

# Workshop on New Trends in Content-based Recommender Systems (CBRecSys 2014)

Toine Bogers  
Department of Communication  
and Psychology  
Aalborg University  
Copenhagen  
2450 Copenhagen, Denmark  
toine@hum.aau.dk

Marijn Koolen  
Institute for Logic, Language  
and Computation  
University of Amsterdam  
Amsterdam, The Netherlands  
marijn.koolen@uva.nl

Iván Cantador  
Escuela Politécnica Superior  
Universidad Autónoma de  
Madrid  
28049 Madrid, Spain  
ivan.cantador@uam.es

## ABSTRACT

While content-based recommendation has been applied successfully in many different domains, it has not seen the same level of attention as collaborative filtering techniques have. However, there are many recommendation domains and applications where content and metadata play a key role, either *in addition to* or *instead of* ratings and implicit usage data. For some domains, such as movies, the relationship between content and usage data has seen thorough investigation already, but for many other domains, such as books, news, scientific articles, and Web pages we still do not know *if* and *how* these data sources should be combined to provided the best recommendation performance. The *CBRecSys 2014* workshop aims to address this by providing a dedicated venue for papers dedicated to all aspects of content-based recommendation.

## Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*information Filtering*; D.2.8 [Software Engineering]: Metrics—*performance measures*

## General Terms

Algorithms, Experimentation, Human Factors, Theory

## Keywords

recommender systems, content-based recommendation, text reviews, user-generated content, implicit feedback, semantics, context

## 1. MOTIVATION AND GOALS

While content-based recommendation has been applied successfully in many different domains [2], it has not seen the same level of attention as collaborative filtering techniques have. In recent years, competitions like the Netflix Prize<sup>1</sup>, CAMRA<sup>2</sup>, and the Ya-

<sup>1</sup><http://www.netflixprize.com/>

<sup>2</sup><http://www.dai-labor.de/camra2010/>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author. Copyright is held by the owner/author(s). *RecSys '14*, Oct 06–10 2014, San Jose or vicinity, CA, USA. ACM 978-1-4503-2668-1/14/10 <http://dx.doi.org/10.1145/2645710.2645784>.

hoo! Music KDD Cup 2011 [1] have spurred on advances in collaborative filtering and how to utilize ratings and usage data. However, there are many recommendation domains and applications where content and metadata play a key role, either *in addition to* or *instead of* ratings and implicit usage data. For some domains, such as movies, the relationship between content and usage data has seen thorough investigation already (e.g. [3]), but for many other domains, such as books, news, scientific articles, and Web pages we still do not know *if* and *how* these data sources should be combined to provided the best recommendation performance.

The *CBRecSys 2014* workshop aims to address this by providing a venue for papers dedicated to all aspects and new trends of content-based recommendation. This would include both recommendation in domains where textual content is abundant (e.g. books, news, scientific articles, jobs, educational resources, and Web pages) as well as dedicated comparisons and combinations of content-based techniques with collaborative filtering approaches.

## 2. TOPICS OF INTEREST

Relevant topics of the workshop include:

- Developing novel recommendation approaches
  - Hybrid strategies combining content-based and collaborative filtering recommendations
  - Content-based approaches to cross-system and cross-domain recommendation
  - Latent factor models for content-based and hybrid recommendation
- Exploiting user-generated content for recommendation
  - Mining microblogging data in content-based recommender systems
  - Social tag-based recommender systems
  - Exploiting Semantic Web and Linked Open Data in content-based recommender systems
- Processing text reviews
  - Estimating (implicit) ratings associated with text reviews
  - Opinion mining and sentiment analysis of text reviews to support content-based recommendation
  - Extracting user personality traits and factors from text reviews for recommendation

- Mining contextual data from content
  - Extraction of contextual signals from textual content for recommendation
  - Incorporating the temporal dimension in content-based recommendation
  - Mood-based recommender systems
- Addressing limitations of recommender systems
  - Addressing the cold-start problem with content-based recommendation approaches
  - Increasing diversity of content-based recommendations
  - Providing novelty in content-based recommendations

In particular, papers submitted to the the workshop have focused on the following topics. Several papers present hybrid systems combining collaborative filtering and content-based recommendation, finding them complementary, with content-based recommendation components especially suitable for tackling the cold-start problem. Other papers investigate how different content features can be used for similarity measures and explore ways to identify which features are the most relevant for a given context. Some papers present approaches to mine user reviews for inferring user preferences on specific attributes of items, essentially deriving more structured feature information from unstructured text. Finally, several papers look at semantic frameworks and Linked Open Data to measure item similarity across different domains.

