data. For items, we focus on their popularity, which corresponds to the amount of interactions received by the items. Since the popularity of items is a signal of the common ratings (or clicks, views, etc.) between users, we aim to explore whether the most common collaborative filtering

algorithms are more prone to suggest popular items. Similarly, for users we focus on the amount of interactions registered by the system from each user, that is, the level of user activity. In this way, we aim to analyze the behavior of algorithms with respect to cold (few interactions) or warm (many interactions) users, as they typically refer to in the literature. Additionally, we interpret the average rating provided by the users as a signal of the level of satisfaction with respect to the system, we name this user feature as *happiness*. In our experiments we aim to investigate whether the recommenders behave fairly for satisfied (happy) and unsatisfied (unhappy) users.

Now, we select two attributes that are more specific to Amazon datasets and that are, to some extent, sensitive for both users and system developers: item price and user helpfulness. The price of an item is indeed an interesting and sensitive attribute, since many users may decide to select or buy a product just because of its price, even when they know that another product might be more beneficial or suitable for them. Hence, by including this attribute we aim to study whether classical recommendation approaches are more (or less) prone to recommend expensive or cheap products – without including such information into the recommendation algorithm – which might be perceived as not fair from the user perspective. The user helpfulness, on the other hand, is a piece of information that is not widely available, but it is becoming a frequent signal in review-based systems, since it allows users to vote on other users' reviews, increasing the confidence on the system. In this way, we aim to analyze if the most helpful users are provided with the best recommendations or not.

Once different user and item attributes are selected, we present how we discretized their values into a small number of classes or clusters. This step is not mandatory since our proposed metric could work with any number of categories or attribute values, however, to make the presentation and discussion of results less cumbersome and confusing, we prefer to limit the number of categories to a maximum of 4 in every case. In general, we decided to create clusters based on quartiles, which are particularly intuitive and allow to be generalizable to datasets of different nature, since the intrinsic distribution of the attributes is taken into account. More specifically, **item price**, **user helpfulness**, and **user interaction** were directly clus-

tered into 4 quartiles according to their original distributions. However, the rest of the attributes presented some problems which made it impossible to apply a standard clustering technique based on the quartiles. First, the **item popularity** showed so many ties for the least popular items that it was not possible to define boundaries for the quartiles. To address this issue, we increased the number of considered quantiles until we obtained 4 distinct clusters; this number corresponds to 30 for Amazon Toys & Games, and 10 for Amazon Video Games and Amazon Electronics. Regarding the last attribute, **user happiness**, we faced a different problem, where most of the average user ratings are very close to 4, hence, the clustering based on quartiles would have lost meaning. For this reason, we decided to set a reasonable threshold equal to 4 (common to the three datasets) to create only two categories: users whose average rating is smaller than 4, and the rest, to separate users according to a predefined level of satisfaction or happiness. Tables 15.4-15.6 present statistics about the resulting clusterings, respectively for Amazon Toys & Games, Amazon Video Games, and Amazon Electronics.

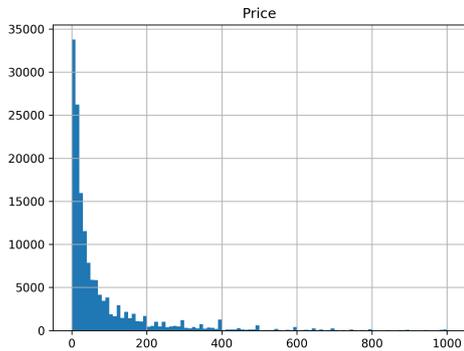
Finally, we note an issue we had to address regarding the computation of quartiles with respect to the availability of side information. First, not all items had associated metadata, whereas this is true for users, information for items is incomplete. Second, items in the training set only correspond to a small fraction of the items in the whole collection; hence, they might not be representative of the entire collection. Because of this, we computed the quartiles (for the item price attribute, which is the only one obtained through the metadata) according to two strategies: either based on the overall metadata information, or based only on the items with metadata that appear in the training set. This information is included in Tables 15.4-15.6 in columns *Price (TS)* for the case where the clustering is computed based on the training set, and in *Price (M)* when the whole metadata is used. Additionally, in Figure 15.2 we present the histograms of the 3 datasets comparing the two strategies to compute the quartiles. In the tables we observe that the resulting item distribution in clusters when using all the metadata is no longer uniform; similarly, in the histograms we see that the distribution is dominated by those very cheap items when using all metadata information, whereas other price values become visible when

only the training items are represented. Hence, because of these issues, we shall work from now on with the strategy based on building the clusters using information from the training set.

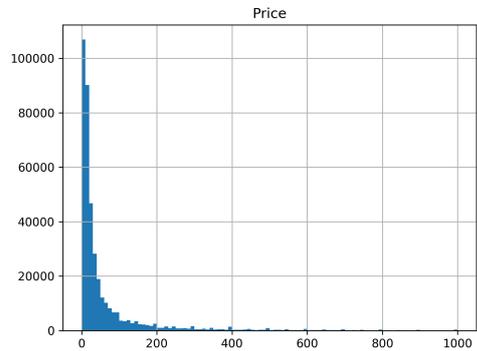
15.4.4 Baseline recommenders

We evaluate several families of Collaborative Filtering recommendation models. Beyond Nearest Neighbors memory-based models, we include latent factors models considering two different kinds of optimization: the minimization of the prediction error, and a pairwise learning-to-rank approach. More specifically, we include:

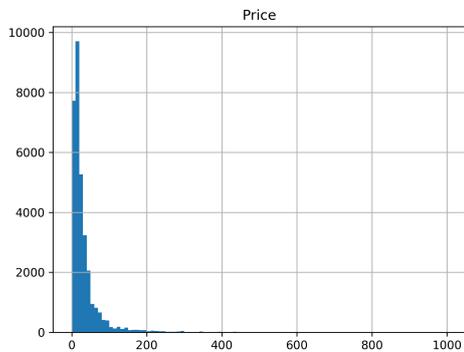
- `Random`, a non-personalized algorithm that produces a random recommendation list for each user. The items are chosen according to a uniform distribution.
- `MostPopular`, a non-personalized algorithm that produces the same recommendation list for all the users. This list is computed by measuring the items' popularity and ordering the items according to that value in descending order. It is acknowledged that popularity ranking typically show very good performance because of statistical biases in the data [49] and it is an important baseline to compare against [108].
- `ItemKNN` [333, 334], an item-based implementation of the K-nearest neighbors algorithm. It finds the K-nearest item neighbors based on a specific similarity function. Usually, as similarity functions, Binarized and standard Cosine Vector Similarity [36, 63, 13], Jaccard Coefficient [129, 307], and Pearson Correlation [177] are considered. The items in neighborhood are then used to predict a score for each user-item pair.
- `UserKNN` [73], a user-based implementation of the K-nearest neighbors algorithm. It finds the K-nearest user neighbors based on a similarity function (usually the same functions as described before for `ItemKNN`). The computed neighbors are then used to predict a score for each user-item pair.



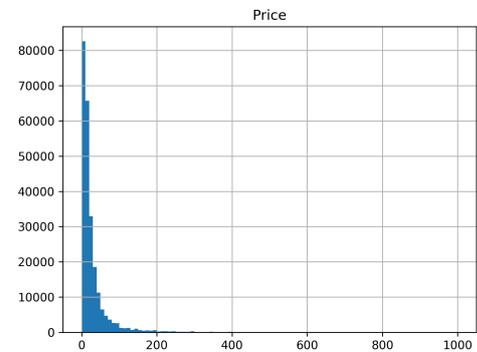
(a) Amazon Electronics: training set



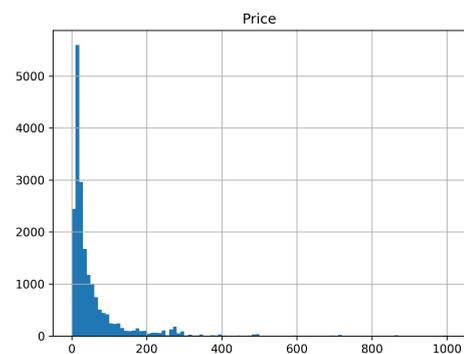
(b) Amazon Electronics: metadata



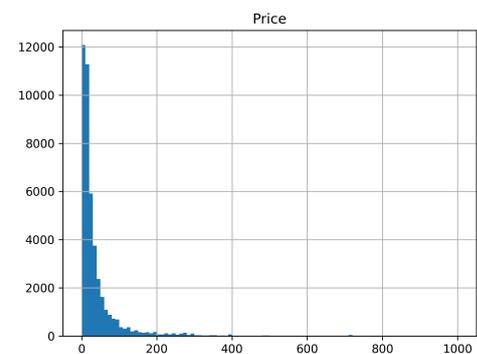
(c) Amazon Toys & Games: training set



(d) Amazon Toys & Games: metadata



(e) Amazon Video Games: training set



(f) Amazon Video Games: metadata

Figure 15.2: Histograms of the item price attribute (considering 100 bins) comparing two strategies to extract the values from (that will be used later to compute the attribute categories): based on items from the training set or based on all the items with associated metadata.

Table 15.4: Statistics about the user and item clustering methods for Amazon Toys & Games, where TS means Training Set, M stands for Metadata, Pop Popularity, Hlpf Helpfulness, Int Interactions, and Hpns Happiness. The rows 25%, 50% and 75% indicate the values of each attribute at that point of the distribution, which correspond to the boundaries between the first and second, second and third, and third and fourth quartiles. Note that, ideally, the number of items (#Items) and users (#Users) in each cluster is expected to be as uniform as possible.

| | Items Clusterings | | | Users Clusterings | | |
|-------------------|-------------------|------------------|---------------|-------------------|---------------|---------------|
| Statistics | Price (TS) | Price (M) | Pop | Hpns | Hlpf | Int |
| count | 19,543 | 19,543 | 24,158 | 1,108 | 1,108 | 1,108 |
| mean | 33.82 | 30.91 | 1.59 | 4.28 | 0.36 | 34.58 |
| std | 57.30 | 56.42 | 2.31 | 0.44 | 0.17 | 40.67 |
| min | 0.01 | 0 | 1 | 1.71 | 0 | 15 |
| 25% | 9.69 | 7.99 | 1 | 4 | 0.24 | 18 |
| 50% | 18.12 | 15.85 | 1 | 4.31 | 0.34 | 23 |
| 75% | 35 | 30.99 | 1 | 4.61 | 0.46 | 35 |
| max | 999.99 | 999.99 | 50 | 5 | 1.00 | 525 |
| Clusters | #Items | #Items | #Items | #Users | #Users | #Users |
| 0 | 4,893 | 3,870 | 22,234 | 264 | 277 | 326 |
| 1 | 4,880 | 4,808 | 729 | 844 | 277 | 246 |
| 2 | 4,911 | 5,213 | 461 | | 277 | 262 |
| 3 | 4,859 | 5,652 | 734 | | 277 | 274 |

Table 15.5: Statistics about the user and item clustering methods for Amazon Video Games, notation as in Table 15.4.

| | Items Clusterings | | | Users Clusterings | | |
|-------------------|-------------------|------------------|---------------|-------------------|---------------|---------------|
| Statistics | Price (TS) | Price (M) | Pop | Hpns | Hlpf | Int |
| count | 8,297 | 8,297 | 8,892 | 479 | 479 | 479 |
| mean | 56.28 | 40.89 | 2.29 | 3.93 | 0.51 | 42.52 |
| std | 85.66 | 67.98 | 2.92 | 0.57 | 0.18 | 65.67 |
| min | 0.01 | 0 | 1 | 1.67 | 0.04 | 15 |
| 25% | 14.99 | 9.99 | 1 | 3.62 | 0.38 | 19 |
| 50% | 28.99 | 19.99 | 1 | 4.01 | 0.51 | 25 |
| 75% | 59.99 | 39.99 | 2 | 4.3 | 0.63 | 43 |
| max | 999.99 | 999.99 | 45 | 5 | 0.98 | 785 |
| Clusters | #Items | #Items | #Items | #Users | #Users | #Users |
| 0 | 2,075 | 1,411 | 6,812 | 231 | 120 | 146 |
| 1 | 2,076 | 2,182 | 736 | 248 | 120 | 99 |
| 2 | 2,143 | 1,829 | 684 | | 119 | 116 |
| 3 | 2,003 | 2,875 | 660 | | 120 | 118 |

Table 15.6: Statistics about the user and item clustering methods for Amazon Electronics, notation as in Table 15.4.

| | Items Clusterings | | | Users Clusterings | | |
|-------------------|-------------------|------------------|---------------|-------------------|---------------|---------------|
| Statistics | Price (TS) | Price (M) | Pop | Hpns | Hlpf | Int |
| count | 47,660 | 47,660 | 56,727 | 5,351 | 5,351 | 5,351 |
| mean | 71.92 | 61.20 | 2.90 | 4.21 | 0.40 | 30.72 |
| std | 124.10 | 118.68 | 6.36 | 0.46 | 0.16 | 25.57 |
| min | 0.01 | 0.01 | 1 | 1.53 | 0 | 15 |
| 25% | 10.06 | 9.95 | 1 | 3.96 | 0.28 | 18 |
| 50% | 24.99 | 19.99 | 1 | 4.27 | 0.39 | 23 |
| 75% | 71 | 51.91 | 2 | 4.54 | 0.51 | 34 |
| max | 999.99 | 999.99 | 275 | 5 | 1.00 | 500 |
| Clusters | #Items | #Items | #Items | #Users | #Users | #Users |
| 0 | 11,915 | 11,058 | 43,674 | 1,434 | 1,338 | 1,607 |
| 1 | 11,940 | 10,042 | 3,743 | 3,917 | 1,338 | 1,230 |
| 2 | 11,893 | 11,562 | 4,345 | | 1,337 | 1,245 |
| 3 | 11,912 | 14,998 | 4,965 | | 1,338 | 1,269 |

Table 15.7: Tuned hyperparameters for each of the tested recommendation methods.

| | Number of Neighbors | Similarity Function |
|----------------|---------------------|---|
| ItemKNN | [50,60,70,90,100] | [Jaccard Coefficient, Binary Cosine, Cosine, Pearson Correlation] |
| UserKNN | [50,60,70,90,100] | [Jaccard Coefficient, Binary Cosine, Cosine, Pearson Correlation] |

| | Number of Latent Factors | Learning Rate | Iterations |
|----------------|--------------------------|----------------------------|--------------|
| SVD++ | [10,20,30,50,100,150] | [0.0005,0.005,0.05] | [30] |
| BPRMF | [10,20,30,50,100,150] | [0.005,0.05,0.5] | [30] |
| BPRSlim | | [0.0005, 0.005, 0.05, 0.5] | [5,10,20,30] |

- SVD++ [221, 224], an algorithm that takes advantage of a simple latent factors model (trained through the stochastic gradient descent method) and it models and compute user and item biases. SVD++ also considers implicit feedback to improve learning.
- BPRMF (Bayesian personalized ranking - Matrix Factorization) [312, 224], a matrix factorization algorithm that exploits the Bayesian Personalized Ranking criterion [312] to minimize the ranking errors.
- BPRSlim (Bayesian personalized ranking - SLIM) [273], an algorithm that produces recommendations using a sparse aggregation coefficient matrix trained with a Sparse Linear method (SLIM), trained maximizing the BPR criterion.

For all these recommenders we have performed a grid search to tune the parameters. We consider the range of values as suggested by the authors or by varying the parameters values around the ones showed in the original papers as the best performing ones; a summary of the considered values is shown in Table 15.7. The optimal values are reported in Table 15.8, and correspond to those that maximize the nDCG metric at cutoff 10.

15.4.5 Evaluation metrics

In our experiments, we compute accuracy metrics as it is standard in the literature [163]. The *top-N* recommendation accuracy metrics we have used are Precision

Table 15.8: Optimal hyperparameters according to nDCG@10.

| | | | |
|---------------------------------|--------------------------------|--------------------------|---------------------------|
| ItemKNN | Amazon Toys & Games | Amazon Videogames | Amazon Electronics |
| Number of Neighbors | 50 | 50 | 50 |
| Similarity Function | Jaccard Coefficient | Jaccard Coefficient | Jaccard Coefficient |
| UserKNN | Amazon Toys & Games | Amazon Videogames | Amazon Electronics |
| Number of Neighbors | 90 | 50 | 100 |
| Similarity Function | Cosine | Jaccard Coefficient | Jaccard Coefficient |
| SVD++ | Amazon Toys & Games | Amazon Videogames | Amazon Electronics |
| Number of Latent Factors | 150 | 100 | 100 |
| Learning Rate | 0.0005 | 0.005 | 0.0005 |
| Iterations | 30 | 30 | 30 |
| BPRMF | Amazon Toys & Games | Amazon Videogames | Amazon Electronics |
| Number of Latent Factors | 10 | 10 | 150 |
| Learning Rate | 0.5 | 0.5 | 0.5 |
| Iterations | 30 | 30 | 30 |
| BPRSlim | Amazon Toys & Games | Amazon Videogames | Amazon Electronics |
| Learning Rate | 5 | 0.05 | 0.05 |
| Iterations | 30 | 30 | 30 |

($P@N$), Recall ($R@N$), and nDCG ($nDCG@N$), all at a cutoff N , which means that we only consider *top-N* items within each recommendation list.

Precision is defined as the proportion of recommended items that are relevant to the user:

$$Precision_u@N = \frac{|Rec_u^N \cap TS_u^+|}{N}$$

where Rec_u^N is the recommendation list up to the N -th element and TS_u^+ is the set of relevant test items for user u . Precision measures the system's ability to reject any non-relevant documents in the recommended set. Recall, on the other hand, is defined as the proportion of relevant items that are actually recommended:

$$Recall_u@N = \frac{|Rec_u^N \cap TS_u^+|}{|TS_u^+|}$$

Hence, recall measures the system's ability to find all the relevant documents. Since these two metrics do not pay attention to whether a relevant item was recommended near the top or closer to the cutoff, in Information Retrieval metrics that explicitly assign a gain to each ranking position are usually considered. The discounted cumulative gain (DCG) is a metric of ranking quality that measures the usefulness of a document based on its position in the result list. Since recommendation results may vary in length depending on the user, it is not possible to compare performance among different users, so the cumulative gain at each position should be normalized across users. Hence, normalized discounted cumulative gain, or nDCG, is defined as:

$$nDCG_u@N = \frac{1}{IDCG@N} \sum_{k=1}^N \frac{2^{r_{uk}} - 1}{\log_2(1 + k)}$$

where k is the position of an item in the recommendation list and $IDCG@N$ indicates the score obtained by an ideal ranking of the recommendation list Rec_u^N that contains only relevant items.

For comparison with the proposed GCE metric, we include two complementary baseline metrics based on the absolute deviation between the mean ratings of different groups as defined in [427]:

$$MAD(R^{(i)}, R^{(j)}) = \left| \frac{\sum R^{(i)}}{|R^{(i)}|} - \frac{\sum R^{(j)}}{|R^{(j)}|} \right| \quad (15.7)$$

where $R^{(i)}$ denotes the predicted ratings for all user-item combinations in group i and $|R^{(i)}|$ is its size. Larger values for MAD mean larger differences between the groups, interpreted as unfairness. Given that our proposed GCE in user-fairness evaluation is based on nDCG, we adapt this definition to also compare between average nDCG for each group. We refer to these two baselines as MAD-rating (or MADr) and MAD-ranking (or MADR). Finally, the reported MAD corresponds to the average MAD between all the pairwise combinations within the groups involved, i.e.,

$$MAD = \text{avg}_{i,j}(MAD(R^{(i)}, R^{(j)})) \quad (15.8)$$

When possible, the metrics are computed on a per-user basis and the reported results are the average of these individual values. We have used the RankSys⁴ framework and adopted the *threshold-based relevant items* condition [83]. The evaluation has been performed considering Top-5, Top-10, and Top-20 recommendations for all datasets and a threshold of 3 with respect to a 1-5 scale for deeming items as relevant for all the three datasets.

15.5 Results

In this section, we discuss the results obtained in our experiments.

15.5.1 Analysis of item fairness results

We show in Table 15.9 a comparison of the item-based GCE using popularity as the item feature for the three tested datasets. Due to space constraints, we only show results for cutoff 10 and the nDCG performance metric, since performance at other cutoffs or based on Precision and Recall were similar. Additionally, we note that larger values of nDCG and GCE (i.e., closer to 0, since it is always negative) correspond to the most accurate or fair systems (under a desired fair distribution p_{f_i}); on the other hand, as explained in Section 15.4.5, recommenders with lower values of MAD are assessed by that metric as more fair.

⁴<http://ranksys.org/>

We observe that accuracy (defined through nDCG) and fairness (either defined as our proposed GCE metric, or using the MAD metric as reference) do not usually match each other, in the sense that the best recommenders for one dimension are different to those for a different dimension; for instance, Random is usually the best recommender based on MADR and MADr, whereas BPRMF and UserKNN are the best in terms of nDCG.

Under equality – i.e., p_{f_0} – UserKNN is the recommender system with highest values of GCE (the most fair) in Amazon Toys & Games and Amazon Video Games, whereas Random and BPRSlim are the most fair ones in Amazon Electronics. As a validation of the proposed metric, the MostPopular recommender always obtains higher (better) values of GCE under p_{f_4} , which is the situation where recommending more popular items is deemed (more) fair by the system designer (they have a larger weight in the probability distribution).

If we now focus on the two extreme non-uniform situations (either very long tail or very popular items, i.e., p_{f_1} or p_{f_4} , respectively), Amazon Electronics and Amazon Video Games show similar results, since BPRSlim has the largest values for popular items and Random for long tail items; on the Amazon Toys & Games dataset, on the other hand, BPRSlim is the most fair regarding long tail items and ItemKNN for popular items, even though BPRSlim also shows good values for popular items, consistent with the results found for the other datasets.

Hence, we conclude that ItemKNN, BPRSlim, and UserKNN are prone to suggest more items from the head of the distribution, although this does not mean they do not recommend tail items, since the values obtained when less popular items are promoted through the fair distribution are not too small either – as it is the case of the MostPopular algorithm, which only recommends items from the fourth category of items, and thus the final GCE value gets distorted by the near-zero probability of the other categories –; SVD++, on the other hand, and BPRMF to a lesser extent (since this depends on the dataset), seems to be tailored to promote mostly popular items, producing similar values as those obtained by the MostPopular algorithm. These results agree with previous observations on the biases evidenced by different algorithms in several datasets [49, 194, 66]. Moreover, if we look at the results

through the lens of which models promote recommendation of long tail items, we can see that BPRSlim and Random are the most capable methods.

