# Evaluating Tag Recommendations for E-Book Annotation Using a Semantic Similarity Metric

**Emanuel Lacic***
Know-Center GmbH
Graz, Austria
elacic@know-center.at

**Dominik Kowald***
Know-Center GmbH
Graz, Austria
dkowald@know-center.at

**Dieter Theiler**
Know-Center GmbH
Graz, Austria
dtheiler@know-center.at

**Matthias Traub**
Know-Center GmbH
Graz, Austria
mtraub@know-center.at

**Lucky Kuffer**
HGV GmbH
Munich, Germany
lucky.kuffer@hgv-online.de

**Stefanie Lindstaedt**
Know-Center GmbH
Graz, Austria
slind@know-center.at

**Elisabeth Lex**
Graz University of Technology
Graz, Austria
elisabeth.lex@tugraz.at

## ABSTRACT

In this paper, we present our work to support publishers and editors in finding descriptive tags for e-books through tag recommendations. We propose a hybrid tag recommendation system for e-books, which leverages search query terms from Amazon users and e-book metadata, which is assigned by publishers and editors. Our idea is to mimic the vocabulary of users in Amazon, who search for and review e-books, and to combine these search terms with editor tags in a hybrid tag recommendation approach. In total, we evaluate 19 tag recommendation algorithms on the review content of Amazon users, which reflects the readers' vocabulary. Our results show that we can improve the performance of tag recommender systems for e-books both concerning tag recommendation accuracy, diversity as well as a novel semantic similarity metric, which we also propose in this paper.

## KEYWORDS

Tag Recommendation; E-Book Annotation; Hybrid Recommendation; Amazon Search Query Terms; Diversity; Semantic Similarity

---

*Both authors contributed equally to this work.

## 1 INTRODUCTION

When people shop for books online in e-book stores such as, e.g., the Amazon Kindle store, they enter search terms with the goal to find e-books that meet their preferences. Such e-books have a variety of metadata such as, e.g., title, author or keywords, which can be used to retrieve e-books that are relevant to the query. As a consequence, from the perspective of e-book publishers and editors, annotating e-books with tags that best describe the content and which meet the vocabulary of users (e.g., when searching and reviewing e-books) is an essential task [22].

**Problem and aim of this work.** Annotating e-books with suitable tags is, however, a complex task as users' vocabulary may differ from the one of editors. Such a vocabulary mismatch yet hinders effective organization and retrieval [21] of e-books. For example, while editors mostly annotate e-books with descriptive tags that reflect the book's content, Amazon users often search for parts of the book title. In the data we use for the present study (see Section 2), we find that around 30% of the Amazon search terms contain parts of e-book titles.

In this paper, we present our work to support editors in the e-book annotation process with tag recommendations [5, 8]. Our idea is to exploit user-generated search query terms in Amazon to mimic the vocabulary of users in Amazon, who search for e-books. We combine these search terms with tags assigned by editors in a hybrid tag recommendation approach. Thus, our aim is to show that we can improve the performance of tag recommender systems for e-books both concerning recommendation accuracy as well as semantic similarity and tag recommendation diversity.

**Related work.** In tag recommender systems, mostly content-based algorithms (e.g., [13, 14]) are used to recommend tags to annotate resources such as e-books. In our work, we incorporate both content features of e-books (i.e., title and description text) as well as Amazon search terms to account for the vocabulary of e-book readers.

Concerning the evaluation of tag recommendation systems, most studies focus on measuring the accuracy of tag recommendations

(e.g., [5]). However, the authors of [2] suggest also to use beyond-accuracy metrics such as diversity to evaluate the quality of tag recommendations. In our work, we measure recommendation diversity in addition to recommendation accuracy and propose a novel metric termed semantic similarity to validate semantic matches of tag recommendations.

**Approach and findings.** We exploit editor tags and user-generated search terms as input for tag recommendation approaches. Our evaluation comprises of a rich set of 19 different algorithms to recommend tags for e-books, which we group into (i) popularity-based, (ii) similarity-based (i.e., using content information), and (iii) hybrid approaches. We evaluate our approaches in terms of accuracy, semantic similarity and diversity on the review content of Amazon users, which reflects the readers' vocabulary. With semantic similarity, we measure how semantically similar (based on learned Doc2Vec [12] embeddings) the list of recommended tags is to the list of relevant tags. We use this additional metric to measure not only exact "hits" of our recommendations but also semantic matches.

