

Influence of tweets and diversification on serendipitous research paper recommender systems

Chifumi Nishioka^{Corresp., Equal first author, 1}, Jörn Hauke^{Equal first author, 2}, Ansgar Scherp³

¹ Kyoto University, Kyoto, Japan

² Christian-Albrechts-Universität Kiel, Kiel, Germany

³ University of Essex, Colchester, United Kingdom

Corresponding Author: Chifumi Nishioka

Email address: nishioka.chifumi.2c@kyoto-u.ac.jp

In recent years, a large body of literature has accumulated around the topic of research paper recommender systems. However, since most studies have focused on the variable of accuracy, they have overlooked the serendipity of recommendations, which is an important determinant of user satisfaction. Serendipity is concerned with the novelty, relevance, and unexpectedness of recommendations, and so serendipitous items are considered those which positively surprise users. The purpose of this article was to examine two key research questions: firstly, whether a user's Tweets can assist in generating more serendipitous recommendations; and secondly, whether the diversification of a list of recommended items further improves serendipity. To investigate these issues, an online experiment was conducted in the domain of computer science with 22 subjects. The results indicate that a user's Tweets do not improve serendipity, but they can reflect recent research interests and are typically heterogeneous. Contrastingly, diversification was found to lead to a greater number of serendipitous research paper recommendations.

Influence of Tweets and Diversification on Serendipitous Research Paper Recommender Systems

Chifumi Nishioka¹, Jörn Hauke², and Ansgar Scherp³

¹Kyoto University Library, Kyoto, Japan

²Kiel University, Kiel, Germany

³University of Essex, Colchester, UK

Corresponding author:

Chifumi Nishioka¹

Email address: nishioka.chifumi.2c@kyoto-u.ac.jp

ABSTRACT

In recent years, a large body of literature has accumulated around the topic of research paper recommender systems. However, since most studies have focused on the variable of accuracy, they have overlooked the serendipity of recommendations, which is an important determinant of user satisfaction. Serendipity is concerned with the novelty, relevance, and unexpectedness of recommendations, and so serendipitous items are considered those which positively surprise users. The purpose of this article was to examine two key research questions: firstly, whether a user's Tweets can assist in generating more serendipitous recommendations; and secondly, whether the diversification of a list of recommended items further improves serendipity. To investigate these issues, an online experiment was conducted in the domain of computer science with 22 subjects. The results indicate that a user's Tweets do not improve serendipity, but they can reflect recent research interests and are typically heterogeneous. Contrastingly, diversification was found to lead to a greater number of serendipitous research paper recommendations.

INTRODUCTION

To help researchers overcome the problem of information overload, various studies have developed recommender systems (Beel et al., 2016; Bai et al., 2019). Recommendations are generated based on considerations such as a user's own papers (Sugiyama and Kan, 2010) or the papers a user has accessed in the past (Nascimento et al., 2011). Most previous studies have focused only on improving the accuracy of recommendations, one example of which is normalised discounted cumulative gain (nDCG). However, several studies on recommender systems conducted in other domains (e.g., movies) have drawn attention to the fact that there are important aspects other than accuracy (McNee et al., 2006; Herlocker et al., 2004). One of these aspects is *serendipity*, which is concerned with the novelty of recommendations and the degree to which recommendations positively surprise users (Ge et al., 2010).

In this article, we study a research paper recommender system focusing on serendipity. Sugiyama and Kan (2015) investigated serendipitous research paper recommendations, focusing on the influence of dissimilar users and the co-author network on recommendation performance. In contrast, this study investigates the following research questions:

- (RQ1) Do a user's Tweets generate serendipitous recommendations?
- (RQ2) Is it possible to improve a recommendation list's serendipity through diversification?

We run an online experiment to facilitate an empirical investigation of these two research questions using three factors. For RQ1, we employ the factor *User Profile Source*, where we compare the two sources of user profiles: firstly, a user's own papers; and secondly, a user's Tweets. The user's own papers are a feature of existing recommender systems, as evidenced by the work conducted by Sugiyama and Kan

43 (2015) and Google Scholar.¹ In this study, we assume that the user's Tweets produce recommendations
44 that cannot be generated based on papers, since researchers Tweet about recent developments and interests
45 that are yet not reflected in their papers (e.g., what they found interesting at a conference or in their social
46 network) (Letierce et al., 2010). In the domain of economics, recommendations based on a user's Tweets
47 received a precision of 60%, which is fairly high (Nishioka and Scherp, 2016).

48 In addition, we analyse the factor *Text Mining Method*, which applies different methods of candidate
49 items (i.e., research papers) for computing profiles, as well as user profiles comprising different content
50 (i.e., Tweets or previous papers).

51 As text mining methods, we compare TF-IDF (Salton and Buckley, 1988) with two of its recent
52 extensions, namely CF-IDF (Goossen et al., 2011) and HCF-IDF (Nishioka et al., 2015). Both have been
53 associated with high levels of performance in recommendation tasks (Goossen et al., 2011; Nishioka
54 et al., 2015). We introduce this factor because text mining methods can have a substantial influence on
55 generating recommendations. For RQ2, we introduce the factor *Ranking Method*, where we compare
56 two ranking methods: firstly, classical cosine similarity; and secondly, the established diversification
57 algorithm IA-Select (Agrawal et al., 2009). Cosine similarity has been widely used in recommender
58 systems (Lops et al., 2011), while IA-Select ranks candidate items with the objective of diversifying
59 recommendations in a list. Since it broadens the coverage of topics in a list, we assume that IA-Select
60 delivers more serendipitous recommendations compared to cosine similarity.

