

Exploiting Publication Contents and Collaboration Networks for Collaborator Recommendation

Elsevier¹

Radarweg 29, Amsterdam

Elsevier Inc^{a,b}, Global Customer Service^{b,}*

^a1600 John F Kennedy Boulevard, Philadelphia

^b360 Park Avenue South, New York

Abstract

In academia, studies show that researchers are usually prolific by well collaboration with collaborators. However, due to the expansion of academic researches in diverse domains, the problem of finding most relevant and potential collaborators from a large volume of big scholarly data has become cumbersome and time-consuming. In this work, we propose an academic collaborators recommendation model called CCRec, an innovative model that combines content-based method with networks-based method. A topic clustering model and a random walk model are adopted to help effectively seek the most potential collaborators for researchers. Using DBLP data sets, we conduct benchmarking experiments to examine the performance of CCRec. Our preliminary experimental results show that CCRec outperforms other state-of-the-art methods in precision, recall and $F1$ score, additionally in addressing the topic drift problems.

Keywords: Collaboration recommendation, publication contents, collaboration networks, topic clustering, random walk.

*Corresponding author

Email address: support@elsevier.com (Global Customer Service)

URL: www.elsevier.com (Elsevier Inc)

¹Since 1880.

1. Introduction

Nowadays, with rapid development of Internet technology, the scale of Internet is beyond the imagination of people. Internet gradually becomes the main carrier of sharing information. Thus, how to obtain the useful one from vast information has become a complex task with the problem of information overload phenomena occurring. Therefore, recommender systems and techniques immensely help people by providing easier access to the specific resources they really need.

In academia, cooperation among researchers is of vital necessary. Studies show that researchers are usually prolific by well collaboration with collaborators [1]. This is to say, the collaborator is a considerable factor well connected with the productivity of a scholar. So researchers tend to get acquainted with the most potential collaborators (MPCs, some influential scholars who are similar in interests and have never collaborator before), or contact more with the most valuable collaborators (MVCs, some influential scholars or colleagues who are active and valuable in adjacent circles). Considering the inherently requirements, there has been a variety of methods proposed recommending collaborators who they have ever collaborated with or never.

In this context, previous studies have exploited mainly three aspects for academic collaboration recommendation, content-based, social network-based and hybrid recommendation. [***references!***]Some traditional content-based methods represent researchers by interests tags, make collaborators recommendation by computing interests similarities.[***references!***] However, interests tags are sometime not accurately to represent researchers' features. Moreover, in general, a researchers shows bias on various academic domains. Such behaviors usually reveal academic features of researchers in different domains. Thus, it would be imperative considering academic domains when recommending collaborators. Our previous work proposed a network-based model, ACRec [2], which solved the problem of recommending MVCs. ACRec make it easier for scientists to collaborate with colleagues in their social network. However, many scientists

also initiate collaborations outside of their social networks. It is burdensome and fraught with risk to initiating collaboration with socially unconnected researchers. In additionally, considering the less value of recommending already known collaborators, unconnected researchers are more deserve to be recommended to seek more MPCs. What’s more, some excellent hybrid models are
35 offered in recent years [3], which provides us many good references.

In this paper, we propose a novel hybrid model exploiting publication contents and collaboration networks for collaborators recommendation (CCRec). Utilizing a topic clustering model [4] [5] and a random walk model, CCR
40 ecs integrates the features of publications content and collaboration networks. We extract the subject terms from all researchers’ publications and cluster these terms into several topics, following distribute researchers to corresponding domains. To represent the feature vectors of each researchers, we run the random walk with restart model (RWR) on each domains, which have been proved to
45 be competent for calculating the rank score of node in social networks. After that, the MPCs recommendation is provided by computing the similarities of researchers’ feature vector.

In summary, we make the following contributions in this paper. 1) To make the most potential collaborators recommendation, we develop a model
50 CCRc, which combines the content-based method and collaboration network-based method. That is more favourable to achieve remarkable personalized collaborators recommendation. 2) To reveal researchers’ academic features in different domains, we present the feature vectors by utilizing a topic clustering model and a random walk model. 3) We conduct extensive experiments on a
55 subset of DBLP data set to evaluate the performance of CCRc in various scenarios as compared against our previous model ACRc and the normal common neighbors-based model (CNRec). Promising results are presented and analyzed.