For the sake of space, from now on we focus our attention to the analysis of the Amazon Toys & Games dataset, the rest are shown in Appendix 15.5.3. Hence, Table 15.10 shows the item-based GCE values obtained using price as the item feature. In this case, and in contrast to the scenario where popularity is used as the item feature, the MostPopular recommender does not show an obvious pattern, since it obtains higher values for p_{f_2} and p_{f_4} ; this is probably due to the inherent biases in the data, indicating that popular items tend to appear in the low-to-medium and high price clusters. These patterns in the data are also evident when checking the results for Random, where the same two clusters (p_{f_2} and p_{f_4}) produce the highest GCE values.

As with the popularity feature, UserKNN is the best method under equality constraints on the same dataset, however, the situation changes drastically when other fair distributions are considered, since the nature of the item features is very different. For instance, now the method that produces more fair recommendations for more expensive items (p_{f_4}) is BPRMF, followed by ItemKNN. Regarding the least expensive items, UserKNN performs the best, followed by BPRSlim, which also obtains good performance values in the two intermediate clusters. These results provide interesting insights into the performance of these models. We can observe that nDCG produces results for most recommendation models that are too close - without providing any transparent/distinguishable difference; under the proposed GCE metric, we obtain more transparent/detailed understanding of how these models behave with respect to promotion of more expensive items against cheap items. This capability might be important for the system designer in order to know which models to choose on different e-commerce settings.

Based on this, we conclude that there are some recommendation techniques more prone to recommend expensive or cheap products, even when this information is not included in the training data of the algorithms. Thus, we observe that ItemKNN, UserKNN, and BPRMF tend to include more expensive items in their recommendations (also considering the results for the other domains as pre-

sented in Table 15.13). However, as discussed before, some of these results might be attributed to the inherent biases of some algorithms to produce more popular items. This effect is, indeed, not negligible, although it depends strongly on the domain: the MostPopular algorithm obtains very good values of ItemGCE in Amazon Electronics for all distributions except when the most expensive items are promoted, a similar situation is found in Amazon Video Games, although the obtained values are larger, whereas the optimal cases in Amazon Toys & Games are found for p_{f_2} and p_{f_4} .

15.5.2 Analysis of user fairness results

In this section, we analyze the user variation of the GCE metric. For this, in Table 15.11 we show the results of the three tested user features (happiness, helpfulness, and interactions) on the Amazon Toys & Games dataset; results for the other datasets are included in the appendix.

The first thing we notice when considering equality as fairness is that the values are much smaller than in the case of item features, even reaching an optimal of 0, for BPRMF in the happiness feature, although the other optimal recommenders in the other datasets for this scenario also obtain values very close to the optimal one.

Let us now analyze the other scenarios, where fairness is not equivalent to equality. In this case, there is no recommender that obtains a perfect value, although again happiness seems to be the easiest feature where something similar to perfect fairness might be achieved, since BPRSlim shows a -0.02 value for p_{f_2} . ItemKNN and BPRSlim tend to obtain good performance values for the three features, in particular ItemKNN is the best recommender for the users with more interactions, together with those users with least helpful reviews; BPRSlim on the other hand is the best one for the users with least interactions and for the happiest users.

In summary, we conclude that ItemKNN is inherently tailored to users with many interactions and less helpful users, whereas BPRSlim seems to provide better and fair recommendations for happy and cold (few interactions) users. Interestingly, although both of these models leverage item-item similarities based on user interactions, we can observe that, under GCE, we find contrasting results with respect to

their performance in the study of item and user fairness.

15.5.3 Discussion

When analyzing the presented approach and reported results from a global point of view, we can finally answer the three research questions posed at the beginning of the chapter. First, regarding (RQ1) *How to define a fairness evaluation metric that considers different notions of fairness (not only equality)?*, we have presented a novel metric that seamlessly work with either user or item features while, at the same time, it is sensitive to different notions of fairness (through the definition of a specific fair distribution): either based on equality (by using a uniform distribution) or favoring some of the attribute values (such as most expensive items or less happy users). This is a critical difference with respect to other metrics proposed in the literature to measure fairness, which should be tailored to either users or items or that implicitly assume equality as fairness (see Section 15.2). In our experiments, this becomes obvious when comparing the results found for the proposed GCE against those found for MAD-based metrics, since the optimal recommender in the latter case is usually Random, mostly because this type of algorithm is unbiased by definition. However, the proposed GCE metric allows to capture other concepts typically considered when evaluating recommender systems such as relevance and ranking.

Second, as an answer to (RQ2) *How do classical recommendation models behave in terms of such an evaluation metric, especially under non-equality definitions of fairness?*, we summarize the results obtained as follows. Recommendation algorithms based on neighbors performs well in general: whereas UserKNN performs well under equality for item attributes, ItemKNN (together with BPRMF) perform well either under equality or non-equality constraints. Additionally, BPRSlim produces fair results under extreme scenarios of fairness (i.e., p_{f_1} or p_{f_4}), again for item attributes. These conclusions also apply, to some extent, to the results not discussed so far which are shown in the appendix. It should be considered that the presented results correspond to the values obtained when optimizing for accuracy (the recommenders were selected according to their nDCG@5 values), hence, a slightly different behavior could have been obtained if each metric was optimized indepen-

dently. We do not include these results because we are more interested in analyzing how state-of-the-art algorithms (typically selected and assessed with respect to accuracy metrics) behave with respect to fairness oriented metrics.

Finally, to answer (RQ3) *Which user and item attributes are more sensitive to different notions of fairness? Which ones impact more on different families of recommendation algorithms?*, we first need to define what we mean by sensitivity to (notions of) fairness. Thus, we assume this can be understood as those cases where results for equality differ too much from results for non-equality. To properly analyze this issue, we compare the rankings obtained for all the tested recommenders (not only the 7 presented which correspond to those with optimal parameters, but the 95 combinations for all parameters) and compute Spearman correlation between the results using the distribution under equality constraints and the other cases (see Table 15.12). We observe that the item popularity is more or less stable, whereas the item price depends heavily on the dataset; on the other hand, the user attributes (helpfulness, interactions, and especially happiness) are the least stable, since their correlations are the lowest ones. This evidences that user attributes are more sensitive to different notions of fairness, since the performance of recommenders change more drastically when equality and non-equality distributions are used.

Full results

In this section we present the results for all the datasets and item and user attributes that were not included in the previous sections for the sake of readability. First, we show in Table 15.13 the item GCE based on the price attribute for the three datasets (instead of only limited to toys, as in Section 15.5.1).

Second, our results on user attributes – that is, interactions, helpfulness, happiness – is presented for the three datasets: Amazon Toys & Games is described in Sections 15.5.1 and 15.5.2, Amazon Electronics in Table 15.14, and Amazon Video Games in Table 15.15.

15.6 Conclusions and future work

In this study, we proposed a flexible, probabilistic framework for evaluating user and item fairness in recommender systems. We conducted extensive experiments on real-world datasets and demonstrate the flexibility of the proposed solution in various settings. In summary, our framework can evaluate fairness beyond equality, can evaluate both user fairness and item fairness, and is designed based on theoretically sound foundations which makes it interpretable. Analyzing the results from the conducted experiments, we observe that an evaluation based on the item fairness as defined in the RecSys Challenge—that is according to the types of users (regular vs. premium)—captures additional nuances about the different submissions.⁵ For instance, the proposed winner system produces balanced recommendations across the two membership classes. This is in contrast to our expectation that premium users should be provided better recommendations. Therefore, even though the winning submission could produce more precise recommendation, it does not meet with our expectation of a fair recommendation for this attribute, which is to recommend better recommendations to premium users.

On the other hand, when exploiting user attributes in a classical recommendation task to evaluate user fairness, we observe interesting insights related to the different recommendation algorithms. So far, we have studied the case where users are clustered according to their activity in the system, but we aim to analyze other scenarios for the final submission. In this case, we have found that algorithms with very similar performance values obtain very different values of user fairness, mostly because the recommendation methods behave strikingly different at each user cluster. Additionally, we compare our proposed metric against baseline metrics defined in the literature (such as MAD [427]), which are extended to be also suitable for ranking scenarios; it becomes evident that these metrics, cannot incorporate other definitions of fairness in its computation, hence their flexibility is very limited.

In the future, we intend to simultaneously incorporate user and item fairness into the generalized cross-entropy computation, in order to evaluate both multiple

⁵In this challenge, the users correspond to the items being recommended.

objectives in a single framework. In addition, conducting user studies to understand the correlation between user satisfactions and fairness computed using GCE is an interesting future direction that we would like to pursue.

Table 15.9: ItemGCE using popularity as feature on the three tested datasets. The fair probability distributions are defined as p_{f_i} so that $p_{f_i}(j) = 0.1$ when $j \neq i$ and 0.7 otherwise – except for p_{f_0} that denotes the uniform distribution – and each column denotes the value obtained by GCE when such probability distribution is used as p_f in Equation 15.3. In bold, highlighted the best values for each metric.

(a) Amazon Electronics

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|---------------|---------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -7.97 | -32.39 | -32.39 | -1.14 | -2.51 | 0.002 | 0.000 |
| MostPopular | 0.008 | -688.97 | -1,874.78 | -1,874.78 | -1,874.78 | -110.06 | 0.029 | 0.751 |
| ItemKNN | 0.004 | -445.87 | -82.79 | -1,778.92 | -1,778.92 | -71.16 | 0.021 | 0.000 |
| UserKNN | 0.014 | -520.00 | -4,074.78 | -85.48 | -85.24 | -83.08 | 0.043 | 0.001 |
| SVD++ | 0.012 | -1,082.10 | -2,944.10 | -2,944.10 | -2,944.10 | -172.96 | 0.045 | 0.029 |
| BPRMF | 0.018 | -1,299.14 | -3,534.45 | -3,534.45 | -3,534.45 | -207.68 | 0.012 | 0.011 |
| BPRSlim | 0.007 | -9.05 | -38.61 | -22.80 | -14.74 | -1.28 | 0.005 | 0.003 |

(b) Amazon Toys & Games

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -15.98 | -2.38 | -44.25 | -44.25 | -44.25 | 0.000 | 0.000 |
| MostPopular | 0.001 | -126.91 | -345.97 | -345.97 | -345.97 | -20.13 | 0.008 | 0.045 |
| ItemKNN | 0.002 | -0.08 | -1.52 | -0.42 | -0.52 | -0.33 | 0.055 | 0.001 |
| UserKNN | 0.004 | -0.01 | -0.81 | -0.38 | -0.54 | -0.53 | 0.028 | 0.001 |
| SVD++ | 0.003 | -305.36 | -831.35 | -831.35 | -831.35 | -48.68 | 0.006 | 0.004 |
| BPRMF | 0.002 | -146.48 | -24.61 | -586.48 | -586.48 | -23.30 | 0.010 | 0.008 |
| BPRSlim | 0.003 | -0.12 | -0.12 | -1.40 | -1.04 | -0.62 | 0.011 | 0.036 |

(c) Amazon Video Games

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -11.85 | -1.87 | -48.40 | -2.11 | -48.40 | 0.001 | 0.001 |
| MostPopular | 0.004 | -490.77 | -1,335.68 | -1,335.68 | -1,335.68 | -78.34 | 0.017 | 0.116 |
| ItemKNN | 0.013 | -1.25 | -3.68 | -1.03 | -7.72 | -0.11 | 0.100 | 0.002 |
| UserKNN | 0.019 | -1.23 | -10.09 | -0.95 | -1.13 | -0.18 | 0.084 | 0.002 |
| SVD++ | 0.005 | -499.25 | -1,358.75 | -1,358.75 | -1,358.75 | -79.70 | 0.023 | 0.017 |
| BPRMF | 0.008 | -3.08 | -2.18 | -17.12 | -8.14 | -0.36 | 0.024 | 0.032 |
| BPRSlim | 0.011 | -1.45 | -3.18 | -8.48 | -2.45 | -0.11 | 0.034 | 0.019 |

Table 15.10: ItemGCE using price as feature on Amazon Toys & Games. Notation as in Table ??.

| Rec | nDCG | P_{f_0} | P_{f_1} | P_{f_2} | P_{f_3} | P_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -11.95 | -48.74 | -1.84 | -48.74 | -2.29 | 0.001 | 0.000 |
| MostPopular | 0.001 | -84.80 | -339.23 | -15.84 | -339.23 | -13.41 | 0.028 | 0.149 |
| ItemKNN | 0.002 | -0.15 | -1.88 | -0.60 | -0.80 | -0.14 | 0.009 | 0.000 |
| UserKNN | 0.004 | -0.07 | -0.43 | -1.17 | -0.92 | -0.22 | 0.025 | 0.002 |
| SVD++ | 0.003 | -203.82 | -33.81 | -815.83 | -815.83 | -32.47 | 0.022 | 0.017 |
| BPRMF | 0.002 | -0.79 | -2.67 | -4.47 | -1.52 | -0.03 | 0.017 | 0.014 |
| BPRSlim | 0.003 | -0.29 | -0.57 | -0.34 | -0.32 | -3.37 | 0.016 | 0.053 |

Table 15.11: UserGCE for Amazon Toys & Games dataset using the three user features considered. Notation as in Table ??, except for the Happiness attribute, where $p_{f_i}(j) = 0.1$ when $j \neq i$ and 0.9 otherwise when used as p_f in Equation ??.

(a) Happiness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | MADR | MADr |
|-------------|--------------|-------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -4.05 | -13.83 | -0.09 | 0.000 | 0.000 |
| MostPopular | 0.001 | -0.01 | -0.44 | -0.23 | 0.000 | 0.121 |
| ItemKNN | 0.002 | -0.25 | -1.44 | -0.04 | 0.002 | 0.008 |
| UserKNN | 0.004 | -0.20 | -1.23 | -0.05 | 0.003 | 0.049 |
| SVD++ | 0.003 | -0.01 | -0.43 | -0.23 | 0.001 | 0.018 |
| BPRMF | 0.002 | 0.00 | -0.33 | -0.31 | 0.000 | 0.104 |
| BPRSlim | 0.003 | -2.01 | -7.20 | -0.02 | 0.003 | 0.939 |

(b) Helpfulness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -6.24 | -0.96 | -25.93 | -1.25 | -25.93 | 0.000 | 0.000 |
| MostPopular | 0.001 | -0.26 | -2.27 | -0.14 | -1.57 | -0.36 | 0.001 | 0.118 |
| ItemKNN | 0.002 | -0.01 | -0.45 | -0.54 | -0.44 | -0.84 | 0.001 | 0.013 |
| UserKNN | 0.004 | -0.10 | -0.91 | -0.14 | -0.60 | -1.32 | 0.003 | 0.031 |
| SVD++ | 0.003 | -0.08 | -1.24 | -0.23 | -0.94 | -0.39 | 0.002 | 0.002 |
| BPRMF | 0.002 | -0.28 | -0.94 | -0.22 | -3.01 | -0.33 | 0.002 | 0.081 |
| BPRSlim | 0.003 | -0.02 | -0.72 | -0.44 | -0.78 | -0.36 | 0.001 | 3.145 |

(c) Interactions

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -9.54 | -26.73 | -26.73 | -26.73 | -1.35 | 0.000 | 0.000 |
| MostPopular | 0.001 | -19.60 | -154.54 | -4.09 | -3.32 | -3.28 | 0.001 | 0.186 |
| ItemKNN | 0.002 | -0.05 | -0.58 | -0.70 | -1.03 | -0.22 | 0.001 | 0.025 |
| UserKNN | 0.004 | -0.02 | -0.84 | -0.55 | -0.63 | -0.31 | 0.001 | 0.038 |
| SVD++ | 0.003 | -0.03 | -0.70 | -0.72 | -0.71 | -0.25 | 0.001 | 0.003 |
| BPRMF | 0.002 | -0.02 | -0.58 | -0.46 | -0.36 | -0.94 | 0.001 | 0.208 |
| BPRSlim | 0.003 | -0.04 | -0.49 | -0.64 | -0.28 | -1.07 | 0.001 | 9.001 |

Table 15.12: Spearman correlation value between recommenders ranked based on GCE values for p_{f_0} and the indicated fair distribution p_f for all the datasets and user and item attributes.

| Attribute | p_f | Amazon Electronics | Amazon Toys & Games | Amazon Video Games |
|--------------|-----------|--------------------|---------------------|--------------------|
| Price | p_{f_1} | 0.10 | 0.85 | 0.85 |
| | p_{f_2} | 0.28 | 0.70 | 0.74 |
| | p_{f_3} | 0.17 | 0.78 | 0.52 |
| | p_{f_4} | 0.73 | 0.63 | 0.83 |
| Popularity | p_{f_1} | 0.77 | 0.84 | 0.91 |
| | p_{f_2} | 0.95 | 0.93 | 0.78 |
| | p_{f_3} | 0.93 | 0.93 | 0.84 |
| | p_{f_4} | 0.99 | 0.92 | 0.94 |
| Happiness | p_{f_1} | 1.00 | 0.63 | 0.59 |
| | p_{f_2} | -1.00 | -0.29 | -0.22 |
| Helpfulness | p_{f_1} | 0.10 | 0.74 | 0.50 |
| | p_{f_2} | 0.25 | 0.35 | 0.11 |
| | p_{f_3} | 0.36 | 0.64 | 0.77 |
| | p_{f_4} | 0.58 | 0.43 | 0.52 |
| Interactions | p_{f_1} | 0.66 | 0.88 | 0.41 |
| | p_{f_2} | 0.55 | 0.73 | 0.39 |
| | p_{f_3} | 0.46 | 0.68 | 0.37 |
| | p_{f_4} | -0.21 | 0.29 | 0.21 |

Table 15.13: ItemGCE using price as feature on the three tested datasets. Notation as in Table 15.9.

| (a) Amazon Electronics | | | | | | | | |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Rec | nDCG | P_{f_0} | P_{f_1} | P_{f_2} | P_{f_3} | P_{f_4} | MADR | MADr |
| Random | 0.000 | -0.09 | -0.29 | -1.15 | -0.28 | -1.14 | 0.000 | 0.000 |
| MostPopular | 0.008 | -0.12 | -0.42 | -0.25 | -0.62 | -1.82 | 0.056 | 1.330 |
| ItemKNN | 0.004 | -0.03 | -0.96 | -0.29 | -0.57 | -0.57 | 0.004 | 0.000 |
| UserKNN | 0.014 | -0.02 | -0.52 | -0.43 | -0.42 | -0.98 | 0.015 | 0.000 |
| SVD++ | 0.012 | -361.26 | -58.07 | -57.78 | -62.81 | -2,829.16 | 0.076 | 0.047 |
| BPRMF | 0.018 | -1.47 | -1.01 | -0.66 | -0.30 | -12.42 | 0.019 | 0.020 |
| BPRSlim | 0.007 | -0.09 | -0.50 | -0.36 | -0.39 | -1.62 | 0.001 | 0.000 |

| (b) Amazon Toys & Games | | | | | | | | |
|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Rec | nDCG | P_{f_0} | P_{f_1} | P_{f_2} | P_{f_3} | P_{f_4} | MADR | MADr |
| Random | 0.000 | -11.95 | -48.74 | -1.84 | -48.74 | -2.29 | 0.001 | 0.000 |
| MostPopular | 0.001 | -84.80 | -339.23 | -15.84 | -339.23 | -13.41 | 0.028 | 0.149 |
| ItemKNN | 0.002 | -0.15 | -1.88 | -0.60 | -0.80 | -0.14 | 0.009 | 0.000 |
| UserKNN | 0.004 | -0.07 | -0.43 | -1.17 | -0.92 | -0.22 | 0.025 | 0.002 |
| SVD++ | 0.003 | -203.82 | -33.81 | -815.83 | -815.83 | -32.47 | 0.022 | 0.017 |
| BPRMF | 0.002 | -0.79 | -2.67 | -4.47 | -1.52 | -0.03 | 0.017 | 0.014 |
| BPRSlim | 0.003 | -0.29 | -0.57 | -0.34 | -0.32 | -3.37 | 0.016 | 0.053 |

| (c) Amazon Video Games | | | | | | | | |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Rec | nDCG | P_{f_0} | P_{f_1} | P_{f_2} | P_{f_3} | P_{f_4} | MADR | MADr |
| Random | 0.000 | -10.57 | -43.16 | -2.21 | -1.60 | -43.16 | 0.002 | 0.001 |
| MostPopular | 0.004 | -164.68 | -28.64 | -27.11 | -26.27 | -1,290.25 | 0.023 | 0.152 |
| ItemKNN | 0.013 | -0.17 | -2.30 | -0.52 | -0.22 | -0.56 | 0.004 | 0.000 |
| UserKNN | 0.019 | -0.12 | -1.80 | -0.75 | -0.33 | -0.28 | 0.063 | 0.004 |
| SVD++ | 0.005 | -0.70 | -3.22 | -1.07 | -0.04 | -3.70 | 0.034 | 0.025 |
| BPRMF | 0.008 | -0.26 | -1.22 | -0.71 | -0.07 | -2.32 | 0.028 | 0.030 |
| BPRSlim | 0.011 | -0.05 | -1.11 | -0.39 | -0.28 | -0.83 | 0.010 | 0.005 |

Table 15.14: UserGCE for Amazon Electronics dataset using the three user features considered. Notation as in Table 15.11.