61 Along with the three factors *User Profile Source*, *Text Mining Method*, and *Ranking Method*, we
62 conduct an online experiment in which 22 subjects receive research paper recommendations in the field of
63 computer science. The results reveal that a user's Tweets do not improve the serendipity of recommender
64 systems. On the other hand, we confirm that the diversification of a recommendation list by IA-Select
65 delivers more serendipitous recommendations to users.

66 The remainder of the paper is organised as follows: firstly, we describe related studies; in turn, we
67 describe the recommender system and the experimental factors and evaluation setup; and finally, before
68 concluding the article, we report on and discuss the experimental results.

69 RELATED WORK

70 Over the last decade, many studies have developed research paper recommender systems (Beel et al.,
71 2016; Bai et al., 2019). According to Beel et al. (2016), more than half of these studies (55%) have applied
72 a content-based approach. Collaborative filtering was applied by 18% and graph-based recommendations,
73 utilising citation networks or co-authorship networks, were applied by 16%. Other researches have
74 employed stereotyping, item-centric recommendations, and hybrid recommendations. In this article, we
75 employ a content-based approach.

76 **Clarifying the notion of serendipity** Most existing studies have evaluated recommender systems by
77 focusing on measures of accuracy, including precision, mean reciprocal rank (MRR), and normalised
78 discounted cumulative gain (nDCG). However, studies that have addressed recommender systems in
79 other domains (e.g., movies) argue that there are important considerations other than accuracy (McNee
80 et al., 2006; Herlocker et al., 2004). One of these considerations is *serendipity*, which is a term that has
81 been defined differently in the literature in the context of recommender systems. For instance, Kotkov
82 et al. (2016) defined serendipity as “a property that indicates how good a recommender system is at
83 suggesting serendipitous items that are relevant, novel, and unexpected for a particular user.” Similarly,
84 Herlocker et al. (2004) defined serendipity as measure of the extent to which the recommended items
85 are both attractive and surprising to the users. Other researchers have offered comparable definitions of
86 serendipity (Shani and Gunawardana, 2011).

87 According to Ge et al. (2010), it is important to recognise two important aspects of serendipity: firstly,
88 a serendipitous item should be unknown to the user and, moreover, should not be expected; and secondly,
89 the item should be interesting, relevant, and useful to the user. Taking these two aspects into account, Ge
90 et al. (2010) proposed a quantitative metric to evaluate the degree to which recommender systems are
91 effective at generating serendipitous recommendations.

92 **Use of social media for serendipitous recommendations** In previous studies addressing content-
93 based research paper recommender systems (Beel et al., 2016; Bai et al., 2019), the authors calculated

¹<https://scholar.google.co.jp/>

Table 1. Experimental factors and design choices

Factor	Possible Design Choices		
<i>User Profile Source</i>	Twitter		Own Papers
<i>Text Mining Method</i>	TF-IDF	CF-IDF	HCF-IDF
<i>Ranking Method</i>	Cosine Similarity		IA-Select

94 recommendations based on a user's own papers (Sugiyama and Kan, 2010) or papers a user has read
 95 in the past (Nascimento et al., 2011). In other domains, several studies have developed content-based
 96 recommender systems (Chen et al., 2010; Orlandi et al., 2012; Shen et al., 2013) that utilise data from a
 97 user's social media accounts, including Twitter and Facebook. Another study proposed research paper
 98 recommendations based on a user's Tweets, which received a relatively high precision of 60% (Nishioka
 99 and Scherp, 2016). However, we hypothesise that because researchers Tweet about recent developments
 100 and interests that are not yet reflected in their papers (Letierce et al., 2010), a user's Tweets will deliver
 101 recommendations that are not generated based on papers.

102 In the context of research paper recommender systems, Sugiyama and Kan (2015) investigated
 103 serendipitous research paper recommendations focusing on the influence of dissimilar users and the co-
 104 author network on the recommendation performance. However, the researchers evaluated their approaches
 105 using accuracy-focused evaluation metrics such as nDCG and MRR. In contrast, this article investigates
 106 serendipitous research paper recommendations from the perspective of Tweets and diversification.

107 **Use of diversification for serendipitous recommendations** As discussed above, novelty is a key con-
 108 cept for serendipity (Ge et al., 2010). One approach that can be used to generate novel recommendations
 109 relates to diversification (Agrawal et al., 2009). This is because diversification leads to the creation of
 110 recommendation lists that include dissimilar items, meaning that users have an opportunity to encounter
 111 items they are unfamiliar with. IA-Select (Agrawal et al., 2009) has been used in the past as a solid
 112 baseline for diversifying lists of recommendations (Vargas and Castells, 2011; Vargas et al., 2011; Wu
 113 et al., 2018). Additionally, MMR (Carbonell and Goldstein, 1998) is a well-known diversification method.
 114 However, since the experimental research conducted by Vargas and Castells (2011) shows that IA-Select
 115 performs better, we employ it in this study's experiment.

116 EXPERIMENTAL FACTORS

117 In this article, we build a content-based recommender system along with the three factors *User Profile*
 118 *Source*, *Text Mining Method*, and *Ranking Method*. It works as follows:

- 119 1. Candidate items of the recommender system (i.e., research papers) are processed by one of the
 120 text mining methods, and paper profiles are generated. A candidate item and a set of candidate
 121 items are referred as d and D , respectively. d 's paper profile P_d is represented by a set of features
 122 F and their weights. Depending on text mining methods, a feature f is either a textual term or a
 123 concept. Formally, paper profiles are described as: $P_d = \{(f, w(f, d)) \mid \forall f \in F\}$. The weighting
 124 function w returns a weight of a feature f for data source I_u . This weight identifies the importance
 125 of the feature f for the user u .
- 126 2. A user profile is generated based on the user profile source (i.e., Tweets or own papers) using the
 127 same text mining method, which is applied to generate paper profiles. I_u is a set of data items i of a
 128 user u . In this article, I_u is either a set of a user's Tweets or a set of a user's own papers. u 's user
 129 profile P_u is represented in a way that it is comparable to P_u as: $P_u = \{(f, w(f, I_u)) \mid \forall f \in F\}$.
- 130 3. One of the ranking methods determines the order of recommended papers.