The remainder of the paper is structured as follows. Section 2 briefly presents the related work. We discuss the details of our model in section 3, which
60 lights the structure of the model. In section 4, we conduct a mass of experiments and analyzes the results. At last, section 5 concludes the paper.

2. RELATED WORK

Social networks have been researched for decades to exploit the relationships and interactions between people. There are also some sample systems helping
65 people find experts on a certain topic [6]. While summarizing the preterit work, researches [3] mainly focus on three parts including content-based method, social networks-based method and minority hybrid method.

Many approaches have been presented to formalize academic collaboration recommendation as a link prediction problem [7] [8]. Link prediction is to pre-
70 dict new links that might be added to the social networks, which is useful for collaborators recommendation, but differs from it in that the latter is context-based. David et al. [9] first defined link prediction: to predict the edges that will be added to the network in the future time t' at the time of t in a social network. Some of these approaches have been applied to large social networks
75 and results show a good performance. Lichtenwalter et al. [10] examined some important factors for link prediction and proposed a general framework.

For content-based method [11] [12], the CollabSeer system recommends col-
laborators for a researcher, which is an open system introduced by Chen et al. considering research interests for collaborators recommendation [13]. Lopes et
80 al. [14] considered researchers' publications area and the vector space model to make collaboration recommendation. Pavlov and Ichise [15] extracted structural attributes from graph of past collaborations and uses them to train a set of predictors using supervised learning algorithms, then these predictors can be used to predict links between existing nodes in the graph. Tin Huynh et al. [16] pro-
85 posed a new approach that uses additional information as new features to make recommendations, i.e., the strength of the relationship between organizations, the importance rating, and the activity scores of researchers. They also propose a new method for evaluating the quality of collaborators recommendations.

Ma et al. [17] analysed how social networks information can benefit rec-
90 ommender systems and proposed a method improving the performance of recommender systems by incorporating social network information. In addition,

Lopes et al. [14] present a novel method for recommending new collaborators and identifying existing collaborations. They took the semantic issues involved in the relationships among researchers in different areas. Barabasi et al. [18]
95 analysed the basic network properties of academic social networks in terms of degree distribution, average separation, clustering coefficient, average degree and so on. In another study, Newman researched a number of statistical properties of scientific collaboration networks like the number of papers, numbers of collaborators and degree of clustering in the networks, etc.

100 Lee et al. [3] exploit how well content-based, social networks-based and hybrid recommendation algorithms predicted coauthor relationship, and the result show that a hybrid algorithm combining content and social networks information outperformed better. T Huynh et al. [16] proposed a method based on combination of probability theory and graph theory for modeling and analysing
105 co-author network. They took the trend information into considering similarity of vertices and experimental results showed that their proposed method TBRSS outperformed other existing methods. Actually, there is not much work for the combination of content and networks.

In summary, there are two deficiency not considered by traditional studies:
110 (1) lack of a better integrated method combining content-based and networks-based, and (2) research fields of researchers are not considered when recommending collaborators.

3. DESIGN OF CCRec

Our proposed design scheme for CCRec is inspired by the reality and truth
115 that a researcher usually desires to know other researchers who have similar research interests and strong influence in academia. As mentioned above, researchers often behave differently across multiple domains of interests. Such behaviors usually reveal academic features of researchers in different domains. Besides, as a social-based model, the RWR model has been proved to be competent for calculating the rank score of node in social networks derived from the
120

