

Novelty Paper Recommendation Using Citation Authority Diffusion

Chun-Han Chen*, Sushilata Devi Mayanglambam*, Fu-Yuan Hsu*, Cheng-Yu Lu†, Hahn-Ming Lee*†, and Jan-Ming Ho†

*Dep. CSIE, National Taiwan University of Science and Technology, Taipei, Taiwan

†Institute of Information Science, Academia Sinica, Taiwan

Corresponding Author Email: cylu@iis.sinica.edu.tw

Abstract—Survey of academic literatures or papers should be considered with both relevance and importance of references. Authors cite related references by considering integrity and novelty. However, the state-of-art publicly academic search engines and services can only recommend related papers of a certain topic. It shows to manually evaluate the novelty of the recommended papers is necessary. In this paper, we propose a citation-network-based methodology, namely Citation Authority Diffusion (*CAD*), to rapidly mine the limited key papers of a topic, and measure the novelty on literature survey. A defined Authority Matrix (*AM*) is used to standardize duplication rate of authors and to describe the authority relation between the citing and the cited papers. Based on *AM*, our *CAD* methodology leverages the Belief Propagation to diffuse the authority among the citation network. Therefore, *CAD* transforms the converged citation network to a novelty paper list to researchers. The experimental results show *CAD* can mine more novelty papers by using real-world cases.

Keywords—Important Paper Recommendation, Citation Network, Diffusion Theory, Belief Propagation

I. INTRODUCTION

When collecting literature survey, researchers should be considered with both relevance and importance of references to their subjects [12]. They may cite the related papers with both integrity and novelty for their research. Currently, public search engines of scholarly papers and some corresponded researches support to recommend the relevant papers toward the target research. However, there exists a bias: these systems and researches only focus on finding the relevant papers rather than the important papers. This may result in the researchers who cannot evaluate the academic value of each relevant paper. Subsequently, the recommended relevant papers require to be manually evaluated the novelty from the researchers [22].

The paper recommendation issue is known as paper recommendation [8], [17], [3], paper summarization [20], [9], [11], [5], and paper suggestion [14], [15]. They are divided into two categories: Some methods discovering the similar material to the current work on a semantic (i.e., content) level and some are in the citation (i.e., structure) level. Citation context analysis [9], [27], [20], [8] is a famous approach for paper recommendation and paper summarization. It uses the viewpoints from other people as the objective measurement to evaluate the target paper. Another famous

approach on recommender system is Collaborative filtering [18], [17], [13], [25] in recent years. It uses the semantic distance among papers to track the papers near the current work. The less distance represents that these papers share the same behavior (i.e., topic). However, these methods cannot ensure the novelty paper recommendation for the target paper.

In this paper, we present a novelty paper recommendation methodology. A citation-network-based methodology named Citation Authority Diffusion (*CAD*) is proposed. The citation network [23] is used to discover potential correlations to figure critical papers. Hence, *CAD* achieves 1) increasing the novelty of survey [12] by filtering out the well-known paper of a specific research domain, and 2) preventing the cold start problem [18], [13], [25] of measuring the importance of each new published papers using citation count.

II. BACKGROUND

Academic paper recommendation is a well-known issue of literature survey. The existing researches toward this issue can be classified as three types: 1) Structure-based Approaches, 2) Content-based Approaches, and 3) Hybrid Approaches.

A. Structure-based Approaches

The structure-based paper recommendation approaches are used to analyze the potential relationship among the papers from the citation network [23]. This is because each paper is profiled via the citation behavior rather the literal content.

The citation behavior can emphasize the common concept among entities, while literature survey has been acquired from multiple sources, especially on the online academic paper corpus like Google Scholar. The Co-reference Resolution Service [7] was a link-structure implementation, which stimulated management of co-reference data, and enabled interaction among multiple Linked Open Data sources.

