# Profiling vs. Time vs. Content: What does Matter for Top-k Publication Recommendation based on Twitter Profiles?

Chifumi Nishioka
Kiel University, Germany
ZBW – Leibniz Information Centre for
Economics, Germany
chni@informatik.uni-kiel.de

Ansgar Scherp
ZBW – Leibniz Information Centre for
Economics, Germany
Kiel University, Germany
a.scherp@zbw.eu

## ABSTRACT

So far it is unclear how different factors of a scientific publication recommender system based on users' tweets have an influence on the recommendation performance. We examine three different factors, namely profiling method, temporal decay, and richness of content. Regarding profiling, we compare CF-IDF that replaces terms in TF-IDF by semantic concepts, HCF-IDF as novel hierarchical variant of CF-IDF, and topic modeling. As temporal decay functions, we apply sliding window and exponential decay. In terms of the richness of content, we compare recommendations using both full-texts and titles of publications and using only titles. Overall, the three factors make twelve recommendation strategies. We have conducted an online experiment with 123 participants and compared the strategies in a within-group design. The best recommendations are achieved by the strategy combining CF-IDF, sliding window, and with full-texts. However, the strategies using the novel HCF-IDF profiling method achieve similar results with just using the titles of the publications. Therefore, HCF-IDF can make recommendations when only short and sparse data is available.

## 1. INTRODUCTION

The social media platform Twitter is popular among scientists to share and discuss their professional thoughts and interests [14]. Thus, they are a natural resource for building up a user's professional profile and using it for recommending scientific publications. Recommending scientific publications based on a user's social media items has several advantages: First, users receive recommendations based on their current and ongoing professional interests. In contrast, systems like Google Scholar and Sugiyama et al. [26] recommend scientific publications based on a user's publication record. It can take up to two years (for conferences) or longer (for journals) until a paper is taken into consideration by the recommender system. Second, content-based profiling from a user's social media items mitigates the well-

known cold-start problem observed in collaborative filtering systems [11]. The cold-start problem refers to the initial situation where a recommender system yet does not know anything about a user's interests. Collaborative filtering systems need to analyze a large amount of user activities in order to provide reasonable recommendations. In contrast, content-based recommender systems like our work make recommendations based on similarity scores between a user profile and candidate items. Therefore, they can generate recommendations based on a single user profile already.

There is various research on user profiling from social media items [5, 21, 24, 29] and recommending scientific publications [15, 26, 28]. However, it is unclear how different profiling methods affect the recommendation performance. In addition, the age of social media items as well as scientific publications has an influence on profiling [24, 21]. But again, it has not been compared. Finally, we investigate whether it is possible to make reasonable recommendations when using only the publications' titles, i. e., when only short and sparse information about the candidate items is available. We have conducted an online experiment to evaluate these three factors of top-$k$ recommendations of scientific publications based on a user's social media profile. In detail, the factors are:

**(i) Profiling Method:** The first factor is the *Profiling Method*, where we use Concept Frequency Inverse Document Frequency (CF-IDF) [7] as baseline. CF-IDF is a modification of TF-IDF where term frequencies are replaced by frequencies of semantic concepts. In an experiment with 19 participants, Goossen et al. have shown that CF-IDF outperforms TF-IDF for news article recommendations [7]. Recently, we have extended the statistical strength of CF-IDF with the semantics provided by a hierarchical knowledge base [19]. The resulting Hierarchical CF-IDF (HCF-IDF) model is capable of revealing semantic concepts that are not explicitly mentioned in texts but still are highly relevant. This is achieved by applying a spreading activation over a hierarchical knowledge base, which is typically provided as domain-specific taxonomy. Please note that we also considered using BM25 and TF-IDF as profiling method. However, our earlier work showed that HCF-IDF performs better for user profiling from social media items [19]. As third method, we apply Latent Dirichlet Allocation (LDA) [2, 1], a state-of-the-art topic modeling method. LDA is a generative machine learning approach and thus does not require any prior information such as a knowledge base.

**(ii) Decay Function:** As the second factor, we investigate two temporal *Decay Function*s. They are based on

the idea that the importance of information declines gradually as time passes. We compare sliding window [24] and exponential decay [21, 26]. Both decay functions have been used in the past for user profiling [24, 21, 26]. But so far no comparative study was carried out.

**(iii) Document Content:** The third factor defines the richness of *Document Content* used for profiling candidate items (i.e., scientific publications). We compare the use of full-texts and titles of scientific publications for profiling versus profiling only based on titles.

We compared twelve recommendation strategies making use of different combinations of the three factors described above. For the experiment, we have recruited $n = 123$ participants who are posting about their professional interests on Twitter. For each strategy, the participants have received recommendations of five publications from a large corpus of $|D| = 279,381$ scientific publications in the broader field of economics. We used rankscore [4] to measure the recommendation performance. We also computed Mean Average Precision (MAP), Precision, Mean Reciprocal Rank (MRR), and normalized Discounted Cumulative Gain (nDCG), which show similar results and documented in the TR [20].

The results are very interesting: The strategy that employs the profiling method CF-IDF and the decay function Sliding window with both titles and full-texts achieves the overall best recommendation performance. Although the strategy using CF-IDF shows the highest performance, it has a drawback that it requires full-texts of scientific publications. Thus, it is remarkable that the strategies with HCF-IDF can achieve comparable results using only titles. We observe no significant difference between the best performing strategy and strategies with HCF-IDF. Thus, we conclude that the use of the spreading activation function over the hierarchical knowledge base enables HCF-IDF to compensate for the sparseness when only titles are available due to e.g., legal reasons to hinder the use of full-texts. Please note, there is no lack in domain-specific hierarchical knowledge bases such as the one used in the experiment for economics. In fact, these knowledge bases are freely available for many domains[1]. Furthermore, they are manually crafted by domain experts and thus are of high quality.

In addition, we have applied a correlation analysis between the recommendation performance and the number of tweets a participant has published, the number of concepts extracted from these tweets, the number of concepts extracted per tweet, and the percentage of tweets containing at least one concept respectively. Our results show no significant correlations in any strategies. Thus, the methods are robust against the amount of tweets.

Subsequently, we review related work in Section 2. Section 3 introduces the problem definition. In Section 4, we describe the three experimental factors used in o4ur recommender system. We present the experiment setup and procedure in Section 5. The results are presented in Section 6 and discussed in Section 7 before we conclude the paper.

## 2. RELATED WORK

Recommender systems are categorized into content-based recommender and collaborative filtering [11]. Collaborative filtering requires analyzing a large amount of user activities in order to predict items to other users [29]. In contrast, we

---

[1]http://www.w3.org/2001/sw/wiki/SKOS/Datasets

focus on content-based recommender, which suggest items based on similarity scores between a user profile and candidate items. A content-based recommender can make recommendations based on data from a single user already. Thus it does not suffer from the cold start problem. Recommender systems for scientific publications mostly employed user profiles based on publications [26, 27] or clicks [15]. Instead, we create user profiles based on social media items.

Many works have extracted user interests from social media platforms [5, 21, 24, 29]. Chen et al. [5] studied a recommender system incorporating Twitter, which recommended URLs based on a user's tweets and follower-followee relationships. In order to find out the best recommendation strategy, they evaluated twelve strategies from three factors: content sources, topic interest models for users, and social popularity. Referring to the factor content sources, Chen et al. showed that profiling based on one's own tweets performed better than based on tweets by one's followees. Hence, we build up user profiles from social media items produced by the users themselves.

In the past years, profiling methods based on semantic concepts (i.e., ontology-based profiling) extraction have been developed [7, 16]. They extract semantic concepts from texts, using a structured knowledge base, e.g., DBpedia. Goossen et al. [7] proposed CF-IDF, as an extention of TF-IDF. CF-IDF counts frequencies of a concept instead of a term. Their news arcticle recommendation experiment with 19 participants demonstrated that CF-IDF outperforms TF-IDF. Lu et al. [16] proposed a recommender system for tweets based on what a user tweeted. They constructed user profiles represented as a set of weighted Wikipedia concepts that correspond to Wikipedia articles. The experiment demonstrated that concept-based approaches outperform TF-IDF. Other works employed a hierarchical structure of a knowledge base for profiling [12, 18, 16] and demonstrated their effectiveness. These approaches can reveal user interests that are not explicitly mentioned in the texts, using a structure of a knowledge base and spreading activation. In particular, Middleton et al. [18] constructed user profiles based on a hierarchical knowledge base using spreading activation for a recommender system of scientific publications. Their user experiment compared a profiling method using the structure of a hierarchical knowledge base and a method not using the structure. The result demonstrated superiority of using the hierarchical knowledge base. Topic modeling such as LDA [2] is one of the most popular profiling methods. It is used in the context of social media [10] but particularly suited for document profiling.

Time-aware user profiles are constructed based on the assumption that the degree of user interests declines as time passes. The decline of user interests is modeled by a decay function. In the past, the decay functions sliding window [24] and exponential decay [21, 26] have been employed for user profiling. However, they have not been compared so far like we do in this work.

## 3. PROBLEM DEFINITION

We address the problem of taking the social media stream as input in order to recommend items such as scientific publications the user might be interested in. The problem can be decomposed into three parts: (1) First, we need to extract the professional interests that a user exposes through his social media stream and represent the interests in a user

**Table 1: Symbol Notation**

| | |
|---|---|
| $u$ | a user |
| $i$ | a social media item |
| $I_u$ | the set of $u$'s social media items |
| $c$ | a concept |
| $C$ | the set of concepts |
| $d$ | a candidate item (scientific publication) |
| $D$ | the set of candidate items |
| $t_i, t_d$ | the time stamp of $i$ and $d$, respectively |
| $P_u$ | $u$'s user profile |
| $P_d$ | $d$'s document profile |
| $\Phi$ | a profiling function |
| $w'$ | a weighting function (not considering temporal decay) |
| $f$ | a decay function |
| $w$ | a weighting function that extends $w'$ with temporal decay |
| $\sigma$ | a similarity function |

**Table 2: Three factors and their choices for the experiment spanning in total** $3 \times 2 \times 2 = 12$ **strategies**

| Factor | Possible Design Choices | | |
|---|---|---|---|
| *Profiling Method* | CF-IDF | HCF-IDF | LDA |
| *Decay Function* | Sliding window | | Exponential decay |
| *Document Content* | All (title + full-text) | | Title |

items (i.e., documents whose similarity scores with $P_u$ are ranked in the top-$k$) are recommended to the user $u$. The similarity functions $\sigma$ are described in Section 4.3.

## 4. EXPERIMENTAL FACTORS

According to the three factors (i)-(iii) stated in the introduction, we form the design space of our experiment. We illustrate the design space in Table 2, where each cell is a possible design choice we can make in one of the three factors. Subsequently, we detail the factor *Profiling Method* in Section 4.1 and the factor *Decay Function* in Section 4.2. Further, we describe similarity functions $\sigma$ in Section 4.3. The factor *Document Content* investigates whether full-texts of scientific publications enhance the recommendation performance compared to using only titles.

### 4.1 Profiling Method

We investigate three methods for user profiling and document profiling. For each method, we define a weighting function $w'$ that gives a certain weight to each concept $c$. The final weighting function $w$ taking temporal decay into account is described in Section 4.2.

**CF-IDF:** Compared to the traditional TF-IDF, CF-IDF (Concept Frequency Inverse Document Frequency) counts frequencies of a semantic concept instead of term frequencies [7]. Semantic concepts or short concepts are stored in an external knowledge base. Each concept has a unique resource identifier (URI) and one or more labels describing the concept [2]. The concept's labels are treated as synonyms. As an example, the concept "clothing industry" has the URI `http://zbw.eu/stw/version/latest/descriptor/13128-2` and is defined in the thesaurus STW, a domain-specific knowledge base for economics (described in Section 5.3). The concept has not only the label "clothing industry" but also the synonymous labels "garment industry" and "apparel industry". We count the label frequency, i.e., the number of times the label appears, in the social media items and candidate items. Subsequently, we calculate the concept frequency, i.e., the number of times the concept appears, by summing up the frequencies of the labels referring to the concept. For instance, if the labels "clothing industry" and "garment industry" appear twice and once in a text, the total frequency of the concept referring to "clothing industry" is three.

For the social media items $I_u$ of the user $u$, CF-IDF is computed along with Equation 2.

$$w'_{cf\text{-}idf}(c, i) = cf(c, i) \cdot \log \frac{|I_u| + |I_r|}{|\{i \in I_u \cup I_r : c \in i\}|}, \quad (2)$$

where $cf(c, i) = \frac{\text{the number of times concept } c \text{ appears in } i}{\text{the number of times all concepts appear in } i}$. The denominator $|\{i \in I_u \cup I_r : c \in i\}|$ counts the number of social media items that contain a concept $c$. $I_r$ is a set of random social media items.

profile. (2) Likewise, we profile candidate items (i.e., scientific publications) and represent them in a way that they are comparable with the user profile. (3) We need a ranking function to compute the top-$k$ items based on similarity scores between the user profile and each candidate item. In the following, we formalize the three steps required to create a recommender system based on a user's professional interests extracted from the social media stream. Symbols used in this paper are summarized in Table 1.

**(1) User profiling from social media items.** We consider $I_u$ as set of social media items $i$ produced by user $u$. A social media item $i \in I_u$ has a certain time stamp $t_i$. Subsequently, $P_u$, the user profile of the user $u$, is created over a set of concepts $C$ by assigning a specific weight for each concept $c \in C$. Generally speaking, a concept $c$ is a key subject in a dedicated field, coming from a given domain-specific knowledge base $C$. For instance, "financial crisis" is a concept in the field of economics. We construct $P_u$ by employing different user profiling functions $\Phi$ and we compare them. Formally, user profiles are defined as:

$$P_u = \Phi(I_u, C) := \{(c, w(c, I_u)) \mid \forall c \in C\} \quad (1)$$

Here, $w$ is an arbitrary weighting function that returns a weight of a concept $c$ in a user's social media stream $I_u$. Thus, it determines how important a concept $c$ is for the user $u$. Profiling functions $\Phi$ and weighting functions $w$ are described in Sections 4.1 and 4.2. Specifically, we describe weighting functions $w'$ that do not consider temporal decay in Section 4.1 and provide weighting functions $w$ which extend $w'$ with temporal decay in Section 4.2.

**(2) Profiling candidate items.** We have a set of candidate items $D$. A candidate item $d \in D$ has a time stamp $t_d$, indicating its published year. To determine the similarity scores between a user profile $P_u$ and each candidate item $d \in D$, we need to construct profiles of candidate items in a way that they are comparable with the user profile. Formally, we represent a candidate item $d$ as a profile $P_d = \Phi(d, C) := \{(c, w(c, d)) \mid \forall c \in C\}$. Since our candidate items are scientific publications, we refer to this process document profiling.

**(3) Ranking candidate items.** We rank candidate items based on similarity scores between the user profile $P_u$ and a document profile $P_d$. A similarity function $\sigma$ takes as input a user profile $P_u$ and document profile $P_d$. It is defined as $\sigma(P_u, P_d) \rightarrow [0, 1]$. The similarity function is applied to all candidate items $d \in D$. Finally, the top-$k$ most relevant

---

[2]https://www.w3.org/DesignIssues/LinkedData.html

We employ a set of random social media items $I_r$, because it allows to better distinguish the relevant concepts in the user's social media items $I_u$, as Chen et al. [5] and Lu et al. [16] did for TF-IDF. For instance, assuming there are two social media items from a user $u$ and both include the concept "currency competition". Although "currency competition" should have a high weight in the user profile, in this case IDF and a final CF-IDF score would be 0 because "currency competition" is common in a user $u$'s social media items. The random social media items are sampled from public microblog postings. In our case, they are obtained from the public Twitter stream using the Twitter API.