131 The experimental design is illustrated in Table 1, where each cell is a possible design choice in each factor.

132 In this section, we first provide a detailed account of the factor *User Profile Source*. In turn, we
 133 show three of the different text mining methods that were applied in the experiment. Finally, we note the
 134 details of the factor *Ranking Method*, which examines whether diversification improves the serendipity of
 135 recommendations.

136 User Profile Source

137 In this factor, we compare the following two data sources that are used to build a user profile.

- 138 • **Research papers:** The research papers written by a user are used as a baseline. This approach is
139 motivated by previous studies that have investigated research paper recommender systems, including
140 Sugiyama and Kan (2010) and Google Scholar.
- 141 • **Twitter:** In contrast to the user’s papers, we assume that using Tweets leads to more serendipitous
142 recommendations. It is common practice among researchers to Tweet about their professional
143 interests (Letierce et al., 2010). Therefore, Tweets can be used to build a user profile in the
144 context of a research paper recommender system. We hypothesise that a user’s Tweets improve
145 the serendipitous nature of recommendations because researchers are likely to Tweet about recent
146 interests and information (e.g., from social networks) that are not yet reflected in their papers.

147 Text Mining Method

148 For each of the two data sources (i.e., the user’s own papers or their Tweets) and the candidate items,
149 we apply a text mining method using one of three text mining methods. Specifically, we compare
150 three methods, namely TF-IDF (Salton and Buckley, 1988), CF-IDF (Goossen et al., 2011), and HCF-
151 IDF (Nishioka et al., 2015), to build paper profiles and a user profile. This factor was introduced because
152 the effectiveness of each text mining method is informed by the type of content that will be analysed (e.g.,
153 Tweets or research papers). For each method, a weighting function w is defined. This weighting function
154 assigns a specific weight to each feature f , which is a term in TF-IDF and a semantic concept in CF-IDF
155 and HCF-IDF.

- **TF-IDF:** Since TF-IDF is frequently used in recommender systems as a baseline (Goossen et al.,
2011), we also use it in this study. Terms are lemmatised and stop words are removed.² In addition,
terms with fewer than three characters are filtered out due to ambiguity. After pre-processing texts,
TF-IDF is computed as:

$$w_{tf-idf}(w,t) = tf(w,t) \cdot \log \frac{|D|}{|\{w \in d : d \in D\}|}. \quad (1)$$

156 tf returns the frequency of a term w in a text t . A text t is either a user profile source I_u or
157 candidate item d . The term frequency acts under the assumption that more frequent terms are
158 more important (Salton and Buckley, 1988). The second term of the equation presents the inverse
159 document frequency, which measures the relative importance of a term w in a corpus D (i.e., a set
160 of candidate items).

- **CF-IDF:** Concept frequency inverse document frequency (CF-IDF) (Goossen et al., 2011) is
161 an extension of TF-IDF, which replaces terms with semantic concepts from a knowledge base.
162 The use of a knowledge base decreases noise in profiles (Abel et al., 2011). In addition, since a
163 knowledge base can store multiple labels for a concept, the method directly supports synonyms.
164 For example, the concept “recommender systems” of the ACM Computing Classification Systems
165 (ACM CCS) has multiple labels, including “recommendation systems”, “recommendation engine”,
166 and “recommendation platforms”.
167

The weighting function w for CF-IDF is defined as:

$$w_{cf-idf}(a,t) = cf(a,t) \cdot \log \frac{|D|}{|\{a \in d : d \in D\}|}. \quad (2)$$

168 cf returns the frequency of a semantic concept a in a text t . The second term presents the IDF,
169 which measures the relative importance of a semantic concept a in a corpus D .

- **HCF-IDF:** Finally, we apply hierarchical concept frequency inverse document frequency (HCF-
IDF) (Nishioka et al., 2015), which is an extension of CF-IDF. HCF-IDF applies a propagation
function (Kapanipathi et al., 2014) over a hierarchical structure of a knowledge base to assign a
weight to concepts at higher levels. In this way, it identifies concepts that are not mentioned in a

²<http://www.nltk.org/book/ch02.html>

text but which are highly relevant. HCF-IDF calculates the weight of a semantic concept a in a text t as follows:

$$w_{hcf-idf}(a,t) = BL(a,t) \cdot \log \frac{|D|}{|\{d \in D : a \in d\}|}. \quad (3)$$

$BL(a,t)$ is the BellLog propagation function (Kapanipathi et al., 2014), which is defined as:

$$BL(a,t) = cf(a,t) + FL(a) \cdot \sum_{a_j \in pc(a)} BL(a_j,t), \quad (4)$$

170 where $cf(a,t)$ is a frequency of a concept a in a text t , and $FL(a) = \frac{1}{\log_{10}(nodes(h(a)+1))}$. The
 171 propagation function underlies the assumption that, in human memory, information is represented
 172 through associations or semantic networks (Collins and Loftus, 1975). The function $h(a)$ returns
 173 the level, where a concept a is located in the knowledge base. Additionally, $nodes$ provides the
 174 number of concepts at a given level in a knowledge base, and $pc(a)$ returns all parent concepts of a
 175 concept a . In this study, we employ HCF-IDF since it has been shown to work effectively for short
 176 pieces of text, including Tweets (Nishioka and Scherp, 2016), in the domain of economics.