(a) Happiness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -1.68 | -6.13 | -0.01 | 0.000 | 0.000 |
| MostPopular | 0.008 | -0.02 | -0.51 | -0.19 | 0.003 | 0.271 |
| ItemKNN | 0.004 | -0.04 | -0.61 | -0.15 | 0.002 | 0.004 |
| UserKNN | 0.014 | -0.04 | -0.63 | -0.15 | 0.007 | 0.003 |
| SVD++ | 0.012 | -0.04 | -0.61 | -0.16 | 0.006 | 0.017 |
| BPRMF | 0.018 | -0.05 | -0.69 | -0.13 | 0.010 | 0.004 |
| BPRSlim | 0.007 | -0.06 | -0.73 | -0.12 | 0.004 | 0.158 |

(b) Helpfulness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -2.53 | -11.04 | -11.04 | -0.35 | -0.81 | 0.000 | 0.000 |
| MostPopular | 0.008 | -0.01 | -0.42 | -0.44 | -0.71 | -0.67 | 0.002 | 0.759 |
| ItemKNN | 0.004 | -0.01 | -0.44 | -0.46 | -0.69 | -0.63 | 0.001 | 0.009 |
| UserKNN | 0.014 | -0.06 | -0.25 | -0.42 | -0.86 | -1.08 | 0.007 | 0.012 |
| SVD++ | 0.012 | -0.03 | -0.28 | -0.52 | -0.82 | -0.75 | 0.004 | 0.003 |
| BPRMF | 0.018 | -0.04 | -0.25 | -0.46 | -0.80 | -1.02 | 0.008 | 0.010 |
| BPRSlim | 0.007 | -0.04 | -0.23 | -0.52 | -0.94 | -0.84 | 0.003 | 0.444 |

(c) Interactions

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|-------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -2.45 | -0.74 | -10.74 | -10.74 | -0.35 | 0.000 | 0.000 |
| MostPopular | 0.008 | -0.01 | -0.67 | -0.59 | -0.57 | -0.39 | 0.001 | 1.805 |
| ItemKNN | 0.004 | -0.16 | -1.24 | -1.67 | -0.41 | -0.15 | 0.003 | 0.027 |
| UserKNN | 0.014 | -0.01 | -0.73 | -0.53 | -0.55 | -0.41 | 0.002 | 0.023 |
| SVD++ | 0.012 | 0.00 | -0.60 | -0.57 | -0.62 | -0.41 | 0.002 | 0.000 |
| BPRMF | 0.018 | 0.00 | -0.54 | -0.61 | -0.58 | -0.45 | 0.002 | 0.013 |
| BPRSlim | 0.007 | -0.05 | -1.11 | -0.74 | -0.48 | -0.25 | 0.003 | 1.302 |

Table 15.15: UserGCE for Amazon Video Games dataset using the three user features considered. Notation as in Table 15.11.

(a) Happiness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | MADR | MADr |
|-------------|--------------|-------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -7.22 | -24.10 | -0.22 | 0.000 | 0.000 |
| MostPopular | 0.004 | -0.05 | -0.66 | -0.14 | 0.002 | 0.373 |
| ItemKNN | 0.013 | 0.00 | -0.37 | -0.27 | 0.002 | 0.006 |
| UserKNN | 0.019 | -0.01 | -0.42 | -0.24 | 0.004 | 0.010 |
| SVD++ | 0.005 | -0.48 | -2.22 | -0.01 | 0.006 | 0.618 |
| BPRMF | 0.008 | -0.12 | -0.97 | -0.08 | 0.007 | 0.038 |
| BPRSlim | 0.011 | 0.00 | -0.26 | -0.39 | 0.002 | 0.099 |

(b) Helpfulness

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -7.58 | -31.08 | -1.95 | -1.10 | -31.08 | 0.000 | 0.000 |
| MostPopular | 0.004 | -0.01 | -0.43 | -0.65 | -0.45 | -0.67 | 0.001 | 0.201 |
| ItemKNN | 0.013 | -0.05 | -0.32 | -0.44 | -1.30 | -0.56 | 0.005 | 0.015 |
| UserKNN | 0.019 | -0.10 | -0.16 | -0.50 | -1.21 | -1.08 | 0.012 | 0.006 |
| SVD++ | 0.005 | -0.44 | -0.42 | -0.56 | -4.51 | -0.29 | 0.003 | 0.103 |
| BPRMF | 0.008 | -0.29 | -0.23 | -0.28 | -2.92 | -1.14 | 0.006 | 0.103 |
| BPRSlim | 0.011 | -0.10 | -0.25 | -0.33 | -1.47 | -0.93 | 0.006 | 0.346 |

(c) Interactions

| Rec | nDCG | p_{f_0} | p_{f_1} | p_{f_2} | p_{f_3} | p_{f_4} | MADR | MADr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random | 0.000 | -6.46 | -0.94 | -26.69 | -26.69 | -1.60 | 0.000 | 0.000 |
| MostPopular | 0.004 | -0.13 | -1.00 | -1.68 | -0.30 | -0.25 | 0.003 | 0.356 |
| ItemKNN | 0.013 | -0.05 | -1.02 | -0.22 | -0.71 | -0.60 | 0.006 | 0.046 |
| UserKNN | 0.019 | -0.04 | -1.13 | -0.41 | -0.51 | -0.40 | 0.006 | 0.014 |
| SVD++ | 0.005 | -0.11 | -1.73 | -0.59 | -0.60 | -0.19 | 0.003 | 0.047 |
| BPRMF | 0.008 | -0.04 | -0.44 | -0.43 | -0.44 | -1.14 | 0.002 | 0.327 |
| BPRSlim | 0.011 | -0.08 | -0.59 | -0.45 | -1.51 | -0.27 | 0.005 | 1.260 |

Part IV

Conclusion and Future Work

Before continuing, some warnings are required. This short final chapter has a quite different tone. Most of this research thesis is written in the first plural person. I have chosen this way because I know that all my achievements are not only mines. My supervisors and all the researchers I have a collaboration with have contributed. Without their effort, my ideas, my experiments, would still be rough. However, this chapter is devoted to my overall conclusions of this research work. For this reason, if you are looking for the results of the research lines, I am sad to say that you are in the wrong place. I have written each specific research line section to be self-conclusive. These lines will be devoted to a broader vision of this research work. In these years, we have analyzed the impact of semantic knowledge on Recommender Systems. We have analyzed the different ways to inject semantic knowledge in the best recommenders. We have proposed several new factorization techniques to factorize user-feature matrices, to fix the meaning of the latent factors, to compute the best and the worst factors. We have written SPARQL queries able to exploit preference theories to produce recommendations. We have proposed new high-performance interpretable models. Then, we also have explored the world of semi-structured knowledge. In this sense, we have analyzed the role of Time, and the Popularity of items. We have exploited Time to build a new diversification algorithm and to define a new personalized popularity algorithm. We have investigated the notion of similarity and we have proposed a new family of similarities. We have explored the evaluation of Recommender Systems proposing a new generalized fairness metric and a technique to analyze the discriminative power of evaluation metrics during the training phase. All the results of the above approaches are competitive or better than the state-of-art. All our findings have been peer-reviewed and some extensions are still under-review. There is still room for improvement for each of these techniques. I have indicated some of these possible extensions at the end of each research line. However, some of the most important research pieces of this research path are not in this work. Failed experiments, wrong settings, and imaginative theories have marked this research path. Without those failures, achievements would not be there. This research path, intended as defeats and victories, has deeply changed me. My point of view is changed. In the beginning, the idea of obtaining

bad results terrified me. Now I am proud of my failures because I have learned to analyze them with the wish of understanding them. I have learned to be very rigorous preparing experiments and formalizing a new theory. When I see the results, I have learned to look further ahead. More important, I have learned that I have a lot of limits. Comparing against research giants forces you to analyze yourself. As for the proposals you have read, also for me there is room for improvement. Writing style, abstracting at a very high level, connecting theoretical points are only some of the aspects I will work on it. In my humble opinion, my research work is only a small rock of a mountain that could be built exploiting knowledge representation and machine learning together. Interpretability of black-boxes, evaluation of fairness, automatic completion of knowledge bases are only some of the topics these technologies can face. On the other side, the knowledge representation background has been crucial also for dealing with semi-structured knowledge, proposing new time-aware and dissimilarity-based algorithms.

There is still much to discover, and failures to understand.

Bibliography

- [1] Title vii of the civil rights act of 1964. <https://www.eeoc.gov/laws/statutes/titlevii.cfm>. Accessed: 2019-07-31.
- [2] Dynamic item-based recommendation algorithm with time decay. In *Sixth International Conference on Natural Computation, ICNC 2010, Yantai, Shandong, China, 10-12 August 2010*, pages 242–247. IEEE, 2010. Withdrawn.
- [3] *Netflix TechBlog*, 2012 (accessed June 5, 2019). <http://techblog.netflix.com/2012/04/netflix-recommendations-beyond-5-stars.html>.
- [4] *Amazon product data*, 2014 (accessed July 31, 2019). <http://jmcauley.ucsd.edu/data/amazon/>.
- [5] *Etsy.com-Shop for anything from creative people everywhere*, 2019 (accessed August 12, 2019). <https://www.etsy.com>.
- [6] B. Abdollahi and O. Nasraoui. Explainable matrix factorization for collaborative filtering. In J. Bourdeau, J. Hendler, R. Nkambou, I. Horrocks, and B. Y. Zhao, editors, *Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11-15, 2016, Companion Volume*, pages 5–6. ACM, 2016.
- [7] B. Abdollahi and O. Nasraoui. Explainable restricted boltzmann machines for collaborative filtering. *CoRR*, abs/1606.07129, 2016.

- [8] H. Abdollahpouri, G. Adomavicius, R. Burke, I. Guy, D. Jannach, T. Kamishima, J. Krasnodebski, and L. A. Pizzato. Beyond personalization: Research directions in multistakeholder recommendation. *CoRR*, abs/1905.01986, 2019.
- [9] H. Abdollahpouri, R. Burke, and B. Mobasher. Recommender systems as multistakeholder environments. In *Proceedings of the 25th Conference on User Modeling, Adaptation and Personalization, UMAP 2017, Bratislava, Slovakia, July 09 - 12, 2017*, pages 347–348, 2017.
- [10] R. Abebe, J. M. Kleinberg, and D. C. Parkes. Fair division via social comparison. In K. Larson, M. Winikoff, S. Das, and E. H. Durfee, editors, *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems, AAMAS 2017, São Paulo, Brazil, May 8-12, 2017*, pages 281–289. ACM, 2017.
- [11] F. Abel, S. Araújo, Q. Gao, and G. Houben. Analyzing cross-system user modeling on the social web. In S. Auer, O. Díaz, and G. A. Papadopoulos, editors, *Web Engineering - 11th International Conference, ICWE 2011, Paphos, Cyprus, June 20-24, 2011*, volume 6757 of *Lecture Notes in Computer Science*, pages 28–43. Springer, 2011.
- [12] F. Abel, E. Herder, G. Houben, N. Henze, and D. Krause. Cross-system user modeling and personalization on the social web. *User Model. User-Adapt. Interact.*, 23(2-3):169–209, 2013.
- [13] S. Abiteboul, P. Buneman, D. Suciu, J. Allan, R. Papka, V. Lavrenko, A. Blum, T. Mitchell, A. Bonifati, S. Ceri, et al. News weeder: Learning to filter netnews. In *Proceedings of the 12th International Conference of Machine Learning (ICML95)*, 1995.
- [14] G. Adomavicius and Y. Kwon. Improving aggregate recommendation diversity using ranking-based techniques. *IEEE Trans. Knowl. Data Eng.*, 24(5):896–911, 2012.

- [15] G. Adomavicius and Y. Kwon. Multi-criteria recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 847–880. Springer, 2015.
- [16] G. Adomavicius and A. Tuzhilin. Multidimensional recommender systems: A data warehousing approach. In L. Fiege, G. Mühl, and U. G. Wilhelm, editors, *Electronic Commerce, Second International Workshop, WELCOM 2001 Heidelberg, Germany, November 16-17, 2001, Proceedings*, volume 2232 of *Lecture Notes in Computer Science*, pages 180–192. Springer, 2001.
- [17] G. Adomavicius and A. Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. Knowl. Data Eng.*, 17(6):734–749, 2005.
- [18] G. Adomavicius and A. Tuzhilin. Context-aware recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 191–226. Springer, 2015.
- [19] F. Aiolli. A preliminary study on a recommender system for the million songs dataset challenge. In R. Basili, F. Sebastiani, and G. Semeraro, editors, *Proceedings of the 4th Italian Information Retrieval Workshop, Pisa, Italy, January 16-17, 2013*, volume 964 of *CEUR Workshop Proceedings*, pages 73–83. CEUR-WS.org, 2013.
- [20] L. Akoglu and C. Faloutsos. Valuepick: Towards a value-oriented dual-goal recommender system. In W. Fan, W. Hsu, G. I. Webb, B. Liu, C. Zhang, D. Gunopulos, and X. Wu, editors, *ICDMW 2010, The 10th IEEE International Conference on Data Mining Workshops, Sydney, Australia, 13 December 2010*, pages 1151–1158. IEEE Computer Society, 2010.
- [21] V. W. Anelli, P. Basile, D. G. Bridge, T. D. Noia, P. Lops, C. Musto, F. Narducci, and M. Zanker. Knowledge-aware and conversational recommender systems. In S. Pera, M. D. Ekstrand, X. Amatriain, and J. O’Donovan, editors, *Proceedings of the 12th ACM Conference on Recommender Systems*,

- RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018*, pages 521–522. ACM, 2018.
- [22] V. W. Anelli, V. Bellini, A. Cali, G. D. Santis, T. D. Noia, and E. D. Sciascio. Querying deep web data sources as linked data. In R. Akerkar, A. Cuzzocrea, J. Cao, and M. Hacid, editors, *Proceedings of the 7th International Conference on Web Intelligence, Mining and Semantics, WIMS 2017, Amantea, Italy, June 19-22, 2017*, pages 32:1–32:7. ACM, 2017.
- [23] V. W. Anelli, V. Bellini, T. D. Noia, W. L. Bruna, P. Tomeo, and E. D. Sciascio. An analysis on time- and session-aware diversification in recommender systems. In M. Bieliková, E. Herder, F. Cena, and M. C. Desmarais, editors, *Proceedings of the 25th Conference on User Modeling, Adaptation and Personalization, UMAP 2017, Bratislava, Slovakia, July 09 - 12, 2017*, pages 270–274. ACM, 2017.
- [24] V. W. Anelli, A. Cali, T. D. Noia, M. Palmonari, and A. Ragone. Exposing open street map in the linked data cloud. In H. Fujita, M. Ali, A. Selamat, J. Sasaki, and M. Kurematsu, editors, *Trends in Applied Knowledge-Based Systems and Data Science - 29th International Conference on Industrial Engineering and Other Applications of Applied Intelligent Systems, IEA/AIE 2016, Morioka, Japan, August 2-4, 2016, Proceedings*, volume 9799 of *Lecture Notes in Computer Science*, pages 344–355. Springer, 2016.
- [25] V. W. Anelli and T. D. Noia. 2nd workshop on knowledge-aware and conversational recommender systems - kars. In W. Zhu, D. Tao, X. Cheng, P. Cui, E. A. Rundensteiner, D. Carmel, Q. He, and J. X. Yu, editors, *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019.*, pages 3001–3002. ACM, 2019.
- [26] V. W. Anelli, T. D. Noia, P. Lops, C. Musto, M. Zanker, P. Basile, D. G. Bridge, and F. Narducci, editors. *Proceedings of the Workshop on Knowledge-aware and Conversational Recommender Systems*

2018 co-located with 12th ACM Conference on Recommender Systems, *KaRS@RecSys 2018, Vancouver, Canada, October 7, 2018*, volume 2290 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2019.

- [27] V. W. Anelli, T. D. Noia, P. Lops, and E. D. Sciascio. Feature factorization for top-n recommendation: From item rating to features relevance. In Y. Zheng, W. Pan, S. S. Sahebi, and I. Fernández, editors, *Proceedings of the 1st Workshop on Intelligent Recommender Systems by Knowledge Transfer & Learning co-located with ACM Conference on Recommender Systems (RecSys 2017), Como, Italy, August 27, 2017.*, volume 1887 of *CEUR Workshop Proceedings*, pages 16–21. CEUR-WS.org, 2017.
- [28] V. W. Anelli, T. D. Noia, E. D. Sciascio, C. Pomo, and A. Ragone. On the discriminative power of hyper-parameters in cross-validation and how to choose them. In T. Bogers, A. Said, P. Brusilovsky, and D. Tikk, editors, *Proceedings of the 13th ACM Conference on Recommender Systems, RecSys 2017, Copenhagen, Denmark, September 16-20, 2019.*, pages 447–451. ACM, 2019.
- [29] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta. Time-aware personalized popularity in top-n recommendation. In *Workshop on Recommendation in Complex Scenarios co-located with 12th ACM Conference on Recommender Systems (RecSys 2018), Vancouver, BC, Canada, October 2-7, 2018*, 2018.
- [30] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta. How to make latent factors interpretable by feeding factorization machines with knowledge graphs. In C. Ghidini, O. Hartig, M. Maleshkova, V. Svátek, I. F. Cruz, A. Hogan, J. Song, M. Lefrançois, and F. Gandon, editors, *The Semantic Web - ISWC 2019 - 18th International Semantic Web Conference, Auckland, New Zealand, October 26-30, 2019, Proceedings, Part I*, volume 11778 of *Lecture Notes in Computer Science*, pages 38–56. Springer, 2019.

- [31] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta. The importance of being dissimilar in recommendation. In C. Hung and G. A. Papadopoulos, editors, *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, SAC 2019, Limassol, Cyprus, April 8-12, 2019*, pages 816–821. ACM, 2019.
- [32] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta. Local popularity and time in top-n recommendation. In L. Azzopardi, B. Stein, N. Fuhr, P. Mayr, C. Hauff, and D. Hiemstra, editors, *Advances in Information Retrieval - 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14-18, 2019, Proceedings, Part I*, volume 11437 of *Lecture Notes in Computer Science*, pages 861–868. Springer, 2019.
- [33] A. Ashkan, B. Kveton, S. Berkovsky, and Z. Wen. Optimal greedy diversity for recommendation. In Q. Yang and M. J. Wooldridge, editors, *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 1742–1748. AAAI Press, 2015.
- [34] A. Azaria, A. Hassidim, S. Kraus, A. Eshkol, O. Weintraub, and I. Netanel. Movie recommender system for profit maximization. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 121–128. ACM, 2013.
- [35] R. A. Baeza-Yates and B. A. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [36] M. Balabanovic and Y. Shoham. Content-based, collaborative recommendation. *Commun. ACM*, 40(3):66–72, 1997.
- [37] L. Baltrunas and F. Ricci. Item weighting techniques for collaborative filtering. In B. Berendt, D. Mladenic, M. de Gemmis, G. Semeraro, M. Spiliopoulou, G. Stumme, V. Svátek, and F. Zelezný, editors, *Knowledge*

Discovery Enhanced with Semantic and Social Information, volume 220 of *Studies in Computational Intelligence*, pages 109–126. 2009.

- [38] H. Bao, Q. Li, S. S. Liao, S. Song, and H. Gao. A new temporal and social pmf-based method to predict users' interests in micro-blogging. *Decision Support Systems*, 55(3):698–709, 2013.
- [39] S. Barocas and A. D. Selbst. Big data's disparate impact. *Cal. L. Rev.*, 104:671, 2016.
- [40] K. Bauman, B. Liu, and A. Tuzhilin. Aspect based recommendations: Recommending items with the most valuable aspects based on user reviews. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 717–725. ACM, 2017.
- [41] I. Bayer. fastfm: A library for factorization machines. *J. Mach. Learn. Res.*, 17:184:1–184:5, 2016.
- [42] F. Belém, R. L. T. Santos, J. M. Almeida, and M. A. Gonçalves. Topic diversity in tag recommendation. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 141–148. ACM, 2013.
- [43] N. J. Belkin and S. E. Robertson. Some ethical and political implications of theoretical research in information science. In *Proceedings of the Association for Information Science, ASIS '76*, pages 597–605, 1976.
- [44] R. M. Bell and Y. Koren. Lessons from the netflix prize challenge. *SIGKDD Explorations*, 9(2):75–79, 2007.
- [45] R. M. Bell and Y. Koren. Scalable collaborative filtering with jointly derived neighborhood interpolation weights. In *Proceedings of the 7th IEEE International Conference on Data Mining (ICDM 2007), October 28-31, 2007, Omaha, Nebraska, USA*, pages 43–52. IEEE Computer Society, 2007.

- [46] R. M. Bell, Y. Koren, and C. Volinsky. Modeling relationships at multiple scales to improve accuracy of large recommender systems. In P. Berkhin, R. Caruana, and X. Wu, editors, *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Jose, California, USA, August 12-15, 2007*, pages 95–104. ACM, 2007.
- [47] A. Bellogín, I. Cantador, and P. Castells. A study of heterogeneity in recommendations for a social music service. In *Proceedings of the 1st International Workshop on Information Heterogeneity and Fusion in Recommender Systems, HetRec '10*, pages 1–8, New York, NY, USA, 2010. ACM.
- [48] A. Bellogín, P. Castells, and I. Cantador. Precision-oriented evaluation of recommender systems: an algorithmic comparison. In B. Mobasher, R. D. Burke, D. Jannach, and G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 333–336. ACM, 2011.
- [49] A. Bellogín, P. Castells, and I. Cantador. Statistical biases in information retrieval metrics for recommender systems. *Inf. Retr. Journal*, 20(6):606–634, 2017.
- [50] A. Bellogín, I. Fernández-Tobías, I. Cantador, and P. Tomeo. Neighbor selection for cold users in collaborative filtering with positive-only feedback. In F. Herrera, S. Damas, R. Montes, S. Alonso, O. Cordón, A. González, and A. Troncoso, editors, *Advances in Artificial Intelligence - 18th Conference of the Spanish Association for Artificial Intelligence, CAEPIA 2018, Granada, Spain, October 23-26, 2018, Proceedings*, volume 11160 of *Lecture Notes in Computer Science*, pages 3–12. Springer, 2018.
- [51] A. Bellogín and P. Sánchez. Revisiting neighbourhood-based recommenders for temporal scenarios. In M. Bieliková, V. Bogina, T. Kuflik, and R. Sasson, editors, *Proceedings of the 1st Workshop on Temporal Reasoning in Recommender Systems co-located with 11th International Conference on Recom-*

- mender Systems (RecSys 2017), Como, Italy, August 27-31, 2017.*, volume 1922 of *CEUR Workshop Proceedings*, pages 40–44. CEUR-WS.org, 2017.
- [52] J. Bennett, S. Lanning, and N. Netflix. The netflix prize. In *In KDD Cup and Workshop in conjunction with KDD*, 2007.
- [53] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyperparameter optimization. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. C. N. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain.*, pages 2546–2554, 2011.
- [54] J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, 13:281–305, 2012.
- [55] J. Bergstra, B. Komer, C. Eliasmith, D. Yamins, and D. D. Cox. Hyperopt: a python library for model selection and hyperparameter optimization. *Computational Science & Discovery*, 8(1):014008, 2015.
- [56] J. Bergstra, D. Yamins, and D. D. Cox. Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms. In *Proc. of the 12th Python in science conf.*, pages 13–20. Citeseer, 2013.
- [57] S. Berkovsky, T. Kuflik, and F. Ricci. Distributed collaborative filtering with domain specialization. In J. A. Konstan, J. Riedl, and B. Smyth, editors, *Proceedings of the 2007 ACM Conference on Recommender Systems, RecSys 2007, Minneapolis, MN, USA, October 19-20, 2007*, pages 33–40. ACM, 2007.
- [58] S. Berkovsky, T. Kuflik, and F. Ricci. Mediation of user models for enhanced personalization in recommender systems. *User Model. User-Adapt. Interact.*, 18(3):245–286, 2008.
- [59] T. Berners-Lee, R. T. Fielding, and L. Masinter. Uniform resource identifiers (URI): generic syntax. *RFC*, 2396:1–40, 1998.