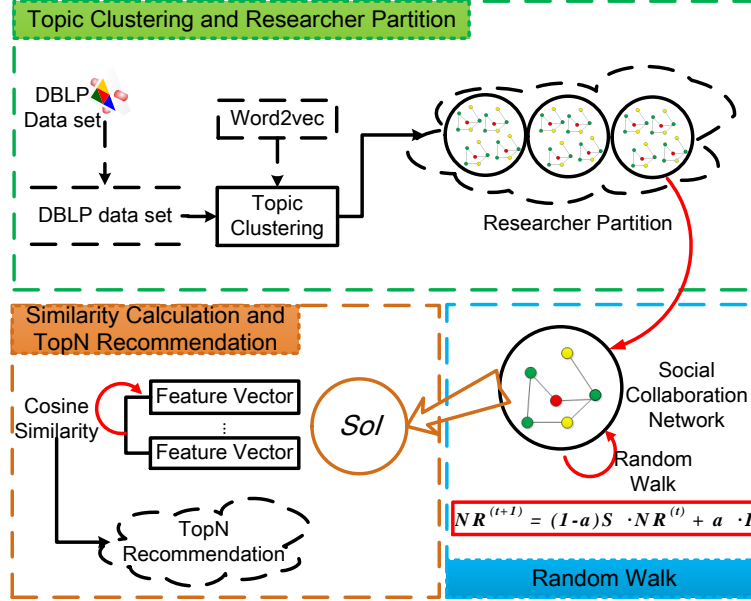


Figure 1: The architecture diagram of CCRec model

co-authorship [2], researchers' strength of influence in specific domains can be well reflected by RWR. In this work, we first adopt a content-based method to acquire multiple domains of interests. and then using the social network-based method of RWR to measure the researchers' strength of influence in different domains. After that, We use the feature vector to evaluate the similarity of researchers and then obtain the recommendation list. The detailed process is described bellow and the corresponding pseudo-code is illustrated in Algorithm 1. Figure 1 depicts three components of CCRec.

3.1. Topic Clustering and Researcher Partition

It is a content-based method for topic clustering and researcher partition, which generates various domains and maps all researchers into these domains. In this work, we use a famous tool of NLP (Natural Language Processing), word2vec, which provides an efficient implementation of the continuous of *bag-of-words* and *skip-gram* architectures for computing vector representations of

135 words. It takes a text corpus as the input and produces the word vectors as the output. The final word vector file can be used as features in many natural language processing and machine learning applications. The word vectors can be also used for deriving word classes from huge data sets. This is achieved by performing K-means clustering on top of the word vectors. The output is
140 a vocabulary file with words and their corresponding domains IDs. In case of CCRRec model, the input data is a set of titles from all the papers created by each researcher. The titles are split in many sequential words. In addition, there is necessary to filter out some irrelevant words, e.g. "of", "the", "and", etc. For extracting from titles, this set of preprocessed words can outline the
145 core contents of papers, which are signified as valuable and reliable corpus to denote a variety of academic topics. With this English corpus, word2vec obtains various domains and clusters the words into the domains.

In addition, CCRRec partitions researchers to specific domains with following method. 1, Extract subject terms from a researcher's publications. 2, Traverse
150 all the terms and check the word vector. The model tags the researcher for particular domains which contains these subject terms. It should be emphasized that one researcher always belongs to several domains and there are also many researchers in one domain. Fig. 2 describe an example. Assuming that CCRRec extracts 12 subject terms from the publications titles of researcher *S1*. After
155 topic clustering, we can see that, three of these subject terms are assigned to domain *A*, and seven in *B*, two in *C*. Thus, researcher *S1* is tagged for domains *A*, *B* and *C*. Through this method, each domains contain numerous researchers.

3.2. Feature Vector Calculation

As mentioned in section 2, in general, researchers devote to several adjacent
160 domains. But in the case of attention and strength of influence in various domains, they are offering some biases. To measure the distribution of researchers' interests, we define the *SoI* (Strength of Influence) to denote the academic values (Rank Score) of researchers in different domains, which can be regarded as the elements of feature vector of researchers. Considering each of the domains,