177 Ranking Method

178 Finally, we rank all the candidate items to determine which items should be recommended to a user. In this
 179 factor, we compare two ranking methods: cosine similarity and diversification with IA-Select (Agrawal
 180 et al., 2009).

- 181 • **Cosine similarity:** As a baseline, we employ a cosine similarity, which has been widely used
 182 in content-based recommender systems. The top- k items with largest cosine similarities are
 183 recommended.
- 184 • **IA-Select:** Following this, we employ IA-Select (Agrawal et al., 2009) to deliver serendipitous
 185 recommendations. IA-Select was originally introduced for information retrieval, but it is also
 186 used in recommender systems to improve serendipity (Vargas et al., 2012). This use case stems
 187 from the algorithm's ability to diversify recommendations in a list, which relies on the avoidance
 188 of recommending similar items (e.g., research papers) together. The basic idea of IA-Select is
 189 that, for those features of a user profile that have been covered by papers already selected for
 190 recommendation, the weights are lowered in an iterative manner. At the outset, the algorithm
 191 computes cosine similarities between a user and each candidate item. In turn, IA-Select adds
 192 the item with the largest cosine similarity to the recommendation list. After selecting the item,
 193 IA-Select decreases the weights of features covered by the selected item in the user profile. These
 194 steps are repeated until k recommendations are determined.
 195 For example, recommendations for the user profile $P_u = ((f_1, 0.1), (f_2, 0.9))$ will contain mostly
 196 those documents that include feature f_2 . However, with IA-Select, the f_2 score is decremented
 197 iteratively in the event that documents contain the f_2 feature. Thus, the probability increases that
 198 documents covering the f_1 feature are included in the list of recommended items.

199 Overall, the three factors with the design choices described above result in $2 \times 3 \times 2 = 12$ available
 200 strategies. The evaluation procedure used to compare the strategies is provided below.

201 EVALUATION

202 To address the two research questions with the three experimental factors described in the previous section,
 203 we conduct an online experiment with 22 subjects. The experiment is based in the field of computer
 204 science, in which an open access culture to research papers exists, and Twitter is chosen as the focal point
 205 because it is an established means by which researchers disseminate their works. The experimental design
 206 adopted in this study is consistent with previous studies (Nishioka and Scherp, 2016; Chen et al., 2010).

207 In this section, the experimental design is described, after which an account of the utilised datasets
 208 (i.e., a corpus of research papers and a knowledge graph of text mining methods) is given. Following this,
 209 descriptive statistics are presented for the research subjects, and finally, the serendipity score is stated.
 210 The purpose of the serendipity score is to evaluate the degree to which each recommender strategy is
 211 effective in generating serendipitous recommendations.

Recommendation (1/12)

Please evaluate the following randomized list of the top five publications "interesting" or "not interesting".
Click on a title to see its abstract in a new window.

Please Note: The list might contain publications which you have already seen, since the system makes recommendations under different, independent strategies.

-	Robin J. Wilson, "Stamps, computing on", Encyclopedia of Computer Science, 2003	<input type="radio"/> interesting <input type="radio"/> not interesting
-	Sven Uebelacker, Susanne Quiel, "The Social Engineering Personality Framework", STAST '14 Proceedings of the 2014 Workshop on Socio-Technical Aspects in Security and Trust, 2014	<input type="radio"/> interesting <input type="radio"/> not interesting
-	Katharina Krombholz, Heidelinde Hobel, Markus Huber, Edgar Weippl, "Social engineering attacks on the knowledge worker", Proceedings of the 6th International Conference on Security of Information and Networks, 2013	<input type="radio"/> interesting <input type="radio"/> not interesting
-	Michael Workman, "Gaining Access with Social Engineering: An Empirical Study of the Threat", Information Systems Security, 2007	<input type="radio"/> interesting <input type="radio"/> not interesting
-	Anker Helms Jørgensen, Brad A. Myers, "User interface history", CHI '08 Extended Abstracts on Human Factors in Computing Systems, 2008	<input type="radio"/> interesting <input type="radio"/> not interesting

Figure 1. Screenshot of the evaluation page. Each subject rated an item as either “interesting” or “not interesting” based on their research interests.

212 Procedure

213 We implemented a web application that enabled the subjects ($n = 22$) to evaluate the twelve recommenda-
 214 tion strategies described above. First, subjects started on the welcome page, which asked for their consent
 215 to collect their data. Thereafter, the subjects were asked to input their Twitter handle and their name,
 216 as recorded in DBLP Persons.³ Based on the user’s name, we retrieved a list of their research papers
 217 and obtained the content of the papers by mapping them to the ACM-Citation-Network V8 dataset (see
 218 below). The top 5 recommendations were computed for each strategy, as shown in Figure 1. Thus, each
 219 subject evaluated $5 \cdot 12 = 60$ items as “interesting” or “not interesting” based on the perceived relevance
 220 to their research interests.

221 A binary evaluation was chosen to minimise the effort of the rating process, consistent with several
 222 previous studies (Nishioka and Scherp, 2016; Chen et al., 2010). As shown in Figure 1, the recommended
 223 items were displayed with bibliographic information such as the authors, title, year, and venue. Finally,
 224 the subjects were provided with the opportunity to access and read the research paper directly by clicking
 225 on a link. In order to avoid bias, the sequence in which the twelve strategies appeared was randomised
 226 for each subject. At the same time, the list of the top 5 items for each strategy was also randomised to
 227 avoid the well-documented phenomenon of ranking bias (Bostandjiev et al., 2012; Chen et al., 2010). The
 228 subjects were informed about the randomised order of the strategies and items on the evaluation page.