- [60] T. Berners-Lee, J. Hendler, and O. Lassila. The semantic web. *Scientific American*, 284(5):34–43, May 2001.
- [61] A. Beutel, E. H. Chi, Z. Cheng, H. Pham, and J. R. Anderson. Beyond globally optimal: Focused learning for improved recommendations. In R. Barrett, R. Cummings, E. Agichtein, and E. Gabrilovich, editors, *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, pages 203–212. ACM, 2017.
- [62] A. J. Biega, K. P. Gummadi, and G. Weikum. Equity of attention: Amortizing individual fairness in rankings. In K. Collins-Thompson, Q. Mei, B. D. Davison, Y. Liu, and E. Yilmaz, editors, *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 405–414. ACM, 2018.
- [63] D. Billsus and M. J. Pazzani. User modeling for adaptive news access. *User Model. User-Adapt. Interact.*, 10(2-3):147–180, 2000.
- [64] C. Bizer, T. Heath, and T. Berners-Lee. Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*, 5(3):1–22, 2009.
- [65] Y. Blanco-Fernández, J. J. Pazos-Arias, A. Gil-Solla, M. R. Cabrer, and M. L. Nores. Providing entertainment by content-based filtering and semantic reasoning in intelligent recommender systems. *IEEE Trans. Consumer Electronics*, 54(2):727–735, 2008.
- [66] L. Boratto, G. Fenu, and M. Marras. The effect of algorithmic bias on recommender systems for massive open online courses. In L. Azzopardi, B. Stein, N. Fuhr, P. Mayr, C. Hauff, and D. Hiemstra, editors, *Advances in Information Retrieval - 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14-18, 2019, Proceedings, Part I*, volume 11437 of *Lecture Notes in Computer Science*, pages 457–472. Springer, 2019.
- [67] S. Börzsönyi, D. Kossmann, and K. Stocker. The skyline operator. In D. Georgakopoulos and A. Buchmann, editors, *Proceedings of the 17th In-*

- ternational Conference on Data Engineering, April 2-6, 2001, Heidelberg, Germany*, pages 421–430. IEEE Computer Society, 2001.
- [68] S. Bostandjiev, J. O’Donovan, and T. Höllerer. Tasteweights: a visual interactive hybrid recommender system. In P. Cunningham, N. J. Hurley, I. Guy, and S. S. Anand, editors, *Sixth ACM Conference on Recommender Systems, RecSys ’12, Dublin, Ireland, September 9-13, 2012*, pages 35–42. ACM, 2012.
- [69] Z. I. Botev and D. P. Kroese. The generalized cross entropy method, with applications to probability density estimation. *Methodology and Computing in Applied Probability*, 13(1):1–27, 2011.
- [70] M. R. Bouadjenek, H. Hacid, and M. Bouzeghoub. Social networks and information retrieval, how are they converging? A survey, a taxonomy and an analysis of social information retrieval approaches and platforms. *Inf. Syst.*, 56:1–18, 2016.
- [71] C. Boutilier, R. I. Brafman, C. Domshlak, H. H. Hoos, and D. Poole. Cp-nets: A tool for representing and reasoning with conditional ceteris paribus preference statements. *J. Artif. Intell. Res.*, 21:135–191, 2004.
- [72] R. I. Brafman and C. Domshlak. Introducing variable importance tradeoffs into cp-nets. In A. Darwiche and N. Friedman, editors, *UAI ’02, Proceedings of the 18th Conference in Uncertainty in Artificial Intelligence, University of Alberta, Edmonton, Alberta, Canada, August 1-4, 2002*, pages 69–76. Morgan Kaufmann, 2002.
- [73] J. S. Breese, D. Heckerman, and C. M. Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In G. F. Cooper and S. Moral, editors, *UAI ’98: Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence, University of Wisconsin Business School, Madison, Wisconsin, USA, July 24-26, 1998*, pages 43–52. Morgan Kaufmann, 1998.
- [74] G. Brewka, S. Benferhat, and D. L. Berre. Qualitative choice logic. *Artif. Intell.*, 157(1-2):203–237, 2004.

- [75] G. Brewka, J. P. Delgrande, J. Romero, and T. Schaub. asprin: Customizing answer set preferences without a headache. In B. Bonet and S. Koenig, editors, *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA.*, pages 1467–1474. AAAI Press, 2015.
- [76] E. Brochu, V. M. Cora, and N. de Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *CoRR*, abs/1012.2599, 2010.
- [77] C. Buckley, G. Salton, and J. Allan. The effect of adding relevance information in a relevance feedback environment. In W. B. Croft and C. J. van Rijsbergen, editors, *Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval. Dublin, Ireland, 3-6 July 1994 (Special Issue of the SIGIR Forum)*, pages 292–300. ACM/Springer, 1994.
- [78] R. Burke. Multisided fairness for recommendation. *CoRR*, abs/1707.00093, 2017.
- [79] R. Burke, N. Sonboli, and A. Ordonez-Gauger. Balanced neighborhoods for multi-sided fairness in recommendation. In S. A. Friedler and C. Wilson, editors, *Conference on Fairness, Accountability and Transparency, FAT 2018, 23-24 February 2018, New York, NY, USA*, volume 81 of *Proceedings of Machine Learning Research*, pages 202–214. PMLR, 2018.
- [80] R. D. Burke. Hybrid recommender systems: Survey and experiments. *User Model. User-Adapt. Interact.*, 12(4):331–370, 2002.
- [81] R. D. Burke. Hybrid web recommender systems. In P. Brusilovsky, A. Kobsa, and W. Nejdl, editors, *The Adaptive Web, Methods and Strategies of Web Personalization*, volume 4321 of *Lecture Notes in Computer Science*, pages 377–408. Springer, 2007.

- [82] R. D. Burke, H. Abdollahpouri, B. Mobasher, and T. Gupta. Towards multi-stakeholder utility evaluation of recommender systems. In F. Cena, M. C. Desmarais, and D. Dicheva, editors, *Late-breaking Results, Posters, Demos, Doctoral Consortium and Workshops Proceedings of the 24th ACM Conference on User Modeling, Adaptation and Personalisation (UMAP 2016), Halifax, Canada, July 13-16, 2016.*, volume 1618 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2016.
- [83] P. G. Campos, F. Díez, and I. Cantador. Time-aware recommender systems: a comprehensive survey and analysis of existing evaluation protocols. *User Model. User-Adapt. Interact.*, 24(1-2):67–119, 2014.
- [84] I. Cantador, P. Brusilovsky, and T. Kuflik. Second workshop on information heterogeneity and fusion in recommender systems (hetrec2011). In B. Mobasher, R. D. Burke, D. Jannach, and G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 387–388. ACM, 2011.
- [85] I. Cantador, I. Fernández-Tobías, and A. Bellogín. Relating personality types with user preferences in multiple entertainment domains. In S. Berkovsky, E. Herder, P. Lops, and O. C. Santos, editors, *Late-Breaking Results, Project Papers and Workshop Proceedings of the 21st Conference on User Modeling, Adaptation, and Personalization.*, Rome, Italy, June 10-14, 2013, volume 997 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2013.
- [86] I. Cantador, I. Fernández-Tobías, S. Berkovsky, and P. Cremonesi. Cross-domain recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 919–959. Springer, 2015.
- [87] B. Cao, N. N. Liu, and Q. Yang. Transfer learning for collective link prediction in multiple heterogenous domains. In J. Fürnkranz and T. Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pages 159–166. Omnipress, 2010.

- [88] Z. Cao, T. Qin, T. Liu, M. Tsai, and H. Li. Learning to rank: from pairwise approach to listwise approach. In Z. Ghahramani, editor, *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, volume 227 of *ACM International Conference Proceeding Series*, pages 129–136. ACM, 2007.
- [89] C. Cappiello, T. D. Noia, B. A. Marcu, and M. Matera. A quality model for linked data exploration. In A. Bozzon, P. Cudré-Mauroux, and C. Pautasso, editors, *Web Engineering - 16th International Conference, ICWE 2016, Lugano, Switzerland, June 6-9, 2016. Proceedings*, volume 9671 of *Lecture Notes in Computer Science*, pages 397–404. Springer, 2016.
- [90] J. G. Carbonell and J. Goldstein. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In W. B. Croft, A. Moffat, C. J. van Rijsbergen, R. Wilkinson, and J. Zobel, editors, *SIGIR '98: Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 24-28 1998, Melbourne, Australia*, pages 335–336. ACM, 1998.
- [91] P. Castells, N. J. Hurley, and S. Vargas. Novelty and diversity in recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 881–918. Springer, 2015.
- [92] S. Chakraborty, R. Tomsett, R. Raghavendra, D. Harborne, M. Alzantot, F. Cerutti, M. B. Srivastava, A. D. Preece, S. Julier, R. M. Rao, T. D. Kelley, D. Braines, M. Sensoy, C. J. Willis, and P. Gurrām. Interpretability of deep learning models: A survey of results. In *2017 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation, SmartWorld/SCALCOM/UIC/ATC/CBD-Com/IOP/SCI 2017, San Francisco, CA, USA, August 4-8, 2017*, pages 1–6. IEEE, 2017.

- [93] S. Chan, P. C. Treleaven, and L. Capra. Continuous hyperparameter optimization for large-scale recommender systems. In X. Hu, T. Y. Lin, V. V. Raghavan, B. W. Wah, R. A. Baeza-Yates, G. C. Fox, C. Shahabi, M. Smith, Q. Yang, R. Ghani, W. Fan, R. Lempel, and R. Nambiar, editors, *Proceedings of the 2013 IEEE International Conference on Big Data, 6-9 October 2013, Santa Clara, CA, USA*, pages 350–358. IEEE Computer Society, 2013.
- [94] O. Chapelle, S. Ji, C. Liao, E. Velipasaoglu, L. Lai, and S. Wu. Intent-based diversification of web search results: metrics and algorithms. *Inf. Retr.*, 14(6):572–592, 2011.
- [95] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan. Expected reciprocal rank for graded relevance. In D. W. Cheung, I. Song, W. W. Chu, X. Hu, and J. J. Lin, editors, *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009, Hong Kong, China, November 2-6, 2009*, pages 621–630. ACM, 2009.
- [96] L. Chen, F. Hsu, M. Chen, and Y. Hsu. Developing recommender systems with the consideration of product profitability for sellers. *Inf. Sci.*, 178(4):1032–1048, 2008.
- [97] X. Chen, Z. Qin, Y. Zhang, and T. Xu. Learning to rank features for recommendation over multiple categories. In R. Perego, F. Sebastiani, J. A. Aslam, I. Ruthven, and J. Zobel, editors, *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval, SIGIR 2016, Pisa, Italy, July 17-21, 2016*, pages 305–314. ACM, 2016.
- [98] J. Chomicki. Preference formulas in relational queries. *ACM Trans. Database Syst.*, 28(4):427–466, 2003.
- [99] J. Chomicki. Logical foundations of preference queries. *IEEE Data Eng. Bull.*, 34(2):3–10, 2011.
- [100] E. Christakopoulou and G. Karypis. Local item-item models for top-n recommendation. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Pro-*

- ceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 67–74. ACM, 2016.
- [101] R. Chung, D. Sundaram, and A. Srinivasan. Integrated personal recommender systems. In M. L. Gini, R. J. Kauffman, D. Sarppo, C. Dellarocas, and F. Dignum, editors, *Proceedings of the 9th International Conference on Electronic Commerce: The Wireless World of Electronic Commerce, 2007, University of Minnesota, Minneapolis, MN, USA, August 19-22, 2007*, volume 258 of *ACM International Conference Proceeding Series*, pages 65–74. ACM, 2007.
- [102] C. Cornelio, A. Loreggia, and V. A. Saraswat. Logical conditional preference theories. *CoRR*, abs/1504.06374, 2015.
- [103] P. Covington, J. Adams, and E. Sargin. Deep neural networks for youtube recommendations. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 191–198. ACM, 2016.
- [104] F. A. Cowell. Measurement of inequality. *Handbook of income distribution*, 1:87–166, 2000.
- [105] F. A. Cowell and K. Kuga. Inequality measurement: an axiomatic approach. *European Economic Review*, 15(3):287–305, 1981.
- [106] H. S. M. Cramer, V. Evers, S. Ramlal, M. van Someren, L. Rutledge, N. Stash, L. Aroyo, and B. J. Wielinga. The effects of transparency on trust in and acceptance of a content-based art recommender. *User Model. User-Adapt. Interact.*, 18(5):455–496, 2008.
- [107] P. Cremonesi, F. Garzotto, S. Negro, A. V. Papadopoulos, and R. Turrin. Looking for ”good” recommendations: A comparative evaluation of recommender systems. In P. F. Campos, T. C. N. Graham, J. A. Jorge, N. J. Nunes, P. A. Palanque, and M. Winckler, editors, *Human-Computer Interaction -*

INTERACT 2011 - 13th IFIP TC 13 International Conference, Lisbon, Portugal, September 5-9, 2011, Proceedings, Part III, volume 6948 of *Lecture Notes in Computer Science*, pages 152–168. Springer, 2011.

- [108] P. Cremonesi, Y. Koren, and R. Turrin. Performance of recommender algorithms on top-n recommendation tasks. In X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, editors, *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, pages 39–46. ACM, 2010.
- [109] P. Cremonesi and M. Quadrana. Cross-domain recommendations without overlapping data: myth or reality? In A. Kobsa, M. X. Zhou, M. Ester, and Y. Koren, editors, *Eighth ACM Conference on Recommender Systems, RecSys '14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014*, pages 297–300. ACM, 2014.
- [110] P. Cremonesi, A. Tripodi, and R. Turrin. Cross-domain recommender systems. In M. Spiliopoulou, H. Wang, D. J. Cook, J. Pei, W. Wang, O. R. Zaïane, and X. Wu, editors, *Data Mining Workshops (ICDMW), 2011 IEEE 11th International Conference on, Vancouver, BC, Canada, December 11, 2011*, pages 496–503. IEEE Computer Society, 2011.
- [111] P. Cremonesi, R. Turrin, E. Lentini, and M. Matteucci. An evaluation methodology for collaborative recommender systems. In *2008 Int. Conf. on Automated Solutions for Cross Media Content and Multi-Channel Distribution*, pages 224–231. IEEE, 2008.
- [112] I. Csiszár. A class of measures of informativity of observation channels. *Periodica Mathematica Hungarica*, 2(1-4):191–213, 1972.
- [113] E. Q. da Silva, C. G. Camilo-Junior, L. M. L. Pascoal, and T. C. Rosa. An evolutionary approach for combining results of recommender systems techniques based on collaborative filtering. *Expert Syst. Appl.*, 53:204–218, 2016.

- [114] M. F. Dacrema, A. Gasparin, and P. Cremonesi. Deriving item features relevance from collaborative domain knowledge. In V. W. Anelli, T. D. Noia, P. Lops, C. Musto, M. Zanker, P. Basile, D. G. Bridge, and F. Narducci, editors, *Proceedings of the Workshop on Knowledge-aware and Conversational Recommender Systems 2018 co-located with 12th ACM Conference on Recommender Systems, KaRS@RecSys 2018, Vancouver, Canada, October 7, 2018.*, volume 2290 of *CEUR Workshop Proceedings*, pages 1–4. CEUR-WS.org, 2018.
- [115] A. Das, C. Mathieu, and D. Ricketts. Maximizing profit using recommender systems. *CoRR*, abs/0908.3633, 2009.
- [116] M. de Gemmis, P. Lops, C. Musto, F. Narducci, and G. Semeraro. Semantics-aware content-based recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 119–159. Springer, 2015.
- [117] B. F. de Souza, A. C. P. de Leon Ferreira de Carvalho, R. Calvo, and R. P. Ishii. Multiclass SVM model selection using particle swarm optimization. In N. Kasabov, M. Köppen, A. König, A. Abraham, and Q. Song, editors, *6th International Conference on Hybrid Intelligent Systems (HIS 2006), 13-15 December 2006, Auckland, New Zealand*, page 31. IEEE Computer Society, 2006.
- [118] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *JASIS*, 41(6):391–407, 1990.
- [119] M. Deshpande and G. Karypis. Item-based top- N recommendation algorithms. *ACM Trans. Inf. Syst.*, 22(1):143–177, 2004.
- [120] A. Dhurandhar, S. Oh, and M. Petrik. Building an interpretable recommender via loss-preserving transformation. *CoRR*, abs/1606.05819, 2016.

- [121] Y. Dimopoulos, L. Michael, and F. Athienitou. Ceteris paribus preference elicitation with predictive guarantees. In C. Boutilier, editor, *IJCAI 2009, Proceedings of the 21st International Joint Conference on Artificial Intelligence, Pasadena, California, USA, July 11-17, 2009*, pages 1890–1895, 2009.
- [122] Y. Ding and X. Li. Time weight collaborative filtering. In O. Herzog, H. Schek, N. Fuhr, A. Chowdhury, and W. Teiken, editors, *Proceedings of the 2005 ACM CIKM International Conference on Information and Knowledge Management, Bremen, Germany, October 31 - November 5, 2005*, pages 485–492. ACM, 2005.
- [123] C. Domshlak, E. Hüllermeier, S. Kaci, and H. Prade. Preferences in AI: an overview. *Artif. Intell.*, 175(7-8):1037–1052, 2011.
- [124] C. Domshlak, E. Hüllermeier, S. Kaci, H. Prade, F. Yaman, T. J. Walsh, M. L. Littman, and M. desJardins. Representing, processing, and learning preferences: Theoretical and practical challenges democratic approximation of lexicographic preference models. *Artif. Intell.*, 175(7):1290–1307, 2011. doi:10.1016/j.artint.2010.11.012.
- [125] C. Domshlak and T. Joachims. Unstructuring user preferences: Efficient non-parametric utility revelation. In *UAI '05, Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence, Edinburgh, Scotland, July 26-29, 2005*, pages 169–177. AUAI Press, 2005.
- [126] C. Domshlak and T. Joachims. Efficient and non-parametric reasoning over user preferences. *User Model. User-Adapt. Interact.*, 17(1-2):41–69, 2007.
- [127] C. Domshlak, S. D. Prestwich, F. Rossi, K. B. Venable, and T. Walsh. Hard and soft constraints for reasoning about qualitative conditional preferences. *J. Heuristics*, 12(4-5):263–285, 2006.
- [128] C. Domshlak, F. Rossi, K. B. Venable, and T. Walsh. Reasoning about soft constraints and conditional preferences: complexity results and approxima-

- tion techniques. In G. Gottlob and T. Walsh, editors, *IJCAI-03, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexico, August 9-15, 2003*, pages 215–220. Morgan Kaufmann, 2003.
- [129] W. Dong, M. Charikar, and K. Li. Efficient k-nearest neighbor graph construction for generic similarity measures. In S. Srinivasan, K. Ramamritham, A. Kumar, M. P. Ravindra, E. Bertino, and R. Kumar, editors, *Proceedings of the 20th International Conference on World Wide Web, WWW 2011, Hyderabad, India, March 28 - April 1, 2011*, pages 577–586. ACM, 2011.
- [130] C. N. dos Santos, L. Barbosa, D. Bogdanova, and B. Zadrozny. Learning hybrid representations to retrieve semantically equivalent questions. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 2: Short Papers*, pages 694–699. The Association for Computer Linguistics, 2015.
- [131] J. Doyle. Prospects for preferences. *Computational Intelligence*, 20(2):111–136, 2004.
- [132] N. Drawel, H. Qu, J. Bentahar, and E. Shakshuki. Specification and automatic verification of trust-based multi-agent systems. *Future Generation Computer Systems*, 2018.
- [133] M. J. Dürst and M. Suignard. Internationalized resource identifiers (iris). *RFC*, 3987:1–46, 2005.
- [134] M. D. Ekstrand, F. M. Harper, M. C. Willemsen, and J. A. Konstan. User perception of differences in recommender algorithms. In A. Kobsa, M. X. Zhou, M. Ester, and Y. Koren, editors, *Eighth ACM Conference on Recommender Systems, RecSys '14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014*, pages 161–168. ACM, 2014.