We have conducted a simple pre-experiment to empirically determine the optimal amount of random tweets to be used in the profiling method in the context of our experiment of recommending economics publications. Given this pre-experiment, we set the size of random social media items to five times of $|I_u|$. In more detail, we applied different sizes of $I_r$, starting from 0 to 1000 random tweets. For 26 Twitter accounts, we computed the IDF scores for user profile over $I_u \cup I_r$ and compared it using cosine similarity with the user profile computed only over $I_u$. The Twitter accounts were taken from a list of famous economists[3] that are frequently tweeting. We ensured that the set of random tweets $I_r$ is disjoint do the user's tweets, i.e.. $I_r \cap I_u = \emptyset$. Particularly, we looked into the changes of the cosine similarity while adding more random tweets. We observed the changes in the IDF scores became stable after about a factor of five w.r.t. to $|I_u|$. The changes indicate the influence of the IDF scores to user profile. Using this technique is effective as the IDF score ensures that too generic concepts do not get too high weights in the user profiling. Those generic concepts are at the upper levels of the hierarchy of the domain-specific knowledge base. In our case those concepts are like "product" and "economics". Please note that the factor may depend on the domain of economics considered in this paper and that a different factor may be chosen for other domains.

Regarding document profiling, CF-IDF is computed as defined in Equation 3. The computation is basically identical with the one for user profiling shown in Equation 2. The difference is that CF is computed over single documents and IDF is computed over the document collection.

$$w'_{cf\text{-}idf}(c,d) = cf(c,d) \cdot \log \frac{|D|}{|\{d \in D \ : c \in d\}|} \qquad (3)$$

**HCF-IDF:** The novel profiling method HCF-IDF (Hierarchical CF-IDF) [19] extends CF-IDF by using a hierarchical knowledge base, where the concepts are hierarchically organized in a taxonomy. HCF-IDF can reveal concepts that are indirectly mentioned in texts by applying a spreading activation over the hierarchical knowledge base. Figure 1 shows an example where a user's profile includes the concept "social recommendation". Due to the hierarchical structure of the knowledge base, also the concepts "web searching" and "world wide web" are activated and obtain non-zero weights even if they are not mentioned. Different from the profiling methods using spreading activation [12, 18], HCF-IDF avoids to provide too high weights to generic concepts like "economy", as it employs IDF. Specifically, HCF-IDF combines the statistical strength of CF-IDF with semantics
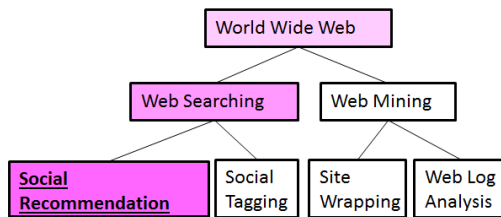
**Figure 1: An example of HCF-IDF**

provided by the hierarchical knowledge base. We compute HCF-IDF along with Equation 4.

$$w'_{hcf\text{-}idf}(c,i) = BL(c,i) \cdot \log \frac{|I_u| + |I_r|}{|\{i \in I_u \cup I_r : c \in i\}|} \qquad (4)$$

$BL(c,i)$ denotes the spreading activation function BellLog from Kapanipathi et al. [12]. It returns a weight of a concept $c$ in a social media item $i$ and is defined below:

$$BL(c,i) = cf(c,i) + FL(i) \cdot \sum_{c_j \in C_l(c)} BL(c_j,i), \qquad (5)$$

where $FL(c) = \frac{1}{\log_{10}(nodes(h(c)+1))}$. $h(c)$ returns the level where a concept $c$ is located in the knowledge base and $nodes$ provides the number of concepts at a given level in a knowledge base. For example, in Figure 1 $h$("web searching") returns 2 and $nodes(h("web searching") + 1)$ returns 4. $C_l(c)$ returns the set of concepts located in one level lower than the concept $c$. In Figure 1 the function $C_l$("world wide web") returns "web searching" and "web mining".

For scientific publications, weights are computed as defined in Equation 6. The computation is basically identical with the one for user profiling as shown in Equation 4. The difference is that $BL$ is applied over single documents and IDF is computed over the document collection.

$$w'_{hcf\text{-}idf}(c,d) = BL(c,d) \cdot \log \frac{|D|}{|d \in D : c \in d|} \qquad (6)$$

**LDA:** As third profiling method, we use LDA [2, 1], an unsupervised topic modeling method. LDA identifies latent topics in a document collection, where each document is represented as a probability distribution over topics, while each topic is again represented as a probability distribution over a number of words. Please note that for user profiling, we treat the set of social media items $I_u$ published by a user $u$ as one *single* social media document in this profiling method. It is known that topic models that treat a user's microblog postings as one combined social media document outperform topic models computed over single postings of a user for recommendation tasks [10]. We first create a topic model for the entire document collection $D$ (using the parameters and tools described in detail in Section 5.3). Subsequently, we run LDA with the given topic model for the document collection $D$ and infer a probability distribution over topics for the user's social media document $I_u$.

Again, we use the same notation of concepts $c$ as introduced above: Each topic generated by LDA is treated as a concept $c \in C$. The weight of a concept $c$ is defined by $w'_{lda}(c,I_u) = p(c \mid I_u)$ for user profiles and $w'_{lda}(c,d) = p(c \mid d)$ for document profiles, where $p(c \mid d)$ and $p(c \mid I_u)$ denote the probability of the concept (i .e., topic) $c$ in the social items $I_u$ and document $d$, respectively.

## 4.2 Decay Function

We compare two decay functions $f$, namely sliding window and exponential decay. In the past, both functions have been used in recommender systems [24, 21, 26]. However, so far they have not been empirically compared. The profiling functions $w'$ described in the previous section are combined with a decay function $f$ in order to obtain a final weight $w$. The final weights are computed by Equation 7 for the set of social media items and Equation 8 for the candidate items.

$$w(c, I_u) = \sum_{c \in i : i \in I_u} f(t_i) \cdot w'(c, i) \qquad (7)$$

$$w(c, d) = f(t_d) \cdot w'(c, d) \qquad (8)$$

Please note that when employing LDA, the decay functions can only be applied on the candidate items, because we treat the user's social media items as one single document.

**Sliding Window:** There are two kinds of sliding window functions, whose window size is defined by (a) the number of items [13] and (b) the period of time [25]. The approach (a) is employed to identify relatively short-term features (e.g., user interests from web browsing histories) [13], while the approach (b) is used to identify long-term features [25]. We aim at extracting a user's professional interests, which are rather long-term. Thus, we take the approach (b) and use only social media items and documents that are younger than a given threshold point in time $thresh$. The sliding window function can be represented as Equation 9.

$$f_{sw}(t) = \begin{cases} 1 & for \; t \geq thresh \\ 0 & for \; t < thresh \end{cases} \qquad (9)$$

For user profiles, we set the threshold based on the work by Orlandi et al. [21]. They found out that the half life time is $thresh_{social} = 250 \; days$. For document profiles, Sangam et al. [22] observed that the half-life time of the scientific publications in the field of social science is $9.04 \; years$. In our experiment, we use a dataset of scientific publications in economics (see Section 5.3), which has a large overlap with social science. Thus, we set $thresh_{doc} = 9.04 \; years$ [22] and remove scientific publications published more than $9.04 \; years$ ago from the candidate items.

**Exponential Decay:** The exponential decay function is defined as shown in Equation 10.

$$f_{exp}(t) = e^{-(t_{current}-t)/\tau}, \qquad (10)$$

where $t_{current}$ denotes the current time and $\tau$ is a positive number presenting mean-life [21]. For user profiles, we set $\tau = 360 \; days$ based on Orlandi [21]. Since Sangam et al. [22] found out that the mean-life of scientific publications in social sciences is $13.05 \; years$, we set $\tau = 13.05 \; years$ for document profiles.

## 4.3 Similarity Functions

We calculate the similarity scores between a user profile $P_u$ and each document profile $P_d$. We cast a user profile $P_u$ and document profiles $P_d$ to a user profile vector $\vec{p}_u$ and document profile vectors $\vec{p}_d$, respectively. Each element in the vectors corresponds to a weight of a concept $c$.

**Temporal Cosine Similarity:** We employ the temporal cosine similarity function described in Equation 11 for the profiling methods CF-IDF and HCF-IDF.

$$\sigma_{tcossim}(P_u, P_d) = f(t_d) \cdot \frac{\vec{p_u} \cdot \vec{p_d}}{||\vec{p_u}|| \cdot ||\vec{p_d}||}, \qquad (11)$$

It extends the cosine similarity by the function $f(t_d)$, which results in higher similarity score to newer documents. $f(t_d)$ is a decay function from Equation 9 or Equation 10. $t_d$ is time stamp of a scientific publication $d$. i.e., the year at which $d$ was published.

**Dot Product:** For LDA, we employ the dot product computed as $\sigma_{dp}(p_u, p_d) = \vec{p}_u \cdot \vec{p}_d$. Since LDA represents documents as probability distribution, it is more reasonable to use Kullback-Leibler divergence (KL divergence). However, the dot product outperforms cosine similarity and Kullback-Leibler divergence (KL divergence) when representing documents using LDA [9].

## 5. EVALUATION

We conducted an online experiment with $n = 123$ participants in order to identify the best strategy for a recommender system along the factors described in Section 4. As social media platform, we choose Twitter as it is widely used in scientific communities [14]. We design our experiment following the experiment setup and procedure of Chen et al. [5]: Each participant obtains top-5 recommendations for each of the twelve strategies formed from the three factors. The recommendation performance of each strategy is measured by the rankscore [4]. Below, we describe the details of our experiment procedure and participants. Subsequently, we explain the dataset and the knowledge base used in the experiment. Finally, we introduce our evaluation metric.

## 5.1 Procedure

The participants are invited to a web application implementing the twelve recommendation strategies. First, participants input their public Twitter handles and e-mail address. Then, the participants' tweets are retrieved from the Twitter API. Subsequently, user profiles are created from the tweets using each of the three profiling methods and two decay functions. Based on the user profiles, personalized top-$k$ recommendations of scientific publications are generated for each of the twelve strategies. We set the number of recommendations per strategy $k = 5$ along with Chen et al. [5]. After computing the recommendations, the participants receive an e-mail invitation to assess the recommendations. The participant go through all of the twelve strategies like as Chen et al. [5]. Thus, we apply a repeated measures design. Each participant obtains $12 \cdot 5 = 60$ recommendations in total throughout the experiment.

Prior to starting the experiment, participants are informed about the task of the experiment, i.e., rating the recommended publications based on relevance to their research interests, and confirmed consent. On each of the subsequent pages, the participants see a list of five recommendations produced by one of the twelve strategies. An example screenshot of the evaluation page is shown in Figure 2.

For each recommended scientific publication, the participants see its bibliographic information, i.e., authors, title, and year of publication. In addition, participants can look into the original PDF files by clicking on a link attached to

**Recommendation (1/12)**

Please evaluate the following randomized list of the top five publications "interesting" or "not interesting".
Clicking a title, you can see the content of a publication.

- Szulc, Elzbieta, "Modelling of the Dependence Between the Space-time Processes," 2008 — ○ interesting ○ not interesting
- Ichino, Andrea; Schwerdt, Guido; Winter-Ebmer, Rudolf; Zweimüller, Josef, "Too old to work, too young to retire?," 2007 — ○ interesting ○ not interesting
- Rodríguez-Pose, Andrés, "Economic geographers and the limelight : the reaction to the 2009 world development report," 2010 — ○ interesting ○ not interesting
- Stöllinger, Roman, "International spillovers in a world of technology clubs," 2012 — ○ interesting ○ not interesting
- den Berg, Gerard J. van; Vikström, Johan, "Monitoring job offer decisions, punishments, exit to work, and job quality," 2009 — ○ interesting ○ not interesting
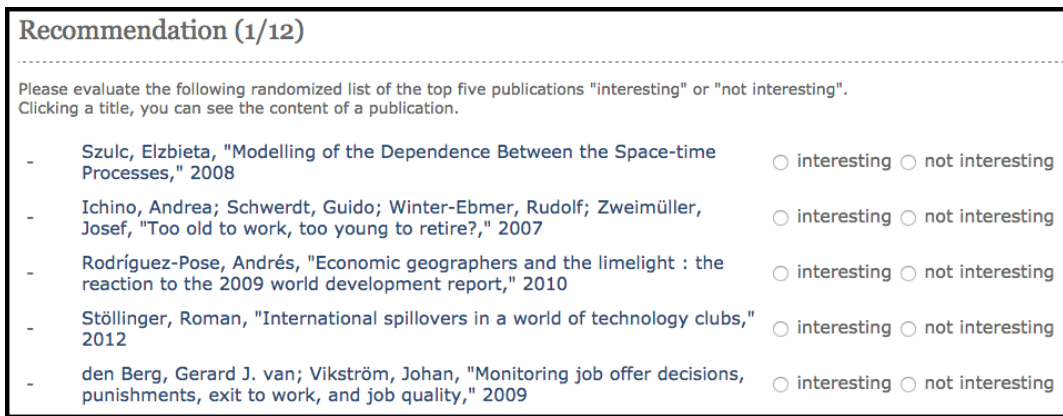
**Figure 2: Screenshot of our experiment web page showing a randomized list of top-5 recommendations for the first of twelve strategies (which again are randomly ordered). For each recommendation the participants could assess the bibliographic record as well as click on the full-text document. The participants rated each recommended publication as "interesting" or "not interesting"' based on their research interests.**

the bibliographic record. In order to avoid bias, the participants go through the twelve strategies in random order. For each strategy, the participants receive one list of five recommendations. The five recommendations in the lists are again shown in random order to the participants to avoid the well-known ranking bias. Typically, participants assume that top-ranked recommendations are essentially more relevant [3, 5]. Thus, again prior to starting the experiment we have explicitly informed the participants that we have randomized the order of the items in the top-5 lists. However, the actual ranks of the recommendations as well as their positions where the recommended items appeared on the participants' screen are stored in the database for later analyses. Participants evaluate each recommendation as "interesting" or "not interesting" by clicking on radio buttons next to the publication records like Chen et al. [5]. Please note, the participants had to evaluate all recommended items.

At the end of the experiment, we collect the demographic information of each participant, including gender, age, highest academic degree, major, years of profession, and current employment status (academia/industry). Finally participants could state free comments regarding the experiment.

## 5.2 Participants

We recruited $n = 123$ participants through mailing lists, tweets, and word-of-mouth on the Internet. Initially 160 participants registered their Twitter handles and email address for our experiment. Among them, 134 participants started the experiment after receiving the e-mail invitation. From these 134 participants, only eleven dropped in the course of assessing the recommendations in the twelve strategies. Thus, finally we obtain evaluations for all strategies from $n = 123$ participants. From these, 27 participants are female. The average age of the participants is 32.83 years (SD: 7.34). Regarding the highest academic degree, we have acquired 21 with a Bachelor, 58 have a Master, 32 a PhD, and 12 are lecturers/professors. While 83 participants work in academia, 40 work in industry. Tweets of the participants were retrieved via Twitter API. We only collected tweets in English as the scientific publications are also in English. The participants published on average 1096.82 En-

glish tweets (SD: 1048.46). The maximum and minimum numbers of tweets are 3192 and 2, respectively. Twitter users who have not produced any tweets in the last 250 *days* could not register and participate in the experiment, since we use a 250 *days* threshold for the decay function Sliding window (see Section 4.2). Five Twitter users could not participate in the experiment for this reason.