229 The actual ranks of the recommended items, as well as their position on the evaluation page, were
 230 stored in a database for later analyses. After evaluating all strategies, the subjects were asked to complete
 231 a questionnaire focusing on demographic information (e.g., age, profession, highest academic degree, and
 232 current employment status). Finally, an opportunity was provided for the subjects to provide qualitative
 233 feedback.

234 Datasets

235 The candidate items for the experiment were computer science articles drawn from a large dataset of
 236 research papers. To analyse and extract semantic concepts from the research papers and Tweets, an
 237 external computer science knowledge base was used. This section describes the research papers and
 238 knowledge graphs used for the experiment.

239 **Research papers** Since the experiment recommended research papers from the field of computer
 240 science, a corpus of research papers and a knowledge base from the same field were used. The ACM
 241 citation network V8 dataset⁴, provided by ArnetMiner (Tang et al., 2008), was used as the corpus of

³<https://dblp.uni-trier.de/pers/>

⁴<https://lfs.aminer.org/lab-datasets/citation/citation-acm-v8.txt.tgz>

242 research papers. From the dataset, 1,669,237 of the available 2,381,688 research papers were included
243 that had a title, author, year of publication, venue, and abstract. Titles and abstracts were used to generate
244 paper profiles.

245 **Knowledge graph** The ACM Computing Classification System (CCS) was used as the knowledge
246 graph for CF-IDF and HCF-IDF.⁵ The knowledge graph, which is freely available, is characterised by its
247 focus on computer science, as well as its hierarchical structure. It consists of 2,299 concepts and 9,054
248 labels, which are organized on six levels. On average, a concept is represented by 3.94 labels (SD: 3.49).

249 For the text mining methods (i.e., CF-IDF and HCF-IDF), we extracted concepts from each user's
250 Tweets and research papers by matching the text with the labels of the concepts in the knowledge graph.
251 As such, we applied what is known in the literature as the gazetteer-based approach. Before processing,
252 we lemmatised both the Tweets and research papers using Stanford Core NLP⁶, and stop words were
253 removed. Regarding Tweets, which often contain hashtags to indicate topics and user mentions, only the
254 # and @ symbols were removed from the Tweets. This decision stemmed from an observation made by
255 Feng and Wang (2014), namely that the combination of Tweets' texts with hashtags and user mentions
256 results in the optimal recommendation performance.

257 Subjects

258 Overall, 22 subjects were recruited through Twitter and mailing lists. 20 were male and 2 were female,
259 and the average age was 36.45 years old (SD: 5.55). Several of the subjects held master's degrees ($n = 2$),
260 while the others held a PhD ($n = 13$) or were lecturers or professors ($n = 7$). In terms of the subjects'
261 employment status, 19 were working in academia and three in industry. On average, the subjects published
262 1256.97 Tweets (SD: 1155.8), with the minimum value being 26 and the maximum value being 3158.

263 An average of 4968.03 terms (SD: 4641.76) was extracted from the Tweets, along with an average of
264 297.91 concepts (SD: 271.88). Thus, on average, 3.95 (SD: 0.54) terms and 0.24 concepts (SD: 0.10) were
265 included per Tweet. Regarding the use of research papers for user profiling, the subjects had published an
266 average of 11.41 papers (SD: 13.53). On average, 687.68 terms (SD: 842.52) and 80.23 concepts (SD:
267 107.73) were identified in their research papers. This led to 60.27 terms (SD: 18.95) and 5.77 concepts
268 (SD: 3.59) per paper.

269 Subjects needed 39 seconds (SD: 43 seconds) on average to evaluate all five recommended items per
270 strategy. Thus, the average length of time needed to complete the experiment was 468 seconds. It is worth
271 noting that this time does not include reading the instructions on the welcome page, inputting the Twitter
272 handle and DBLP record, and completing the questionnaire.

273 Evaluation Metric

To evaluate the serendipity of recommendations, we used the serendipity score (SRDP) (Ge et al., 2010).
This evaluation metric takes into account both the unexpectedness and usefulness of candidate items,
which is defined as:

$$274 \quad SRDP = \sum_{d \in UE} \frac{rate(d)}{|UE|}. \quad (5)$$

275 UE denotes a set of unexpected items that are recommended to a user. An item is regarded as unexpected
276 if it is not included in a recommendation list computed by the primitive strategy. We used the strategy
277 Own Papers \times TF-IDF \times Cosine Similarity as a primitive strategy since it is a combination of baselines.
278 The function $rate(d)$ returns an evaluation rate of an item d given by a subject. As such, if a subject
evaluated an item as "interesting", the function would return 1, otherwise 0.

279 RESULTS

280 The purpose of this section is to present the results of the experiment. At the outset, the quantitative
281 analysis is examined, which shows the optimal strategy in terms of SRDP. In turn, the impact of each of
282 the three experimental factors is analysed.

⁵<https://www.acm.org/publications/class-2012>

⁶<https://stanfordnlp.github.io/CoreNLP/>

283 Comparison of the Twelve Strategies

284 The results of the twelve strategies in terms of their SRDP values are presented in Table 2. As previously
 285 noted, this study drew on Own Papers \times TF-IDF \times Cosine Similarity as a primitive strategy. Thus, for
 286 this particular strategy, the mean and standard deviation are .00.