- [135] M. D. Ekstrand, M. Tian, I. M. Azpiazu, J. D. Ekstrand, O. Anuyah, D. McNeill, and M. S. Pera. All the cool kids, how do they fit in?: Popularity and demographic biases in recommender evaluation and effectiveness. In S. A. Friedler and C. Wilson, editors, *Conference on Fairness, Accountability and Transparency, FAT 2018, 23-24 February 2018, New York, NY, USA*, volume 81 of *Proceedings of Machine Learning Research*, pages 172–186. PMLR, 2018.
- [136] A. Elbadrawy and G. Karypis. User-specific feature-based similarity models for top- n recommendation of new items. *ACM TIST*, 6(3):33:1–33:20, 2015.
- [137] M. Enrich, M. Braunhofer, and F. Ricci. Cold-start management with cross-domain collaborative filtering and tags. In C. Huemer and P. Lops, editors, *E-Commerce and Web Technologies - 14th International Conference, EC-Web 2013, Prague, Czech Republic, August 27-28, 2013. Proceedings*, volume 152 of *Lecture Notes in Business Information Processing*, pages 101–112. Springer, 2013.
- [138] R. Fagin. Combining fuzzy information from multiple systems. *J. Comput. Syst. Sci.*, 58(1):83–99, 1999.
- [139] R. Falcone, A. Sapienza, and C. Castelfranchi. The relevance of categories for trusting information sources. *ACM Trans. Internet Techn.*, 15(4):13:1–13:21, 2015.
- [140] I. Fernández-Tobías, M. Braunhofer, M. Elahi, F. Ricci, and I. Cantador. Alleviating the new user problem in collaborative filtering by exploiting personality information. *User Model. User-Adapt. Interact.*, 26(2-3):221–255, 2016.
- [141] I. Fernández-Tobías and I. Cantador. Exploiting social tags in matrix factorization models for cross-domain collaborative filtering. In T. Bogers, M. Koolen, and I. Cantador, editors, *Proceedings of the 1st Workshop on New Trends in Content-based Recommender Systems co-located with the 8th*

- ACM Conference on Recommender Systems, CBRecSys@RecSys 2014, Foster City, Silicon Valley, California, USA, October 6, 2014.*, volume 1245 of *CEUR Workshop Proceedings*, pages 34–41. CEUR-WS.org, 2014.
- [142] I. Fernández-Tobías, I. Cantador, M. Kaminskas, and F. Ricci. A generic semantic-based framework for cross-domain recommendation. In *Proceedings of the 2Nd International Workshop on Information Heterogeneity and Fusion in Recommender Systems, HetRec '11*, pages 25–32, New York, NY, USA, 2011. ACM.
- [143] I. Fernández-Tobías, I. Cantador, and L. Plaza. An emotion dimensional model based on social tags: Crossing folksonomies and enhancing recommendations. In C. Huemer and P. Lops, editors, *E-Commerce and Web Technologies - 14th International Conference, EC-Web 2013, Prague, Czech Republic, August 27-28, 2013. Proceedings*, volume 152 of *Lecture Notes in Business Information Processing*, pages 88–100. Springer, 2013.
- [144] I. Fernández-Tobías, I. Cantador, P. Tomeo, V. W. Anelli, and T. D. Noia. Addressing the user cold start with cross-domain collaborative filtering: exploiting item metadata in matrix factorization. *User Model. User-Adapt. Interact.*, 29(2):443–486, 2019.
- [145] I. Fernández-Tobías, P. Tomeo, I. Cantador, T. D. Noia, and E. D. Sciascio. Accuracy and diversity in cross-domain recommendations for cold-start users with positive-only feedback. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 119–122. ACM, 2016.
- [146] V. Fionda and G. Pirrò. Querying graphs with preferences. In Q. He, A. Iyengar, W. Nejdl, J. Pei, and R. Rastogi, editors, *22nd ACM International Conference on Information and Knowledge Management, CIKM'13, San Francisco, CA, USA, October 27 - November 1, 2013*, pages 929–938. ACM, 2013.

- [147] P. C. Fishburn. *Utility theory for decision making*. Publications in operations research. J. Wiley, 1970.
- [148] D. M. Fleder and K. Hosanagar. Recommender systems and their impact on sales diversity. In J. K. MacKie-Mason, D. C. Parkes, and P. Resnick, editors, *Proceedings 8th ACM Conference on Electronic Commerce (EC-2007)*, San Diego, California, USA, June 11-15, 2007, pages 192–199. ACM, 2007.
- [149] S. Flesca and S. Greco. Partially ordered regular languages for graph queries. *J. Comput. Syst. Sci.*, 70(1):1–25, 2005.
- [150] F. Friedrichs and C. Igel. Evolutionary tuning of multiple SVM parameters. *Neurocomputing*, 64:107–117, 2005.
- [151] S. Funk. Netflix update: Try this at home, 2006.
- [152] J. Fürnkranz and E. Hüllermeier, editors. *Preference Learning*. Springer, 2010.
- [153] E. Gabrilovich and S. Markovitch. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In M. M. Veloso, editor, *IJCAI 2007, Proceedings of the 20th International Joint Conference on Artificial Intelligence, Hyderabad, India, January 6-12, 2007*, pages 1606–1611, 2007.
- [154] Z. Gantner, S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme. Mymedialite: a free recommender system library. In B. Mobasher, R. D. Burke, D. Jannach, and G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 305–308. ACM, 2011.
- [155] S. Gao, H. Luo, D. Chen, S. Li, P. Gallinari, and J. Guo. Cross-domain recommendation via cluster-level latent factor model. In H. Blockeel, K. Kersting, S. Nijssen, and F. Zelezný, editors, *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2013*,

- Prague, Czech Republic, September 23-27, 2013, Proceedings, Part II*, volume 8189 of *Lecture Notes in Computer Science*, pages 161–176. Springer, 2013.
- [156] S. Givon and V. Lavrenko. Predicting social-tags for cold start book recommendations. In L. D. Bergman, A. Tuzhilin, R. D. Burke, A. Felfernig, and L. Schmidt-Thieme, editors, *Proceedings of the 2009 ACM Conference on Recommender Systems, RecSys 2009, New York, NY, USA, October 23-25, 2009*, pages 333–336. ACM, 2009.
- [157] O. Goga, H. Lei, S. H. K. Parthasarathi, G. Friedland, R. Sommer, and R. Teixeira. Exploiting innocuous activity for correlating users across sites. In D. Schwabe, V. A. F. Almeida, H. Glaser, R. A. Baeza-Yates, and S. B. Moon, editors, *22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013*, pages 447–458. International World Wide Web Conferences Steering Committee / ACM, 2013.
- [158] J. Goldsmith, J. Lang, M. Truszczynski, and N. Wilson. The computational complexity of dominance and consistency in cp-nets. *J. Artif. Intell. Res.*, 33:403–432, 2008.
- [159] G. Grahne, A. Thomo, and W. W. Wadge. Preferential regular path queries. *Fundam. Inform.*, 89(2-3):259–288, 2008.
- [160] N. Grgic-Hlaca, M. B. Zafar, K. P. Gummadi, and A. Weller. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. In S. A. McIlraith and K. Q. Weinberger, editors, *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 51–60. AAAI Press, 2018.
- [161] J. T. Guerin, T. E. Allen, and J. Goldsmith. Learning cp-net preferences online from user queries. In P. Perny, M. Pirlot, and A. Tsoukiàs, editors,

- Algorithmic Decision Theory - Third International Conference, ADT 2013, Bruxelles, Belgium, November 12-14, 2013, Proceedings*, volume 8176 of *Lecture Notes in Computer Science*, pages 208–220. Springer, 2013.
- [162] M. Gueroussova, A. Polleres, and S. A. McIlraith. SPARQL with qualitative and quantitative preferences. In I. Celino, E. D. Valle, M. Krötzsch, and S. Schlobach, editors, *Proceedings of the 2nd International Workshop on Ordering and Reasoning, OrdRing 2013, Co-located with the 12th International Semantic Web Conference (ISWC 2013), Sydney, Australia, October 22nd, 2013*, volume 1059 of *CEUR Workshop Proceedings*, pages 2–8. CEUR-WS.org, 2013.
- [163] A. Gunawardana and G. Shani. Evaluating recommender systems. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 265–308. Springer, 2015.
- [164] G. Guo, J. Zhang, Z. Sun, and N. Yorke-Smith. Librec: A java library for recommender systems. In A. I. Cristea, J. Masthoff, A. Said, and N. Tintarev, editors, *Posters, Demos, Late-breaking Results and Workshop Proceedings of the 23rd Conference on User Modeling, Adaptation, and Personalization (UMAP 2015), Dublin, Ireland, June 29 - July 3, 2015.*, volume 1388 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2015.
- [165] E. Han and G. Karypis. Feature-based recommendation system. In O. Herzog, H. Schek, N. Fuhr, A. Chowdhury, and W. Teiken, editors, *Proceedings of the 2005 ACM CIKM International Conference on Information and Knowledge Management, Bremen, Germany, October 31 - November 5, 2005*, pages 446–452. ACM, 2005.
- [166] S. O. Hansson. What is ceteris paribus preference? *J. Philosophical Logic*, 25(3):307–332, 1996.
- [167] M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual*

Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pages 3315–3323, 2016.

- [168] T. H. Haveliwala. Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *IEEE Trans. Knowl. Data Eng.*, 15(4):784–796, 2003.
- [169] J. Havrda and F. Charvát. Quantification method of classification processes. concept of structural α -entropy. *Kybernetika*, 3(1):30–35, 1967.
- [170] L. He, N. N. Liu, and Q. Yang. Active dual collaborative filtering with both item and attribute feedback. In W. Burgard and D. Roth, editors, *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2011, San Francisco, California, USA, August 7-11, 2011*. AAAI Press, 2011.
- [171] M. He, Q. Meng, and S. Zhang. Collaborative additional variational autoencoder for top-n recommender systems. *IEEE Access*, 7:5707–5713, 2019.
- [172] M. He, J. Zhang, P. Yang, and K. Yao. Robust transfer learning for cross-domain collaborative filtering using multiple rating patterns approximation. In Y. Chang, C. Zhai, Y. Liu, and Y. Maarek, editors, *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, pages 225–233. ACM, 2018.
- [173] R. He and J. J. McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In J. Bourdeau, J. Hendler, R. Nkambou, I. Horrocks, and B. Y. Zhao, editors, *Proceedings of the 25th International Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11 - 15, 2016*, pages 507–517. ACM, 2016.
- [174] X. He and T. Chua. Neural factorization machines for sparse predictive analytics. In N. Kando, T. Sakai, H. Joho, H. Li, A. P. de Vries, and R. W. White, editors, *Proceedings of the 40th International ACM SIGIR Conference on Re-*

search and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017, pages 355–364. ACM, 2017.

- [175] J. L. Herlocker, J. A. Konstan, A. Borchers, and J. Riedl. An algorithmic framework for performing collaborative filtering. In F. C. Gey, M. A. Hearst, and R. M. Tong, editors, *SIGIR '99: Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 15-19, 1999, Berkeley, CA, USA*, pages 230–237. ACM, 1999.
- [176] J. L. Herlocker, J. A. Konstan, and J. Riedl. Explaining collaborative filtering recommendations. In W. A. Kellogg and S. Whittaker, editors, *CSCW 2000, Proceeding on the ACM 2000 Conference on Computer Supported Cooperative Work, Philadelphia, PA, USA, December 2-6, 2000*, pages 241–250. ACM, 2000.
- [177] J. L. Herlocker, J. A. Konstan, and J. Riedl. An empirical analysis of design choices in neighborhood-based collaborative filtering algorithms. *Inf. Retr.*, 5(4):287–310, 2002.
- [178] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. Riedl. Evaluating collaborative filtering recommender systems. *ACM Trans. Inf. Syst.*, 22(1):5–53, 2004.
- [179] B. Hidasi, M. Quadrana, A. Karatzoglou, and D. Tikk. Parallel recurrent neural network architectures for feature-rich session-based recommendations. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 241–248. ACM, 2016.
- [180] B. Hidasi and D. Tikk. Initializing matrix factorization methods on implicit feedback databases. *J. UCS*, 19(12):1834–1853, 2013.
- [181] W. C. Hill, L. Stead, M. Rosenstein, and G. W. Furnas. Recommending and evaluating choices in a virtual community of use. In I. R. Katz, R. L.

- Mack, L. Marks, M. B. Rosson, and J. Nielsen, editors, *Human Factors in Computing Systems, CHI '95 Conference Proceedings, Denver, Colorado, USA, May 7-11, 1995.*, pages 194–201. ACM/Addison-Wesley, 1995.
- [182] L. Hu, J. Cao, G. Xu, L. Cao, Z. Gu, and C. Zhu. Personalized recommendation via cross-domain triadic factorization. In D. Schwabe, V. A. F. Almeida, H. Glaser, R. A. Baeza-Yates, and S. B. Moon, editors, *22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013*, pages 595–606. International World Wide Web Conferences Steering Committee / ACM, 2013.
- [183] L. Hu, S. Jian, L. Cao, and Q. Chen. Interpretable recommendation via attraction modeling: Learning multilevel attractiveness over multimodal movie contents. In J. Lang, editor, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 3400–3406. ijcai.org, 2018.
- [184] Y. Hu, Y. Koren, and C. Volinsky. Collaborative filtering for implicit feedback datasets. In *Proceedings of the 8th IEEE International Conference on Data Mining (ICDM 2008), December 15-19, 2008, Pisa, Italy*, pages 263–272. IEEE Computer Society, 2008.
- [185] I. Hulpus, N. Prangnawarat, and C. Hayes. Path-based semantic relatedness on linked data and its use to word and entity disambiguation. In M. Arenas, Ó. Corcho, E. Simperl, M. Strohmaier, M. d’Aquin, K. Srinivas, P. T. Groth, M. Dumontier, J. Heflin, K. Thirunarayan, and S. Staab, editors, *The Semantic Web - ISWC 2015 - 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part I*, volume 9366 of *Lecture Notes in Computer Science*, pages 442–457. Springer, 2015.
- [186] N. Hurley and M. Zhang. Novelty and diversity in top-n recommendation - analysis and evaluation. *ACM Trans. Internet Techn.*, 10(4):14:1–14:30, 2011.

- [187] F. Hutter, H. H. Hoos, and K. Leyton-Brown. An efficient approach for assessing hyperparameter importance. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, volume 32 of *JMLR Workshop and Conference Proceedings*, pages 754–762. JMLR.org, 2014.
- [188] F. Hutter, J. Lücke, and L. Schmidt-Thieme. Beyond manual tuning of hyperparameters. *KI*, 29(4):329–337, 2015.
- [189] M. Jahrer, A. Töscher, and R. A. Legenstein. Combining predictions for accurate recommender systems. In B. Rao, B. Krishnapuram, A. Tomkins, and Q. Yang, editors, *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, July 25-28, 2010*, pages 693–702. ACM, 2010.
- [190] P. Jain, P. Kumaraguru, and A. Joshi. @i seek 'fb.me': identifying users across multiple online social networks. In L. Carr, A. H. F. Laender, B. F. Lóscio, I. King, M. Fontoura, D. Vrandecic, L. Aroyo, J. P. M. de Oliveira, F. Lima, and E. Wilde, editors, *22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013, Companion Volume*, pages 1259–1268. International World Wide Web Conferences Steering Committee / ACM, 2013.
- [191] T. Jambor and J. Wang. Optimizing multiple objectives in collaborative filtering. In X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, editors, *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, pages 55–62. ACM, 2010.
- [192] D. Jannach and G. Adomavicius. Price and profit awareness in recommender systems. *CoRR*, abs/1707.08029, 2017.
- [193] D. Jannach, L. Lerche, F. Gedikli, and G. Bonnin. What recommenders recommend - an analysis of accuracy, popularity, and sales diversity effects.

- In S. Carberry, S. Weibelzahl, A. Micarelli, and G. Semeraro, editors, *User Modeling, Adaptation, and Personalization - 21th International Conference, UMAP 2013, Rome, Italy, June 10-14, 2013, Proceedings*, volume 7899 of *Lecture Notes in Computer Science*, pages 25–37. Springer, 2013.
- [194] D. Jannach, L. Lerche, I. Kamehkhosh, and M. Jugovac. What recommenders recommend: an analysis of recommendation biases and possible countermeasures. *User Model. User-Adapt. Interact.*, 25(5):427–491, 2015.
- [195] D. Jannach and M. Ludewig. When recurrent neural networks meet the neighborhood for session-based recommendation. In P. Cremonesi, F. Ricci, S. Berkovsky, and A. Tuzhilin, editors, *Proceedings of the Eleventh ACM Conference on Recommender Systems, RecSys 2017, Como, Italy, August 27-31, 2017*, pages 306–310. ACM, 2017.
- [196] D. Jannach, P. Resnick, A. Tuzhilin, and M. Zanker. Recommender systems - : beyond matrix completion. *Commun. ACM*, 59(11):94–102, 2016.
- [197] D. Jannach, M. Zanker, A. Felfernig, and G. Friedrich. *Recommender Systems - An Introduction*. Cambridge University Press, 2010.
- [198] K. Järvelin and J. Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf. Syst.*, 20(4):422–446, 2002.
- [199] G. Jawaheer, M. Szomszor, and P. Kostkova. Comparison of implicit and explicit feedback from an online music recommendation service. In *Proceedings of the 1st International Workshop on Information Heterogeneity and Fusion in Recommender Systems, HetRec '10*, pages 47–51, New York, NY, USA, 2010. ACM.
- [200] K. Ji, R. Sun, W. Shu, and X. Li. Next-song recommendation with temporal dynamics. *Knowl.-Based Syst.*, 88:134–143, 2015.
- [201] K. S. Jones, S. Walker, and S. E. Robertson. A probabilistic model of information retrieval: development and comparative experiments - part 1. *Inf. Process. Manage.*, 36(6):779–808, 2000.

- [202] K. S. Jones, S. Walker, and S. E. Robertson. A probabilistic model of information retrieval: development and comparative experiments - part 2. *Inf. Process. Manage.*, 36(6):809–840, 2000.
- [203] Y. Juan, Y. Zhuang, W. Chin, and C. Lin. Field-aware factorization machines for CTR prediction. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 43–50. ACM, 2016.
- [204] M. Jugovac, D. Jannach, and L. Lerche. Efficient optimization of multiple recommendation quality factors according to individual user tendencies. *Expert Syst. Appl.*, 81:321–331, 2017.
- [205] S. Kabbur, X. Ning, and G. Karypis. FISM: factored item similarity models for top-n recommender systems. In I. S. Dhillon, Y. Koren, R. Ghani, T. E. Senator, P. Bradley, R. Parekh, J. He, R. L. Grossman, and R. Uthurusamy, editors, *The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11-14, 2013*, pages 659–667. ACM, 2013.
- [206] S. Kalloori and F. Ricci. Improving cold start recommendation by mapping feature-based preferences to item comparisons. In M. Bieliková, E. Herder, F. Cena, and M. C. Desmarais, editors, *Proceedings of the 25th Conference on User Modeling, Adaptation and Personalization, UMAP 2017, Bratislava, Slovakia, July 09 - 12, 2017*, pages 289–293. ACM, 2017.
- [207] M. Kaminskis, I. Fernández-Tobías, I. Cantador, and F. Ricci. Ontology-based identification of music for places. In L. Cantoni and Z. P. Xiang, editors, *Information and Communication Technologies in Tourism 2013, ENTER 3013, Proceedings of the International Conference in Innsbruck, Austria, January 22-25, 2013.*, pages 436–447. Springer, 2013.
- [208] J. N. Kapur and H. K. Kesavan. *The generalized maximum entropy principle (with applications)*. Sandford Educational Press Waterloo, Ontario, 1987.

- [209] G. Karypis. Evaluation of item-based top-n recommendation algorithms. In *Proceedings of the 2001 ACM CIKM International Conference on Information and Knowledge Management, Atlanta, Georgia, USA, November 5-10, 2001*, pages 247–254. ACM, 2001.
- [210] R. Katarya and O. P. Verma. Effective collaborative movie recommender system using asymmetric user similarity and matrix factorization. In *Computing, Communication and Automation (ICCCA), 2016 International Conference on*, pages 71–75. IEEE, 2016.
- [211] M. Kendall and J. Gibbons. Rank correlation methods, trans. *JD Gibbons (5th edn ed.)*. Edward Arnold: London, 1990.
- [212] H. Khrouf and R. Troncy. Hybrid event recommendation using linked data and user diversity. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 185–192. ACM, 2013.
- [213] W. Kießling. Foundations of preferences in database systems. In *Proceedings of 28th International Conference on Very Large Data Bases, VLDB 2002, Hong Kong, August 20-23, 2002*, pages 311–322. Morgan Kaufmann, 2002.
- [214] W. Kießling, M. Endres, and F. Wenzel. The preference SQL system - an overview. *IEEE Data Eng. Bull.*, 34(2):11–18, 2011.
- [215] D. Kim and B. Yum. Collaborative filtering based on iterative principal component analysis. *Expert Syst. Appl.*, 28(4):823–830, 2005.
- [216] Y. Kim, K. Stratos, and R. Sarikaya. Frustratingly easy neural domain adaptation. In N. Calzolari, Y. Matsumoto, and R. Prasad, editors, *COLING 2016, 26th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, December 11-16, 2016, Osaka, Japan*, pages 387–396. ACL, 2016.
- [217] D. Kluver and J. A. Konstan. Evaluating recommender behavior for new users. In A. Kobsa, M. X. Zhou, M. Ester, and Y. Koren, editors, *Eighth*

ACM Conference on Recommender Systems, RecSys '14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014, pages 121–128. ACM, 2014.

- [218] N. Koenigstein, G. Dror, and Y. Koren. Yahoo! music recommendations: modeling music ratings with temporal dynamics and item taxonomy. In B. Mobasher, R. D. Burke, D. Jannach, and G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 165–172. ACM, 2011.
- [219] R. Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence, IJCAI 95, Montréal Québec, Canada, August 20-25 1995, 2 Volumes*, pages 1137–1145. Morgan Kaufmann, 1995.
- [220] J. A. Konstan, S. M. McNee, C. Ziegler, R. Torres, N. Kapoor, and J. Riedl. Lessons on applying automated recommender systems to information-seeking tasks. In *Proceedings, The Twenty-First National Conference on Artificial Intelligence and the Eighteenth Innovative Applications of Artificial Intelligence Conference, July 16-20, 2006, Boston, Massachusetts, USA*, pages 1630–1633. AAAI Press, 2006.
- [221] Y. Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In Y. Li, B. Liu, and S. Sarawagi, editors, *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Las Vegas, Nevada, USA, August 24-27, 2008*, pages 426–434. ACM, 2008.
- [222] Y. Koren. Collaborative filtering with temporal dynamics. In J. F. E. IV, F. Fogelman-Soulié, P. A. Flach, and M. J. Zaki, editors, *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Paris, France, June 28 - July 1, 2009*, pages 447–456. ACM, 2009.