The participants spent on average 517.54 seconds to complete the assessment of the $5 \times 12 = 60$ recommendations (SD: 376.72). This does not include the time spent to register for the experiment, read the instructions, and filling out the final questionnaire. As incentive, each participant received the information about his most similar economist among 26 famous economists[4] and the top-5 dominant semantic concepts in their tweets after the experiment. In addition, the participants could opt-in to a raffle for one of two Amazon vouchers worth of 50 €.

## 5.3 Dataset Preparation

We use a large-scale dataset of scientific publications in the field of economics as candidate items and a high-quality taxonomy as a knowledge base for profiling methods.

**Dataset of Scientific Publications.** We collaborate with the providers of EconBiz[5], a portal for scientific publications in economics managed by ZBW, the German National Library of Economics. From this portal, we obtained 1 million URLs of open access publications and extracted full-texts and metadata (i.e., authors, title, year of publication) of 413,098 scientific publications. Finally, we determined the document language[6] and got 279,381 scientific publications in English, which were used in this experiment.

**Knowledge Base in Economics.** The ZBW also maintains and further develops the hierarchical knowledge base STW[7], a thesaurus specialized for the field of economics. The STW is freely available and is of high quality due to its manual maintenance by domain experts. The knowledge

---

[4]http://www.huffingtonpost.com/2012/11/13/economists-twitter_n_2122781.html
[5]http://www.econbiz.de/
[6]https://code.google.com/p/language-detection/
[7]http://zbw.eu/stw/version/8.12/about.en.html

base is poly-hierarchically organized with six levels. It contains $6,335$ semantic concepts and $11,679$ labels. The hierarchically organized concepts are connected with each other via $14,875$ edges. In order to extract as many labels as possible, we enhanced the original STW with DBpedia redirects[8]. From DBpedia redirects we can retrieve the synonymous labels for a concept. STW contains $2,692$ concepts that have both a DBpedia mapping and one or more DBpedia redirects. As an example, for the concept "Telecommunications industry" in the thesaurus, we obtain the DBpedia redirects "Telecommunications operator" and "Telephone companies" and use them as synonymous labels referring to the concept "Telecommunications industry". Finally, our extended STW contains $6,335$ concepts and $37,733$ labels. This extended STW is used for the profiling methods CF-IDF and HCF-IDF. For CF-IDF, we ignore the edges between concepts.

**Processing of the tweets and publications.** For the profiling methods CF-IDF and HCF-IDF, we extract semantic concepts from the participants' tweets and the scientific publications by matching the texts with the labels from the extended STW (i.e., a gazetteer-based approach). Before processing, we lemmatize both the tweets and the scientific publications using Stanford Core NLP[9] and remove stop words. Regarding the tweets, some of them contain hashtags indicating topics (e.g., #election) and user mentions (e.g., @UNICEF). We remove only the symbols # and @ from the tweets as Feng et al. [6] observed that the combination of the tweets' textual content with the hashtags and user mentions made the highest performance for tag recommendation.

This process extracts only the users' professional interests from tweets and helps to avoid noise (i.e., topics not relevant to professional interests in economics). A participant has published on average 1096.82 tweets (SD: 1048.46). On average $1,214.93$ concepts (SD: 1181.43) are contained in a participant's tweets and 1.07 concepts (SD: 0.31) are contained per tweet. Regarding CF-IDF and HCF-IDF, we calculate the ratio of the number of tweets containing at least one concept and the total number of tweets the user has published. This indicates the percentage of tweets that have contributed to creating the user profile. On average, 62.24% of the tweets (SD: 13.55) that a participant has published contain at least one concept in economics. These tweets are assumed to be relevant to the professional interests.

**LDA.** For constructing profiles by LDA, we use JGibbLDA[10]. We first run LDA to generate the topic model based on the given document set $D$. Following Blei et al. [1], we lemmatize the scientific publications using Stanford NLP Core. Subsequently, we remove stop words and words that appear in fewer than 25 scientific publications. We optimized the number of topics $K$ regarding the maximum mean log likelihood of words given topics as suggested by Griffiths et al. [8]. We experimented with $K = 20, 50, 100, 200, 500, 1000,$ and $5000$ and obtained the highest log likelihood for $K = 100$. All topic models were computed over 500 iterations. Regarding the further parameters for LDA, we set $\alpha = 0.5$ and $\beta = 0.1$ as suggested by Griffiths et al. [8]. To infer a topic distribution over a user's tweets, we run LDA again using the topic model for the document set $D$ with 200

---

[8] http://oldwiki.dbpedia.org/Downloads39#redirects
[9] http://nlp.stanford.edu/software/corenlp.shtml
[10] http://jgibblda.sourceforge.net/

---

iterations. Prior to this, we prepare the tweets of a user $u$ in a single social media document as described in Section 4.1.

## 5.4 Evaluation Metric

In order to assess the recommendation performance, we compute the rankscore [4] as used by Bostandjiev et al. [3] and introduced by Jannach et al [11]. Rankscore posits that each successive item in a list is less likely to be viewed by users with an exponential decay, as defined in Equation 12.

$$rankscore' = \sum_{d \in hits} \frac{1}{2^{\frac{rank_d - 1}{\theta - 1}}} \qquad (12)$$

$\theta$ denotes a viewing halflife parameter controlling the speed of the exponential decay. As suggested by Breese et al. [4], we set $\theta = 5$. $hits$ refers to the set of documents $d$ evaluated as "interesting" and $rank_d$ denotes the rank of a recommended item $d$ in a list. Please note $rank_d$ denotes the actual rank stored in the database different from the position where a item $d$ appears in the list (cf. Section 5.1). The normalized rankscore is computed by $rankscore = rankscore'$ $/rankscore_{max}$, where the maximum rankscore $rankscore_{max} = \sum_{j=1}^{k} \frac{1}{2^{\frac{j-1}{\theta - 1}}}$. Here, $k$ is the number of the recommended items. We set $k = 5$. We also computed Mean Average Precision (MAP), Precision@5, Mean Reciprocal Rank (MRR), and normalized Discounted Cumulative Gain (nDCG). Overall, the results are similar to the rankscore and thus omitted for reasons of brevity. The interested reader may refer to the details in the appendix [20].

## 6. RESULTS

In this section, we document the results of the experiment[11] and conduct the statistical analyses. We set a significance level of $\alpha = 5\%$ for all statistical tests (please do not confuse with $\alpha$ for LDA in Section 5.3).

### 6.1 Quantitative Analyses

We first report the best performing strategy among the twelve strategies. Subsequently, we analyze the influence by the experimental factors followed by investigating the correlations between the recommendation performance and the numbers of tweets written by a user. Finally, we analyze the performance related to the number of times the participants clicked on the full-text of a publication.

**Best performing strategy.** Table 3 documents the average rankscores of the twelve strategies sorted in decreasing order. Overall, the best performing strategy is the strategy CF-IDF $\times$ Sliding window $\times$ All. We apply a one-way repeated-measure ANOVA in order to identify if there are significant differences between the strategies. For using ANOVA, we first need to verify whether the variances of the rankscores of the twelve strategies are equal. This is done by using Mauchly's test, which reveals a violation of sphericity in the strategies ($\chi^2(65) = 435.90$, $p = .00$). It may lead to positively biased F-statistics and increases the risk of false positives. To reduce this risk, we apply a Greenhouse-Geisser correction of $\epsilon = .61$ and run the one-way repeated-measure ANOVA. It reveals a significant difference in the rankscores of the strategies ($F(6.60, 805.33) = 21.98$, $p = .00$). To assess the pair-wise significant differences between

---

[11] The anonymized experimental data is available from: http://dx.doi.org/10.7802/1224

**Table 4: Post-hoc analysis with pairwise p-values over the twelve strategies using Shaffer's MSRB procedure. The p-values are marked in bold font if $p < .05$, which indicates a significant difference between the two strategies. Strategies are sorted by rankscores as shown in Table 3.**

| | | | | All / Sliding window / HCF-IDF 2. | Title / Sliding window / HCF-IDF 3. | Title / Exponential decay / HCF-IDF 4. | All / Exponential decay / CF-IDF 5. | All / Exponential decay / HCF-IDF 6. | Title / Exponential decay / CF-IDF 7. | Title / Sliding window / CF-IDF 8. | Title / Exponential decay / LDA 9. | Title / Sliding window / LDA 10. | All / Exponential decay / LDA 11. | All / Sliding window / LDA 12. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. | CF-IDF | Sliding window | All | .99 | .97 | .72 | .22 | .12 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 2. | HCF-IDF | Sliding window | All | | .99 | .99 | .99 | .99 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 3. | HCF-IDF | Sliding window | Title | | | .99 | .99 | .99 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 4. | HCF-IDF | Exponential decay | Title | | | | .99 | .99 | **.01** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 5. | CF-IDF | Exponential decay | All | | | | | .99 | **.04** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 6. | HCF-IDF | Exponential decay | All | | | | | | .12 | **.02** | **.00** | **.00** | **.00** | **.00** |
| 7. | CF-IDF | Exponential decay | Title | | | | | | | .99 | .99 | .41 | .28 | **.01** |
| 8. | CF-IDF | Sliding window | Title | | | | | | | | .99 | .84 | .61 | **.03** |
| 9. | LDA | Exponential decay | Title | | | | | | | | | .99 | .99 | .72 |
| 10. | LDA | Sliding window | Title | | | | | | | | | | .99 | .99 |
| 11. | LDA | Exponential decay | All | | | | | | | | | | | .88 |

the twelve strategies, a post-hoc analysis is conducted. We have applied Shaffer's modified sequentially rejective Bonferroni procedure (Shaffer's MSRB procedure) [23] that takes into account the number of different experiment conditions, i.e., the number of recommendation strategies. The result of the post-hoc analysis is presented in Table 4. The vertical and horizontal dimensions of the Table 4 show the eleven-by-eleven comparison of the twelve strategies. As one can see, we observe various significant differences between the strategies ($p < .05$, marked in bold font). For example, while we observe a significant difference between the strategies CF-IDF × Sliding window × Title and HCF-IDF × Sliding window × All ($t(122) = 4.77$, $p = .00$), there is no significant difference between the strategies CF-IDF × Exponential decay × Title and LDA × Sliding window × Title ($t(122) = 2.43$, n.s., $p = .41$).

**Table 3: Rankscores of the strategies in decreasing order. M and SD denote mean and standard deviation, respectively.**

| | Strategy | | | Rankscore |
|---|---|---|---|---|
| | **Profiling Method** | **Decay Function** | **Content** | **M (SD)** |
| 1. | CF-IDF | Sliding window | All | .59 (.33) |
| 2. | HCF-IDF | Sliding window | All | .56 (.34) |
| 3. | HCF-IDF | Sliding window | Title | .55 (.33) |
| 4. | HCF-IDF | Exponential decay | Title | .52 (.30) |
| 5. | CF-IDF | Exponential decay | All | .51 (.32) |
| 6. | HCF-IDF | Exponential decay | All | .49 (.30) |
| 7. | CF-IDF | Exponential decay | Title | .41 (.29) |
| 8. | CF-IDF | Sliding window | Title | .39 (.27) |
| 9. | LDA | Exponential decay | Title | .35 (.31) |
| 10. | LDA | Sliding window | Title | .33 (.31) |
| 11. | LDA | Exponential decay | All | .32 (.30) |
| 12. | LDA | Sliding window | All | .27 (.33) |

**Difference in experiment factors.** Subsequently, we analyze the results with respect to each experimental fac-

**Table 5: Three-way repeated-measure ANOVA with Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* | 58.40 | .48 | **.00** |
| *Decay Function* | 1.17 | .01 | .28 |
| *Document Content* | 5.18 | .04 | **.02** |
| *Profiling Method × Decay Function* | 4.63 | .04 | **.01** |
| *Profiling Method × Document Content* | 17.09 | .14 | **.00** |
| *Decay Function × Document Content* | 4.69 | .04 | **.03** |
| *Profiling Method × Decay Function × Document Content* | 3.35 | .03 | **.04** |

tor. To this end, we first apply Mendoza's test [17] to check for violations of sphericity against the factors. Mendoza's test is an extension of Mauchly's test to adopt to multi-way repeated-measure ANOVA. It shows significances with the global ($\chi^2(65) = 435.90$, $p = .00$) and the factors *Profiling Method* ($\chi^2(2) = 12.21$, $p = .00$), *Profiling Method × Decay Function* ($\chi^2(2) = 20.02$, $p = .00$), and *Profiling Method × Document Content* ($\chi^2(2) = 8.61$, $p = .01$). Subsequently, we run a three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$ for the global and $\epsilon = .91$ for the factors *Profiling Method*, $\epsilon = .87$ for *Profiling Method × Decay Function*, and $\epsilon = .93$ for *Profiling Method × Document Content*. Table 5 shows the results of the ANOVA with F-ratio, effect size $\eta^2$, and p-value. The effect size is small when $\eta^2 > .02$, medium when $\eta^2 > .13$, and large when $\eta^2 > .26$. The analyses reveal significant differences in all three factors and their contributions except the factor *Decay Function*. For all factors with significant differences, we apply again a post-hoc analysis using Shaffer's MSRB procedure with respect to each factor. In terms of the factor *Profiling Method*, the post-hoc analysis reveals significant differences between all pairs of HCF-IDF, CF-IDF, and LDA (details of the post-hoc analysis are omitted

for the reasons of brevity and documented in our TR [20]). Although the strategy CF-IDF $\times$ Sliding window $\times$ All performs best as shown in Table 3, the best *Profiling Method* is HCF-IDF as it performs under all other factors better than CF-IDF and LDA. Regarding the factor *Document Content*, "All" outperforms "Title" ($F(1, 122) = 5.18$, $p = .02$). Regarding the factor *Profiling Method $\times$ Decay Function*, the result suggests that the strategies with the Exponential decay function perform better than those with the Sliding window function when LDA is employed. In addition, there are significant differences among the three profiling methods when a decay function is fixed. In both decay functions, HCF-IDF performs best, followed by CF-IDF, and LDA. Referring to the factor *Profiling Method $\times$ Document Content*, the result indicates that All is a better choice than Title, when CF-IDF is employed. In profiling methods HCF-IDF and LDA, the factor *Document Content* makes no significant difference. It indicates that HCF-IDF does perform well when only titles of candidate items are available. In addition there are significant differences among the profiling methods when a choice of *Document Content* is fixed. In those cases, HCF-IDF always outperforms others. In terms of the factor *Decay Function $\times$ Document Content*, All is a better choice than Title, when Sliding window is used.