B. Content-based Approaches

The content-based approaches use comparing method between the concepts in each pair of papers on the semantic level [9]. The notion is the set of Bag-of-Words [4] extracted

from the paper which expresses the semantic meaning of it. Bag-of-Word is a set of single keywords which represent some textual entities such as patents, technical reports, and academic papers. These approaches profile each paper through the literal content rather than the structure among the network.

A research named Citation Semantic Link Network (C-SLN) is used to profile the semantic level information over the citation networks. Huang et al. [9] used several Natural Language Processing methods as a framework to construct a C-SLN. Based on C-SLN, they could find papers of high importance effectively. In addition, they also integrated the citation functions and detected opinion communities among papers.

C. Hybrid Approaches

This kind of approaches combines the both ideas mentioned above. These approaches profile each paper via not only its citation behavior, but also its literal information. A good example is Citation Context [21]. Citation Context is a set of sentences surrounding the reference holders (i.e., the format: [1]) used to refer to other papers. Based on Citation Context, researchers develop many methods to extract the essential information for key paper finding.

He et al. [8] proposed a novel non-parametric probabilistic approach named CRM, which could measure the context-based relevance between a target paper and its Citation Context. They presented a context-aware citation recommender system, which could recommend a small number of good papers for literature survey. According to the study, a non-parametric probabilistic model and its scalable closed form solutions were used. In addition, it could suggest a set of citations for a target paper with good quality.

III. CITATION AUTHORITY DIFFUSION

In order to solve the novelty paper recommendation, a methodology named Citation Authority Diffusion (*CAD*) is proposed. This approach is composed of three modules: 1) Information Collection; 2) Information Organization; and 3) Information Presentation, as shown in Figure 1.

A. Information Collection

The Information Collection module is used to extract the concept from the target research, and collect some of the survey materials corresponding to the concept efficiently. The concept of this module is to identify the academic domains of a target research. This module consists of two units: Key Concept Extractor and Survey Material Finder.

1) *Key Concept Extractor*: The main purpose of Key Concept Extractor is to discover the potential concept of a target research or a specific research and is the pre-processing stage of extracting key concept to build the bird-eye map of the target research. The input data is a

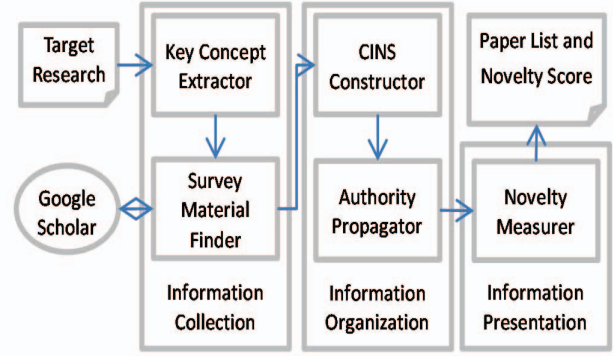


Figure 1. The system architecture of *CAD*.

target research (*TR*). *TR* includes the titles, abstracts, summaries, or even full texts. This information may represent the domain expertise. The title of target research most likely contains the basic concept of the entire research and the remaining sections may play a support role to extend concepts. This unit not only filters out the noises (e.g., prepositions, conjunctions, and so on), but also keeps the essential information from a target research.

Our proposed methodology achieves this unit by the C-value method because it enhances the common statistical measure of frequency of occurrence for term extraction, making it sensitive to a particular type of multi-word terms, the nested terms [6], [26]. There are three steps in this method: 1) POS (part-of-speech) tagging, 2) linguistic filtering, and 3) statistical ranking. In this unit, the Stanford POS tagger [24] and Porter stemmer [10] are adopted. The output of this unit is the concepts of a target research (*CTR*) containing key concept of target research and the corresponding C-value.

2) *Survey Material Finder*: The idea of Survey Material Finder is to collect the relevant survey materials (*SM*) corresponding to the target research. It simulates the human behavior on gathering information through web sites using an online digital library. In this unit, Survey Material Finder gathers a *CTR* which extracted from the target research and queries them in Google Scholar. It also allocates each of the concept result with the different ratio according to the C-value generated by Key Concept Extractor.