- [223] Y. Koren and R. M. Bell. Advances in collaborative filtering. In F. Ricci, L. Rokach, and B. Shapira, editors, *Recommender Systems Handbook*, pages 77–118. Springer, 2015.
- [224] Y. Koren, R. M. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [225] Y. Koren and J. Sill. Ordrec: an ordinal model for predicting personalized item rating distributions. In B. Mobasher, R. D. Burke, D. Jannach, and G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 117–124. ACM, 2011.
- [226] F. Koriche and B. Zanuttini. Learning conditional preference networks. *Artif. Intell.*, 174(11):685–703, 2010.
- [227] A. B. Kouki. *Performance prediction and evaluation in recommender systems: An information retrieval perspective*. PhD thesis, Universidad Autónoma de Madrid, 2012.
- [228] O. Küçükünç, E. Saule, K. Kaya, and Ü. V. Çatalyürek. Diversifying citation recommendations. *ACM TIST*, 5(4):55:1–55:21, 2014.
- [229] N. Lathia, S. Hailes, and L. Capra. Temporal collaborative filtering with adaptive neighbourhoods. In J. Allan, J. A. Aslam, M. Sanderson, C. Zhai, and J. Zobel, editors, *Proceedings of the 32nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2009, Boston, MA, USA, July 19-23, 2009*, pages 796–797. ACM, 2009.
- [230] J. Lees-Miller, F. Anderson, B. Hoehn, and R. Greiner. Does wikipedia information help netflix predictions? In M. A. Wani, X. Chen, D. Casasent, L. A. Kurgan, T. Hu, and K. Hafeez, editors, *Seventh International Conference on Machine Learning and Applications, ICMLA 2008, San Diego, California, USA, 11-13 December 2008*, pages 337–343. IEEE Computer Society, 2008.

- [231] A. Leonardo and F. Vasco. OWLPref: Uma representação declarativa de preferências para web semântica. *Anais do XXVII Congresso da SBC*, pages 1411–1420, 2007.
- [232] B. Li, Q. Yang, and X. Xue. Can movies and books collaborate? cross-domain collaborative filtering for sparsity reduction. In C. Boutilier, editor, *IJCAI 2009, Proceedings of the 21st International Joint Conference on Artificial Intelligence, Pasadena, California, USA, July 11-17, 2009*, pages 2052–2057, 2009.
- [233] B. Li, Q. Yang, and X. Xue. Transfer learning for collaborative filtering via a rating-matrix generative model. In A. P. Danyluk, L. Bottou, and M. L. Littman, editors, *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, volume 382 of *ACM International Conference Proceeding Series*, pages 617–624. ACM, 2009.
- [234] C. Li, M. A. Soliman, K. C. Chang, and I. F. Ilyas. Ranksql: Supporting ranking queries in relational database management systems. In K. Böhm, C. S. Jensen, L. M. Haas, M. L. Kersten, P. Larson, and B. C. Ooi, editors, *Proceedings of the 31st International Conference on Very Large Data Bases, Trondheim, Norway, August 30 - September 2, 2005*, pages 1342–1345. ACM, 2005.
- [235] X. Li, H. Fang, Q. Yang, and J. Zhang. Who is your best friend?: Ranking social network friends according to trust relationship. In T. Mitrovic, J. Zhang, L. Chen, and D. Chin, editors, *Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization, UMAP 2018, Singapore, July 08-11, 2018*, pages 301–309. ACM, 2018.
- [236] X. Li, E. Gavves, C. G. M. Snoek, M. Worring, and A. W. M. Smeulders. Personalizing automated image annotation using cross-entropy. In K. S. Candan, S. Panchanathan, B. Prabhakaran, H. Sundaram, W. Feng, and N. Sebe, editors, *Proceedings of the 19th International Conference on Multimedia 2011*,

Scottsdale, AZ, USA, November 28 - December 1, 2011, pages 233–242. ACM, 2011.

- [237] X. Li and J. She. Collaborative variational autoencoder for recommender systems. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 305–314. ACM, 2017.
- [238] X. Li, C. G. M. Snoek, and M. Worring. Unsupervised multi-feature tag relevance learning for social image retrieval. In S. Li, X. Gao, and N. Sebe, editors, *Proceedings of the 9th ACM International Conference on Image and Video Retrieval, CIVR 2010, Xi'an, China, July 5-7, 2010*, pages 10–17. ACM, 2010.
- [239] G. Linden, B. Smith, and J. York. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 7(1):76–80, 2003.
- [240] H. Liu, X. He, F. Feng, L. Nie, R. Liu, and H. Zhang. Discrete factorization machines for fast feature-based recommendation. In J. Lang, editor, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 3449–3455. ijcai.org, 2018.
- [241] J. Liu, Z. Li, J. Tang, Y. Jiang, and H. Lu. Personalized geo-specific tag recommendation for photos on social websites. *IEEE Trans. Multimedia*, 16(3):588–600, 2014.
- [242] J. Liu, J. Shi, W. Cai, B. Liu, W. Pan, Q. Yang, and Z. Ming. Transfer learning from APP domain to news domain for dual cold-start recommendation. In Y. Zheng, W. Pan, S. S. Sahebi, and I. Fernández, editors, *Proceedings of the 1st Workshop on Intelligent Recommender Systems by Knowledge Transfer & Learning co-located with ACM Conference on Recommender Systems (RecSys 2017), Como, Italy, August 27, 2017.*, volume 1887 of *CEUR Workshop Proceedings*, pages 38–41. CEUR-WS.org, 2017.

- [243] J. Liu, Y. Xiong, C. Wu, Z. Yao, and W. Liu. Learning conditional preference networks from inconsistent examples. *IEEE Trans. Knowl. Data Eng.*, 26(2):376–390, 2014.
- [244] J. Liu, Z. Yao, Y. Xiong, W. Liu, and C. Wu. Learning conditional preference network from noisy samples using hypothesis testing. *Knowl.-Based Syst.*, 40:7–16, 2013.
- [245] N. N. Liu and Q. Yang. Eigenrank: a ranking-oriented approach to collaborative filtering. In S. Myaeng, D. W. Oard, F. Sebastiani, T. Chua, and M. Leong, editors, *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2008, Singapore, July 20-24, 2008*, pages 83–90. ACM, 2008.
- [246] N. N. Liu, M. Zhao, E. W. Xiang, and Q. Yang. Online evolutionary collaborative filtering. In X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, editors, *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, pages 95–102. ACM, 2010.
- [247] N. N. Liu, M. Zhao, and Q. Yang. Probabilistic latent preference analysis for collaborative filtering. In D. W. Cheung, I. Song, W. W. Chu, X. Hu, and J. J. Lin, editors, *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009, Hong Kong, China, November 2-6, 2009*, pages 759–766. ACM, 2009.
- [248] W. Liu and R. Burke. Personalizing fairness-aware re-ranking. *CoRR*, abs/1809.02921, 2018.
- [249] W. Liu, C. Wu, B. Feng, and J. Liu. Conditional preference in recommender systems. *Expert Syst. Appl.*, 42(2):774–788, 2015.
- [250] B. Loni, Y. Shi, M. Larson, and A. Hanjalic. Cross-domain collaborative filtering with factorization machines. In M. de Rijke, T. Kenter, A. P. de Vries,

- C. Zhai, F. de Jong, K. Radinsky, and K. Hofmann, editors, *Advances in Information Retrieval - 36th European Conference on IR Research, ECIR 2014, Amsterdam, The Netherlands, April 13-16, 2014. Proceedings*, volume 8416 of *Lecture Notes in Computer Science*, pages 656–661. Springer, 2014.
- [251] X. Lu, A. Moffat, and J. S. Culpepper. The effect of pooling and evaluation depth on IR metrics. *Inf. Retr. Journal*, 19(4):416–445, 2016.
- [252] X. Luo, M. Zhou, S. Li, Z. You, Y. Xia, and Q. Zhu. A nonnegative latent factor model for large-scale sparse matrices in recommender systems via alternating direction method. *IEEE Trans. Neural Netw. Learning Syst.*, 27(3):579–592, 2016.
- [253] G. S. Manku, S. Rajagopalan, and B. G. Lindsay. Random sampling techniques for space efficient online computation of order statistics of large datasets. In A. Delis, C. Faloutsos, and S. Ghandeharizadeh, editors, *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA.*, pages 251–262. ACM Press, 1999.
- [254] J. J. McAuley and J. Leskovec. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In D. Schwabe, V. A. F. Almeida, H. Glaser, R. A. Baeza-Yates, and S. B. Moon, editors, *22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013*, pages 897–908. International World Wide Web Conferences Steering Committee / ACM, 2013.
- [255] J. J. McAuley, C. Targett, Q. Shi, and A. van den Hengel. Image-based recommendations on styles and substitutes. In R. A. Baeza-Yates, M. Lalmas, A. Moffat, and B. A. Ribeiro-Neto, editors, *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, Santiago, Chile, August 9-13, 2015*, pages 43–52. ACM, 2015.
- [256] S. M. McNee, J. Riedl, and J. A. Konstan. Being accurate is not enough: how accuracy metrics have hurt recommender systems. In G. M. Olson and

- R. Jeffries, editors, *Extended Abstracts Proceedings of the 2006 Conference on Human Factors in Computing Systems, CHI 2006, Montréal, Québec, Canada, April 22-27, 2006*, pages 1097–1101. ACM, 2006.
- [257] R. Mehrotra, A. Anderson, F. Diaz, A. Sharma, H. M. Wallach, and E. Yilmaz. Auditing search engines for differential satisfaction across demographics. In R. Barrett, R. Cummings, E. Agichtein, and E. Gabrilovich, editors, *Proceedings of the 26th International Conference on World Wide Web Companion, Perth, Australia, April 3-7, 2017*, pages 626–633. ACM, 2017.
- [258] R. Mehrotra, J. McInerney, H. Bouchard, M. Lalmas, and F. Diaz. Towards a fair marketplace: Counterfactual evaluation of the trade-off between relevance, fairness & satisfaction in recommendation systems. In A. Cuzzocrea, J. Allan, N. W. Paton, D. Srivastava, R. Agrawal, A. Z. Broder, M. J. Zaki, K. S. Candan, A. Labrinidis, A. Schuster, and H. Wang, editors, *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018*, pages 2243–2251. ACM, 2018.
- [259] R. Meymandpour and J. G. Davis. Enhancing recommender systems using linked open data-based semantic analysis of items. In J. G. Davis and A. Bozzon, editors, *3rd Australasian Web Conference, AWC 2015, Sydney, Australia, January 2015*, volume 166 of *CRPIT*, pages 11–17. Australian Computer Society, 2015.
- [260] S. E. Middleton, N. Shadbolt, and D. D. Roure. Ontological user profiling in recommender systems. *ACM Trans. Inf. Syst.*, 22(1):54–88, 2004.
- [261] M. Millán, M. F. Trujillo, and E. Ortiz. A collaborative recommender system based on asymmetric user similarity. In H. Yin, P. Tiño, E. Corchado, W. Byrne, and X. Yao, editors, *Intelligent Data Engineering and Automated Learning - IDEAL 2007, 8th International Conference, Birmingham, UK, December 16-19, 2007, Proceedings*, volume 4881 of *Lecture Notes in Computer Science*, pages 663–672. Springer, 2007.

- [262] D. Milne and I. H. Witten. An effective, low-cost measure of semantic relatedness obtained from wikipedia links. In *Proceedings of AAAI Workshop on Wikipedia and Artificial Intelligence: an Evolving Synergy*, pages 25–30. AAAI Press, 2008.
- [263] N. Mirbakhsh and C. X. Ling. Improving top-n recommendation for cold-start users via cross-domain information. *TKDD*, 9(4):33:1–33:19, 2015.
- [264] M. Montaner, B. López, and J. L. de la Rosa. A taxonomy of recommender agents on the internet. *Artif. Intell. Rev.*, 19(4):285–330, 2003.
- [265] C. Musto, P. Basile, P. Lops, M. de Gemmis, and G. Semeraro. Linked open data-enabled strategies for top-n recommendations. In T. Bogers, M. Koolen, and I. Cantador, editors, *Proceedings of the 1st Workshop on New Trends in Content-based Recommender Systems co-located with the 8th ACM Conference on Recommender Systems, CBRecSys@RecSys 2014, Foster City, Silicon Valley, California, USA, October 6, 2014.*, volume 1245 of *CEUR Workshop Proceedings*, pages 49–56. CEUR-WS.org, 2014.
- [266] C. Musto, P. Basile, P. Lops, M. de Gemmis, and G. Semeraro. Introducing linked open data in graph-based recommender systems. *Inf. Process. Manage.*, 53(2):405–435, 2017.
- [267] C. Musto, P. Lops, P. Basile, M. de Gemmis, and G. Semeraro. Semantics-aware graph-based recommender systems exploiting linked open data. In J. Vassileva, J. Blustein, L. Aroyo, and S. K. D’Mello, editors, *Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization, UMAP 2016, Halifax, NS, Canada, July 13 - 17, 2016*, pages 229–237. ACM, 2016.
- [268] C. Musto, P. Lops, M. de Gemmis, and G. Semeraro. Semantics-aware recommender systems exploiting linked open data and graph-based features. In P. Champin, F. L. Gandon, M. Lalmas, and P. G. Ipeirotis, editors, *Companion of the The Web Conference 2018 on The Web Conference 2018, WWW 2018, Lyon , France, April 23-27, 2018*, pages 457–460. ACM, 2018.

- [269] C. Musto, F. Narducci, P. Lops, M. de Gemmis, and G. Semeraro. Explod: A framework for explaining recommendations based on the linked open data cloud. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 151–154. ACM, 2016.
- [270] C. Musto, G. Semeraro, P. Lops, M. de Gemmis, and F. Narducci. Leveraging social media sources to generate personalized music playlists. In C. Huemer and P. Lops, editors, *E-Commerce and Web Technologies - 13th International Conference, EC-Web 2012, Vienna, Austria, September 4-5, 2012. Proceedings*, volume 123 of *Lecture Notes in Business Information Processing*, pages 112–123. Springer, 2012.
- [271] M. Naseri, M. Elahi, and P. Cremonesi. Investigating the decision making process of users based on the polimovie dataset. In M. Ge and F. Ricci, editors, *Proceedings of the 2nd International Workshop on Decision Making and Recommender Systems, Bolzano, Italy, October 22-23, 2015.*, volume 1533 of *CEUR Workshop Proceedings*, pages 41–44. CEUR-WS.org, 2015.
- [272] M. Nasery, M. Braunhofer, and F. Ricci. Recommendations with optimal combination of feature-based and item-based preferences. In J. Vassileva, J. Blustein, L. Aroyo, and S. K. D’Mello, editors, *Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization, UMAP 2016, Halifax, NS, Canada, July 13 - 17, 2016*, pages 269–273. ACM, 2016.
- [273] X. Ning and G. Karypis. SLIM: sparse linear methods for top-n recommender systems. In D. J. Cook, J. Pei, W. Wang, O. R. Zaïane, and X. Wu, editors, *11th IEEE International Conference on Data Mining, ICDM 2011, Vancouver, BC, Canada, December 11-14, 2011*, pages 497–506. IEEE Computer Society, 2011.
- [274] X. Ning and G. Karypis. Sparse linear methods with side information for top-n recommendations. In P. Cunningham, N. J. Hurley, I. Guy, and S. S.

- Anand, editors, *Sixth ACM Conference on Recommender Systems, RecSys '12, Dublin, Ireland, September 9-13, 2012*, pages 155–162. ACM, 2012.
- [275] T. D. Noia. LOSM: a SPARQL endpoint to query open street map. In S. Villata, J. Z. Pan, and M. Dragoni, editors, *Proceedings of the ISWC 2015 Posters & Demonstrations Track co-located with the 14th International Semantic Web Conference (ISWC-2015), Bethlehem, PA, USA, October 11, 2015.*, volume 1486 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2015.
- [276] T. D. Noia, T. Lukasiewicz, M. V. Martínez, G. I. Simari, and O. Tifrea-Marcuska. Combining existential rules with the power of cp-theories. In Q. Yang and M. J. Wooldridge, editors, *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 2918–2925. AAAI Press, 2015.
- [277] T. D. Noia, C. Magarelli, A. Maurino, M. Palmonari, and A. Rula. Using ontology-based data summarization to develop semantics-aware recommender systems. In A. Gangemi, R. Navigli, M. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai, and M. Alam, editors, *The Semantic Web - 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3-7, 2018, Proceedings*, volume 10843 of *Lecture Notes in Computer Science*, pages 128–144. Springer, 2018.
- [278] T. D. Noia, R. Mirizzi, V. C. Ostuni, and D. Romito. Exploiting the web of data in model-based recommender systems. In P. Cunningham, N. J. Hurley, I. Guy, and S. S. Anand, editors, *Sixth ACM Conference on Recommender Systems, RecSys '12, Dublin, Ireland, September 9-13, 2012*, pages 253–256. ACM, 2012.
- [279] T. D. Noia, R. Mirizzi, V. C. Ostuni, D. Romito, and M. Zanker. Linked open data to support content-based recommender systems. In V. Presutti and H. S. Pinto, editors, *I-SEMANTICS 2012 - 8th International Conference on Semantic Systems, I-SEMANTICS '12, Graz, Austria, September 5-7, 2012*, pages 1–8. ACM, 2012.

- [280] T. D. Noia and V. C. Ostuni. Recommender systems and linked open data. In W. Faber and A. Paschke, editors, *Reasoning Web. Web Logic Rules - 11th International Summer School 2015, Berlin, Germany, July 31 - August 4, 2015, Tutorial Lectures*, volume 9203 of *Lecture Notes in Computer Science*, pages 88–113. Springer, 2015.
- [281] T. D. Noia, V. C. Ostuni, J. Rosati, P. Tomeo, E. D. Sciascio, R. Mirizzi, and C. Bartolini. Building a relatedness graph from linked open data: A case study in the IT domain. *Expert Syst. Appl.*, 44:354–366, 2016.
- [282] T. D. Noia, V. C. Ostuni, P. Tomeo, and E. D. Sciascio. Sprank: Semantic path-based ranking for top- N recommendations using linked open data. *ACM TIST*, 8(1):9:1–9:34, 2016.
- [283] T. D. Noia, J. Rosati, P. Tomeo, and E. D. Sciascio. Adaptive multi-attribute diversity for recommender systems. *Inf. Sci.*, 382-383:234–253, 2017.
- [284] I. Nunes, S. Miles, M. Luck, and C. J. P. de Lucena. An introduction to reasoning over qualitative multi-attribute preferences. *Knowledge Eng. Review*, 30(3):342–372, 2015.
- [285] J. Oh, S. Park, H. Yu, M. Song, and S. Park. Novel recommendation based on personal popularity tendency. In D. J. Cook, J. Pei, W. Wang, O. R. Zaïane, and X. Wu, editors, *11th IEEE International Conference on Data Mining, ICDM 2011, Vancouver, BC, Canada, December 11-14, 2011*, pages 507–516. IEEE Computer Society, 2011.
- [286] S. Oramas, V. C. Ostuni, T. D. Noia, X. Serra, and E. D. Sciascio. Sound and music recommendation with knowledge graphs. *ACM TIST*, 8(2):21:1–21:21, 2017.
- [287] F. Ortega, J. L. Sánchez, J. Bobadilla, and A. Gutiérrez. Improving collaborative filtering-based recommender systems results using pareto dominance. *Inf. Sci.*, 239:50–61, 2013.

- [288] V. C. Ostuni, T. D. Noia, E. D. Sciascio, and R. Mirizzi. Top-n recommendations from implicit feedback leveraging linked open data. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 85–92. ACM, 2013.
- [289] E. Palumbo, G. Rizzo, and R. Troncy. entity2rec: Learning user-item relatedness from knowledge graphs for top-n item recommendation. In P. Cremonesi, F. Ricci, S. Berkovsky, and A. Tuzhilin, editors, *Proceedings of the Eleventh ACM Conference on Recommender Systems, RecSys 2017, Como, Italy, August 27-31, 2017*, pages 32–36. ACM, 2017.
- [290] W. Pan. A survey of transfer learning for collaborative recommendation with auxiliary data. *Neurocomputing*, 177:447–453, 2016.
- [291] W. Pan, N. N. Liu, E. W. Xiang, and Q. Yang. Transfer learning to predict missing ratings via heterogeneous user feedbacks. In T. Walsh, editor, *IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011*, pages 2318–2323. IJCAI/AAAI, 2011.
- [292] W. Pan, E. W. Xiang, N. N. Liu, and Q. Yang. Transfer learning in collaborative filtering for sparsity reduction. In M. Fox and D. Poole, editors, *Proceedings of the Twenty-Fourth AAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*. AAAI Press, 2010.
- [293] U. Panniello, M. Gorgoglione, S. Hill, and K. Hosanagar. Incorporating profit margins into recommender systems: A randomized field experiment of purchasing behavior and consumer trust. 2014.
- [294] E. Pantaleo, V. W. Anelli, T. D. Noia, and G. Sérasset. Etytree: A graphical and interactive etymology dictionary based on wiktionary. In R. Barrett, R. Cummings, E. Agichtein, and E. Gabrilovich, editors, *Proceedings of the 26th International Conference on World Wide Web Companion, Perth, Australia, April 3-7, 2017*, pages 1635–1640. ACM, 2017.