**Correlation of recommendation performance with the number of tweets, the number of concepts, the number of concepts per tweet, and the percentage of tweets containing at least one concept.** We computed Pearson's $r$ and Kendall's $\tau$ between the users' mean rankscores and each of the number of tweets, concepts, concepts per tweet and the percentage of tweets containing at least one concept. A correlation may show a dependency that could influence the recommendation performance. The results show no significant correlation: As stated in Section 5.3, a participant has published on average 1096.82 tweets (SD: 1048.46). There is no significant correlation with the rankscores ($r(121) = .04$, n.s., $p = .62$ and $\tau = .00$, n.s., $p = .98$). Referring to the number of concepts, on average $1,214.93$ concepts (SD: 1181.43) are contained in a participant's Twitter stream. The correlation coefficients are non-significant ($r(121) = .05$, n.s., $p = .60$ and $\tau = -.01$, n.s., $p = .94$). Regarding the number of concepts per tweet, a participant's tweet contains on average 1.07 concepts (SD: 0.31) with again no significant correlation to the rankscores ($r(121) = -.05$, n.s., $p = .59$ and $\tau = -.02$, n.s., $p = .71$). Regarding the tweets that contribute in computing the user profiles for the methods with CF-IDF and HCF-IDF, we calculate the percentage of the number of tweets containing at least one concept and the number of tweets for each user. On average, 62.24% of the tweets (SD: 13.55) that a participant has published contain at least one concept, with no significant correlation ($r(121) = -.04$, n.s., $p = .67$ and $\tau = -.03$, n.s., $p = .73$)

## 6.2 Questionnaire Feedback

At the end of the experiment, the participants were asked to rate: "How easy it was to make the decisions whether a recommended publication is interesting". Using a 5-point Likert scale, where values between 1 and 5 refer to very difficult to very easy, the result is fairly high with an average of 3.68 (SD: 0.88). Regarding question "Whether the participants noticed a difference among the twelve strategies", the result is similarly high with an average of 3.46 (SD: 1.20).

In the free text feedback, one participant denoted that the recommender system failed to pick up his primary field despite having tweeted about that field. Apart from this, we received many positive comments (e.g., interesting, useful).

## 7. DISCUSSION

The strategies with HCF-IDF perform almost equally well compared to the best performing strategy CF-IDF $\times$ Sliding window $\times$ All. There is no significant difference between them as described in Table 4. The strong advantage of HCF-IDF is that it reaches its performance already when using only the titles of the scientific publications. The reason is that spreading activation over the hierarchical knowledge base used in HCF-IDF successfully reveals concepts that are not explicitly mentioned in the texts. CF-IDF works well when full-texts are available. Referring to LDA, the recommendation performance of the strategies with LDA is overall low, even if full-texts are available. A possible reason is that LDA cannot construct accurate user profiles because of the shortness and sparseness of social media items. Without accurate user profiles it is impossible to make good recommendations, even if full-texts are available. In fact, a slight correlation between the rankscores of LDA and the number of tweets is observed [20]. It indicates that participants with more tweets receive better recommendations. Please note as documented in [20], rankscores are almost exact same values with Precision@5 and nDCG. Although rankscores are slightly different with MAP and MRR, the order of performance of strategies are almost identical. Thus, the arguments described in this paper do not be influenced by differences among those evaluation metrics.

Our dataset covers scientific publications in the broader field of economics. Thus, although the dataset is obtained from a portal of economics literature, it contains scientific publications from various fields including, e.g., social sciences, political sciences, and information sciences. In the experiment, 31 of 123 participants do not have a major in economics. We have conducted an ANOVA test to identify whether the recommendation performance is significantly different for participants from economics and those not in economics. The result shows that majors make no significant difference ($F(1, 121) = 0.01$, n.s., $p = .94$). Thus, we assume that our approach may be transferred to other domains. Furthermore, there are a lot of domain-specific hierarchical knowledge bases in other domains freely available such as Medical Subject Headings (MeSH) for medicine and ACM Computing Classification System (ACM CCS) for computer science. An overview of freely available hierarchical knowledge bases is maintained by the W3C as cited in the introduction. The knowledge bases are of similar structure to the STW used in this paper. They are of high quality as they are manually crafted by domain experts. Therefore, HCF-IDF can be easily applied to other fields. Our approach could be integrated with other social media platforms (e.g., Facebook, LinkedIn), where users generate short and sparse texts. In addition, HCF-IDF is robust against the number of tweets a user published, because there is no correlation between the number of tweets and the rankscores of the strategies with HCF-IDF.

Our results may potentially be influenced by the amount of time that each participant spent for evaluating the $5 \times 12 = 60$ recommended publications by the twelve strategies in the experiment. However, they spent on average 517.54

seconds (SD: 376.72) to complete the evaluation of the 60 recommendations. In addition, we randomized the order of the strategies presented to the participants to counterbalance any influence on the order of the strategies. Thus, we think that our results are not influenced by it. Another potential threat to the validity of our results could be the procedure how we recruited the participants. We believe that the risk is low since we collected enough participants regarding each demographic factor (as shown in Section 5.2). Regarding the demographic factors, we found significant differences only for the participants' highest academic degree and participants' gender (details are documented in the TR [20]). However, they do not affect the order of the recommendation performance of the different strategies.

# 8. CONCLUSIONS

This paper contributes to content-based recommender systems for scientific publications based on user profiles extracted from social media platforms. We have constructed twelve different recommendation strategies along three factors, namely profiling method, decay function, and document content. The online experiment revealed that titles of scientific publications are sufficient to achieve competitive recommendation results when employing the profiling method HCF-IDF. Thus, the spreading activation over the hierarchical knowledge base enables HCF-IDF to extract a sufficient number of concepts from titles to compute competitive recommendations. This is an important result as full-texts are not always available, e. g., due to legal reasons.

# 9. REFERENCES

[1] D. M. Blei and J. D. Lafferty. Dynamic topic models. In *ICML*. ACM, 2006.

[2] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *JMLR*, 3, 2003.

[3] S. Bostandjiev, J. O'Donovan, and T. Höllerer. Taste-Weights: a visual interactive hybrid recommender system. In *RecSys*. ACM, 2012.

[4] J. S. Breese, D. Heckerman, and C. Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In *UAI*. Morgan Kaufmann, 1998.

[5] J. Chen, R. Nairn, L. Nelson, M. Bernstein, and E. Chi. Short and tweet: experiments on recommending content from information streams. In *CHI*. ACM, 2010.

[6] W. Feng and J. Wang. We can learn your# hashtags: Connecting tweets to explicit topics. In *ICDE*. IEEE, 2014.

[7] F. Goossen, W. IJntema, F. Frasincar, F. Hogenboom, and U. Kaymak. News personalization using the CF--IDF semantic recommender. In *WIMS*. ACM, 2011.

[8] T. L. Griffiths and M. Steyvers. Finding scientific topics. *NAS*, 101, 2004.

[9] T. J. Hazen. Direct and latent modeling techniques for computing spoken document similarity. In *the Spoken Language Technology*. IEEE, 2010.

[10] L. Hong and B. D. Davison. Empirical study of topic modeling in Twitter. In *SOMA*. ACM, 2010.

[11] D. Jannach, M. Zanker, A. Felfernig, and G. Friedrich. *Recommender systems: an introduction*. Cambridge University Press, 2010.

[12] P. Kapanipathi, P. Jain, C. Venkataramani, and A. Sheth. User interests identification on Twitter using a hierarchical knowledge base. In *ESWC*. Springer, 2014.

[13] M. K. Khribi, M. Jemni, and O. Nasraoui. Automatic recommendations for e-learning personalization based on web usage mining techniques and information retrieval. In *ICALT*. IEEE, 2008.

[14] J. Letierce, A. Passant, J. G. Breslin, and S. Decker. Understanding how twitter is used to spread scientific messages. In *WebSci*. Web Science Trust, 2010.

[15] Y. Li, M. Yang, and Z. M. Zhang. Scientific articles recommendation. In *CIKM*. ACM, 2013.

[16] C. Lu, W. Lam, and Y. Zhang. Twitter user modeling and tweets recommendation based on Wikipedia concept graph. In *AAAI Workshops*, 2012.

[17] J. L. Mendoza. A significance test for multisample sphericity. *Psychometrika*, 45(4), 1980.

[18] S. E. Middleton, D. C. De Roure, and N. R. Shadbolt. Capturing knowledge of user preferences: ontologies in recommender systems. In *K-CAP*. ACM, 2001.

[19] C. Nishioka, G. Große-Bölting, and A. Scherp. Influence of time on user profiling and recommending researchers in social media. In *i-KNOW*. ACM, 2015.

[20] C. Nishioka and A. Scherp. Profiling vs. time vs. content: What does matter for top-k publication recommendation based on twitter profiles? - an extended technical report. http://arxiv.org/abs/1603.07016.

[21] F. Orlandi, J. Breslin, and A. Passant. Aggregated, interoperable and multi-domain user profiles for the social web. In *I-SEMANTICS*. ACM, 2012.

[22] S. L. Sangam and S. S. Mogali. Obsolescence of literature in the field of social sciences. *PEARL*, 7(3), 2013.

[23] J. P. Shaffer. Modified sequentially rejective multiple test procedures. *J. of the ASA*, 81(395), 1986.

[24] W. Shen, J. Wang, P. Luo, and M. Wang. Linking named entities in tweets with knowledge base via user interest modeling. In *KDD*. ACM, 2013.

[25] S. J. Soltysiak and I. B. Crabtree. Automatic learning of user profiles - towards the personalisation of agent services. *BT Tech. J.*, 16(3), 1998.

[26] K. Sugiyama and M.-Y. Kan. Scholarly paper recommendation via user's recent research interests. In *JCDL*. ACM, 2010.

[27] K. Sugiyama and M.-Y. Kan. Exploiting potential citation papers in scholarly paper recommendation. In *JCDL*, pages 153–162. ACM, 2013.

[28] C. Wang and D. M. Blei. Collaborative topic modeling for recommending scientific articles. In *KDD*. ACM, 2011.

[29] Z. Zhao, Z. Cheng, L. Hong, and E. H. Chi. Improving user topic interest profiles by behavior factorization. In *WWW*. IW3C2, 2015.

# APPENDIX

## A. THREE-WAY REPEATED ANOVA FOR RANKSCORES

We describe the details of the three-way repeated ANOVA for rankscores conducted in Section 6.1. Specifically, we provide the analyses with respect to each factor that shows a significant difference as documented in Table 5. For all statistical analyses, we use $\alpha = .05$ as significance level. Regarding the effect size $\eta^2$, the effect size is small when $\eta^2 > .02$, medium when $\eta^2 > .13$, and large when $\eta^2 > .26$. In terms of the effect size $d$ measured by Cohen's $d$, the effect size is interpreted small when $d = .20$, medium when $d = .50$, and large when $d = .80$.

**The factor _Profiling Method_.** Tables 6(a), (b) and (c) show the rankscores, the post-hoc analysis for the factor _Profiling Method_, and the effect size, respectively. Table 6(a) presents the means and standard deviations of the three profiling methods. Table 6(b) shows p-values of each pair. Since Table 5 shows that this factor has the largest effect size, we further compute the effect size using Cohen's d for each pair shown in Table 6(c).

**Table 6: Rankscores, Post-hoc analysis for the factor _Profiling Method_ using Shaffer's MSRB procedure, and effect size.**

a) Rankscores

| Choice | M | SD |
|--------|------|------|
| HCF-IDF | .53 | .32 |
| CF-IDF | .48 | .31 |
| LDA | .32 | .31 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|--------|---------|------|
| CF-IDF | **.00** | **.00** |
| HCF-IDF | | **.00** |

c) Effect size using Cohen's d

| | HCF-IDF | LDA |
|--------|---------|------|
| CF-IDF | .17 | .50 |
| HCF-IDF | | .67 |

**The factor _Document Content_.** Table 7 shows the post-hoc analysis for the factor _Document Content_. The result shows that the recommender systems perform better when using both titles and full-texts.

**Table 7: Rankscores and Post-hoc analysis for the factor _Document Content_ using Shaffer's MSRB procedure.**

a) Rankscores

| Choice | M | SD |
|--------|------|------|
| All | .46 | .21 |
| Title | .43 | .20 |

b) Post-hoc analysis p-values

| | Title |
|-----|-------|
| All | **.02** |

**The factor _Profiling Method_ × _Decay Function_.** Table 8 shows the results of ANOVA regarding the factor _Profiling Method_ when a choice of the factor _Decay Function_ is

fixed and vice versa. Mendoza's test found a violation of sphericity in the factor _Profiling Method_ when Sliding window is used ($\chi^2(2) = 9.26$, $p = .01$) and Exponential decay is used ($\chi^2(2) = 11.16$, $p = .00$). Thus, we run a one-way repeated-measure ANOVA with Greenhouse-Geisser correction of $\eta = .93$ for the first row in Table 12 and $\eta = .92$ for the second row in Table 12. We observe significant differences when a choice of the factor _Decay Function_ is fixed and when LDA is employed. The post-hoc analyses of them are shown in Table 9, Table 10, and Table 11, respectively. In Table 9 and Table 10, a choice of the factor _Decay Function_ is fixed. The results demonstrate that HCF-IDF performs best, followed by CF-IDF and LDA. Table 11 shows the post-hoc analysis of the factor _Decay Function_ when LDA is employed. It indicates Exponential decay performs better than Sliding window for LDA.

**Table 8: ANOVA for _Profiling Method_ × _Decay Function_ interaction**

| Factor | F | $\eta^2$ | p |
|--------|------|------|------|
| _Profiling Method_ at Sliding window | 52.71 | .43 | **.00** |
| _Profiling Method_ at Exponential decay | 26.89 | .22 | **.00** |
| _Decay Function_ at CF-IDF | 3.69 | .03 | .06 |
| _Decay Function_ at HCF-IDF | 2.33 | .02 | .12 |
| _Decay Function_ at LDA | 5.26 | .04 | **.02** |

**Table 9: Rankscores and Post-hoc analysis for the factor _Profiling Method_ at Sliding window using Shaffer's MSRB procedure.**

a) Rankscores

| Choice | M | SD |
|--------|------|------|
| HCF-IDF | .55 | .33 |
| CF-IDF | .49 | .32 |
| LDA | .30 | .32 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|--------|---------|------|
| CF-IDF | **.01** | **.00** |
| HCF-IDF | | **.00** |

**Table 10: Rankscores and Post-hoc analysis for the factor _Profiling Method_ at Exponential decay using Shaffer's MSRB procedure.**

a) Rankscores

| Choice | M | SD |
|--------|------|------|
| HCF-IDF | .51 | .30 |
| CF-IDF | .46 | .31 |
| LDA | .34 | .31 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|--------|---------|------|
| CF-IDF | **.02** | **.00** |
| HCF-IDF | | **.00** |

**The factor _Profiling Method_ × _Document Content_.** Table 12 shows the results of ANOVA regarding the factor _Profiling Method_ when a choice of the factor _Document Content_ is fixed and vice versa. We observe there is a significant difference when a choice of the factor _Document Content_ is fixed and CF-IDF is employed. Mendoza's test found a violation of sphericity in the factor _Profiling Method_ when All

**Table 11: Rankscores and Post-hoc analysis for the factor *Decay Function* at LDA using Shaffer's MSRB procedure.**

**a) Rankscores**

| Choice | M | SD |
|---|---|---|
| Exponential decay | .34 | .31 |
| Sliding window | .30 | .32 |

**b) Post-hoc analysis p-value**

| | Exponential decay |
|---|---|
| Sliding window | **.02** |

(i.e., titles and full-texts) is used ($\chi^2(2) = 25.24$, $p = .00$). Thus, we run a one-way repeated-measure ANOVA with Greenhouse-Geisser correction of $\eta = .84$ for the second row in Table 12. Table 13 presents the post-hoc analysis when Title is selected for the factor *Document Profiling*. We see that HCF-IDF outperforms others with significant differences. On the other hand, Table 14 shows the post-hoc analysis when All is chosen for the factor *Document Profiling*. There is no significant difference between CF-IDF and HCF-IDF. Table 15 shows the post-hoc analysis of the factor *Document Content* when CF-IDF is employed. It indicates that the strategies with CF-IDF and All significantly outperforms those with CF-IDF and Title.