Our evaluation results show that combining both data sources enhances the quality of tag recommendations for annotating e-books. Furthermore, approaches that solely train on Amazon search terms provide poor performance in terms of accuracy but deliver good results in terms of semantic similarity and recommendation diversity.

## 2 METHOD

In this section, we describe our dataset as well as our tag recommendation approaches we propose to annotate e-books.

### 2.1 Dataset

Our dataset contains two sources of data, one to generate tag recommendations and another one to evaluate tag recommendations. HGV GmbH has collected all data sources[1] and we provide the dataset statistics in Table 1.

**Data used to generate recommendations.** We employ two sources of e-book annotation data: (i) editor tags, and (ii) Amazon search terms. For editor tags, we collect data of 48,705 e-books from 13 publishers, namely Kunstmann, Delius-Klasnig, VUR, HJR, Diogenes, Campus, Kiwi, Beltz, Chbeck, Rowohlt, Droemer, Fischer and Neopubli. Apart from the editor tags, this data contains metadata fields of e-books such as the ISBN, the title, a description text, the author and a list of BISACs, which are identifiers for book categories.

For the Amazon search terms, we collect search query logs of 21,243 e-books for 12 months (i.e., November 2017 to October 2018). Apart from the search terms, this data contains the e-books' ISBNs, titles and description texts.

Table 1 shows that the overlap of e-books that have editor tags and Amazon search terms is small (i.e., only 497). Furthermore, author and BISAC (i.e., the book category identifier) information are primarily available for e-books that contain editor tags. Consequently, both data sources provide complementary information,

---

[1]Currently the data used in this study cannot be made publicly available because of copyright issues, but we will try to provide a public version of it soon in the future.

| Data used to generate recommendations | # |
|---|---|
| Number of e-books | 69,451 |
| thereof with editor tags | 48,705 |
| thereof with Amazon search terms | 21,243 |
| thereof with editor tags and Amazon search terms | 497 |
| Number of distinct authors | 25,086 |
| Number of distinct BISACs (= category IDs) | 1,448 |
| Number of distinct editor tags | 114,707 |
| Number of distinct Amazon search terms | 8,240 |
| **Data used to evaluate recommendations** | **#** |
| Number of e-books with Amazon review keywords | 2,896 |
| Number of distinct Amazon review keywords | 33,663 |
| Avg. number of distinct Amazon review keywords per e-book | 30 |

**Table 1: Statistics of our e-book annotation dataset used to generate and evaluate tag recommendations.**

which underpins the intention of this work, i.e., to evaluate tag recommendation approaches using annotation sources from different contexts.

**Data used to evaluate recommendations.** For evaluation, we use a third set of e-book annotations, namely Amazon review keywords. These review keywords are extracted from the Amazon review texts and are typically provided in the review section of books on Amazon. Our idea is to not favor one or the other data source (i.e., editor tags and Amazon search terms) when evaluating our approaches against expected tags. At the same time, we consider Amazon review keywords to be a good mixture of editor tags and search terms as they describe both the content and the users' opinions on the e-books (i.e., the readers' vocabulary). As shown in Table 1, we collect Amazon review keywords for 2,896 e-books (publishers: Kiwi, Rowohlt, Fischer, and Droemer), which leads to 33,663 distinct review keywords and on average 30 keyword assignments per e-book.

### 2.2 Tag Recommendation Approaches

We implement three types of tag recommendation approaches, i.e., (i) popularity-based, (ii) similarity-based (i.e., using content information), and (iii) hybrid approaches. Due to the lack of personalized tags (i.e., we do not know which user has assigned a tag), we do not implement other types of algorithms such as collaborative filtering [15]. In total, we evaluate 19 different algorithms to recommend tags for annotating e-books.

**Popularity-based approaches.** We recommend the most frequently used tags in the dataset, which is a common strategy for tag recommendations [10]. That is, a most popular $MP_{Editor}$ approach for editor tags and a most popular $MP_{Amazon}$ approach for Amazon search terms. For e-books, for which we also have author (= $MP_{Editor}^{Author}$ and $MP_{Amazon}^{Author}$) or BISAC (=$MP_{Editor}^{BISAC}$ and $MP_{Amazon}^{BISAC}$) information, we use these features to further filter the recommended tags, i.e., to only recommend tags that were used to annotate e-books of a specific author or a specific BISAC.