- [295] U. Paquet and N. Koenigstein. One-class collaborative filtering with random graphs. In D. Schwabe, V. A. F. Almeida, H. Glaser, R. A. Baeza-Yates, and S. B. Moon, editors, *22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013*, pages 999–1008. International World Wide Web Conferences Steering Committee / ACM, 2013.
- [296] D. Parra and X. Amatriain. Walk the talk - analyzing the relation between implicit and explicit feedback for preference elicitation. In J. A. Konstan, R. Conejo, J. Marzo, and N. Oliver, editors, *User Modeling, Adaption and Personalization - 19th International Conference, UMAP 2011, Girona, Spain, July 11-15, 2011. Proceedings*, volume 6787 of *Lecture Notes in Computer Science*, pages 255–268. Springer, 2011.
- [297] A. Passant. dbrec - music recommendations using dbpedia. In P. F. Patel-Schneider, Y. Pan, P. Hitzler, P. Mika, L. Zhang, J. Z. Pan, I. Horrocks, and B. Glimm, editors, *The Semantic Web - ISWC 2010 - 9th International Semantic Web Conference, ISWC 2010, Shanghai, China, November 7-11, 2010, Revised Selected Papers, Part II*, volume 6497 of *Lecture Notes in Computer Science*, pages 209–224. Springer, 2010.
- [298] A. Paterek. Improving regularized singular value decomposition for collaborative filtering. In *Proceedings of KDD cup and workshop*, volume 2007, pages 5–8, 2007.
- [299] H. Paulheim and J. Fürnkranz. Unsupervised generation of data mining features from linked open data. In D. D. Burdescu, R. Akerkar, and C. Badica, editors, *2nd International Conference on Web Intelligence, Mining and Semantics, WIMS '12, Craiova, Romania, June 6-8, 2012*, pages 31:1–31:12. ACM, 2012.
- [300] M. J. Pazzani and D. Billsus. Content-based recommendation systems. In P. Brusilovsky, A. Kobsa, and W. Nejdl, editors, *The Adaptive Web, Methods and Strategies of Web Personalization*, volume 4321 of *Lecture Notes in Computer Science*, pages 325–341. Springer, 2007.

- [301] G. Peake and J. Wang. Explanation mining: Post hoc interpretability of latent factor models for recommendation systems. In Y. Guo and F. Farooq, editors, *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, pages 2060–2069. ACM, 2018.
- [302] G. Piao and J. G. Breslin. Measuring semantic distance for linked open data-enabled recommender systems. In S. Ossowski, editor, *Proceedings of the 31st Annual ACM Symposium on Applied Computing, Pisa, Italy, April 4-8, 2016*, pages 315–320. ACM, 2016.
- [303] G. Pigozzi, A. Tsoukiàs, and P. Viappiani. Preferences in artificial intelligence. *Ann. Math. Artif. Intell.*, 77(3-4):361–401, 2016.
- [304] I. Pilászy, D. Zibriczky, and D. Tikk. Fast als-based matrix factorization for explicit and implicit feedback datasets. In X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, editors, *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, pages 71–78. ACM, 2010.
- [305] N. Pinto, D. Doukhan, J. J. DiCarlo, and D. D. Cox. A high-throughput screening approach to discovering good forms of biologically inspired visual representation. *PLoS Computational Biology*, 5(11), 2009.
- [306] P. Pirasteh, D. Hwang, and J. J. Jung. Exploiting matrix factorization to asymmetric user similarities in recommendation systems. *Knowl.-Based Syst.*, 83:51–57, 2015.
- [307] A. M. Qamar, É. Gaussier, J. Chevallet, and J. Lim. Similarity learning for nearest neighbor classification. In *Proceedings of the 8th IEEE International Conference on Data Mining (ICDM 2008), December 15-19, 2008, Pisa, Italy*, pages 983–988. IEEE Computer Society, 2008.

- [308] Z. Qian, P. Zhong, and R. Wang. Tag refinement for user-contributed images via graph learning and nonnegative tensor factorization. *IEEE Signal Process. Lett.*, 22(9):1302–1305, 2015.
- [309] A. Rana and D. Bridge. Explanation chains: Recommendations by explanation. In D. Tikk and P. Pu, editors, *Proceedings of the Poster Track of the 11th ACM Conference on Recommender Systems (RecSys 2017), Como, Italy, August 28, 2017.*, volume 1905 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2017.
- [310] S. Rendle. Factorization machines. In G. I. Webb, B. Liu, C. Zhang, D. Gunopulos, and X. Wu, editors, *ICDM 2010, The 10th IEEE International Conference on Data Mining, Sydney, Australia, 14-17 December 2010*, pages 995–1000. IEEE Computer Society, 2010.
- [311] S. Rendle. *Context-Aware Ranking with Factorization Models*, volume 330 of *Studies in Computational Intelligence*. Springer, 2011.
- [312] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme. BPR: bayesian personalized ranking from implicit feedback. In J. A. Bilmes and A. Y. Ng, editors, *UAI 2009, Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, Montreal, QC, Canada, June 18-21, 2009*, pages 452–461. AUAI Press, 2009.
- [313] S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme. Factorizing personalized markov chains for next-basket recommendation. In M. Rappa, P. Jones, J. Freire, and S. Chakrabarti, editors, *Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010*, pages 811–820. ACM, 2010.
- [314] S. Rendle, Z. Gantner, C. Freudenthaler, and L. Schmidt-Thieme. Fast context-aware recommendations with factorization machines. In W. Ma, J. Nie, R. A. Baeza-Yates, T. Chua, and W. B. Croft, editors, *Proceeding*

of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011, pages 635–644. ACM, 2011.

- [315] S. Rendle and L. Schmidt-Thieme. Pairwise interaction tensor factorization for personalized tag recommendation. In B. D. Davison, T. Suel, N. Craswell, and B. Liu, editors, *Proceedings of the Third International Conference on Web Search and Web Data Mining, WSDM 2010, New York, NY, USA, February 4-6, 2010*, pages 81–90. ACM, 2010.
- [316] P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl. GroupLens: An open architecture for collaborative filtering of netnews. In *Proceedings of the 1994 ACM Conference on Computer Supported Cooperative Work, CSCW '94*, pages 175–186, New York, NY, USA, 1994. ACM.
- [317] F. Ricci, L. Rokach, and B. Shapira, editors. *Recommender Systems Handbook*. Springer, 2015.
- [318] P. Ristoski, J. Rosati, T. D. Noia, R. D. Leone, and H. Paulheim. Rdf2vec: RDF graph embeddings and their applications. *Semantic Web*, 10(4):721–752, 2019.
- [319] S. E. Robertson and H. Zaragoza. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4):333–389, 2009.
- [320] J. Rosati, T. D. Noia, T. Lukasiewicz, R. D. Leone, and A. Maurino. Preference queries with ceteris paribus semantics for linked data. In C. Debruyne, H. Panetto, R. Meersman, T. S. Dillon, G. Weichhart, Y. An, and C. A. Ardagna, editors, *On the Move to Meaningful Internet Systems: OTM 2015 Conferences - Confederated International Conferences: CoopIS, ODBASE, and C&TC 2015, Rhodes, Greece, October 26-30, 2015, Proceedings*, volume 9415 of *Lecture Notes in Computer Science*, pages 423–442. Springer, 2015.

- [321] M. Rowe. Semanticsvd++: Incorporating semantic taste evolution for predicting ratings. In *2014 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT), Warsaw, Poland, August 11-14, 2014 - Volume II*, pages 213–220. IEEE Computer Society, 2014.
- [322] N. Sachdeva, K. Gupta, and V. Pudi. Attentive neural architecture incorporating song features for music recommendation. In S. Pera, M. D. Ekstrand, X. Amatriain, and J. O’Donovan, editors, *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018*, pages 417–421. ACM, 2018.
- [323] S. Sahebi and P. Brusilovsky. Cross-domain collaborative recommendation in a cold-start context: The impact of user profile size on the quality of recommendation. In S. Carberry, S. Weibelzahl, A. Micarelli, and G. Semeraro, editors, *User Modeling, Adaptation, and Personalization - 21th International Conference, UMAP 2013, Rome, Italy, June 10-14, 2013, Proceedings*, volume 7899 of *Lecture Notes in Computer Science*, pages 289–295. Springer, 2013.
- [324] S. Sahebi, P. Brusilovsky, and V. Bobrov. Cross-domain recommendation for large-scale data. In Y. Zheng, W. Pan, S. S. Sahebi, and I. Fernández, editors, *Proceedings of the 1st Workshop on Intelligent Recommender Systems by Knowledge Transfer & Learning co-located with ACM Conference on Recommender Systems (RecSys 2017), Como, Italy, August 27, 2017.*, volume 1887 of *CEUR Workshop Proceedings*, pages 9–15. CEUR-WS.org, 2017.
- [325] A. Said and A. Bellogín. Comparative recommender system evaluation: benchmarking recommendation frameworks. In A. Kobsa, M. X. Zhou, M. Ester, and Y. Koren, editors, *Eighth ACM Conference on Recommender Systems, RecSys ’14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014*, pages 129–136. ACM, 2014.

- [326] T. Sakai. Evaluating evaluation metrics based on the bootstrap. In E. N. Efthimiadis, S. T. Dumais, D. Hawking, and K. Järvelin, editors, *SIGIR 2006: Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Seattle, Washington, USA, August 6-11, 2006*, pages 525–532. ACM, 2006.
- [327] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pages 1257–1264. Curran Associates, Inc., 2007.
- [328] G. Salton. *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*. Addison-Wesley, 1989.
- [329] J. Sang, C. Xu, and J. Liu. User-aware image tag refinement via ternary semantic analysis. *IEEE Trans. Multimedia*, 14(3-2):883–895, 2012.
- [330] R. L. T. Santos, C. Macdonald, and I. Ounis. Exploiting query reformulations for web search result diversification. In M. Rappa, P. Jones, J. Freire, and S. Chakrabarti, editors, *Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010*, pages 881–890. ACM, 2010.
- [331] P. Sapiezynski, V. Kassarnig, and C. Wilson. Academic performance prediction in a gender-imbalanced environment. In *1st FATREC Workshop on Responsible Recommendation*, 2017.
- [332] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl. Application of dimensionality reduction in recommender system - a case study. In *ACM WebKDD 2000 Workshop*. ACM SIGKDD, 2000.
- [333] B. M. Sarwar, G. Karypis, J. A. Konstan, and J. Riedl. Analysis of recommendation algorithms for e-commerce. In *EC*, pages 158–167, 2000.

- [334] B. M. Sarwar, G. Karypis, J. A. Konstan, and J. Riedl. Item-based collaborative filtering recommendation algorithms. In V. Y. Shen, N. Saito, M. R. Lyu, and M. E. Zurko, editors, *Proceedings of the Tenth International World Wide Web Conference, WWW 10, Hong Kong, China, May 1-5, 2001*, pages 285–295. ACM, 2001.
- [335] N. Sawant, R. Datta, J. Li, and J. Z. Wang. Quest for relevant tags using local interaction networks and visual content. In J. Z. Wang, N. Bouje-maa, N. O. Ramirez, and A. Natsev, editors, *Proceedings of the 11th ACM SIGMM International Conference on Multimedia Information Retrieval, MIR 2010, Philadelphia, Pennsylvania, USA, March 29-31, 2010*, pages 231–240. ACM, 2010.
- [336] J. B. Schafer, D. Frankowski, J. L. Herlocker, and S. Sen. Collaborative filtering recommender systems. In P. Brusilovsky, A. Kobsa, and W. Nejdl, editors, *The Adaptive Web, Methods and Strategies of Web Personalization*, volume 4321 of *Lecture Notes in Computer Science*, pages 291–324. Springer, 2007.
- [337] M. Schmachtenberg, T. Strufe, and H. Paulheim. Enhancing a location-based recommendation system by enrichment with structured data from the web. In R. Akerkar, N. Bassiliades, J. Davies, and V. Ermolayev, editors, *4th International Conference on Web Intelligence, Mining and Semantics (WIMS 14), WIMS '14, Thessaloniki, Greece, June 2-4, 2014*, pages 17:1–17:12. ACM, 2014.
- [338] N. Shadbolt, T. Berners-Lee, and W. Hall. The semantic web revisited. *IEEE Intelligent Systems*, 21(3):96–101, 2006.
- [339] G. Shani and A. Gunawardana. Evaluating recommendation systems. In F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, editors, *Recommender Systems Handbook*, pages 257–297. Springer, 2011.
- [340] G. Shani, D. Heckerman, and R. I. Brafman. An mdp-based recommender system. *J. Mach. Learn. Res.*, 6:1265–1295, 2005.

- [341] B. Shapira, L. Rokach, and S. Freilikhman. Facebook single and cross domain data for recommendation systems. *User Model. User-Adapt. Interact.*, 23(2-3):211–247, 2013.
- [342] U. Shardanand and P. Maes. Social information filtering: Algorithms for automating "word of mouth". In I. R. Katz, R. L. Mack, L. Marks, M. B. Rosson, and J. Nielsen, editors, *Human Factors in Computing Systems, CHI '95 Conference Proceedings, Denver, Colorado, USA, May 7-11, 1995.*, pages 210–217. ACM/Addison-Wesley, 1995.
- [343] Y. Shi, A. Karatzoglou, L. Baltrunas, M. Larson, and A. Hanjalic. xclimf: optimizing expected reciprocal rank for data with multiple levels of relevance. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 431–434. ACM, 2013.
- [344] Y. Shi, M. Larson, and A. Hanjalic. List-wise learning to rank with matrix factorization for collaborative filtering. In X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, editors, *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, pages 269–272. ACM, 2010.
- [345] Y. Shi, M. Larson, and A. Hanjalic. Tags as bridges between domains: Improving recommendation with tag-induced cross-domain collaborative filtering. In J. A. Konstan, R. Conejo, J. Marzo, and N. Oliver, editors, *User Modeling, Adaption and Personalization - 19th International Conference, UMAP 2011, Girona, Spain, July 11-15, 2011. Proceedings*, volume 6787 of *Lecture Notes in Computer Science*, pages 305–316. Springer, 2011.
- [346] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating activation differences. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3145–3153. PMLR, 2017.

- [347] W. Siberski, J. Z. Pan, and U. Thaden. Querying the semantic web with preferences. In I. F. Cruz, S. Decker, D. Allemang, C. Preist, D. Schwabe, P. Mika, M. Uschold, and L. Aroyo, editors, *The Semantic Web - ISWC 2006, 5th International Semantic Web Conference, ISWC 2006, Athens, GA, USA, November 5-9, 2006, Proceedings*, volume 4273 of *Lecture Notes in Computer Science*, pages 612–624. Springer, 2006.
- [348] A. Singh and T. Joachims. Fairness of exposure in rankings. In Y. Guo and F. Farooq, editors, *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, pages 2219–2228. ACM, 2018.
- [349] R. R. Sinha and K. Swearingen. The role of transparency in recommender systems. In L. G. Terveen and D. R. Wixon, editors, *Extended abstracts of the 2002 Conference on Human Factors in Computing Systems, CHI 2002, Minneapolis, Minnesota, USA, April 20-25, 2002*, pages 830–831. ACM, 2002.
- [350] M. R. Smith, L. Mitchell, C. G. Giraud-Carrier, and T. R. Martinez. Recommending learning algorithms and their associated hyperparameters. In J. Vanschoren, P. Brazdil, C. Soares, and L. Kotthoff, editors, *Proceedings of the International Workshop on Meta-learning and Algorithm Selection co-located with 21st European Conference on Artificial Intelligence, MetaSel@ECAI 2014, Prague, Czech Republic, August 19, 2014.*, volume 1201 of *CEUR Workshop Proceedings*, pages 39–40. CEUR-WS.org, 2014.
- [351] M. D. Smucker, J. Allan, and B. Carterette. A comparison of statistical significance tests for information retrieval evaluation. In M. J. Silva, A. H. F. Laender, R. A. Baeza-Yates, D. L. McGuinness, B. Olstad, Ø. H. Olsen, and A. O. Falcão, editors, *Proceedings of the Sixteenth ACM Conference on Information and Knowledge Management, CIKM 2007, Lisbon, Portugal, November 6-10, 2007*, pages 623–632. ACM, 2007.

- [352] E. Smyth. *Would the Internet widen public participation?* PhD thesis, University of Leeds, United Kingdom, 2012.
- [353] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pages 2960–2968, 2012.
- [354] T. Sørensen. A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Biol. Skr.*, 5:1–34, 1948.
- [355] T. Speicher, H. Heidari, N. Grgic-Hlaca, K. P. Gummadi, A. Singla, A. Weller, and M. B. Zafar. A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices. In Y. Guo and F. Farooq, editors, *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, pages 2239–2248. ACM, 2018.
- [356] H. Steck. Evaluation of recommendations: rating-prediction and ranking. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 213–220. ACM, 2013.
- [357] A. Stewart, E. Diaz-Aviles, W. Nejdl, L. B. Marinho, A. Nanopoulos, and L. Schmidt-Thieme. Cross-tagging for personalized open social networking. In C. Cattuto, G. Ruffo, and F. Menczer, editors, *HYPertext 2009, Proceedings of the 20th ACM Conference on Hypertext and Hypermedia, Torino, Italy, June 29 - July 1, 2009*, pages 271–278. ACM, 2009.
- [358] Z. Sun, J. Yang, J. Zhang, A. Bozzon, L. Huang, and C. Xu. Recurrent knowledge graph embedding for effective recommendation. In S. Pera, M. D.

- Ekstrand, X. Amatriain, and J. O'Donovan, editors, *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018*, pages 297–305. ACM, 2018.
- [359] Ö. Süreer, R. Burke, and E. C. Malthouse. Multistakeholder recommendation with provider constraints. In *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018*, pages 54–62, 2018.
- [360] P. Symeonidis, A. Nanopoulos, and Y. Manolopoulos. Justified recommendations based on content and rating data. In *WebKDD Workshop on Web Mining and Web Usage Analysis*, 2008.
- [361] M. Szomszor, H. Alani, I. Cantador, K. O'Hara, and N. Shadbolt. Semantic modelling of user interests based on cross-folksonomy analysis. In A. P. Sheth, S. Staab, M. Dean, M. Paolucci, D. Maynard, T. W. Finin, and K. Thirunarayan, editors, *The Semantic Web - ISWC 2008, 7th International Semantic Web Conference, ISWC 2008, Karlsruhe, Germany, October 26-30, 2008. Proceedings*, volume 5318 of *Lecture Notes in Computer Science*, pages 632–648. Springer, 2008.
- [362] M. Szomszor, I. Cantador, and H. Alani. Correlating user profiles from multiple folksonomies. In P. Brusilovsky and H. C. Davis, editors, *HYPertext 2008, Proceedings of the 19th ACM Conference on Hypertext and Hypermedia, Pittsburgh, PA, USA, June 19-21, 2008*, pages 33–42. ACM, 2008.
- [363] G. Takács, I. Pilászy, B. Németh, and D. Tikk. Matrix factorization and neighbor based algorithms for the netflix prize problem. In P. Pu, D. G. Bridge, B. Mobasher, and F. Ricci, editors, *Proceedings of the 2008 ACM Conference on Recommender Systems, RecSys 2008, Lausanne, Switzerland, October 23-25, 2008*, pages 267–274. ACM, 2008.
- [364] G. Takács, I. Pilászy, B. Németh, and D. Tikk. Scalable collaborative filtering approaches for large recommender systems. *J. Mach. Learn. Res.*, 10:623–656, 2009.

- [365] G. Takács and D. Tikk. Alternating least squares for personalized ranking. In P. Cunningham, N. J. Hurley, I. Guy, and S. S. Anand, editors, *Sixth ACM Conference on Recommender Systems, RecSys '12, Dublin, Ireland, September 9-13, 2012*, pages 83–90. ACM, 2012.
- [366] A. Taneja and A. Arora. Cross domain recommendation using multidimensional tensor factorization. *Expert Syst. Appl.*, 92:304–316, 2018.
- [367] N. Tintarev and J. Masthoff. A survey of explanations in recommender systems. In *Proceedings of the 23rd International Conference on Data Engineering Workshops, ICDE 2007, 15-20 April 2007, Istanbul, Turkey*, pages 801–810. IEEE Computer Society, 2007.
- [368] N. Tintarev and J. Masthoff. Designing and evaluating explanations for recommender systems. In F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, editors, *Recommender Systems Handbook*, pages 479–510. Springer, 2011.
- [369] A. Tiroshi, S. Berkovsky, M. A. Kâafar, T. Chen, and T. Kuflik. Cross social networks interests predictions based on graph features. In Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, editors, *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*, pages 319–322. ACM, 2013.
- [370] A. Tiroshi and T. Kuflik. Domain ranking for cross domain collaborative filtering. In J. Masthoff, B. Mobasher, M. C. Desmarais, and R. Nkambou, editors, *User Modeling, Adaptation, and Personalization - 20th International Conference, UMAP 2012, Montreal, Canada, July 16-20, 2012. Proceedings*, volume 7379 of *Lecture Notes in Computer Science*, pages 328–333. Springer, 2012.
- [371] A. Troumpoukis, S. Konstantopoulos, and A. Charalambidis. An extension of SPARQL for expressing qualitative preferences. In C. d'Amato, M. Fernández, V. A. M. Tamma, F. Lécué, P. Cudré-Mauroux, J. F. Sequeda, C. Lange, and J. Heflin, editors, *The Semantic Web - ISWC 2017 - 16th International Semantic Web Conference, Vienna, Austria, October 21-25, 2017*,

Proceedings, Part I, volume 10587 of *Lecture Notes in Computer Science*, pages 711–727. Springer, 2017.