**Table 12: ANOVA for *Profiling Method* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* at Title | 26.15 | .21 | **.00** |
| *Profiling Method* at All | 55.28 | .45 | **.00** |
| *Document Content* at CF-IDF | 32.95 | .27 | **.00** |
| *Document Content* at HCF-IDF | 0.43 | .00 | .51 |
| *Document Content* at LDA | 2.06 | .02 | .15 |

**Table 13: Rankscores and Post-hoc analysis for the factor *Profiling Method* at Title using Shaffer's MSRB procedure.**

**a) Rankscores**

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .54 | .31 |
| CF-IDF | .40 | .28 |
| LDA | .34 | .31 |

**b) Post-hoc analysis p-values**

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.00** | **.04** |
| HCF-IDF | | **.00** |

**The factor *Decay Function* × *Document Content*.** Table 16 shows the results of ANOVA regarding the factor *Decay Function* when a choice of the factor *Document Content* is fixed and vice versa. According to Table 16, there is a significant difference among the factor *Document Content*, when Sliding window is used. The rankscores and post-hoc analysis of it are shown in Tables 17(a) and (b). It indicates that All significantly enhances the performance of the recommender system when Sliding window is used.

## B. MEAN AVERAGE PRECISION

**Table 14: Rankscores and Post-hoc analysis for the factor *Profiling Method* at All using Shaffer's MSRB procedure.**

**a) Rankscores**

| Choice | M | SD |
|---|---|---|
| CF-IDF | .55 | .33 |
| HCF-IDF | .53 | .32 |
| LDA | .30 | .32 |

**b) Post-hoc analysis p-values**

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .20 | **.00** |
| HCF-IDF | | **.00** |

**Table 15: Rankscores and Post-hoc analysis for the factor *Document Content* at CF-IDF using Shaffer's MSRB procedure.**

**a) Rankscores**

| Choice | M | SD |
|---|---|---|
| All | .55 | .33 |
| Title | .40 | .28 |

**b) Post-hoc analysis p-value**

| | All |
|---|---|
| Title | **.00** |

**Table 16: ANOVA for *Decay Function* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Decay Function* at Title | 0.04 | .00 | .85 |
| *Decay Function* at All | 3.16 | .03 | .08 |
| *Document Content* at Sliding window | 9.44 | .08 | **.00** |
| *Document Content* at Exponential decay | 0.56 | .00 | .46 |

**Table 17: Rankscores and Post-hoc analysis for the factor *Document Content* at Sliding window using Shaffer's MSRB procedure.**

**a) Rankscores**

| Choice | M | SD |
|---|---|---|
| All | .48 | .36 |
| Title | .42 | .32 |

**b) Post-hoc analysis p-value**

| | All |
|---|---|
| Title | **.00** |

We look into the recommendation performance by computing Mean Average Precision (MAP). Average Precision (AP) is calculated as described in Equation 13.

$$AP = \frac{1}{|hits|} \sum_{d \in hits} Precision@rank_d, \quad (13)$$

where $hits$ and $rank_d$ stand for the set of relevant items and the rank of the item $d$, respectively. $|hits|$ is the number of relevant items in the recommendation list. $Precision@rank_d$ denotes the precision at cut off $rank_d$ in the recommendation list. Mean Average Precision (MAP) is the mean of the Average Precision scores for each participant. In this section, we evaluate the recommendation performance using MAP. Particularly, we first compare the twelve different strategies. Subsequently, we investigate the difference among the different experiment factors.

## B.1 Best performing strategy

Table 18 shows the Mean Average Precisions (MAP) of the twelve strategies. The order of the strategies is almost same with rankscores shown in Table 3. In order to investigate significant differences among strategies, we first apply Mauchly's test and found a violation of sphericity in the strategies ($\chi^2(65) = 353.51$, $p = .00$). Subsequently, we run a one-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .65$. It reveals a significant difference of the strategies' MAPs ($F(7.17, 875.15) = 15.59$, $p = .00$). To assess the statistical significance of pair-wise differences between the twelve strategies, a post-hoc analysis is performed using Shaffer's MSRB procedure [23]. The result of the post-hoc analysis is presented in Table 19. The vertical and horizontal dimensions of the Table 19 show the eleven-by-eleven comparison of the twelve strategies. As one can see, we observe various significant differences between strategies (marked in bold font).

**Table 18: Mean Average Precision (MAP) of the strategies in decreasing order.**

|     | Strategy | | | MAP |
| --- | --- | --- | --- | --- |
|     | Profiling Method | Decay Function | Content | M (SD) |
| 1.  | CF-IDF | Sliding window | All | .71 (.32) |
| 2.  | HCF-IDF | Exponential decay | All | .65 (.33) |
| 3.  | HCF-IDF | Exponential decay | Title | .65 (.32) |
| 4.  | CF-IDF | Exponential decay | All | .65 (.35) |
| 5.  | HCF-IDF | Sliding window | Title | .65 (.34) |
| 6.  | HCF-IDF | Sliding window | All | .65 (.34) |
| 7.  | CF-IDF | Exponential decay | Title | .58 (.35) |
| 8.  | CF-IDF | Sliding window | Title | .55 (.34) |
| 9.  | LDA | Exponential decay | All | .47 (.39) |
| 10. | LDA | Exponential decay | Title | .44 (.34) |
| 11. | LDA | Sliding window | Title | .43 (.35) |
| 12. | LDA | Sliding window | All | .40 (.42) |

## B.2 Difference in experiment factors

Subsequently, we analyze the results of MAPs with respect to each factor. First, we apply Mendoza's test [17] which shows violations of sphericity against the factors *Profiling Method* $\times$ *Decay Function* ($\chi^2(2) = 10.30$, $p = .01$), and *Profiling Method* $\times$ *Document Content* ($\chi^2(2) = 13.18$, $p = .00$). Thus, we run three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .92$ for the factor *Profiling Method* $\times$ *Decay Function*, and $\epsilon = .91$ for the factor *Profiling Method* $\times$ *Document Content*. Table 20 shows the results of applying an ANOVA. $\eta^2$ indicates the effect size of each factor. For all the factors that make significant difference, we conduct a post-hoc analysis using Shaffer's MSRB Procedure.

**The factor *Profiling Method*.** Tables 21(a), (b) and (c) show the MAPs with respect to each profiling method, the post-hoc analysis for the factor *Profiling Method*, and the effect size, respectively. Table 21(a) presents the means and standard deviations of the three profiling methods. Table 21(b) shows p-values of each pair. Since Table 20 shows that the factor *Profiling Method* has the largest effect size, we further compute the effect size using Cohen's d for each pair shown in Table 21(c). The result shows that CF-IDF and HCF-IDF are superior to LDA. In contrast, there is

**Table 20: Three-way repeated-measure ANOVA with Greenhouse-Geisser correction with F-ratio, $\eta^2$, and p-value for MAP.**

| Factor | F | $\eta^2$ | p |
| --- | --- | --- | --- |
| *Profiling Method* | 51.79 | .42 | **.00** |
| *Decay Function* | 0.33 | .00 | .57 |
| *Document Content* | 5.16 | .04 | **.02** |
| *Profiling Method* $\times$ *Decay Function* | 1.66 | .01 | .20 |
| *Profiling Method* $\times$ *Document Content* | 4.76 | .02 | **.01** |
| *Decay Function* $\times$ *Document Content* | 0.02 | .00 | .90 |
| *Profiling Method* $\times$ *Decay Function* $\times$ *Document Content* | 3.19 | .03 | **.04** |

no significant difference between CF-IDF and HCF-IDF, although MAP of HCF-IDF is slightly higher than CF-IDF.

**Table 21: MAPs, Post-hoc analysis for the factor *Profiling Method* using Shaffer's MSRB procedure, and effect size.**

a) MAPs

| Choice | M | SD |
| --- | --- | --- |
| HCF-IDF | .65 | .33 |
| CF-IDF | .62 | .35 |
| LDA | .43 | .38 |

b) Post-hoc analysis p-values

|     | HCF-IDF | LDA |
| --- | --- | --- |
| CF-IDF | .15 | **.00** |
| HCF-IDF | | **.00** |

c) Effect size using Cohen's d

|     | HCF-IDF | LDA |
| --- | --- | --- |
| CF-IDF | .09 | .52 |
| HCF-IDF | | .62 |

**The factor *Document Content*.** Table 22 shows the post-hoc analysis for the factor *Document Content*. It indicates that the recommender system works better when All (i.e., full-texts and titles) is taken into consideration for computing recommendations.

**Table 22: MAPs and Post-hoc analysis for the factor *Document Content* using Shaffer's MSRB procedure.**

a) MAPs

| Choice | M | SD |
| --- | --- | --- |
| All | .59 | .38 |
| Title | .55 | .35 |

b) Post-hoc analysis p-values

|     | Title |
| --- | --- |
| All | **.02** |

**The factor *Profiling Method* $\times$ *Document Content*.** Table 23 shows the results of ANOVA regarding the factor *Profiling Method* when a choice of the factor *Document Content* is fixed and vice versa. We observe there are significant differences when a choice of the factor *Document Content* is fixed and CF-IDF is employed. Mendoza's test found a violation of sphericity in the factor *Profiling Method* when All is taken ($\chi^2(2) = 31.35$, $p = .00$). Thus, we run a one-way repeated-measure ANOVA with Greenhouse-Geisser correc-

**Table 19: Post-hoc analysis of Mean Average Precision (MAP) with pairwise p-values over the strategies using Shaffer's MSRB procedure. The p-values are marked in bold font if $p < .05$, which indicates a significant difference between the two strategies. Strategies are sorted by Precision@5 as shown in Table 18.**

| | | | All / Exponential decay / HCF-IDF / 2. | Title / Exponential decay / HCF-IDF / 3. | All / Exponential decay / CF-IDF / 4. | Title / Sliding window / HCF-IDF / 5. | All / Sliding window / HCF-IDF / 6. | Title / Exponential decay / CF-IDF / 7. | Title / Sliding window / CF-IDF / 8. | All / Exponential decay / LDA / 9. | Title / Exponential decay / LDA / 10. | Title / Sliding window / LDA / 11. | All / Sliding window / LDA / 12. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. | CF-IDF | Sliding window | All | .99 | .99 | .99 | .99 | .99 | **.02** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 2. | HCF-IDF | Exponential decay | All | | .99 | .99 | .99 | .99 | .99 | .25 | **.00** | **.00** | **.00** | **.00** |
| 3. | HCF-IDF | Exponential decay | Title | | | .99 | .99 | .99 | .99 | .99 | **.00** | **.00** | **.00** | **.00** |
| 4. | CF-IDF | Exponential decay | All | | | | .99 | .99 | .99 | .33 | **.01** | **.00** | **.00** | **.00** |
| 5. | HCF-IDF | Sliding window | Title | | | | | .99 | .99 | .54 | **.00** | **.00** | **.00** | **.00** |
| 6. | HCF-IDF | Sliding window | All | | | | | | .99 | .60 | **.00** | **.00** | **.00** | **.00** |
| 7. | CF-IDF | Exponential decay | Title | | | | | | | .99 | .67 | **.02** | **.01** | **.00** |
| 8. | CF-IDF | Sliding window | Title | | | | | | | | .99 | .06 | **.03** | **.02** |
| 9. | LDA | Exponential decay | All | | | | | | | | | .99 | .99 | .64 |
| 10. | LDA | Exponential decay | Title | | | | | | | | | | .99 | .99 |
| 11. | LDA | Sliding window | Title | | | | | | | | | | | .99 |

tion of $\eta = .81$ for the second row in Table 23. Subsequently, we conduct the post-hoc analyses for each factor in which shows a significant difference. Table 24 presents the post-hoc analysis when Title is employed. We see that HCF-IDF outperforms others with significant differences. Table 25 shows the post-hoc analysis when All is chosen for the factor *Document Content*. Different from the result shown in Table 24, CF-IDF performs slightly better than HCF-IDF, although there is no significant difference between them. Both CF-IDF and HCF-IDF demonstrate better recommendation performance than LDA. Table 26 shows the post-hoc analysis of the factor *Document Content* when CF-IDF is employed. It indicates that the strategies with CF-IDF and All significantly outperforms those with CF-IDF and Title.

**Table 23: ANOVA for *Profiling Method* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* at Title | 23.99 | .20 | **.00** |
| *Profiling Method* at All | 36.35 | .30 | **.00** |
| *Document Content* at CF-IDF | 14.69 | .12 | **.00** |
| *Document Content* at HCF-IDF | 0.00 | .00 | .95 |
| *Document Content* at LDA | 0.01 | .00 | .93 |

## C. PRECISION

In this section, we evaluate the recommendation performance by computing Precision, especially Precision@5 (P@5). Precision is computed as described in Equation 14.

$$Precision@k = \frac{1}{k} \sum_{i=1}^{k} rel(i), \qquad (14)$$

where $rel(k)$ returns 1 if the item ranked at $i$ is relevant and 0 if irrelevant. In this paper, we set $k = 5$, since five items are recommended by each strategy in the experiment. Using

**Table 24: MAPs and Post-hoc analysis for the factor *Profiling Method* at Title using Shaffer's MSRB procedure.**

a) MAPs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .65 | .33 |
| CF-IDF | .56 | .35 |
| LDA | .43 | .35 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.01** | **.00** |
| HCF-IDF | | **.00** |

**Table 25: MAPs and Post-hoc analysis for the factor *Profiling Method* at All using Shaffer's MSRB procedure.**

a) MAPs

| Choice | M | SD |
|---|---|---|
| CF-IDF | .68 | .34 |
| HCF-IDF | .65 | .34 |
| LDA | .44 | .41 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .21 | **.00** |
| HCF-IDF | | **.00** |

Precision@5, we first compare the twelve different strategies. Subsequently, we investigate the difference among the different experiment factors.

### C.1 Best performing strategy

Table 27 shows Precision@5 of each strategy. For the statistical analyses, we first applied Mauchly's test and found a violation of sphericity in the strategies ($\chi^2(65) = 421.32$, $p = .00$). Subsequently, we run a one-way repeated-measure

**Table 26: MAPs and Post-hoc analysis for the factor *Document Content* at CF-IDF using Shaffer's MSRB procedure.**

**a) MAPs**

| Choice | M | SD |
|---|---|---|
| All | .68 | .34 |
| Title | .56 | .35 |

**b) Post-hoc analysis p-values**

| | All |
|---|---|
| Title | **.00** |

ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. It reveals a significant difference of the strategies' MAPs ($F(6.62, 808.00) = 21.85$, $p = .00$). To assess the statistical significance of pair-wise differences between the twelve strategies, a post-hoc analysis is performed using Shaffer's MSRB procedure [23]. The result of the post-hoc analysis is presented in Table 28. The vertical and horizontal dimensions of the Table 28 show the eleven-by-eleven comparison of the twelve strategies. As one can see, we observe various significant differences between strategies (marked in bold font).