Survey Material Finder receives the key concept set *CTR* of a target research and generates the survey materials corresponding to target research. At first, Key Concept Extractor generates *CTR* containing two characteristics in each concept: the key phrase and the corresponding C-value. The key phrase is the combination of keywords which identifies the meaning of a concept and the C-value is a normalized score which identifies the significance of the phrase. Then, Survey Material Finder sends each of the concepts from *CTR* into Google Scholar. It receives some

results sorted by importance from Google Scholar and keeps them in a result buffer. In other words, the data are flowing bi-directionally i.e., from Survey Material Finder to Google Scholar and vice versa as indicated in the Figure 1. Finally, this unit uses each C-value from *CTR* to generate the survey materials.

B. Information Organization

The Information Organization module is used to explore the key information of a specific research. Its notion is to discover the potential paper and relationships among the survey materials. This module comprises two units: Citation Network of Survey (*CINS*) Constructor and Authority Propagator.

1) *CINS Constructor*: The goal of *CINS* Constructor is to build the *CINS* with the current survey materials. The main purpose is to find the potential papers and to integrate the relations among all the resources. Hence, this unit is used to collect the high relevance papers based on each material received from the previous module. There are two kinds of resources for relevance paper finding: the papers which cite the current materials (i.e., citers) and the papers which are cited by the current materials (i.e., references). Thus, to find the two resources is the significant task of this unit.

The citers and the references can be found via the citation network easily. Most of the corresponding bibliography could easily be detected based on the Equation (1).

$$Rel(s, d) = \frac{InDeg(s|d)}{InDeg(d)} \quad (1)$$

In (1), d is one of the current material of survey materials, s is the candidate paper of d for *CINS* patching, $InDeg(d)$ is the citer count of d , and $InDeg(s|d)$ is the citer count only from d 's citer to the paper i.e., the paper s is cited by d 's citer. The output $Rel(s, d)$ is the relatedness score from d to the corresponding s . In addition, another task is to filter out the irrelevant candidates which may be considered by the relatedness score. Thus, a threshold should be identified for relatedness as shown in Equation (2).

$$Fil(d) = \frac{|Sib(d)|}{\sum_{s \in Sib(d)} Rel(s, d)^{-1}} \quad (2)$$

In (2), $Sib(d)$ is candidate paper set of d , and $|Sib(d)|$ is the count number of $Sib(d)$. The output $Fil(d)$ is the threshold score for $Rel(s, d)$ detection. If $Rel(s, d)$ is greater than $Fil(d)$, the candidate s will add into the *CINS*. Similarly, in (1), $OutDeg(d)$ could be used as the reference count of d .

In this unit, we attempt to find most of the relevant papers toward the survey materials and would like to build the *CINS*. Each current materials d and track the siblings (e.g., citers) $Sib(d)$ of d is profiled. Generally, the more common references they shared, the higher correlation they gained.

Thus, the citation network provides the platform for the patching papers finding.

2) *Authority Propagator*: The Authority Propagator is aimed to generate the key information (*KI*) corresponding to the topics of the target research. It is used to identify the key information from the citation network. The key information is a model which consists of the popular papers corresponding to the academic fields of the *CINS* and a target research. The Authority Propagator not only tries to calculate the duplication rate (i.e., authority), but also discovers the authority of papers. For one thing, the duplication rate of each author related to the papers will be calculated. For another, the Authority Matrix (*AM*) will be produced during the process. Finally, the model of target research can be generated according to the Authority Matrix of each citation link and the duplication rate of each author. It gives us a standard for literature survey novelty measurement.

This unit is used to traverse the paths from the *CINS* and calculates the duplication rate of every corresponding author. Each author in the *CINS* represents a authority entity for paper quality evaluation. The paper and each link represent a citation behavior. In addition, this unit is also used to dynamically generate the Authority Matrix for authority diffusion. At last, Authority Propagator updates the authority of each paper in the citation network and filters out the low authority papers. At the same time, essential information of *CINS* (i.e., *KI*) is generated.