We combine both data sources (i.e., editor tags and Amazon search terms) using a round-robin combination strategy, which ensures an equal weight for both sources. This gives us three additional popularity-based algorithms (= $MP_{Combined}$, $MP_{Combined}^{Author}$ and $MP_{Combined}^{BISAC}$).

**Similarity-based approaches.** We exploit the textual content of e-books (i.e., description or title) to recommend relevant tags [4]. For this, we first employ a content-based filtering approach [1] based on TF-IDF [18] to find top-$N$ similar e-books[2]. For each of the similar e-books, we then either extract the assigned editor tags (= $SIM_{Editor}^{Description}$ and $SIM_{Editor}^{Title}$) or the Amazon search terms (= $SIM_{Amazon}^{Description}$ and $SIM_{Amazon}^{Title}$). To combine the tags of the top-$N$ similar e-books, we use the cross-source algorithm [3], which favors tags that were used to annotate more than one similar e-book (i.e., tags that come from multiple recommendation sources). The final tag relevancy is calculated as:

$$W_{t_i} = |S_{t_i}| \cdot \sum_{s_{t_i} \in S} W_{s_{t_i}} \qquad (1)$$

where $|S_{t_i}|$ denotes the number of distinct e-books, which yielded the recommendation of tag $t_i$, to favor tags that come from multiple sources and $W_{s_{t_i}}$ is the similarity score of the corresponding e-book. We again use a round-robin strategy to combine both data sources (= $SIM_{Combined}^{Description}$ and $SIM_{Combined}^{Title}$).

**Hybrid approaches.** We use the previously mentioned cross-source algorithm [3] to construct four hybrid recommendation approaches. In this case, tags are favored that are recommended by more than one algorithm.

Hence, to create a popularity-based hybrid (= $HYB^{MP}$), we combine the best three performing popularity-based approaches from the ones (i) without any contextual signal, (ii) with the author as context, and (iii) with BISAC as context. In the case of the similarity-based hybrid (= $HYB^{SIM}$), we utilize the two best performing similarity-based approaches from the ones (i) which use the title, and (ii) which use the description text. We further define $HYB^{All}$, a hybrid approach that combines the three popularity-based methods of $HYB^{MP}$ and the two similarity-based approaches of $HYB^{SIM}$. Finally, we define $HYB^{Best}$ as a hybrid approach that uses the best performing popularity-based and the best performing similarity-based approach (see Figure 1 in Section 4 for more details about the particular algorithm combinations).

## 3 EXPERIMENTAL SETUP

In this section, we describe our evaluation protocol as well as the measures we use to evaluate and compare our tag recommendation approaches.

### 3.1 Evaluation Protocol

For evaluation, we use the third set of e-book annotations, namely Amazon review keywords. As described in Section 2.1, these review keywords are extracted from the Amazon review texts and thus, reflect the users' vocabulary. We evaluate our approaches for the 2,896 e-books, for whom we got review keywords. To follow common practice for tag recommendation evaluation [9], we predict the assigned review keywords (= our test set) for respective e-books.

### 3.2 Evaluation Metrics

In this work, we measure (i) recommendation accuracy, (ii) semantic similarity, and (iii) recommendation diversity to evaluate the quality of our approaches from different perspectives.

**Recommendation accuracy.** We use *Normalized Discounted Cumulative Gain* (nDCG) [17] to measure the accuracy of the tag recommendation approaches. The nDCG measure is a standard ranking-dependent metric that not only measures how many tags can be correctly predicted but also takes into account their position in the recommendation list with length of $k$. It is based on the *Discounted Cummulative Gain*, which is given by:

$$DCG@k = \sum_{k=1}^{|r_b^k|} \left( \frac{2^{T(k)} - 1}{log_2(1 + k)} \right) \qquad (2)$$

where $T(k)$ is a function that returns 1 if the recommended tag at position $i$ in the recommended list is relevant. We then calculate DCG@$k$ for every evaluated e-book by dividing DCG@$k$ with the ideal DCG value iDCG@$k$, which is the highest possible DCG value that can be achieved if all the relevant tags would be recommended in the correct order. It is given by the following formula [17]:

$$nDCG@k = \frac{1}{|B|} \sum_{b \in B} \left( \frac{DCG@k}{iDCG@k} \right) \qquad (3)$$

**Semantic similarity.** One precondition of standard recommendation accuracy measures is that to generate a "hit", the recommended tag needs to be an exact syntactical match to the one from the test set. When tags are recommended from one data source and compared to tags from another source, this can be problematic. For example, if we recommend the tag "victim" but expect the tag "prey", we would mark this as a mismatch, therefore being a bad recommendation. But if we know that the corresponding e-book is a crime novel, the recommended tag would be (semantically) descriptive to reflect the book's content. Hence, in this paper, we propose to additionally measure the semantic similarity between recommended tags and tags from the test set (i.e., the Amazon review keywords).