**Table 27: Precision@5 (P@5) of the strategies in decreasing order.**

| | Strategy | | | P@5 |
|---|---|---|---|---|
| | Profiling Method | Decay Function | Content | M (SD) |
| 1. | CF-IDF | Sliding window | All | .59 (.33) |
| 2. | HCF-IDF | Sliding window | All | .56 (.33) |
| 3. | HCF-IDF | Sliding window | Title | .55 (.33) |
| 4. | HCF-IDF | Exponential decay | Title | .52 (.30) |
| 5. | CF-IDF | Exponential decay | All | .50 (.32) |
| 6. | HCF-IDF | Exponential decay | All | .48 (.30) |
| 7. | CF-IDF | Exponential decay | Title | .40 (.29) |
| 8. | CF-IDF | Sliding window | Title | .39 (.27) |
| 9. | LDA | Exponential decay | Title | .37 (.31) |
| 10. | LDA | Sliding window | Title | .34 (.31) |
| 11. | LDA | Exponential decay | All | .31 (.30) |
| 12. | LDA | Sliding window | All | .27 (.33) |

## C.2 Difference in experiment factors

Subsequently, we analyze the results of Precision@5 with respect to each factor. First, we apply Mendoza's test [17] which showed violations of sphericity against the factors *Profiling Method* ($\chi^2(2) = 13.92$, $p = .00$), *Profiling Method* $\times$ *Decay Function* ($\chi^2(2) = 19.64$, $p = .00$), and *Profiling Method* $\times$ *Document Content* ($\chi^2(2) = 7.23$, $p = .03$). Thus, we run three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .90$ for the factor *Profiling Method*, $\epsilon = .87$ for the factor *Profiling Method* $\times$ *Decay Function*, and $\epsilon = .95$ for the factor *Profiling Method* $\times$ *Document Content*. Table 29 shows the result of an ANOVA with F-ratio, $\eta^2$ and p-value. $\eta^2$ indicates the effect size of each factor. The effect size is small when $\eta^2 > .02$, medium when $\eta^2 > .13$, and large when $\eta^2 > .26$. For all factors that make significant difference, we conduct a post-hoc analysis using Shaffer's MSRB Procedure.

**The factor *Profiling Method*.** Tables 30(a), (b) and (c) show the Precision@5, the post-hoc analysis for the fac-

**Table 29: Three-way repeated-measure ANOVA with Greenhouse-Geisser correction with F-ratio, $\eta^2$ and p-value for Precision@5.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* | 54.24 | .42 | **.00** |
| *Decay Function* | 1.75 | .00 | .19 |
| *Document Content* | 3.23 | .04 | .08 |
| *Profiling Method* $\times$ *Decay Function* | 6.32 | .01 | **.00** |
| *Profiling Method* $\times$ *Document Content* | 20.53 | .02 | **.00** |
| *Decay Function* $\times$ *Document Content* | 7.13 | .00 | **.01** |
| *Profiling Method* $\times$ *Decay Function* $\times$ *Document Content* | 2.61 | .03 | .07 |

tor *Profiling Method*, and the effect size, respectively. Table 30(a) presents the means and standard deviations of the three profiling methods. Table 30(b) shows p-values of each pair. Since Table 29 shows that this factor has the largest effect size, we further compute the effect size using Cohen's d for each pair shown in Table 21(c). There are significant differences between all pairs of the three profiling methods and among the three profiling methods HCF-IDF performs best significantly.

**Table 30: Precision@5, Post-hoc analysis for the factor *Profiling Method* using Shaffer's MSRB procedure, and effect size.**

**a) Precision@5**

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .53 | .31 |
| CF-IDF | .47 | .31 |
| LDA | .32 | .31 |

**b) Post-hoc analysis p-values**

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.00** | **.00** |
| HCF-IDF | | **.00** |

**c) Effect size using Cohen's d**

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .09 | .52 |
| HCF-IDF | | .62 |

**The factor *Profiling Method* $\times$ *Decay Function*.** Table 31 shows the results of ANOVA regarding the *profiling method* when a choice of the *Decay Function* is fixed and vice versa. There are significant differences when the choice of the factor *Decay Function* is fixed. In both decay functions, all pairs of the three profiling methods show significant differences. Specifically, HCF-IDF performs best, followed by CF-IDF and LDA. When CF-IDF is employed, Sliding window makes significantly better recommendations than Exponential decay ($F(1, 122) = 5.44$, $p = .02$). In contrast, when LDA is employed, Exponential decay performs significantly better than Sliding window ($F(1, 122) = 6.75$, $p = .01$). The factor *Decay Function* does not make difference on the recommendation performance when HCF-IDF is employed.

**The factor *Profiling Method* $\times$ *Document Content*.** Table 32 shows the results of ANOVA regarding the factor *Profiling Method* when a choice of the *Document Content* is fixed and vice versa. When the choice of the *Document Content* is Title, HCF-IDF performs best and significantly

**Table 28: Post-hoc analysis of Precision@5 (P@5) with pairwise p-values over the strategies using Shaffer's MSRB procedure. The p-values are marked in bold font if $p < .05$, which indicates a significant difference between the two strategies. Strategies are sorted by Precision@5 as shown in Table 27.**

| | | | | All | Title | Title | All | All | Title | Title | Title | Title | All | All |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Sliding window | Sliding window | Exponential decay | Exponential decay | Exponential decay | Exponential decay | Sliding window | Exponential decay | Sliding window | Exponential decay | Sliding window |
| | | | | HCF-IDF | HCF-IDF | HCF-IDF | CF-IDF | HCF-IDF | CF-IDF | CF-IDF | LDA | LDA | LDA | LDA |
| | | | | 2. | 3. | 4. | 5. | 6. | 7. | 8. | 9. | 10. | 11. | 12. |
| 1. | CF-IDF | Sliding window | All | .99 | .99 | .60 | .09 | **.05** | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 2. | HCF-IDF | Sliding window | All | | .99 | .99 | .99 | .50 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 3. | HCF-IDF | Sliding window | Title | | | .99 | .99 | .99 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 4. | HCF-IDF | Exponential decay | Title | | | | .99 | .99 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 5. | CF-IDF | Exponential decay | All | | | | | .99 | **.03** | **.01** | **.02** | **.00** | **.00** | **.00** |
| 6. | HCF-IDF | Exponential decay | All | | | | | | .09 | **.03** | **.03** | **.00** | **.00** | **.00** |
| 7. | CF-IDF | Exponential decay | Title | | | | | | | .99 | .99 | .99 | .26 | **.01** |
| 8. | CF-IDF | Sliding window | Title | | | | | | | | .99 | .99 | .34 | **.02** |
| 9. | LDA | Exponential decay | Title | | | | | | | | | .99 | .99 | .09 |
| 10. | LDA | Sliding window | Title | | | | | | | | | | .99 | .82 |
| 11. | LDA | Exponential decay | All | | | | | | | | | | | .99 |

**Table 31: ANOVA for *Profiling Method* × *Decay Function* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* at Sliding window | 52.98 | .20 | **.00** |
| *Profiling Method* at Exponential decay | 22.52 | .30 | **.00** |
| *Decay Function* at CF-IDF | 5.44 | .12 | **.02** |
| *Decay Function* at HCF-IDF | 3.25 | .00 | .07 |
| *Decay Function* at LDA | 6.75 | .00 | **.01** |

better than both CF-IDF and LDA. There is no significant difference between CF-IDF and LDA. When the choice of the *Document Content* is All, HCF-IDF performs best. But, there is no significant difference between CF-IDF and HCF-IDF and both profiling methods are significantly superior to LDA. When CF-IDF is employed, All is the better choice than Title. In contrast, Title performs better than All, when LDA is employed.

**Table 32: ANOVA for *Profiling Method* × *Document Content* (*Document Content*) interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* at Title | 23.37 | .20 | **.00** |
| *Profiling Method* at All | 56.54 | .30 | **.00** |
| *Document Content* at CF-IDF | 33.39 | .12 | **.00** |
| *Document Content* at HCF-IDF | 0.44 | .00 | .51 |
| *Document Content* at LDA | 4.68 | .00 | **.03** |

**The factor *Decay Function* × *Document Content*.** Table 33 shows the results of ANOVA regarding the factor *Decay Function* when a choice of the factor *Document Content* is fixed and vice versa. When All is choosen for the factor *Document Content*, Sliding window is the better decay function. When Sliding window is employed in the strate-

gies, the strategies with All is significantly better than those with Title.

**Table 33: ANOVA for *Decay Function* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Decay Function* at Title | 0.08 | .20 | .78 |
| *Decay Function* at All | 4.99 | .30 | **.03** |
| *Document Content* at Sliding window | 8.74 | .12 | **.00** |
| *Document Content* at Exponential decay | 0.00 | .00 | .97 |

## D. MEAN RECIPROCAL RANK

In this section, we evaluate the recommendation performance by computing Mean Reciprocal Rank (MRR). Reciprocal Rank is defined as Equation 15.

$$RR = \frac{1}{rank_{first}}, \qquad (15)$$

where $rank_{first}$ indicates the rank position of the first item which is evaluated as interesting. Mean Reciprocal Rank (MRR) is the mean of the Reciprocal Rank scores for each participant. If there is no relevant items in the recommendation list, RR outputs 0. Using MRR, we first compare the twelve different strategies. Subsequently, we investigate the difference among the different experiment factors.

### D.1 Best Performing Strategy

Table 34 shows the Mean Reciprocal Ranks (MRR) of each strategies. The order of the strategies are different from rankscores shown in Table 3. For the statistical analyses, we first applied Mauchly's test and found a violation of sphericity in the strategies ($\chi^2(65) = 308.70$, $p = .00$). Subsequently, we ran a one-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .67$. It revealed a significant difference of the strategies' MRRs ($F(0.18, 2.53) =$

**Table 35: Post-hoc analysis of Mean Reciprocal Rank (MRR) with pairwise p-values over the strategies using Shaffer's MSRB procedure. The p-values are marked in bold font if $p < .05$, which indicates a significant difference between the two strategies. Strategies are sorted by MRR as shown in Table 34.**

| | | | | Exponential decay / All / CF-IDF / 2. | Exponential decay / All / HCF-IDF / 3. | Exponential decay / Title / HCF-IDF / 4. | Sliding window / Title / HCF-IDF / 5. | Sliding window / All / HCF-IDF / 6. | Exponential decay / Title / CF-IDF / 7. | Sliding window / Title / CF-IDF / 8. | Exponential decay / All / LDA / 9. | Exponential decay / Title / LDA / 10. | Sliding window / Title / LDA / 11. | Sliding window / All / LDA / 12. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. | CF-IDF | Sliding window | All | .99 | .99 | .99 | .99 | .99 | .12 | **.03** | **.00** | **.00** | **.00** | **.00** |
| 2. | CF-IDF | Exponential decay | All | | .99 | .99 | .99 | .99 | .99 | .99 | **.01** | **.00** | **.00** | **.00** |
| 3. | HCF-IDF | Exponential decay | All | | | .99 | .99 | .99 | .99 | .99 | **.01** | **.00** | **.00** | **.01** |
| 4. | HCF-IDF | Exponential decay | Title | | | | .99 | .99 | .99 | .99 | **.01** | **.00** | **.00** | **.00** |
| 5. | HCF-IDF | Sliding window | Title | | | | | .99 | .99 | .99 | **.01** | **.00** | **.00** | **.00** |
| 6. | HCF-IDF | Sliding window | All | | | | | | .99 | .99 | **.01** | **.00** | **.00** | **.01** |
| 7. | CF-IDF | Exponential decay | Title | | | | | | | .99 | .94 | **.00** | **.00** | **.00** |
| 8. | CF-IDF | Sliding window | Title | | | | | | | | .99 | **.00** | **.00** | **.02** |
| 9. | LDA | Exponential decay | All | | | | | | | | | .99 | .99 | .58 |
| 10. | LDA | Exponential decay | Title | | | | | | | | | | .99 | .99 |
| 11. | LDA | Sliding window | Title | | | | | | | | | | | .99 |

14.40, $p = .00$). To assess the statistical significance of pair-wise differences between the twelve strategies, a post-hoc analysis was performed using Shaffer's MSRB procedure [23]. The result of the post-hoc analysis is presented in Table 35. The vertical and horizontal dimensions of the Table 35 show the eleven-by-eleven comparison of the twelve strategies. As one can see, we observe various significant differences between strategies (marked in bold font).

**Table 34: Mean Reciprocal Rank (MRR) of the strategies in decreasing order.**

| | Strategy | | | MRR |
|---|---|---|---|---|
| | **Profiling Method** | **Decay Function** | **Content** | **M (SD)** |
| 1 | CF-IDF | Sliding window | All | .73 (.35) |
| 2 | CF-IDF | Exponential decay | All | .69 (.39) |
| 3 | HCF-IDF | Exponential decay | All | .68 (.37) |
| 4 | HCF-IDF | Exponential decay | Title | .68 (.37) |
| 5 | HCF-IDF | Sliding window | Title | .67 (.38) |
| 6 | HCF-IDF | Sliding window | All | .67 (.37) |
| 7 | CF-IDF | Exponential decay | Title | .61 (.39) |
| 8 | CF-IDF | Sliding window | Title | .59 (.39) |
| 9 | LDA | Exponential decay | All | .50 (.43) |
| 10 | LDA | Exponential decay | Title | .43 (.37) |
| 11 | LDA | Sliding window | Title | .42 (.38) |
| 12 | LDA | Sliding window | All | .41 (.44) |

## D.2 Difference in experiment factors

Subsequently, we look into the results of MRRs with respect to each factor. First, we applied Mendoza's test [17] which showed violations of sphericity against the factors *Profiling Method × Decay Function* ($\chi^2(2) = 8.16$, $p = .02$), and *Profiling Method × Document Content* ($\chi^2(2) = 8.85$, $p = .01$). Thus, we ran three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .94$ for *Profiling Method × Decay Function*, and $\epsilon = .93$ for *profiling method × document content*. Table 36 shows the results of an ANOVA with F-ratio, $\eta^2$ and p-value. The analysis revealed significant differences only in the two factors *Profiling Method* and *Document Content*.

**Table 36: Three-way repeated-measure ANOVA with Greenhouse-Geisser correction with F-ratio, $\eta^2$ and p-value for MRR.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* | 50.65 | .42 | **.00** |
| *Decay Function* | 0.56 | .00 | .45 |
| *Document Content* | 5.10 | .04 | **.03** |
| *Profiling Method × Decay Function* | 1.28 | .01 | .28 |
| *Profiling Method × Document Content* | 2.83 | .02 | .06 |
| *Decay Function × Document Content* | 0.13 | .00 | .72 |
| *Profiling Method × Decay Function × Document Content* | 2.33 | .02 | .10 |

**The factor _Profiling Method_.** Tables 37(a), (b) and (c) show the MRRs, the post-hoc analysis for the factor *Profiling Method*, and the effect size, respectively. Table 37(a) presents the means and standard deviations of the three profiling methods. Table 37(b) shows p-values of each pair. Since Table 36 shows that this factor has the largest effect size, we further compute the effect size using Cohen's d for each pair shown in Table 37(c).

**The factor *Document Content*.** Table 38 shows the post-hoc analysis for the factor *Document Content*. It indicates that generally the recommender system work better when full-texts are available.

## E. NORMALIZED DISCOUNTED CUMULATIVE GAIN

In this section, we evaluate the recommendation performance by Normalized Discounted Cumulative Gain (nDCG).