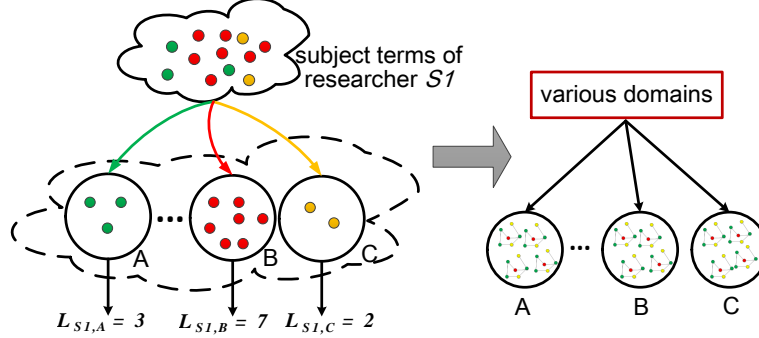


Figure 2: Researcher Partition

165 there are numerous researchers with similar research interests. Their co-author relationships can be modeled by a social network. Thus, there are many co-author networks corresponding to different domains. The *SoI* is measured by RWR model based on the co-author networks. The core equation of the RWR model is shown below:

$$R_d^{(t+1)} = \alpha \mathbf{S} R_d^{(t)} + (1 - \alpha)q \quad (1)$$

170 where R_d represents the rank score vector of all researchers in domain d , q is the initial vector R^0 , and α denotes the damping coefficient. RWR is an iterative process. After limited iterations, the vector R will be convergent. In this scenario, $SoI_s = R_{d,s}$. That is, the final value of the vector item $R_{d,s}$ is the *SoI* of researcher s .

175 In addition, with the help of RWR, The *SoI* in various domains is quantified for each researcher. To measure researchers academic feature, we define the vector F with *SoI*.

3.3. Collaboration Recommendation by Feature Vector Similarity

CCRec recommends collaborators for researchers based on their similarities. To measure the academic features similarities of researchers, we borrow a standard method, *cosine similarity* (CS). CS is employed to define the similarity

between two users s_1 and s_2 based on their feature vectors F_{s_1} and F_{s_2} .

$$Sim(s_1, s_2) = \frac{\sum_{i=1}^n (F_{s_1,i} * F_{s_2,i})}{\sqrt{\sum_{i=1}^n F_{s_1,i}^2} * \sqrt{\sum_{i=1}^n F_{s_2,i}^2}} \quad (2)$$

Finally, we consider that researchers with high similarities have common interests, they should be recommended to each other as potential academic collaborators. Hence, CCRec provides a TopN recommendation list for each researcher.

4. Evaluation and Analysis

We conduct various experiments using data from DBLP [19], a computer science bibliography website hosted at University Trier. We extracted the subsets of the entire data using the required information, which are all in the field of data mining involving 34 journals and 49 conferences. The data was modeled by an academic social network, which contains 59659 nodes (authors) and 90282 edges (coauthor relations). We divided the data set into two parts: the data before year 2011 as a training set, and others as a testing set.

We embarked on benchmarking experiments on CCRec. To evaluate the performance of CCRec model in a better way, we use three metrics which are widely used in the recommender systems, *Precision*, *Recall* and *F1* [20]. We compared CCRec with the two following approaches. ACRec: a random walk recommendation model based on collaboration networks [2]. CNRec: a common neighbors based recommendation model [14]. Four groups of experiments were conducted. 1, Find the most valuable collaborators, who may have known each other before, or be active in adjacent circles. 2, Recommend most potential collaborators, who have never cooperated with the target researcher before. 3, Evaluate how domains clustering impact the performance of CCRec. 4, Exploit the impact of clustered domains number on CCRec. For each experiment, we randomly chose 100 constant researchers who are at least somewhat active in academic activities, that is they have co-authored more than 30 person-time with others. We make collaborators recommendation for these 100 researchers, moreover compute the average of precision, recall and *F1*.