287 The purpose of an analysis of variance (ANOVA) is to detect significant differences between variables.
 288 Therefore, in this study, ANOVA was used to identify whether any of the strategies were significantly
 289 different. The significance level was set to $\alpha = .05$. Mauchly's test revealed a violation of sphericity
 290 ($\chi^2(54) = 80.912, p = .01$), which could lead to positively biased F-statistics and, consequently, an
 291 increase in the risk of false positives. Therefore, a Greenhouse-Geisser correction with $\epsilon = 0.58$ was
 292 applied.

293 The results of the ANOVA test revealed that significant differences existed between the strategies
 294 ($F(5.85, 122.75) = 3.51, p = .00$). Therefore, Shaffer's modified sequentially rejective Bonferroni
 295 procedure was undertaken to compute the pairwise differences between the strategies (Shaffer, 1986). We
 296 observed significant differences between the primitive strategy and one of the other strategies.

Table 2. SRDP and the number of unexpected items included in the twelve strategies. The values are ordered by SRDP. M and SD denote mean and standard deviation, respectively.

	Strategy			SRDP	UE
	Text Mining Method	Profiling Source	Ranking Method	M (SD)	M (SD)
1.	TF-IDF	Own Papers	IA-Select	.45 (.38)	2.95 (1.05)
2.	CF-IDF	Twitter	CosSim	.39 (.31)	4.91 (0.29)
3.	TF-IDF	Twitter	IA-Select	.36 (.29)	4.91 (0.43)
4.	CF-IDF	Twitter	IA-Select	.31 (.22)	4.95 (0.21)
5.	CF-IDF	Own Papers	CosSim	.26 (.28)	4.91 (0.29)
6.	CF-IDF	Own Papers	IA-Select	.25 (.28)	4.91 (0.29)
7.	HCF-IDF	Own Papers	IA-Select	.24 (.22)	4.95 (0.21)
8.	HCF-IDF	Twitter	CosSim	.22 (.28)	5.00 (0.00)
9.	TF-IDF	Twitter	CosSim	.20 (.24)	4.95 (0.21)
10.	HCF-IDF	Twitter	IA-Select	.18 (.21)	5.00 (0.00)
11.	HCF-IDF	Own Papers	CosSim	.16 (.18)	5.00 (0.00)
12.	TF-IDF	Own Papers	CosSim	.00 (.00)	0.00 (0.00)

297 Impact of Experimental Factors

298 In order to analyse the impact of each experimental factor, a three-way repeated measures ANOVA was
 299 conducted. The Mendoza test identified violations of sphericity for the following factors: firstly, *User*
 300 *Profile Source* \times *Text Mining Method* \times *Ranking Method* ($\chi^2(65) = 101.83, p = .0039$); and secondly,
 301 *Text Mining Method* \times *Ranking Method* ($\chi^2(2) = 12.01, p = .0025$) (Mendoza, 1980). Thus, a three-way
 302 repeated measures ANOVA was applied with a Greenhouse-Geisser correction of $\epsilon = .54$ for the factors
 303 *User Profile Source* \times *Text Mining Method* \times *Ranking Method* and $\epsilon = .69$ for the factor *Text Mining*
 304 *Method* \times *Ranking Method*. Table 3 shows the results with the F-Ratio, effect size η^2 , and p -value.

305 Regarding the single factors, *Ranking Method* had the largest impact on SRDP, as the effect size η^2
 306 indicates. For all the factors with significant differences, we applied a post-hoc analysis using Shaffer's
 307 MSRB procedure. The results of the post-hoc analysis revealed that the strategies using IA-Select resulted
 308 in higher SRDP values when compared to those using cosine similarity. In addition, we observed a
 309 significant difference in the factors *User Profile Source* \times *Ranking Method* and *Text Mining Method* \times
 310 *Ranking Method*. For both factors, post-hoc analyses revealed significant differences when a baseline was
 311 used in either of the two factors. When a baseline was used in one factor, $|UE|$ became small unless a
 312 method other than a baseline was used in the other factor.

313 DISCUSSION

314 This section discusses the study's results in relation to the two research questions. In turn, we review the re-
 315 sults for the *Text Mining Method* factor, which was found to have the largest influence on recommendation

Table 3. Three-way repeated measures ANOVA for SRDP with Greenhouse-Geisser correction and F-ratio, effect size η^2 , and p-value.

Factor	F	η^2	p
<i>User Profile Source</i>	2.21	.11	.15
<i>Text Mining Method</i>	3.02	.14	.06
<i>Ranking Method</i>	14.06	.67	.00
<i>User Profile Source</i> \times <i>Text Mining Method</i>	0.98	.05	.38
<i>User Profile Source</i> \times <i>Ranking Method</i>	18.20	.87	.00
<i>Text Mining Method</i> \times <i>Ranking Method</i>	17.80	.85	.00
<i>User Profile Source</i> \times <i>Text Mining M.</i> \times <i>Ranking M.</i>	2.39	.11	.11

316 performance among the three factors.

317 **RQ1** : Do a user's Tweets generate serendipitous recommendations?

318 Regarding RQ1, the results of the experiment indicate that a user's Tweets do not improve the
 319 serendipity of recommendations. As shown in the rightmost column of Table 2, Tweets deliver unexpected
 320 recommendations to users, but only a small fraction of these are interesting to the users. One way to
 321 account for this result is by drawing attention to the high probability that the users employed their Twitter
 322 accounts for purposes other than professional, research-related ones. In particular, the users are likely
 323 to have used their Twitter accounts to express private interests. We presume that taking private interests
 324 into consideration delivers serendipitous recommendations. This is because the recommender system
 325 will then suggest research papers that include both professional interests and private interests, and which
 326 are thus likely to be serendipitous. In the future, it may be helpful to introduce explanation interfaces
 327 for recommender systems (Herlocker et al., 2000; Tintarev and Masthoff, 2007). The purpose of these
 328 explanation interfaces is to show why a specific item is being recommended to users, thereby enabling
 329 users to find a connection between a recommended paper and their interests.