**Table 37: MRRs, Post-hoc analysis for the factor *Profiling Method* using Shaffer's MSRB procedure, and effect size.**

a) MRRs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .68 | .38 |
| CF-IDF | .66 | .37 |
| LDA | .44 | .41 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .34 | **.00** |
| HCF-IDF | | **.00** |

c) Effect size using Cohen's d

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .05 | .56 |
| HCF-IDF | | .61 |

**Table 38: MRRs and Post-hoc analysis for the factor *Document Content* using Shaffer's MSRB procedure.**

a) MRRs

| Choice | M | SD |
|---|---|---|
| All | .61 | .41 |
| Title | .57 | .39 |

b) Post-hoc analysis p-values

| | All |
|---|---|
| Title | **.03** |

Discounted Cumulative Gain (DCG) is calculated as Equation 16.

$$DCG = \sum_{i=1}^{k} \frac{2^{rel(i)} - 1}{\log_2 i},  \qquad (16)$$

where $rel(k)$ returns 1 if the item ranked at $i$ is relevant and 0 if irrelevant. Similar to rankscore, the items ranked at higher positions have a larger influence on output score. First, we compare the twelve different strategies using the metric. Subsequently, we investigate the difference among the different experiment factors.

## E.1 Best performing strategy

Table 39 shows the Normalized Discounted Cumulative Gain (nDCG) of the twelve strategies. The order of the strategies is identical with rankscores shown in Table 3. For the statistical analyses, we first apply Mauchly's test and found a violation of sphericity in the strategies ($\chi^2(65) = 424.00$, $p = .00$). Subsequently, we run a one-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .61$. It reveals a significant difference of the strategies' nDCG ($F(6.69, 816.37) = 21.16$, $p = .00$). To assess the statistical significance of pair-wise differences between the twelve strategies, a post-hoc analysis is performed using Shaffer's MSRB procedure [23]. The result of the post-hoc analysis is presented in Table 40. The vertical and horizontal dimensions of the Table 40 show the eleven-by-eleven comparison of the twelve strategies. As one can see, we observe various significant differences between strategies (marked in bold font).

## E.2 Difference in experiment factors

Subsequently, we analyze the results of nDCGs with respect to each factor. First, we apply Mendoza's test [17]

**Table 39: nDCGs of the strategies in decreasing order. M and SD denote mean and standard deviation, respectively.**

| | Strategy | | | nDCG |
|---|---|---|---|---|
| | Profiling Method | Decay Function | Content | M (SD) |
| 1. | CF-IDF | Sliding window | All | .59 (.33) |
| 2. | HCF-IDF | Sliding window | All | .56 (.34) |
| 3. | HCF-IDF | Sliding window | Title | .55 (.33) |
| 4. | HCF-IDF | Exponential decay | Title | .52 (.30) |
| 5. | CF-IDF | Exponential decay | All | .52 (.32) |
| 6. | HCF-IDF | Exponential decay | All | .50 (.30) |
| 7. | CF-IDF | Exponential decay | Title | .41 (.30) |
| 8. | CF-IDF | Sliding window | Title | .40 (.27) |
| 9. | LDA | Exponential decay | Title | .34 (.31) |
| 10. | LDA | Sliding window | Title | .32 (.31) |
| 11. | LDA | Exponential decay | All | .32 (.31) |
| 12. | LDA | Sliding window | All | .28 (.33) |

which shows violations of sphericity against the factors *Profiling Method* ($\chi^2(2) = 11.29$, $p = .00$), *Profiling Method $\times$ Decay Function* ($\chi^2(2) = 18.90$, $p = .00$), and *Profiling Method $\times$ Document Content* ($\chi^2(2) = 8.61$, $p = .01$). Thus, we run three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .92$ for the factor *Profiling Method*, $\epsilon = .87$ for the factor *Profiling Method $\times$ Decay Function*, and $\epsilon = .94$ for the factor *Profiling Method $\times$ Document Content*. Table 41 shows the results of applying an ANOVA. $\eta^2$ indicates the effect size of each factor. For all the factors that make significant difference, we conduct a post-hoc analysis using Shaffer's MSRB Procedure.

**Table 41: Three-way repeated-measure ANOVA with Greenhouse-Geisser correction with F-ratio, $\eta^2$, and p-value for nDCG.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* | 58.42 | . | **.00** |
| *Decay Function* | 0.80 | . | .37 |
| *Document Content* | 6.33 | . | **.01** |
| *Profiling Method $\times$ Decay Function* | 3.81 | . | **.03** |
| *Profiling Method $\times$ Document Content* | 14.54 | . | **.00** |
| *Decay Function $\times$ Document Content* | 3.57 | . | .06 |
| *Profiling Method $\times$ Decay Function $\times$ Document Content* | 3.09 | . | **.05** |

**The factor *Profiling Method*.** Tables 42(a), (b) and (c) show the MAPs, the post-hoc analysis for the factor *Profiling Method*, and the effect size, respectively. Table 42(a) presents the means and standard deviations of the three profiling methods. Table 42(b) shows p-values of each pair. Since Table 41 shows that the factor *Profiling Method* has the largest effect size, we further compute the effect size using Cohen's d for each pair shown in Table 42(c). The result shows that CF-IDF and HCF-IDF are superior to LDA. In contrast, there is no significant difference between CF-IDF and HCF-IDF, although MAP of HCF-IDF is slightly better than CF-IDF.

**The factor *Document Content*.** Table 43 shows the post-hoc analysis for the factor *Document Content*. It indicates that the recommender system works better when All (i.e.,

**Table 40: Post-hoc analysis of normalized Discounted Cumulative Gain (nDCG) with pairwise p-values over the twelve strategies using Shaffer's MSRB procedure. The p-values are marked in bold font if $p < .05$, which indicates a significant difference between the two strategies. Strategies are sorted by nDCG as shown in Table 39.**

| | | | | HCF-IDF Sliding window All | HCF-IDF Sliding window Title | HCF-IDF Exponential decay Title | CF-IDF Exponential decay All | HCF-IDF Exponential decay All | CF-IDF Exponential decay Title | CF-IDF Sliding window Title | LDA Exponential decay Title | LDA Sliding window Title | LDA Exponential decay All | LDA Sliding window All |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 2. | 3. | 4. | 5. | 6. | 7. | 8. | 9. | 10. | 11. | 12. |
| 1. | CF-IDF | Sliding window | All | .99 | .99 | .85 | .41 | .22 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 2. | HCF-IDF | Sliding window | All | | .99 | .99 | .99 | .99 | **.00** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 3. | HCF-IDF | Sliding window | Title | | | .99 | .99 | .99 | **.01** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 4. | HCF-IDF | Exponential decay | Title | | | | .99 | .99 | **.01** | **.00** | **.00** | **.00** | **.00** | **.00** |
| 5. | CF-IDF | Exponential decay | All | | | | | .99 | .05 | **.01** | **.00** | **.00** | **.00** | **.00** |
| 6. | HCF-IDF | Exponential decay | All | | | | | | .17 | **.03** | **.00** | **.00** | **.00** | **.00** |
| 7. | CF-IDF | Exponential decay | Title | | | | | | | .99 | .85 | .18 | .34 | **.01** |
| 8. | CF-IDF | Sliding window | Title | | | | | | | | .99 | .41 | .77 | **.05** |
| 9. | LDA | Exponential decay | Title | | | | | | | | | .99 | .99 | .99 |
| 10. | LDA | Sliding window | Title | | | | | | | | | | .99 | .99 |
| 11. | LDA | Exponential decay | All | | | | | | | | | | | .98 |

**Table 42: nDCGs, Post-hoc analysis for the factor _Profiling Method_ using Shaffer's MSRB procedure, and effect size.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .53 | .32 |
| CF-IDF | .48 | .32 |
| LDA | .32 | .32 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.00** | **.00** |
| HCF-IDF | | **.00** |

c) Effect size using Cohen's d

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | . | . |
| HCF-IDF | | . |

full-texts and titles) is taken into consideration for computing recommendations.

**Table 43: nDCGs and Post-hoc analysis for the factor _Document Content_ using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| All | .46 | .34 |
| Title | .42 | .31 |

b) Post-hoc analysis p-values

| | Title |
|---|---|
| All | **.01** |

**The factor _Profiling Method_ × _Decay Function_.** Table 44 shows the results of ANOVA regarding the factor _Profiling Method_ when a choice of the factor _Decay Function_ is fixed and vice versa. Mendoza's test found a violation of sphericity in the factor _Profiling Method_ when Sliding window is used ($\chi^2(2) = 7.55$, $p = .02$) and Exponential decay is used ($\chi^2(2) = 10.74$, $p = .00$). Thus, we run a one-way repeated-measure ANOVA with Greenhouse-Geisser correction of $\eta = .94$ for the first row in Table 48 and $\eta = .92$ for the second row in Table 48. We observe significant differences when a choice of the factor _Decay Function_ is fixed and when LDA is employed. The post-hoc analyses of them are shown in Table 45, Table 46, and Table 47, respectively. In Table 45 and Table 46, a choice of the factor _Decay Function_ is fixed. Table 47 shows the post-hoc analysis of the factor _Decay Function_ when LDA is employed. It indicates Exponential decay performs better than Sliding window for LDA.

**Table 44: ANOVA for _Profiling Method_ × _Decay Function_ interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| _Profiling Method_ at Sliding window | 50.59 | . | **.00** |
| _Profiling Method_ at Exponential decay | 27.92 | . | **.00** |
| _Decay Function_ at CF-IDF | 2.79 | . | .10 |
| _Decay Function_ at HCF-IDF | 1.78 | . | .18 |
| _Decay Function_ at LDA | 4.90 | . | **.03** |

**The factor _Profiling Method_ × _Document Content_.** Table 48 shows the results of ANOVA regarding the factor _Profiling Method_ when a choice of the factor _Document Content_ is fixed and vice versa. We observe there are significant differences when a choice of the factor _Document Content_ is fixed and CF-IDF is employed. Mendoza's test found a violation of sphericity in the factor _Profiling Method_ when All is taken ($\chi^2(2) = 24.64$, $p = .00$). Thus, we run a one-way repeated-measure ANOVA with Greenhouse-Geisser correction of $\eta = .84$ for the second row in Table 48. Table 24

**Table 45: nDCGs and Post-hoc analysis for the factor *Profiling Method* at Sliding window using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .55 | .33 |
| CF-IDF | .50 | .32 |
| LDA | .30 | .32 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.02** | **.00** |
| HCF-IDF | | **.00** |

**Table 46: nDCGs and Post-hoc analysis for the factor *Profiling Method* at Exponential decay using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .51 | .30 |
| CF-IDF | .46 | .31 |
| LDA | .34 | .31 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.03** | **.00** |
| HCF-IDF | | **.00** |

**Table 47: nDCGs and Post-hoc analysis for the factor *Decay Function* at LDA using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| Exponential decay | .34 | .31 |
| Sliding window | .30 | .32 |

b) Post-hoc analysis p-value

| | Exponential decay |
|---|---|
| Sliding window | **.03** |

presents the post-hoc analysis when Title is employed. We see that HCF-IDF outperforms others with significant differences. Table 25 shows the post-hoc analysis when All is chosen for the factor *Document Content*. There is no significant difference between CF-IDF and HCF-IDF. Table 26 shows the post-hoc analysis of the factor *Document Content* when CF-IDF is employed. It indicates that the strategies with CF-IDF and All significantly outperforms those with CF-IDF and Title.

**Table 48: ANOVA for *Profiling Method* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* at Title | 26.61 | . | **.00** |
| *Profiling Method* at All | 52.51 | . | **.00** |
| *Document Content* at CF-IDF | 30.81 | . | **.00** |
| *Document Content* at HCF-IDF | 0.31 | . | .58 |
| *Document Content* at LDA | 0.94 | . | .33 |

**The factor *Decay Function* × *Document Content*.** Table 51 shows the results of ANOVA regarding the factor *Decay Function* when a choice of the factor *Document Content* is fixed and vice versa. According to Table 51, there is a

**Table 49: nDCGs and Post-hoc analysis for the factor *Profiling Method* at Title using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .53 | .32 |
| CF-IDF | .41 | .29 |
| LDA | .33 | .31 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | **.00** | **.01** |
| HCF-IDF | | **.00** |

**Table 50: nDCG and Post-hoc analysis for the factor *Profiling Method* at All using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| CF-IDF | .56 | .33 |
| HCF-IDF | .53 | .34 |
| LDA | .30 | .32 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .18 | **.00** |
| HCF-IDF | | **.00** |

significant difference among the factor *Document Content*, when Sliding window is used. The nDCGs and post-hoc analysis of it are shown in Tables 52(a) and (b). It indicates that All significantly enhances the performance of the recommender system when Sliding window is used.

**Table 51: ANOVA for *Decay Function* × *Document Content* interaction**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Decay Function* at Title | 0.06 | . | .81 |
| *Decay Function* at All | 2.28 | . | .13 |
| *Document Content* at Sliding window | 9.96 | . | **.00** |
| *Document Content* at Exponential decay | 1.19 | . | .28 |

**Table 52: nDCGs and Post-hoc analysis for the factor *Document Content* at Sliding window using Shaffer's MSRB procedure.**

a) nDCGs

| Choice | M | SD |
|---|---|---|
| All | .48 | .36 |
| Title | .42 | .32 |

b) Post-hoc analysis p-value

| | All |
|---|---|
| Title | **.00** |

# F. CORRELATION WITH THE NUMBER OF TWEETS, THE NUMBER OF CONCEPTS, THE NUMBER OF CONCEPTS PER TWEET, AND THE PERCENTAGE OF TWEETS CONTAINING AT LEAST ONE CONCEPT

Table 53: Correlation coefficients of each of twelve strategies with the number of tweets, the number of concepts, the number of concepts per tweet, and the percentage of tweets containing at least one concept. We calculate correlation coefficients using Pearson product-moment correlation coefficient and Kendall rank correlation coefficient. Strategies are sorted by rankscores as shown in Table 3. In parentheses, p-values are given and marked in bold font, if < .05. Strategies are sorted by rankscores as shown in Table 3.

| | Strategy | | | # of tweets | | # of concepts | | # of concepts per tweet | |
|---|---|---|---|---|---|---|---|---|---|
| | **Profiling Method** | **Decay Function** | **Document Content** | **Pearson** | **Kendall** | **Pearson** | **Kendall** | **Pearson** | **Kendall** |
| 1. | CF-IDF | Sliding window | All | -.02 (.82) | -.01 (.84) | .00 (.98) | -.01 (.91) | -.06 (.54) | -.01 (.82) |
| 2. | HCF-IDF | Sliding window | All | -.02 (.85) | -.02 (.77) | .00 (.99) | .00 (.96) | .08 (.38) | .05 (.46) |
| 3. | HCF-IDF | Sliding window | Title | -.12 (.20) | -.07 (.26) | -.07 (.41) | -.04 (.51) | .04 (.67) | .05 (.38) |
| 4. | HCF-IDF | Exponential decay | Title | -.01 (.94) | -.03 (.68) | -.03 (.74) | -.04 (.57) | -.10 (.26) | -.09 (.15) |
| 5. | CF-IDF | Exponential decay | All | .02 (.83) | .02 (.80) | .00 (.93) | .01 (.92) | -.09 (.35) | -.08 (.23) |
| 6. | HCF-IDF | Exponential decay | All | .03 (.72) | .02 (.71) | .04 (.63) | .03 (.68) | -.06 (.49) | -.06 (.32) |
| 7. | CF-IDF | Exponential decay | Title | .17 (.07) | .08 (.19) | .13 (.14) | .07 (.25) | -.13 (.17) | -.05 (40) |
| 8. | CF-IDF | Sliding window | Title | .14 (.12) | .09 (.17) | .13 (.15) | .06 (.32) | -.10 (.26) | -.03 (.65) |
| 9. | LDA | Exponential decay | Title | .15 (.09) | .12 (**.05**) | .14 (.12) | .11 (.07) | .02 (.78) | .00 (.99) |
| 10. | LDA | Sliding window | Title | .18 (.05) | .08 (.24) | .14 (.12) | .06 (.34) | -.03 (.73) | -.01 (.89) |
| 11. | LDA | Exponential decay | All | -.10 (.28) | -.13 (**.03**) | -.07 (.46) | -.12 (**.05**) | .06 (.48) | .03 (.65) |
| 12. | LDA | Sliding window | All | -.06 (.50) | -.08 (.25) | -.04 (.66) | -.07 (.28) | -.03 (.76) | .01 (.93) |

| | Strategy | | | percentage of tweets with concepts | |
|---|---|---|---|---|---|
| | **Profiling Method** | **Decay Function** | **Document Content** | **Pearson** | **Kendall** |
| 1. | CF-IDF | Sliding window | All | -.02 (.79) | .01 (.92) |
| 2. | HCF-IDF | Sliding window | All | .11 (.22) | .06 (.32) |
| 3. | HCF-IDF | Sliding window | Title | .06 (.49) | .06 (.31) |
| 4. | HCF-IDF | Exponential decay | Title | -.11 (.24) | -.10 (.11) |
| 5. | CF-IDF | Exponential decay | All | -.08 (.38) | -.07 (.24) |
| 6. | HCF-IDF | Exponential decay | All | -.05 (.62) | -.07 (.24) |
| 7. | CF-IDF | Exponential decay | Title | -.14 (.12) | -.08 (.20) |
| 8. | CF-IDF | Sliding window | Title | -.13 (.16) | -.04 (.50) |
| 9. | LDA | Exponential decay | Title | .03 (.76) | .00 (.97) |
| 10. | LDA | Sliding window | Title | -.03 (.72) | -.01 (.82) |
| 11. | LDA | Exponential decay | All | .07 (.42) | .04 (.54) |
| 12. | LDA | Sliding window | All | -.03 (.71) | -.01 (.92) |

We examine whether the recommendation performance of the twelve strategies measured by the rankscore have correlations with the number of tweets, the number of concepts, the number of concepts per tweet, and the percentage of tweets containing at least one concept. We compute Pearson product-moment correlation coefficient and Kendall rank correlation coefficient between the rankscores of each of the twelve strategies and each of the number of tweets, the number of concepts, the number of concepts per tweet, and the percentage of tweets containing at least one concept. Table 53 provides the results.