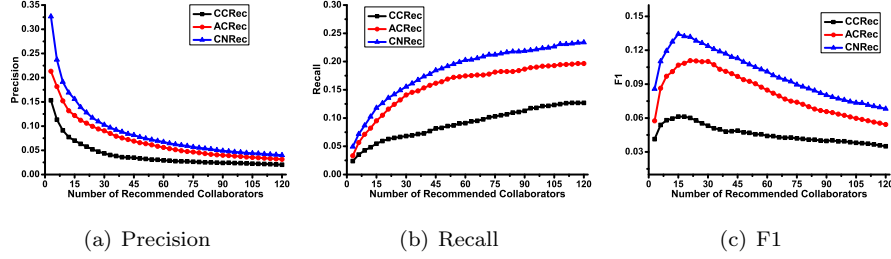


Figure 3: Performance of CCRec and ACRec on most valuable collaborators recommendation

All experiments were performed on a 64-bit Linux-based operation system,
 210 Ubuntu 12.04 with a 4-duo and 32GHz Intel CPU, 4-G Bytes memory. All the
 programs were implemented with Python.

4.1. Most Valuable Collaborators Recommendation

In our previous work [2], We proposed the ACRec model to make the most
 valuable collaborators recommendation for researchers. In this section, we ex-
 215 plored the performance of CCRec and ACRec on making most valuable collab-
 orators recommendation. The comparative results are shown in Fig. 2.

As shown in Fig. 2, The number of recommended collaborators has an
 obvious influence on the metrics with a clear trend. In the case of CCRec,
 as shown in Fig. 2(a), the precision drops when the number of recommended
 220 collaborators is increasing. At the same time, the recall in Fig. 2(b) rises
 with the increase of recommendation list, which approximates to 20% in the
 end. As for ACRec, it has the same trend with CCRec on precision and recall.
 Thus it can be seen, the precision and recall are a pair of contradictory metrics.
 Weighing the two metrics to maximum the profit, G. Shani et al. [20] adopt
 225 the metric $F1$. Fig. 2(c) describes the performance of CCRec and ACRec on
 $F1$. In case of CCRec model, the $F1$ generally increases until the number of
 recommended collaborators is over 15, and then decreases gradually. Since the
 point 15 is exactly the peak of $F1$. We can see that, CCRec performs best when
 recommend 15 collaborators to each researcher, and the $F1$ can reach 6.13%.
 230 However, in this scenario, ACRec gets its' highest $F1$ score 11.01% at point 30.

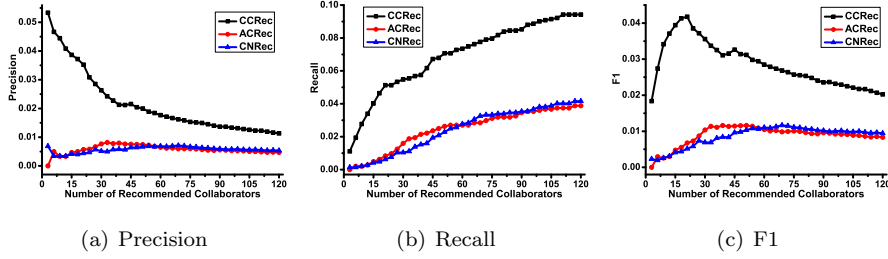


Figure 4: Performance of CCRec, ACRec and CNRec on most potential collaborators recommendation

In terms of Fig. 2, It is much in evidence that ACRec model outperforms CCRec model on making most valuable collaborators recommendation. This is because, ACRec based on the link-importance guiding random walk, considering the walk distance and rank score, seeks the most valuable collaborators who may have known each other before, or active in adjacent circles. Thus, compared with ACRec, there is no obvious superiority for CCRec to find the most valuable collaborators in adjacent circles.

4.2. Most Potential Collaborators Recommendation

We define the Most Potential Collaborators as collaborators who are worthy of being recommended and have never cooperated with the target researcher. Making the most potential collaborators recommendation is of great significance as the new collaborators are more meaningful and practical in academia reality. In this section, we explored the performance of CCRec, ACRec and CNRec on making most valuable collaborators recommendation.