In this unit, we leverage the Belief Propagation with our potential function named Authority Matrix. The Belief Propagation [19] based on the Equation (3) is used to infer the probabilities about maximum likelihood state from each paper of the citation network. For implementing the algorithm of Belief Propagation, our method achieves this unit by the libDAI library [16].

$$m_{ij} = \sum_{\sigma'} \Psi(\sigma', \sigma) \prod_{n \in N(i) \setminus j} m_{ni}(\sigma') \quad (3)$$

In (3), m_{ij} is the message vector sent by paper i to j and $N(i)$ is the set of papers citing i . Meanwhile, an Authority Matrix $\Psi(\sigma', \sigma)$ is proposed, which includes the prior state assignment for each citation pair. This state assignment could be represented in Table I. The paper authority value is the mean of all the author authority values and the author authority value is the standardized author duplication rate among the citation network from 0 to 1. In Table I, $Auth(Citer)$ denoted as σ' represents the authority value of citer paper and $Auth(Cited)$ denoted as σ represents the authority value of cited paper. Authority Propagator uses the citation network as the platform for paper authority passing. There are two states of a node in a citation network: key paper and normal paper. The high authority papers usually cite the papers with a certain level of authority basically. Based on this philosophy, Authority Propagator uses *AM* as

Table I
THE AUTHORITY MATRIX STATEMENTS

Citer \ Cited	Key paper	Normal Paper
Key paper	$Auth(Citer) \times Auth(Cited)$	$Auth(Citer) \times (1 - Auth(Cited))$
Normal Paper	$(1 - Auth(Citer)) \times Auth(Cited)$	$(1 - Auth(Citer)) \times (1 - Auth(Cited))$

a diffusion factor $\Psi(\sigma', \sigma)$ and dynamically update the paper authority. Finally, this unit generates a converged network with high authority and this network KI is the key paper list.

C. Information Presentation

The Information Presentation module is intended to provide an evaluation method for the target research to the KI . The meaning of this module is to make the relation between the generated survey model and the target research comprehensible. Thus, this module is used to generate a survey novelty score toward the target research. People could easily find out what they have done or what they have to do.

For achieving the goal of Information Presentation, this module comprises only one unit: Novelty Measurer. The main purpose of Novelty Measurer is to generate a score about the survey novelty of the target research. The scheme of this unit is to use the novelty as the quality feature. It gives us a bird-eye view of our current research.

After introducing the idea of this unit, this paragraph shows the steps in detail. First, the Novelty Measurer receives the key information from the previous unit. The key information contains the key paper titles. At the same time, the original bibliography of the target research which written by the researcher could be the comparison data to the key information. For extracting the titles from the bibliography in original research, our method achieves this unit by the FreeCite Citation Parser [2]. Each citation of the original bibliography can be parsed to paper title, authors, and so on. Second, this unit calculates the novelty score between the target research's own comparison data and the key information. The novelty score can be implemented in Recall [8] and Co-cited Probability [8], [9], [27]. The detail of these standards of measurement will be presented in next section. Finally, the researcher will obtain the novelty paper list with the novelty score to their current literature survey of the target research.

IV. EXPERIMENTS

For implementing *CAD*, a Survey Importance Measurement (*SIM*) system is developed as an online web service¹. *SIM* receives the title, abstract, keywords, and bibliography of a target research, and then generates a novelty score with a survey paper list toward it.

¹<http://140.118.155.22:8080/SIM/>

A. Dataset

For the experimental purpose [8], we selected all the academic papers published before year 2008 from CiteseerX dataset [1] as the experiment corpus. There were 456,787 unique papers after removing the duplications and 1,612 papers with quality references were chosen [8] as the testing dataset.