## G. USER CLICK RATES

In this section, we describe the details of the analysis regarding click rates.

### G.1 Three-way repeated ANOVA for click rates

Table 54 shows average click rates of each strategy. Click rates differ depending on each strategy. In order to reveal effects on click rates from factors, we run three-way repeated ANOVA. First, we apply Mendoza's test [17] which showed violations of sphericity against the factors *Profiling Method* ($\chi^2(2) = 12.24$, $p = .00$), *Profiling Method $\times$ Decay Function* ($\chi^2(2) = 8.44$, $p = .01$), and *Profiling Method $\times$ Document Content* ($\chi^2(2) = 13.57$, $p = .00$). Thus, we run three-way repeated-measure ANOVA with a Greenhouse-Geisser correction of $\epsilon = .91$ for the factor *Profiling Method*,

$\epsilon = .94$ for the factor *Profiling Method $\times$ Decay Function*, and $\epsilon = .90$ for the factor *Profiling Method $\times$ Document Content*. Table 55 shows the results of an ANOVA with F-ratio, $\eta^2$ and p-value. $\eta^2$ indicates the effect size of each factor. The effect size is small when $\eta^2 > .02$, medium when $\eta^2 > .13$, and large when $\eta^2 > .26$.

Table 54: Average click rates on the PDF files. In parentheses, the standard deviations are shown. Strategies are sorted by rankscores as shown in Table 3.

| | Strategy | | | Click rate |
|---|---|---|---|---|
| | **Profiling Method** | **Decay Function** | **Content** | **Rate** |
| 1. | CF-IDF | Sliding window | All | 10.73% (24.73) |
| 2. | HCF-IDF | Sliding window | All | 10.08% (23.94) |
| 3. | HCF-IDF | Sliding window | Title | 9.11% (23.22) |
| 4. | HCF-IDF | Exponential decay | Title | 7.64% (17.28) |
| 5. | CF-IDF | Exponential decay | All | 9.11% (22.21) |
| 6. | HCF-IDF | Exponential decay | All | 8.29% (20.31) |
| 7. | CF-IDF | Exponential decay | Title | 8.94% (20.03) |
| 8. | CF-IDF | Sliding window | Title | 9.59% (22.81) |
| 9. | LDA | Exponential decay | Title | 4.23% (13.12) |
| 10. | LDA | Sliding window | Title | 4.72% (15.38) |
| 11. | LDA | Exponential decay | All | 9.27% (21.47) |
| 12. | LDA | Sliding window | All | 5.37% (16.41) |

According to 55, the factor *Profiling Method* makes significant difference on click rates. We applied again a post-

**Table 57: Correlation coefficients of each of twelve strategies between click rates and rankscores. We calculate correlation coefficients using Pearson product-moment correlation coefficient and Kendall rank correlation coefficient. Strategies are sorted by rankscores as shown in Table 3. In parentheses, p-values are given and marked in bold font, if < .05. Strategies are sorted by rankscores as shown in Table 3.**

| | Strategy | | | Correlation with click rates | |
|---|---|---|---|---|---|
| | Profiling Method | Decay Function | Con-tent | Pearson | Kendall |
| 1. | CF-IDF | Sliding window | All | .02 (.82) | -.03 (.68) |
| 2. | HCF-IDF | Sliding window | All | .07 (.42) | .05 (.47) |
| 3. | HCF-IDF | Sliding window | Title | -.06 (.52) | -.07 (.34) |
| 4. | HCF-IDF | Exponential decay | Title | .05 (.61) | .05 (.51) |
| 5. | CF-IDF | Exponential decay | All | .12 (.19) | .05 (.51) |
| 6. | HCF-IDF | Exponential decay | All | -.02 (.86) | -.02 (.78) |
| 7. | CF-IDF | Exponential decay | Title | .04 (.67) | .01 (.85) |
| 8. | CF-IDF | Sliding window | Title | .05 (.57) | .00 (.96) |
| 9. | LDA | Exponential decay | Title | .02 (.82) | -.02 (.80) |
| 10. | LDA | Sliding window | Title | -.01 (.90) | .02 (.79) |
| 11. | LDA | Exponential decay | All | .13 (.16) | .18 (**.02**) |
| 12. | LDA | Sliding window | All | .22 (**.01**) | .17 (**.03**) |

**Table 55: Three-way repeated-measure ANOVA for click rates with Greenhouse-Geisser correction with F-ratio, $\eta^2$ and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| *Profiling Method* | 5.23 | .04 | **.01** |
| *Decay Function* | 0.23 | .00 | .64 |
| *Document Content* | 2.60 | .02 | .11 |
| *Profiling Method × Decay Function* | 2.71 | .02 | .07 |
| *Profiling Method × Document Content* | 1.16 | .01 | .32 |
| *Decay Function × Document Content* | 1.00 | .01 | .10 |
| *Profiling Method × Decay Function × Document Content* | 2.31 | .02 | **.04** |

hoc using Shaffer's MSRB Procedure for the factor *Profiling Method*. Table 56 provides the result of the post-hoc analysis for the factor *Profiling Method*. It indicates that click rates of CF-IDF and HCF-IDF are significantly higher than those of LDA.

**Table 56: Click rates, Post-hoc analysis for the factor *Profiling Method* using Shaffer's MSRB procedure, and effect size.**

a) Click rates

| Choice | M | SD |
|---|---|---|
| HCF-IDF | .09 | .21 |
| CF-IDF | .10 | .22 |
| LDA | .06 | .17 |

b) Post-hoc analysis p-values

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .42 | **.01** |
| HCF-IDF | | **.03** |

c) Effect size using Cohen's d

| | HCF-IDF | LDA |
|---|---|---|
| CF-IDF | .05 | .20 |
| HCF-IDF | | .16 |

## G.2  Correlation with rankscores

In addition, we investigate the correlation between click rates and rankscores with respect to each strategy. We use Pearson product-moment correlation coefficient as well as Kendall rank correlation coefficient. Table 57 reports the result of the analysis.

## G.3  Precision for clicked PDF files

Table 58 shows the average precision for clicked PDF files. It is equal to the probability that a participant evaluate a recommended item "interesting" when clicking it. We observe that the precisions for the strategies HCF-IDF it are high even if recommendations are made with only titles.

**Table 58: Precision for clicked PDF files.**

| | Strategy | | | Click rate |
|---|---|---|---|---|
| | Profiling Method | Decay Function | Con-tent | |
| 1. | CF-IDF | Sliding window | All | 71.21% |
| 2. | HCF-IDF | Sliding window | All | 64.52% |
| 3. | HCF-IDF | Sliding window | Title | 55.36% |
| 4. | HCF-IDF | Exponential decay | Title | 68.09% |
| 5. | CF-IDF | Exponential decay | All | 71.43% |
| 6. | HCF-IDF | Exponential decay | All | 60.78% |
| 7. | CF-IDF | Exponential decay | Title | 47.27% |
| 8. | CF-IDF | Sliding window | Title | 49.15% |
| 9. | LDA | Exponential decay | Title | 42.31% |
| 10. | LDA | Sliding window | Title | 37.93% |
| 11. | LDA | Exponential decay | All | 43.86% |
| 12. | LDA | Sliding window | All | 48.48% |

## H.  DEMOGRAPHIC FACTOR

For each demographic factor, we first apply Mendoza's test. Subsequently, we conduct a mixed ANOVA test with one between subject factor (i.e., demographic factor) and one within subject factor (i.e., strategy), adjusted by Greenhouse-Geisser's epsilon. In addition, we provide the post-hoc analyses. However, we omit the post-hoc analysis of the factor strategy for the sake of brevity, because it is not so different from the result of the one-way repeated-measure ANOVA shown in Table 4.

**Gender.** Mendoza's test found a violation of sphericity in the factor strategy ($\chi^2(131) = 489.39$, $p = .00$) when comparing the male ($n = 96$) and female ($n = 27$) participants. Table 59 shows the result of an ANOVA with a Greenhouse-

Geisser correction of $\epsilon = .60$. According to Table 59, we see a significant difference between males and females. The post-hoc analysis is shown in Table 60. We observe female participants are more likely to evaluate recommended items as interesting than males. However, the factor gender does no make any difference about how each of the twelve strategies performs compared to the other strategies, because there is no significant difference in the factor gender $\times$ strategy.

**Table 59: Mixed ANOVA with a between subject factor *Gender* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Gender | 9.69 | .08 | **.00** |
| Strategy | 16.58 | .14 | **.00** |
| Gender $\times$ Strategy | 1.11 | .01 | .36 |

**Table 60: Rankscores and Post-hoc analysis for the factor *Gender* using Shaffer's MSRB procedure.**

a) Rankscores

| Degree | M | SD |
|---|---|---|
| male | .42 | .32 |
| female | .54 | .35 |

b) Post-hoc analysis p-values

| | female |
|---|---|
| male | **.00** |

**Age.** On average, participants are 32.90 years old (SD: 7.36). We divide participants into three groups for an ANOVA (group 1: participants who are $> 29$ years old ($n = 42$), group 2: $<= 29$ and $> 38$ years old ($n = 49$), group 3: $<= 38$ years old ($n = 32$)). We set those thresholds to make three groups have the almost same number of participants. Mendoza's test found a violation of sphericity in the strategies ($\chi^2(197) = 504.35$, $p = .00$). Table 61 shows the result of an ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. It indicates that the age of participants has no effect on the performance of the different strategies.

**Table 61: Mixed ANOVA with a between subject factor *Age* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Age | 2.06 | .03 | .13 |
| Strategy | 14.82 | .12 | **.00** |
| Age $\times$ Strategy | 0.69 | .01 | .77 |

**Highest Academic Degree.** We have participants whose highest academic degree is Bachelor ($n = 21$), Master ($n = 58$), PhD ($n = 32$), and lecturer/professor ($n = 12$). Mendoza's test found a violation of sphericity in the strategies when comparing the distributions among the factors ($\chi^2(263) = 653.03$, $p = .00$). Table 62 shows the result of applying an ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. According to Table 62, we see a significant difference among participants grouped by their highest academic

degrees. The post-hoc analysis is shown in Table 63. We observe that participants whose highest academic degree is Bachelor are more likely to evaluate recommended items as interesting than those whose highest academic degree is lecturer/professor.

**Table 62: Mixed ANOVA with a between subject factor *Highest Academic Degree* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Highest Academic Degree | 3.38 | .09 | **.02** |
| Strategy | 16.02 | .13 | **.00** |
| Highest Academic Degree $\times$ Strategy | 0.77 | .02 | .75 |

**Table 63: Rankscores and Post-hoc analysis for the factor *Highest Academic Degree* using Shaffer's MSRB procedure.**

a) Rankscores

| Degree | M | SD |
|---|---|---|
| Bachelor | .53 | .30 |
| Master | .43 | .33 |
| PhD | .44 | .33 |
| lecturer/professor | .32 | .28 |

b) Post-hoc analysis p-values

| | Master | PhD | lecturer/professor |
|---|---|---|---|
| Bachelor | .20 | .21 | **.01** |
| Master | | .72 | .21 |
| PhD | | | .09 |

**Major.** In the experiment, participants provide information about their majors. We manually classify participants into the two groups: participants whose major is economics ($n = 92$) and others ($n = 31$). Mendoza's test found a violation of sphericity in the strategies for these two groups ($\chi^2(131) = 466.90$, $p = .00$). Table 64 shows the result of an ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. It indicates that the major of participants has no effect on the performance of the different strategies.

**Table 64: Mixed ANOVA with a between subject factor *Major* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Major | 0.01 | .00 | .94 |
| Strategy | 16.41 | .14 | **.00** |
| Major $\times$ Strategy | 1.73 | .01 | .10 |

**Years of Profession.** On average, participants work in their fields for 7.85 years (SD: 6.85). We divide participants into three groups for an ANOVA (group 1: participants who work for $> 5$ years ($n = 44$), group 2: $<= 5$ and $> 10$ years ($n = 34$), group 3: $<= 10$ years ($n = 44$)). We set those thresholds to make three groups have the almost same number of participants. Mendoza's test found a violation of sphericity in the strategies ($\chi^2(197) = 541.67$, $p = .00$). Table 66 shows the result of an ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. It indicates that how long

participants have worked in their fields has no effect on the performance of the different strategies.

**Table 65: Mixed ANOVA with a between subject factor *Years of Profession* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Years of Profession | 0.13 | .00 | .88 |
| Strategy | 21.70 | .18 | **.00** |
| Years of Profession × Strategy | 0.80 | .01 | .66 |

**Employment Type.** We have participants who work in academia ($n = 83$) and industry ($n = 40$). Mendoza's test found a violation of sphericity in the strategies ($\chi^2(131) = 472.14$, $p = .00$). Table 66 shows the result of an ANOVA with a Greenhouse-Geisser correction of $\epsilon = .60$. It indicates that the employment type of participants has no effect on the performance of the different strategies.

**Table 66: Mixed ANOVA with a between subject factor *Employment Type* and a within subject factor *Strategy* Greenhouse-Geisser correction with F-ratio, effect size $\eta^2$, and p-value.**

| Factor | F | $\eta^2$ | p |
|---|---|---|---|
| Employment Type | 0.35 | .00 | .55 |
| Strategy | 18.05 | .15 | **.00** |
| Employment Type × Strategy | 0.97 | .01 | .45 |