Figure 3 shows the performance of CCRec, ACRec and CNRec in terms of precision, recall and $F1$ with the number of recommended collaborators increasing. It can be observed that CCRec significantly outperforms ACRec and CNRec all the time on these three metrics. CCRec shows a downtrend for precision and an uptrend for recall rate. In the case of $F1$, it reaches the peak 4.18% when recommending 21 researchers. We also see the evidence that when making most potential collaborators recommendation, ACRec performs similarly to

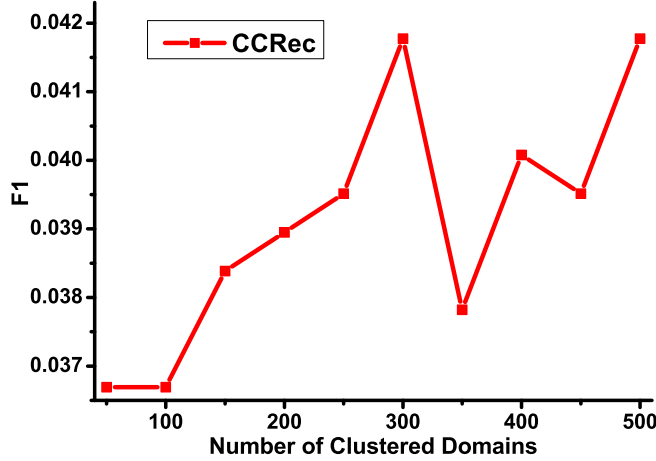


Figure 5: The impact of clustered domains number on CCRec

CNRec.

In a nutshell, CCRec outperforms ACRec and CNRec with higher precision, recall and $F1$ on making the most potential collaborators. We analysed the theory. Each researcher is represented by the feature vector, as well as CCRec model combines publications contents and collaboration networks to define the vector, which has distinct advantages (e.g. rich information, more accurately to represent researchers' feature) in recommending new collaborators.

4.3. Impact of clustered Domains number

In this section, we exploit the impact of clustered domains number on the performance of CCRec. We adopted the following experiment settings. 1, Evaluating how the $F1$ changes with the number of collaborators recommended. 2, Make the most potential collaborators recommendation for those 100 researchers selected above. 3, Recommend 21 potential collaborators for each researcher. Fig. 4 shows the experimental results.

In terms of Fig. 4, the number of clustered domains does have certain effect on the performance of CCRec. If the number of clustered domains is appropriate, the $F1$ score can get some enhancement. In this situation, when

clustering the data mining academia into 300 or 500 domains, CCRec performs
270 best for $F1$, which reaches 4.18%.

In summary, we can still claim that the model combining content-based
method and social networks-based method is really effective. Moreover, in terms
of precision, recall and $F1$, CCRec outperforms ACRec and CNRec on making
most potential collaborators recommendation for academic researchers.

275 5. Conclusions

In this paper, we focused on how to find researchers' MPCs based on big
scholarly data which is necessary in current academia. To this end, we proposed
a novel model named CCRec, by combining the features of publications content
and collaboration networks. A topic clustering model and A random walk model
280 are adopted to obtain the scholars features, and make MPCs recommendation
for researchers. We conducted extensive experiments on a subset of DBLP data
set to evaluate the performance of CCRec. We also conducted the ACRec and
CNRec on the data set as comparisons. The experimental results show that,
CCRec outperforms ACRec and CNRec in precision, recall and $F1$ score. With
285 employing topic clustering model, the problem of topic drift has been solved to
some extent.

Our research on CCRec reveals that the combination of content-based method
and networks-based method can improve the generation of effective academic
collaborations. Nonetheless, there is still room for future study in this direc-
290 tion. We extracted the titles of publications as the corpus of topic clustering
model, which are not more comprehensive than the abstract and main body of
publications. Besides, an exactly metric should be confirmed to evaluate the
topic drift problem. As future work, more experiments and studies should be
conducted.