- [372] V. Tsintzou, E. Pitoura, and P. Tsaparas. Bias disparity in recommendation systems. *CoRR*, abs/1811.01461, 2018.
- [373] K. H. L. Tso and L. Schmidt-Thieme. Evaluation of attribute-aware recommender system algorithms on data with varying characteristics. In W. K. Ng, M. Kitsuregawa, J. Li, and K. Chang, editors, *Advances in Knowledge Discovery and Data Mining, 10th Pacific-Asia Conference, PAKDD 2006, Singapore, April 9-12, 2006, Proceedings*, volume 3918 of *Lecture Notes in Computer Science*, pages 831–840. Springer, 2006.
- [374] T. X. Tuan and T. M. Phuong. 3d convolutional networks for session-based recommendation with content features. In P. Cremonesi, F. Ricci, S. Berkovsky, and A. Tuzhilin, editors, *Proceedings of the Eleventh ACM Conference on Recommender Systems, RecSys 2017, Como, Italy, August 27-31, 2017*, pages 138–146. ACM, 2017.
- [375] D. Valcarce, A. Bellogín, J. Parapar, and P. Castells. On the robustness and discriminative power of information retrieval metrics for top-n recommendation. In S. Pera, M. D. Ekstrand, X. Amatriain, and J. O’Donovan, editors, *Proceedings of the 12th ACM Conference on Recommender Systems, RecSys 2018, Vancouver, BC, Canada, October 2-7, 2018*, pages 260–268. ACM, 2018.
- [376] S. Vargas, L. Baltrunas, A. Karatzoglou, and P. Castells. Coverage, redundancy and size-awareness in genre diversity for recommender systems. In A. Kobsa, M. X. Zhou, M. Ester, and Y. Koren, editors, *Eighth ACM Conference on Recommender Systems, RecSys ’14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014*, pages 209–216. ACM, 2014.
- [377] S. Vargas and P. Castells. Rank and relevance in novelty and diversity metrics for recommender systems. In B. Mobasher, R. D. Burke, D. Jannach, and

- G. Adomavicius, editors, *Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, October 23-27, 2011*, pages 109–116. ACM, 2011.
- [378] S. Vargas and P. Castells. Exploiting the diversity of user preferences for recommendation. In J. Ferreira, J. Magalhães, and P. Calado, editors, *Open research Areas in Information Retrieval, OAIR '13, Lisbon, Portugal, May 15-17, 2013*, pages 129–136. ACM, 2013.
- [379] S. Vargas, P. Castells, and D. Vallet. Intent-oriented diversity in recommender systems. In W. Ma, J. Nie, R. A. Baeza-Yates, T. Chua, and W. B. Croft, editors, *Proceeding of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011*, pages 1211–1212. ACM, 2011.
- [380] H. R. Varian. Economics and search. *SIGIR Forum*, 33(1):1–5, 1999.
- [381] H. R. Varian. The economics of search. In F. C. Gey, M. A. Hearst, and R. M. Tong, editors, *SIGIR '99: Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 15-19, 1999, Berkeley, CA, USA*, page 1. ACM, 1999.
- [382] S. Verma and J. Rubin. Fairness definitions explained. In Y. Brun, B. Johnson, and A. Meliou, editors, *Proceedings of the International Workshop on Software Fairness, FairWare@ICSE 2018, Gothenburg, Sweden, May 29, 2018*, pages 1–7. ACM, 2018.
- [383] J. Vig, S. Sen, and J. Riedl. Tagsplanations: explaining recommendations using tags. In C. Conati, M. Bauer, N. Oliver, and D. S. Weld, editors, *Proceedings of the 14th International Conference on Intelligent User Interfaces, IUI 2009, Sanibel Island, Florida, USA, February 8-11, 2009*, pages 47–56. ACM, 2009.
- [384] M. Vlachos, C. Dünner, R. Heckel, V. G. Vassiliadis, T. P. Parnell, and K. Atasu. Addressing interpretability and cold-start in matrix factorization

- for recommender systems. *IEEE Trans. Knowl. Data Eng.*, 31(7):1253–1266, 2019.
- [385] G. H. Von Wright. *The Logic of Preference*. Edinburgh University Press, 1963.
- [386] V. Walter Anelli, R. De Leone, T. Di Noia, T. Lukasiewicz, and J. Rosati. Combining rdf and sparql with cp-theories to reason about preferences in a linked data setting. *Semantic Web*, pages 1–29, 12 2018.
- [387] H. Wang and C. Wu. A mathematical model for product selection strategies in a recommender system. *Expert Syst. Appl.*, 36(3):7299–7308, 2009.
- [388] H. Wang, F. Zhang, M. Zhao, W. Li, X. Xie, and M. Guo. Multi-task feature learning for knowledge graph enhanced recommendation. In L. Liu, R. W. White, A. Mantrach, F. Silvestri, J. J. McAuley, R. S. Baeza-Yates, and L. Zia, editors, *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, pages 2000–2010. ACM, 2019.
- [389] T. Wang, Y. Cai, H. Leung, Z. Cai, and H. Min. Entropy-based term weighting schemes for text categorization in VSM. In *27th IEEE International Conference on Tools with Artificial Intelligence, ICTAI 2015, Vietri sul Mare, Italy, November 9-11, 2015*, pages 325–332. IEEE Computer Society, 2015.
- [390] X. Wang, X. He, F. Feng, L. Nie, and T. Chua. TEM: tree-enhanced embedding model for explainable recommendation. In P. Champin, F. L. Gandon, M. Lalmas, and P. G. Ipeirotis, editors, *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*, pages 1543–1552. ACM, 2018.
- [391] J. Wasilewski and N. Hurley. Intent-aware diversification using a constrained PLSA. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 39–42. ACM, 2016.

- [392] M. Weimer, A. Karatzoglou, and A. J. Smola. Improving maximum margin matrix factorization. *Machine Learning*, 72(3):263–276, 2008.
- [393] S. Weng and M. Liu. Feature-based recommendations for one-to-one marketing. *Expert Syst. Appl.*, 26(4):493–508, 2004.
- [394] N. Wilson. Extending cp-nets with stronger conditional preference statements. In D. L. McGuinness and G. Ferguson, editors, *Proceedings of the Nineteenth National Conference on Artificial Intelligence, Sixteenth Conference on Innovative Applications of Artificial Intelligence, July 25-29, 2004, San Jose, California, USA*, pages 735–741. AAAI Press / The MIT Press, 2004.
- [395] N. Wilson. Computational techniques for a simple theory of conditional preferences. *Artif. Intell.*, 175(7-8):1053–1091, 2011.
- [396] P. Winoto and T. Y. Tang. If you like the devil wears prada the book, will you also enjoy the devil wears prada the movie? A study of cross-domain recommendations. *New Generation Comput.*, 26(3):209–225, 2008.
- [397] W. S. Wong, R. W. P. Luk, H. V. Leong, L. K. Ho, and D. L. Lee. Re-examining the effects of adding relevance information in a relevance feedback environment. *Inf. Process. Manage.*, 44(3):1086–1116, 2008.
- [398] C. Wongchokprasitti, J. Peltonen, T. Ruotsalo, P. Bandyopadhyay, G. Jacucci, and P. Brusilovsky. User model in a box: Cross-system user model transfer for resolving cold start problems. In F. Ricci, K. Bontcheva, O. Conlan, and S. Lawless, editors, *User Modeling, Adaptation and Personalization - 23rd International Conference, UMAP 2015, Dublin, Ireland, June 29 - July 3, 2015. Proceedings*, volume 9146 of *Lecture Notes in Computer Science*, pages 289–301. Springer, 2015.
- [399] C. Wu, A. Ahmed, A. Beutel, A. J. Smola, and H. Jing. Recurrent recommender networks. In M. de Rijke, M. Shokouhi, A. Tomkins, and M. Zhang, editors, *Proceedings of the Tenth ACM International Conference on Web*

Search and Data Mining, WSDM 2017, Cambridge, United Kingdom, February 6-10, 2017, pages 495–503. ACM, 2017.

- [400] Q. Wu, C. J. C. Burges, K. M. Svore, and J. Gao. Adapting boosting for information retrieval measures. *Inf. Retr.*, 13(3):254–270, 2010.
- [401] J. Xiao, H. Ye, X. He, H. Zhang, F. Wu, and T. Chua. Attentional factorization machines: Learning the weight of feature interactions via attention networks. In C. Sierra, editor, *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 3119–3125. ijcai.org, 2017.
- [402] J. Xu and H. Li. Adarank: a boosting algorithm for information retrieval. In W. Kraaij, A. P. de Vries, C. L. A. Clarke, N. Fuhr, and N. Kando, editors, *SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007*, pages 391–398. ACM, 2007.
- [403] J. Yang, Z. Sun, A. Bozzon, and J. Zhang. Learning hierarchical feature influence for recommendation by recursive regularization. In S. Sen, W. Geyer, J. Freyne, and P. Castells, editors, *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 51–58. ACM, 2016.
- [404] S. Yao and B. Huang. Beyond parity: Fairness objectives for collaborative filtering. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 2921–2930, 2017.
- [405] X. Yu, X. Ren, Y. Sun, Q. Gu, B. Sturt, U. Khandelwal, B. Norick, and J. Han. Personalized entity recommendation: a heterogeneous information network approach. In B. Carterette, F. Diaz, C. Castillo, and D. Metzler, editors,

Seventh ACM International Conference on Web Search and Data Mining, WSDM 2014, New York, NY, USA, February 24-28, 2014, pages 283–292. ACM, 2014.

- [406] F. Yuan, G. Guo, J. M. Jose, L. Chen, H. Yu, and W. Zhang. Boostfm: Boosted factorization machines for top-n feature-based recommendation. In G. A. Papadopoulos, T. Kuflik, F. Chen, C. Duarte, and W. Fu, editors, *Proceedings of the 22nd International Conference on Intelligent User Interfaces, IUI 2017, Limassol, Cyprus, March 13-16, 2017*, pages 45–54. ACM, 2017.
- [407] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In R. Barrett, R. Cummings, E. Agichtein, and E. Gabrilovich, editors, *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, pages 1171–1180. ACM, 2017.
- [408] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi. Fairness constraints: A flexible approach for fair classification. *J. Mach. Learn. Res.*, 20:75:1–75:42, 2019.
- [409] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, K. P. Gummadi, and A. Weller. From parity to preference-based notions of fairness in classification. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 229–239, 2017.
- [410] M. Zanker. The influence of knowledgeable explanations on users’ perception of a recommender system. In P. Cunningham, N. J. Hurley, I. Guy, and S. S. Anand, editors, *Sixth ACM Conference on Recommender Systems, RecSys ’12, Dublin, Ireland, September 9-13, 2012*, pages 269–272. ACM, 2012.

- [411] A. Zaveri, A. Rula, A. Maurino, R. Pietrobon, J. Lehmann, and S. Auer. Quality assessment for linked data: A survey. *Semantic Web*, 7(1):63–93, 2016.
- [412] M. Zehlike, F. Bonchi, C. Castillo, S. Hajian, M. Megahed, and R. A. Baeza-Yates. Fa*ir: A fair top-k ranking algorithm. In E. Lim, M. Winslett, M. Sanderson, A. W. Fu, J. Sun, J. S. Culpepper, E. Lo, J. C. Ho, D. Donato, R. Agrawal, Y. Zheng, C. Castillo, A. Sun, V. S. Tseng, and C. Li, editors, *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, Singapore, November 06 - 10, 2017*, pages 1569–1578. ACM, 2017.
- [413] C. Zhai and J. D. Lafferty. A study of smoothing methods for language models applied to ad hoc information retrieval. In W. B. Croft, D. J. Harper, D. H. Kraft, and J. Zobel, editors, *SIGIR 2001: Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, September 9-13, 2001, New Orleans, Louisiana, USA*, pages 334–342. ACM, 2001.
- [414] F. Zhang, N. J. Yuan, D. Lian, X. Xie, and W. Ma. Collaborative knowledge base embedding for recommender systems. In B. Krishnapuram, M. Shah, A. J. Smola, C. C. Aggarwal, D. Shen, and R. Rastogi, editors, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 353–362. ACM, 2016.
- [415] Q. Zhang, D. Wu, J. Lu, F. Liu, and G. Zhang. A cross-domain recommender system with consistent information transfer. *Decision Support Systems*, 104:49–63, 2017.
- [416] Y. Zhang. Explainable recommendation: Theory and applications. *CoRR*, abs/1708.06409, 2017.
- [417] Y. Zhang, B. Cao, and D. Yeung. Multi-domain collaborative filtering. In P. Grünwald and P. Spirtes, editors, *UAI 2010, Proceedings of the Twenty-*

Sixth Conference on Uncertainty in Artificial Intelligence, Catalina Island, CA, USA, July 8-11, 2010, pages 725–732. AUAI Press, 2010.

- [418] Y. Zhang and X. Chen. Explainable recommendation: A survey and new perspectives. *CoRR*, abs/1804.11192, 2018.
- [419] Y. Zhang, G. Lai, M. Zhang, Y. Zhang, Y. Liu, and S. Ma. Explicit factor models for explainable recommendation based on phrase-level sentiment analysis. In *The 37th Int. Conf. on Research and Development in Information Retrieval, SIGIR '14, Gold Coast , QLD, Australia*, pages 83–92, 2014.
- [420] Y. Zhang, M. Zhang, Y. Zhang, G. Lai, Y. Liu, H. Zhang, and S. Ma. Daily-aware personalized recommendation based on feature-level time series analysis. In A. Gangemi, S. Leonardi, and A. Panconesi, editors, *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015*, pages 1373–1383. ACM, 2015.
- [421] L. Zhao, S. J. Pan, and Q. Yang. A unified framework of active transfer learning for cross-system recommendation. *Artif. Intell.*, 245:38–55, 2017.
- [422] W. X. Zhao, S. Li, Y. He, L. Wang, J. Wen, and X. Li. Exploring demographic information in social media for product recommendation. *Knowl. Inf. Syst.*, 49(1):61–89, 2016.
- [423] Y. Zheng, N. Ghane, and M. Sabouri. Personalized educational learning with multi-stakeholder optimizations. In G. A. Papadopoulos, G. Samaras, S. Weibelzahl, D. Jannach, and O. C. Santos, editors, *Adjunct Publication of the 27th Conference on User Modeling, Adaptation and Personalization, UMAP 2019, Larnaca, Cyprus, June 09-12, 2019*, pages 283–289. ACM, 2019.
- [424] Y. Zheng, B. Mobasher, and R. D. Burke. Incorporating context correlation into context-aware matrix factorization. In D. Jannach, J. Mengin, B. Mobasher, A. Passerini, and P. Viappiani, editors, *Proceedings of the IJ-CAI 2015 Joint Workshop on Constraints and Preferences for Configuration*

and Recommendation and Intelligent Techniques for Web Personalization co-located with the 24th International Joint Conference on Artificial Intelligence (IJCAI 2015), Buenos Aires, Argentina, July 27, 2015., volume 1440 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2015.

- [425] Y. Zhou, D. M. Wilkinson, R. Schreiber, and R. Pan. Large-scale parallel collaborative filtering for the netflix prize. In R. Fleischer and J. Xu, editors, *Algorithmic Aspects in Information and Management, 4th International Conference, AAIM 2008, Shanghai, China, June 23-25, 2008. Proceedings*, volume 5034 of *Lecture Notes in Computer Science*, pages 337–348. Springer, 2008.
- [426] F. Zhu, Y. Wang, C. Chen, G. Liu, M. A. Orgun, and J. Wu. A deep framework for cross-domain and cross-system recommendations. In J. Lang, editor, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 3711–3717. ijcai.org, 2018.
- [427] Z. Zhu, X. Hu, and J. Caverlee. Fairness-aware tensor-based recommendation. In A. Cuzzocrea, J. Allan, N. W. Paton, D. Srivastava, R. Agrawal, A. Z. Broder, M. J. Zaki, K. S. Candan, A. Labrinidis, A. Schuster, and H. Wang, editors, *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018*, pages 1153–1162. ACM, 2018.
- [428] F. Zhuang, P. Luo, H. Xiong, Y. Xiong, Q. He, and Z. Shi. Cross-domain learning from multiple sources: A consensus regularization perspective. *IEEE Trans. Knowl. Data Eng.*, 22(12):1664–1678, 2010.
- [429] C. Ziegler, S. M. McNee, J. A. Konstan, and G. Lausen. Improving recommendation lists through topic diversification. In A. Ellis and T. Hagino, editors, *Proceedings of the 14th international conference on World Wide Web, WWW 2005, Chiba, Japan, May 10-14, 2005*, pages 22–32. ACM, 2005.

- [430] A. Zimdars, D. M. Chickering, and C. Meek. Using temporal data for making recommendations. In J. S. Breese and D. Koller, editors, *UAI '01: Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, University of Washington, Seattle, Washington, USA, August 2-5, 2001*, pages 580–588. Morgan Kaufmann, 2001.

Arrivati a questo punto, sento che qualcosa dovrei scriverla. Non perché debba, ma perché tante di queste cose sono rimaste chiuse in me per tanto tempo.

Voglio ringraziare la mia famiglia. Sinceramente e totalmente. Innanzi tutto, voglio ringraziare Filippo e Flora. Ammettiamolo, non hanno mai capito veramente cosa stessi facendo. Ciononostante, mi hanno supportato in ogni modo, imparando ad esserci senza essere invadenti. Sono consapevole di quanto sia stato difficile per loro. Sono sempre stati un faro per me, e sempre lo saranno. Li amo profondamente, come amo i miei fratelli. Sami ed Ale, così diversi eppure così simili, fragili e fortissimi allo stesso tempo, tentano ogni giorno di essere la migliore versione di se stessi. Sono profondamente orgoglioso di loro, e mi piacerebbe essere il loro scudo nei confronti delle preoccupazioni della vita.

Vorrei ringraziare Tommaso. E' difficile trovare le parole giuste. Proviamoci. Delle profonde trasformazioni a cui ti obbliga il dottorato, lui ne è il principale artefice (o colpevole, questione di punti di vista). Non potrei mai ringraziarlo abbastanza per tutte le cose che ho imparato, e questo è il meno in un dottorato. Mi ha mostrato nuovi livelli di sacrificio, eppure mi ha anche spinto a guardare le cose nella giusta prospettiva, privando di valore ed importanza le cose inutili. Ha avuto la pazienza dell'acqua nello scavare la roccia, dove la roccia è il mio ripetere ostinatamente gli stessi errori ed inventarne di nuovi. Grazie di tutto.

Vorrei ringraziare Eugenio. E' stata una costante nel mio percorso, ed in ogni passaggio importante c'è sempre stata la sua mano. Osserva, analizza, valuta, decide. Mi ha lasciato commettere i miei errori e costruire il mio percorso, senza mai lodare i primi approssimativi tentativi. E' esigente con chi gli sta vicino come lo è con se stesso. Gli sono immensamente grato per questo. Mi ha insegnato anche che c'è un momento per le cose futili ed un momento per parlare delle cose importanti. Quando parla delle seconde, le parole sono poche ed hanno la densità di una stella di neutroni.

Vorrei ringraziare Francesca. E' la bellissima costante della mia vita. Ancora non si rende conto di quanto sia importante il suo sostegno e la sua presenza nella mia vita. Sono sempre stato convinto che in una coppia entrambi debbano essere delle persone compiute ed interessanti e lo credo ancora. Tuttavia sono sicuro che

senza di lei la mia vita sarebbe piú spenta e grigia. Penso di aver scoperto nuovi colori e nuovi suoni con lei, e spero di scoprirne tanti altri.

Vorrei ringraziare i miei quattro nonni. A chiunque vi abbia conosciuti é evidente che io sia specchio, nel bene e nel male, dei vostri caratteri.

Vorrei ringraziare Nicola e Grazia. Ci siete stati sempre e la vostra presenza e le vostre parole sono sempre state una guida per me.

Vorrei ringraziare Mike, un amico fraterno che tenta sempre di esserci e di capire cosa ti turba. Spero che la vita gli restituisca in soddisfazioni almeno una parte di ciò che lui dona agli altri.

Vorrei ringraziare Stefano, Fabia, Marco, Rossana e Giorgio. Sono da sempre una parte di me, e spero che lo siano sempre.

Vorrei ringraziare Anna Lucia. Lei e Claudia sono state per me ciò che di piú simile ci possa essere a delle figlie o a delle sorelline. Mi auguro di avervi donato almeno un po' della mia forza e della mia ostinazione nell'affrontare il mondo a testa alta. E spero che siate felici, anche se vi lamenterete sempre.

Vorrei ringraziare Domenico, l'ho visto crescere e soffrire, sbagliare e correggersi. Nonostante il percorso sia ancora lungo, spero di poter essere un buon fratello maggiore.

Deseo agradecer a Alejandro, Ivan y Pablo, que me recibieron en su laboratorio como si fuera parte de su familia. Espero devolverles su hospitalidad. Alejandro fue un buen compañero de investigación y me enseñó mucho. Gracias!

Vorrei ringraziare Fabio, Stefano, Pasquale, Pierpaolo, Alessandro. Se penso a svagarmi, a staccare da tutto, uscire e divertirmi penso a loro. Spero di riuscire in futuro a restituire loro il bene che mi vogliono.

Vorrei ringraziare Irene. Sei una buona amica e mi fai sorridere sempre, spero di riuscire a farti capire la differenza tra una carta d'identità ed una patente di guida. Vorrei ringraziare Margherita ed Emanuele, Mario e Roberta. Mi avete sempre coinvolto e fatto sentire bene. Vi apprezzo tantissimo sia individualmente che insieme.

Vorrei ringraziare Barbara, per me sei come una sorellina acquisita.

Vorrei ringraziare Michele, Andrea, Dimitri, Manuel, Mona, Alberto, Totó, Gae-

tano, Martina, Carmine, El mudo, Francesco, Nika, Marta, Eva e tutti gli amici di Madrid. Essere a Madrid con voi stata un'esperienza mistico-musicale-culinaria. Vi penso ogni giorno e spero di rivedervi presto.

Vorrei ringraziare Eliana, un'amica vittima del fuoco incrociato. E di solito il fuoco é il suo. Spero che tu possa tornare il prima possibile nelle nostre vite.

Vorrei ringraziare i ragazzi del laboratorio. Yes! Voglio ringraziare Vito, Claudio, Antonio, Felice, GianMaria, Yashar, Lucio, Carmelo, Gaetano, Daniele, Giulia e Giulio. Il mondo della ricerca vi chiede molto ogni giorno e voi rispondete pleasantly, tra un problema con pip all'abbondanza di dati. Resistete stoicamente a labirinti di discussion, interazioni, iterazioni ed esperimenti, subendo call e toss a coin. In attesa della mattina dell'eleganza, vi potete consolare con qualche dolce croccante ed un po' di pulled pork, o con pasta ed insalata per gli amanti del brivido.

Vorrei ringraziare tutti gli amici di Villa Edda, ed in primis Michele e Alessandro, siete dei fratelli. Michele é una parte fondamentale di me stesso. Parlare con lui é come parlare guardandosi allo specchio. Mi riuscirebbe impossibile immaginare una versione alternativa della mia vita in cui lui non c'è. Sarei, semplicemente, un'altra persona. Marco, Daniele, Gabriele, Bucci, Marcolino, Ciccarelli, Matteo, Davide, Flavio, Giorgio, Massimo, Mimmi, Nicola, Paolo, Pierpaolo, Stefano, Giuseppe, Carlo, Adriano, Campanale, Pellecchia. Siete tantissimi e sarebbe dura elencarvi tutti, ma non si può fare un ordinamento, voglio bene a tutti voi.

Vorrei ringraziare gli amici di Tins ed in primis Dario, Pasquale, Gianluca, Mauro e Mariano. Siete tutti sempre nel mio cuore anche se siamo lontani.