Over the last four years, there have been several notable publications in the area of applying deep learning to uncover semantic relationships between textual content (e.g., by learning word embeddings with Word2Vec [7, 16]). Based on this, we propose an alternative measure of recommendation quality by learning the semantic relationships from both vocabularies and then using it to compare how semantically similar the recommended tags are to the expected review keywords. For this, we first extract the textual content in the form of the description text, title, editor tags and Amazon search terms of e-books from our dataset[3]. We then train a Doc2Vec [12] model[4] on the content. Then, we use the model to infer the latent representation for both the complete list of recommended tags as well as the list of expected tags from the test

---

[2]In our experiments, we set $N = 20$, the minimum document frequency to 10 and the minimum word length to 5.

---

[3]We also pre-process the extracted text by removing bad characters, stop words and changing all words to lowercase.
[4]We use the DBOW approach with a size of 50 for the latent vector representation, negative sampling of 10, a learning rate of 0.025 and train it for 10 epochs.
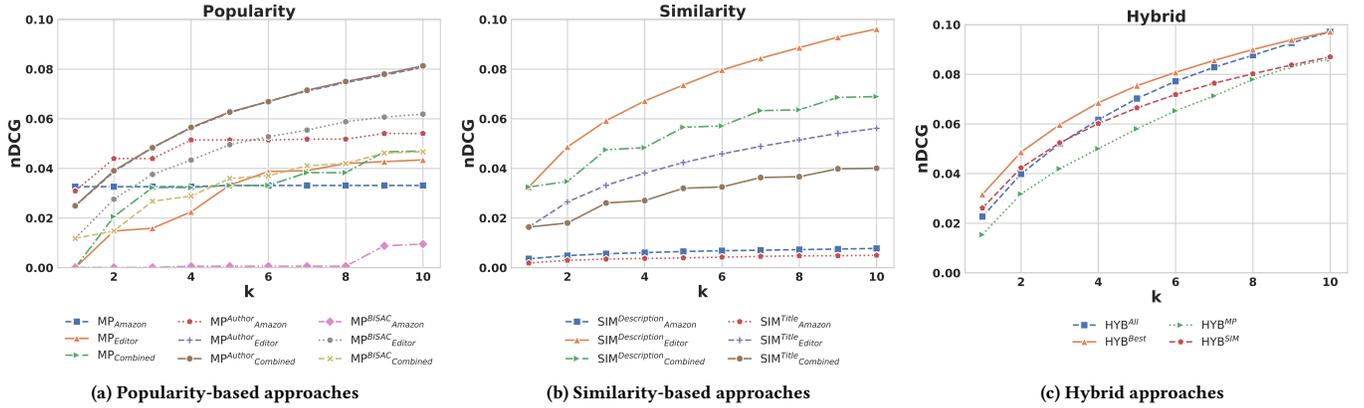
**Figure 1: Accuracy results with respect to $nDCG$ for (a) popularity-based, (b) similarity-based and (c) hybrid tag recommendation approaches. All results are reported for different numbers of recommended tags (i.e., $k \in [1, 10]$).**

set. Finally, we use the cosine similarity measure to calculate how semantically similar these two lists are.

**Recommendation diversity.** As defined in [19], we calculate recommendation diversity as the average dissimilarity of all pairs of tags in the list of recommended tags. Thus, given a distance function $d(t_i, t_j)$ that corresponds to the dissimilarity between two tags $t_i$ and $t_j$ in the list of recommended tags, $D$ is given as the average dissimilarity of all pairs of tags:

$$D@k = \frac{1}{|B|} \sum_{b \in B} \left( \frac{1}{k \cdot (k-1)} \sum_{i \in R} \sum_{j \in r_b^k, j \neq i} d(t_i, t_j) \right) \quad (4)$$

where $|B|$ is the number of evaluated e-books and the dissimilarity function is defined as $d(t_i, t_j) = 1 - sim(t_i, t_j)$. In our experiments, we use the previously trained Doc2Vec model to extract the latent representation of a specific tag. The similarity of two tags $sim(t_i, t_j)$ is then calculated with the Cosine similarity measure using the latent vector representations of respective tags $t_i$ and $t_j$.