295 References

- [1] S. Lee, B. Bozeman, The impact of research collaboration on scientific productivity, *Social studies of science* 35 (5) (2005) 673–702.
- [2] J. Li, F. Xia, W. Wang, Z. Chen, N. Y. Asabere, H. Jiang, Acrec: a co-authorship based random walk model for academic collaboration recommendation, in: *Proceedings of the companion publication of the 23rd international conference on World wide web companion*, International World Wide Web Conferences Steering Committee, 2014, pp. 1209–1214.
- [3] D. H. Lee, P. Brusilovsky, T. Schleyer, Recommending collaborators using social features and mesh terms, *Proceedings of the American Society for Information Science and Technology* 48 (1) (2011) 1–10.
- [4] C. Pan, W. Li, Research paper recommendation with topic analysis, in: *Computer Design and Applications (ICCD)*, 2010 International Conference on, Vol. 4, IEEE, 2010, pp. V4–264.
- [5] M. C. Pham, Y. Cao, R. Klamka, M. Jarke, A clustering approach for collaborative filtering recommendation using social network analysis., *J. UCS* 17 (4) (2011) 583–604.
- [6] J. Freyne, S. Berkovsky, E. M. Daly, W. Geyer, Social networking feeds: recommending items of interest, in: *Proceedings of the fourth ACM conference on Recommender systems*, ACM, 2010, pp. 277–280.
- [7] H.-H. Chen, L. Gou, X. L. Zhang, C. L. Giles, Discovering missing links in networks using vertex similarity measures, in: *Proceedings of the 27th Annual ACM Symposium on Applied Computing*, ACM, 2012, pp. 138–143.
- [8] Y. Sun, R. Barber, M. Gupta, C. C. Aggarwal, J. Han, Co-author relationship prediction in heterogeneous bibliographic networks, in: *Advances in Social Networks Analysis and Mining (ASONAM)*, 2011 International Conference on, IEEE, 2011, pp. 121–128.

- [9] D. Liben-Nowell, J. Kleinberg, The link-prediction problem for social networks, *Journal of the American society for information science and technology* 58 (7) (2007) 1019–1031.
- 325 [10] R. N. Lichtenwalter, J. T. Lussier, N. V. Chawla, New perspectives and methods in link prediction, in: *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2010, pp. 243–252.
- 330 [11] K. Balog, M. de Rijke, Finding similar experts, in: *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, ACM, 2007, pp. 821–822.
- [12] S. D. Gollapalli, P. Mitra, C. L. Giles, Similar researcher search in academic environments, in: *Proceedings of the 12th ACM/IEEE-CS joint conference on Digital Libraries*, ACM, 2012, pp. 167–170.
- 335 [13] H.-H. Chen, L. Gou, X. Zhang, C. L. Giles, Collabseer: a search engine for collaboration discovery, in: *Proceedings of the 11th annual international ACM/IEEE joint conference on Digital libraries*, ACM, 2011, pp. 231–240.
- [14] G. R. Lopes, M. M. Moro, L. K. Wives, J. P. M. De Oliveira, Collaboration recommendation on academic social networks, in: *Advances in Conceptual Modeling—Applications and Challenges*, Springer, 2010, pp. 190–199.
- 340 [15] M. Pavlov, R. Ichise, Finding experts by link prediction in co-authorship networks., *FEWS* 290 (2007) 42–55.
- [16] T. Huynh, K. Hoang, D. Lam, Trend based vertex similarity for academic collaboration recommendation, in: *Computational Collective Intelligence. Technologies and Applications*, Springer, 2013, pp. 11–20.
- 345 [17] H. Ma, D. Zhou, C. Liu, M. R. Lyu, I. King, Recommender systems with social regularization, in: *Proceedings of the fourth ACM international conference on Web search and data mining*, ACM, 2011, pp. 287–296.

- [18] A.-L. Barabási, H. Jeong, Z. Néda, E. Ravasz, A. Schubert, T. Vicsek,
350 Evolution of the social network of scientific collaborations, *Physica A: Statistical mechanics and its applications* 311 (3) (2002) 590–614.
- [19] M. Ley, Dblp: some lessons learned, *Proceedings of the VLDB Endowment* 2 (2) (2009) 1493–1500.
- [20] G. Shani, A. Gunawardana, Evaluating recommendation systems, in: *Recommender systems handbook*, Springer, 2011, pp. 257–297.
355