B. Evaluation Methods

The novelty of recommendation could be evaluated by two standards of measurement: 1) Recall [8], and 2) Co-cited Probability [8], [9], [27]. These measurements are shown on Equation (4) to Equation (6).

$$Recall(TR, KI) = \frac{|TR \cap KI|}{|TR|} \quad (4)$$

$$CP(TR, KI) = \frac{\sum_{i \in TR} \sum_{j \in KI} P(d_i, d_j)}{|TR| \cdot |KI|} \quad (5)$$

$$P(d_i, d_j) = \frac{|Citers(d_i) \cap Citers(d_j)|}{|Citers(d_i) \cup Citers(d_j)|} \quad (6)$$

In (4), Recall is defined as the percentage of the top K recommended papers in KI that appear in the original bibliography from the TR . However, some important or even better papers other than those original ones in the bibliography need to be captured. Thus, the Co-cited Probability is used to solve this problem. In (5), CP (i.e., Co-cited Probability) is the average number of over all $|TR| \cdot |KI|$ unique paper pairs for the TR and the top K of KI , where $|TR|$ and $|KI|$ is the number of TR and KI . Considering (6), for each pair $\langle d_i, d_j \rangle$, we calculated the probability which these two papers have been co-cited by the prestige in the past. Here, d_i is an original paper from TR and the d_j is a recommended paper from KI .

We compare our *CAD* with 3 other baselines using Recall and Co-cited Probability listed in [8]: 1) CRM, a content-based method which use a probabilistic model. It can evaluate the context similarity between a citation context and a paper. CRM can recommend paper for a citation context or a paper list for a literature survey paper with high quality; 2) g-count, a method based on CRM. In this approach, the candidate papers are ranked according to the citation count in the whole corpus; 3) textsim, the other method based on CRM. In this approach, the candidate papers are ranked according to semantic distance with the query information using only title and abstract.

C. Experiment Results and Discussion

The experimental results about Recall represent the novelty of paper recommendation, which is shown in Figure 2. The trends of the Recall for each method are growing because all the methods keep recommending papers to match the original bibliography.

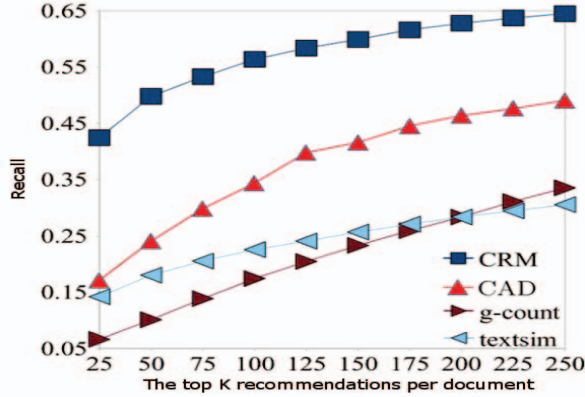


Figure 2. Recall comparison for paper recommendations.

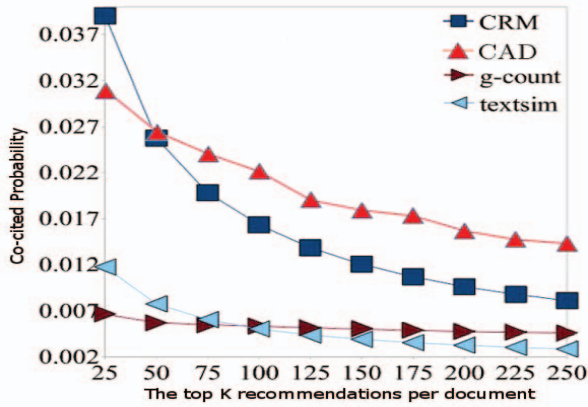


Figure 3. Co-cited Probability comparison for paper recommendations.

Considering Figure 2, the Recall of *CAD* is from 0.17 to 0.49 when *K* is from 25 to 250. The speed of growth of *CAD* is faster than others when *K* is from 25 to 125. This result reflects that *CAD* can keep recommending important papers in a conscious range of recommendations.