## 4 RESULTS

Concerning tag recommendation accuracy, in this section, we report results for different values of $k$ (i.e., number of recommended tags). For the beyond-accuracy experiment, we use the full list of recommended tags (i.e., $k = 10$).

### 4.1 Recommendation Accuracy Evaluation

Figure 1 shows the results of the accuracy experiment for the (i) popularity-based, (ii) similarity-based, and (iii) hybrid tag recommendation approaches.

**Popularity-based approaches.** In Figure 1a, we see that popularity-based approaches based on editor tags tend to perform better than if trained on Amazon search terms. If we take into account contextual information like BISAC or author, we can further improve accuracy in terms of $nDCG$. That is, we find that using popular tags from e-books of a specific author leads to the best accuracy of the popularity-based approaches. This suggests that editors and readers do seem to reuse tags for e-books of same authors. If we

use both editor tags and Amazon search terms, we can further increase accuracy, especially for higher values of $k$ like in the case of $MP_{Combined}$. This is, however, not the case for $MP_{Combined}^{BISAC}$ as the accuracy of the integrated $MP_{Amazon}^{BISAC}$ approach is low. The reason for this is the limited amount of e-books from within the Amazon search query logs that have BISAC information (i.e., only 2.38%).

**Similarity-based approaches.** We further improve accuracy if we first find similar e-books and then extract their top-$k$ tags in a cross-source manner as described in Section 2.2.

As shown in Figure 1b, using the description text to find similar e-books results in more accurate tag recommendations than using the title (i.e., $nDCG@10 = 0.0961$ for $SIM_{Editor}^{Description}$). This is somehow expected as the description text consists of a bigger corpus of words (i.e., multiple sentences) than the title. Concerning the collected Amazon search query logs, extracting and then recommending tags from this source results in a much lower accuracy performance. Thus, these results also suggest to investigate beyond-accuracy metrics as done in Section 4.2.

**Hybrid approaches.** Figure 1c shows the accuracy results of the four hybrid approaches. By combining the best three popularity-based approaches, we outperform all of the initially evaluated popularity algorithms (i.e., $nDCG@10 = 0.0862$ for $HYB^{MP}$). On the contrary, the combination of the two best performing similarity-based approaches $SIM_{Editor}^{Description}$ and $SIM_{Editor}^{Title}$ does not yield better accuracy. The negative impact of using a lower-performing approach such as $SIM_{Editor}^{Title}$ within a hybrid combination can also be observed in $HYB^{All}$ for lower values of $k$. Overall, this confirms our initial intuition that combining the best performing popularity-based approach with the best similarity-based approach should result in the highest accuracy (i.e., $nDCG@10 = 0.0972$ for $HYB^{Best}$). Moreover, our goal, namely to exploit editor tags in combination with search terms used by readers to increase the metadata quality of e-books, is shown to be best supported by applying hybrid approaches as they provide the best prediction results.
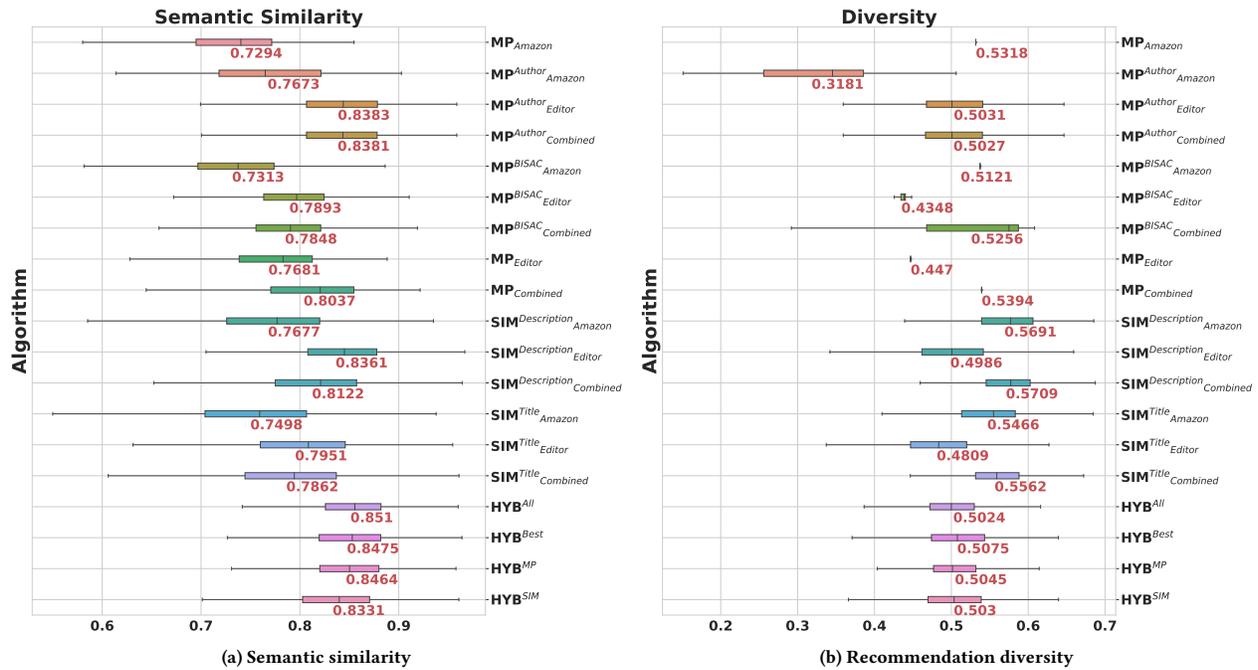
**Figure 2: Beyond-accuracy evaluation results of our tag recommendation approaches. We use the list of 10 recommended tags to calculate the (a) semantic similarity, and (b) recommendation diversity. We provide the boxplots and the mean values for the approaches.**