### 3. CHALLENGE ON BOOK RECOMMENDATION

To facilitate exploration of the above mentioned topics, the workshop features an in-workshop challenge on book recommendation. This challenge focuses on recommending new, interesting books to LibraryThing users based on usage data (which books they have added to their collection) and content-based information about the books available in LibraryThing. The rich textual nature of the task makes the challenge an excellent venue to revisit questions about the benefits of content-based filtering vs. collaborative filtering, and metadata versus ratings information. At the workshop the evaluation results of the challenge were presented.

#### 3.1 Dataset

For this challenge, a large dataset containing user profiles with book ratings and tags, and 2.8 million book descriptions with library metadata, user ratings, tags, and reviews from Amazon and LibraryThing was made available.

The dataset for the book recommendation challenge is comprised of two parts: usage data and book metadata. The first part of the dataset for book recommendation is a log of usage data: who added which books to their collection at what point in time. In addition to this, ratings and tags assigned to books are also included in the usage dataset (where available). The user profiles in this data set contain a total number of 1,830,958 unique books added by 78,633 different LibraryThing users, anonymized through different privacy-preserving measures. This usage data serves as the main data source for evaluating our challenge, which is described in more detail in Section 3.2.

The second part of our challenge dataset for book recommendation is a collection of metadata records for 2.8 million books cataloged on LibraryThing. This collection was crawled from Amazon and LibraryThing by the University of Duisburg-Essen in early

2009. From Amazon, there is formal metadata like book title, author, publisher, publication year, library classification codes, Amazon categories, and similar product information, as well as user-generated content in the form of user ratings and reviews. From LibraryThing, there are user tags and user-provided metadata on awards, book characters and locations, and blurbs. This part of the challenge data has been used successfully for more retrieval-oriented challenges at the INEX 2011–2014 Social Book Search tracks<sup>3</sup>.

#### 3.2 Evaluation

The evaluation of the book recommendation challenge follows the familiar *backtesting* paradigm, where a small number of randomly selected books is withheld for each user with the remaining data to be used as training material. If a user's withheld items are predicted at the top of the ranked result list, i.e., if the algorithm is able to correctly predict the user's interest in those withheld items, then the algorithm has performed well. The main evaluation metric used in the challenge is the ranking-based metric NDCG@10, as ratings information is sparse in the usage data.

To make evaluation easier for challenge participants and to reduce the possibility of over-fitting, we divided the usage data into a training and a validation set. The training material contains 90% of the user profiles and is meant for training the participants' recommendation algorithms. To aid in the learning and parameter optimization phase on the training material, we performed 10-fold cross-validation on the training set. This resulted in 10 training and test sets, one for each of the folds. We encouraged all participants to use these 10 folds to train their system, so the results of the training phase are directly comparable among participants.

The validation set contains the remaining 10% of the user profiles (or 7863 users) and was released at a later date its aim was to produce the final comparison of the different submitted approaches. This validation set contains a set of users that are not included in the original training material to avoid over-fitting. To allow participants to train their systems on these new users as well, we also released a small amount of extra training material corresponding to the users near the end of the challenge.

### 4. WEBSITE AND PROCEEDINGS

The workshop material (list of accepted papers, invited talk, and the workshop schedule) can be found on the CBRecSys 2014 workshop website at <http://ir.ii.uam.es/cbreccsys2014/>. The proceedings are published as a CEUR Workshop Proceedings volume.

### 5. REFERENCES

- [1] G. Dror, N. Koenigstein, Y. Koren, and M. Weimer. The Yahoo! Music Dataset and KDD-Cup '11. In *JMLR Workshop and Conference Proceedings*, volume 18 of *Proceedings of KDD Cup 2011*, pages 3–18, 2012.
- [2] P. Lops, M. de Gemmis, and G. Semeraro. Content-based Recommender Systems: State of the Art and Trends. In *Recommender Systems Handbook*, pages 73–105. 2011.
- [3] I. Pilászy and D. Tikk. Recommending New Movies: Even a Few Ratings Are More Valuable Than Metadata. In *RecSys '09: Proceedings of the Third ACM Conference on Recommender Systems*, pages 93–100. ACM, 2009.

<sup>3</sup><https://inex.mmci.uni-saarland.de/tracks/books/>