330 **RQ2** : Is it possible to improve a recommendation list's serendipity through diversification?

331 In terms of RQ2, the results indicate that the diversification of a recommendation list using the
 332 IA-Select algorithm delivers serendipitous recommendations. This confirms results published elsewhere
 333 in the literature, which have found that IA-Select improves serendipity. Additionally, the iterative decrease
 334 of covered interests was associated with greater variety in recommender systems for scientific publications.
 335 Furthermore, the experiment demonstrated that diversified recommendations are likely to be associated
 336 with greater utility for users.

337 **Text Mining Methods** Among the three factors, the *Text Mining Method* factor was associated with
 338 the most substantial impact on recommender system performance. In contrast to observations made in
 339 previous literature (Goossen et al., 2011; Nishioka and Scherp, 2016), CF-IDF and HCF-IDF did not yield
 340 effective results. It is worth emphasising that this result could have been influenced by the quality of the
 341 knowledge graph used in this study (i.e., ACM CCS), particularly in view of the fact that the performance
 342 of many text mining methods is directly informed by the quality of the knowledge graph (Nishioka et al.,
 343 2015).

344 Another way to account for the poor outcomes relates to the variable of the knowledge graphs' age.
 345 In particular, ACM CCS has not been updated since 2012, despite the fact that computer science is
 346 a rapidly changing field of inquiry. Furthermore, relatively few concepts and labels were included in
 347 the knowledge base, which contrasts with the large number included in the knowledge graphs used in
 348 previous studies. For example, the STW Thesaurus for Economics used 6335 concepts and 37,773 labels,
 349 respectively (Nishioka and Scherp, 2016). Hence, the number of concepts and labels could have influenced
 350 the quality of the knowledge graph and, in turn, the recommender system's performance.

351 In addition, while a previous study that used HCF-IDF (Nishioka and Scherp, 2016) only drew on the
 352 titles of research papers, our study used both titles and abstracts to construct paper profiles and user profiles
 353 when a user's own papers were selected as the user profile source. Furthermore, since our study used
 354 sufficient information when mining research papers, we did not observe any differences among TF-IDF,

355 CF-IDF, and HCF-IDF, which can include related concepts. Finally, as with any empirical experiment,
356 data triangulation is needed before generalising any of the conclusions drawn in this paper. Therefore,
357 further studies of recommender systems in other domains and similar settings should be conducted.

358 CONCLUSION

359 The purpose of this study's online experiment was to determine whether Tweets and the IA-Select
360 algorithm have the capability to deliver serendipitous research paper recommendations. The results
361 revealed that Tweets do not improve the serendipity of recommendations, but IA-Select does. We
362 anticipate that this insight will contribute to the development of future recommender systems, principally
363 because service providers and platform administrators can use the data presented here to make more
364 informed design choices for the systems and services developed. The data from this experiment are
365 publicly available for further study and reuse.⁷

366 REFERENCES

- 367 Abel, F., Herder, E., and Krause, D. (2011). Extraction of professional interests from social web profiles.
368 In *UMAP*.
- 369 Agrawal, R., Gollapudi, S., Halverson, A., and Ieong, S. (2009). Diversifying search results. In
370 *Proceedings of the second ACM international conference on Web Search and Data Mining*, pages 5–14.
371 ACM.
- 372 Bai, X., Wang, M., Lee, I., Yang, Z., Kong, X., and Xia, F. (2019). Scientific paper recommendation: A
373 survey. *IEEE Access*, 7:9324–9339.
- 374 Beel, J., Gipp, B., Langer, S., and Breitingner, C. (2016). Research-paper recommender systems: A
375 literature survey. *International Journal on Digital Libraries*, 17(4):305–338.
- 376 Bostandjiev, S., O'Donovan, J., and Höllerer, T. (2012). Taste-Weights: A visual interactive hybrid
377 recommender system. In *Proceedings of the sixth ACM conference on Recommender Systems*. ACM.
- 378 Carbonell, J. G. and Goldstein, J. (1998). The use of mmr and diversity-based reranking for reordering
379 documents and producing summaries.
- 380 Chen, J., Nairn, R., Nelson, L., Bernstein, M., and Chi, E. (2010). Short and tweet: experiments on
381 recommending content from information streams. In *Proceedings of the SIGCHI conference on Human*
382 *Factors in Computing Systems*. ACM.
- 383 Collins, A. M. and Loftus, E. F. (1975). A spreading-activation theory of semantic processing. *Psycholog-*
384 *ical review*, 82(6):407.
- 385 Feng, W. and Wang, J. (2014). We can learn your# hashtags: Connecting tweets to explicit topics. In
386 *2014 IEEE 30th International Conference on Data Engineering*, pages 856–867. IEEE.
- 387 Ge, M., Delgado-Battenfeld, C., and Jannach, D. (2010). Beyond accuracy: Evaluating recommender
388 systems by coverage and serendipity. In *Proceedings of the fourth ACM conference on Recommender*
389 *systems*, pages 257–260. ACM.
- 390 Goossen, F., IJntema, W., Frasinca, F., Hogenboom, F., and Kaymak, U. (2011). News personalization
391 using the CF-IDF semantic recommender. In *Proceedings of the International Conference on Web*
392 *Intelligence, Mining and Semantics*. ACM.
- 393 Herlocker, J. L., Konstan, J. A., and Riedl, J. (2000). Explaining collaborative filtering recommendations.
394 In *Proceedings of the 2000 ACM Conference on Computer supported cooperative work*, pages 241–250.
395 ACM.
- 396 Herlocker, J. L., Konstan, J. A., Terveen, L. G., and Riedl, J. (2004). Evaluating collaborative filtering
397 recommender systems. *ACM Transactions on Information Systems (TOIS)*, 22(1):5–53.
- 398 Kapanipathi, P., Jain, P., Venkataramani, C., and Sheth, A. (2014). User interests identification on Twitter
399 using a hierarchical knowledge base. In *European Semantic Web Conference*. Springer.
- 400 Kotkov, D., Wang, S., and Veijalainen, J. (2016). A survey of serendipity in recommender systems.
401 *Knowledge-Based Systems*, 111:180–192.
- 402 Letierce, J., Passant, A., Breslin, J. G., and Decker, S. (2010). Understanding how Twitter is used to
403 spread scientific messages. In *WebSci*. Web Science Trust.
- 404 Lops, P., De Gemmis, M., and Semeraro, G. (2011). Content-based recommender systems: State of the
405 art and trends. In *Recommender systems handbook*, pages 73–105. Springer.