## 4.2 Beyond-Accuracy Evaluation

Figure 2 illustrates the results of the experiments, which measure the recommendation impact beyond-accuracy.

**Semantic similarity.** Figure 2a illustrates the results of our proposed semantic similarity measure. To compare our proposed measure to standard accuracy measures such as *nDCG*, we use Kendall's Tau rank correlation [6] as suggested by [11] for automatic evaluation of information-ordering tasks. From that, we rank our recommendation approaches according to both accuracy and semantic similarity and calculate the relation between both rankings. This results in $\tau = 0.743$ with a p-value < 0.00001, which suggests a high correlation between the semantic similarity and the standard accuracy measure.

Therefore, the semantic similarity measure helps us interpret the recommendation quality. For instance, we achieve the lowest *nDCG* values with the similarity-based approaches that recommend Amazon search terms (i.e., $SIM^{Description}_{Amazon}$ and $SIM^{Title}_{Amazon}$). When comparing these results with others from Figure 1b, a conclusion could be quickly drawn that the recommended tags are merely unusable. However, by looking at Figure 2a, we see that, although these approaches do not provide the highest recommendation accuracy, they still result in tag recommendations that are semantically related at a high degree[5] to the expected annotations from the test set. Overall, this suggests that approaches, which provide a poor

accuracy performance concerning *nDCG* but provide a good performance regarding semantic similarity could still be helpful for annotating e-books.

**Recommendation diversity.** Figure 2b shows the diversity of the tag recommendation approaches. We achieve the highest diversity with the similarity-based approaches, which extract Amazon search terms. Their accuracy is, however, very low. Thus, the combination of the two vocabularies can provide a good trade-off between recommendation accuracy and diversity.

## 5 CONCLUSION AND FUTURE WORK

In this paper, we present our work to support editors in the e-book annotation process. Specifically, we aim to provide tag recommendations that incorporate both the vocabulary of the editors and e-book readers. Therefore, we train various configurations of tag recommender approaches on editors' tags and Amazon search terms and evaluate them on a dataset containing Amazon review keywords. We find that combining both data sources enhances the quality of tag recommendations for annotating e-books. Furthermore, while approaches that train only on Amazon search terms provide poor performance concerning recommendation accuracy, we show that they still offer helpful annotations concerning recommendation diversity as well as our novel semantic similarity metric.

**Future work.** For future work, we plan to validate our findings using another dataset, e.g., by recommending tags for scientific articles and books in BibSonomy. With this, we aim to demonstrate

---

[5]A semantic similarity of 1.0 would denote a semantically (and syntactically) perfect fit of tag recommendations to the test set.

the usefulness of the proposed approach in a similar domain and to enhance the reproducibility of our results by using an open dataset.