⁷<https://doi.org/10.5281/zenodo.3367795>

- 406 McNee, S. M., Riedl, J., and Konstan, J. A. (2006). Being accurate is not enough: How accuracy metrics
407 have hurt recommender systems. In *CHI'06 Extended abstracts on Human Factors in Computing*
408 *Systems*, pages 1097–1101. ACM.
- 409 Mendoza, J. L. (1980). A significance test for multisample sphericity. *Psychometrika*, 45(4).
- 410 Nascimento, C., Laender, A. H. F., da Silva, A. S., and Gonçalves, M. A. (2011). A source independent
411 framework for research paper recommendation. In *Proceedings of the 11th Annual International*
412 *ACM/IEEE Joint Conference on Digital Libraries*, pages 297–306. ACM.
- 413 Nishioka, C., Große-Bölting, G., and Scherp, A. (2015). Influence of time on user profiling and rec-
414 ommending researchers in social media. In *Proceedings of the 15th International Conference on*
415 *Knowledge Technologies and Data-driven Business*. ACM.
- 416 Nishioka, C. and Scherp, A. (2016). Profiling vs. time vs. content: What does matter for top-k publication
417 recommendation based on Twitter profiles? In *Proceedings of 2016 IEEE/ACM Joint Conference on*
418 *Digital Libraries*, pages 171–180. ACM.
- 419 Orlandi, F., Breslin, J., and Passant, A. (2012). Aggregated, interoperable and multi-domain user profiles
420 for the social web. In *Proceedings of the 8th International Conference on Semantic Systems*. ACM.
- 421 Salton, G. and Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information*
422 *Processing & Management*, 24(5).
- 423 Shaffer, J. P. (1986). Modified sequentially rejective multiple test procedures. *Journal of the American*
424 *Statistical Association*, 81(395).
- 425 Shani, G. and Gunawardana, A. (2011). Evaluating recommendation systems. In *Recommender systems*
426 *handbook*, pages 257–297. Springer.
- 427 Shen, W., Wang, J., Luo, P., and Wang, M. (2013). Linking named entities in tweets with knowledge
428 base via user interest modeling. In *Proceedings of the 19th ACM SIGKDD international conference on*
429 *Knowledge Discovery and Data Mining*. ACM.
- 430 Sugiyama, K. and Kan, M.-Y. (2010). Scholarly paper recommendation via user’s recent research interests.
431 In *Proceedings of the 10th annual joint conference on Digital Libraries*. ACM.
- 432 Sugiyama, K. and Kan, M.-Y. (2015). Towards higher relevance and serendipity in scholarly paper
433 recommendation. *ACM SIGWEB Newsletter*, (Winter):4.
- 434 Tang, J., Zhang, J., Yao, L., Li, J., Zhang, L., and Su, Z. (2008). ArnetMiner: Extraction and mining
435 of academic social networks. In *Proceedings of the 14th ACM SIGKDD international conference on*
436 *Knowledge Discovery and Data Mining*, pages 990–998. ACM.
- 437 Tintarev, N. and Masthoff, J. (2007). Effective explanations of recommendations: user-centered design.
438 In *Proceedings of the 2007 ACM Conference on Recommender Systems*, pages 153–156. ACM.
- 439 Vargas, S. and Castells, P. (2011). Rank and relevance in novelty and diversity metrics for recommender
440 systems. In *Proceedings of the fifth ACM conference on Recommender systems*, pages 109–116. ACM.
- 441 Vargas, S., Castells, P., and Vallet, D. (2011). Intent-oriented diversity in recommender systems. In *Pro-*
442 *ceedings of the 34th international ACM SIGIR conference on Research and development in Information*
443 *Retrieval*, pages 1211–1212. ACM.
- 444 Vargas, S., Castells, P., and Vallet, D. (2012). Explicit relevance models in intent-oriented information
445 retrieval diversification. In *Proceedings of the 35th International ACM SIGIR Conference on Research*
446 *and Development in Information Retrieval*, pages 75–84. ACM.
- 447 Wu, Y., Liu, Y., Chen, F., Zhang, M., and Ma, S. (2018). Beyond greedy search: Pruned exhaustive search
448 for diversified result ranking. In *Proceedings of the 2018 ACM SIGIR International Conference on*
449 *Theory of Information Retrieval*, pages 99–106. ACM.