Moreover, we plan to evaluate our tag recommendation approaches in a study with domain users. Also, we want to improve our similarity-based approaches by integrating novel embedding approaches [7, 16] as we did, for example, with our proposed semantic similarity evaluation metric. Finally, we aim to incorporate explanations for recommended tags so that editors of e-book annotations receive additional support in annotating e-books [20]. By making the underlying (semantic) reasoning visible to the editor who is in charge of tailoring annotations, we aim to support two goals: (i) allowing readers to discover e-books more efficiently, and (ii) enabling publishers to leverage semi-automatic categorization processes for e-books. In turn, providing explanations fosters control over which vocabulary to choose when tagging e-books for different application contexts.

## REFERENCES

[1] M. Balabanovic and Y. Shoham. Fab: content-based, collaborative recommendation. *Communications of the ACM*, 40(3):66–73, 1997.
[2] F. Belém, R. Santos, J. Almeida, and M. Gonçalves. Topic diversity in tag recommendation. In *Proc. of RecSys'13*, pages 141–148. ACM, 2013.
[3] S. Bostandjiev, J. O'Donovan, and T. Höllerer. Tasteweights: a visual interactive hybrid recommender system. In *Proc. of RecSys'12*, pages 35–42. ACM, 2012.
[4] I. Cantador, A. Bellogín, and D. Vallet. Content-based recommendation in social tagging systems. In *Proc. of RecSys'10*, pages 237–240. ACM, 2010.
[5] R. Jäschke, L. Marinho, A. Hotho, L. Schmidt-Thieme, and G. Stumme. Tag recommendations in folksonomies. In *European Conference on Principles of Data Mining and Knowledge Discovery*, pages 506–514. Springer, 2007.
[6] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2), 1938.
[7] T. Kenter and M. De Rijke. Short text similarity with word embeddings. In *Proc. of CIKM'15*, pages 1411–1420. ACM, 2015.
[8] D. Kowald, S. Kopeinik, and E. Lex. The tagrec framework as a toolkit for the development of tag-based recommender systems. In *Adjunct Publication of the 25th Conference on User Modeling, Adaptation and Personalization*, pages 23–28. ACM, 2017.
[9] D. Kowald and E. Lex. Evaluating tag recommender algorithms in real-world folksonomies: A comparative study. In *Proc. of RecSys'15*. ACM, 2015.
[10] D. Kowald and E. Lex. The influence of frequency, recency and semantic context on the reuse of tags in social tagging systems. In *Proc. of ACM HT'16*, 2016.
[11] M. Lapata. Automatic evaluation of information ordering: Kendall's tau. *Computational Linguistics*, 32(4):471–484, 2006.
[12] Q. Le and T. Mikolov. Distributed representations of sentences and documents. In *Proc. of ICML'14*.
[13] P. Lops, M. De Gemmis, G. Semeraro, C. Musto, and F. Narducci. Content-based and collaborative techniques for tag recommendation: an empirical evaluation. *Journal of Intelligent Information Systems*, 40(1):41–61, 2013.
[14] Y.-T. Lu, S.-I. Yu, T.-C. Chang, and J. Y.-j. Hsu. A content-based method to enhance tag recommendation. In *Twenty-First International Joint Conference on Artificial Intelligence*, 2009.
[15] L. B. Marinho and L. Schmidt-Thieme. Collaborative tag recommendations. In *Data Analysis, Machine Learning and Applications*. Springer, 2008.
[16] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
[17] D. Parra and S. Sahebi. Recommender systems : Sources of knowledge and evaluation metrics. In *Advanced Techniques in Web Intelligence-2: Web User Browsing Behaviour and Preference Analysis*, pages 149–175. Springer, 2013.
[18] J. Ramos et al. Using tf-idf to determine word relevance in document queries. In *Proc. of first instructional conference on machine learning*, volume 242, 2003.
[19] B. Smyth and P. McClave. Similarity vs. diversity. In *Proc. of ICCBR '01*, pages 347–361. Springer-Verlag, 2001.
[20] J. Vig, S. Sen, and J. Riedl. Tagsplanations: explaining recommendations using tags. In *Proc. of IUI'09*, pages 47–56. ACM, 2009.
[21] L. Zhao and J. Callan. Term necessity prediction. In *Proc. of CIKM'2010)*, 2010.
[22] A. Zubiaga, C. Körner, and M. Strohmaier. Tags vs shelves: from social tagging to social classification. In *Proc. of HT'11*, pages 93–102. ACM, 2011.