Recall uses only the references of target paper as the correct answers, but there is still a lot of articles not referenced. Therefore, it would loss the novelty assessment. Although the Recall curve of *CAD* is much higher than *g-count* and *textsim*, there still exists a gap between *CAD* and *CRM* method. Considering (4) again, the main purpose of Recall is to set the bibliography of the target research as the ground truth. However, what researchers want to evaluate is not only their bibliography, but also the important or even better papers other than their bibliography. In order to fulfill the need of the evaluation, the Co-cited Probability also represents the novelty of paper recommendation, which is shown in Figure 3. The trends of the Co-cited Probability for each method are decayed because all the methods keep recommending papers to narrow down the probability.

In Figure 3, the Co-cited Probability of *CAD* is from 0.029 to 0.013 when *K* is from 25 to 250. The speed of

decay of *CAD* is slower than others. But we could not find the reason why *CAD* gives the lower Co-cited Probability when *K* is lower than 50. Possibly, when *K* is smaller, the generated papers from the key information may include some low authority papers which will reduce $P(d_i, d_j)$. However, we intend to figure it out more details in our future work. Although the Co-cited Probability of *CAD* is lower than *CRM* at the beginning, it is quickly going beyond the *CRM* with the slow decay. This result shows that *CAD* can keep recommending novelty papers in a conscious range of recommendations.

The Co-cited Probability solves the problem with novelty evaluation occurred by Recall. It measures not only the recommended papers which appear in the bibliography, but also the papers other than the original citations.

V. CONCLUSION AND FURTHER WORK

Conducting a literature survey should be considered with both relevance and importance of references to their subjects. In this paper, we proposed a Citation Authority Diffusion (*CAD*) methodology to facilitate the identification of key papers and improve the novelty of literature survey. A Survey Integrity Measurement (*SIM*) system, as an online web service, is also developed to automatically provide an authoritative paper list to the target research. *SIM* can compare with the research bibliography and generate a novelty score to the literature survey as well. In our experiment, we chose CiteseerX [1] as our dataset. The results showed that *CAD* can not only recommend papers to match the original bibliography, but also provides the important papers even excluded from the original citations. Thus, researchers could easily figure out how to ignore the familiar papers when collecting literature survey. Furthermore, they can focus on the important papers among the recent relevant researches.

Although our proposed *CAD* could help researchers to choose novelty papers, there are still some works that need further consideration. Several open questions, such as further dataset and concept extraction are needed to be investigated.

Further dataset. In the experimental dataset, different dataset may reflect in different search results. In this paper, we use the more credible CiteseerX computer science dataset as our experimental one. However, it cannot be confirmed that such a dataset is suitable. This may need more experiments to ensure the accuracy of the novelty.

Concept extraction. *CAD* is still needed to extract the concept from the target research because it is a hybrid methodology. Hence, *CAD* needs to extract the irrelevant or even empty concept if researchers provide limited research information. It needs more feature extracting methods to ensure the quantity of the information.

Our proposed *CAD* provides a new basis for discovering research survey, and measuring novelty of literature survey. We are convinced that such *CAD* can ensure more efficiency and would provide higher novelty level of literature survey.

ACKNOWLEDGMENT

This research was partly founded by the National Science Council of Taiwan under grant number NSC96-2628-E-011-084-MY3, NSC98-2221-E-001-010-MY3, and NSC99-2221-E-011-075-MY3.

REFERENCES

- [1] Citeseerx. <http://citeseerx.ist.psu.edu/>.
- [2] Freecite. <http://freecite.library.brown.edu/>.
- [3] ADOMAVICIUS, G., AND TUZHILIN, A. Context-aware recommender systems. *Recommender Systems Handbook* (2011), 217–253.
- [4] BERRY, M., AND CASTELLANOS, M. *Survey of text mining II*. Springer, 2008.
- [5] ELKISS, A., SHEN, S., FADER, A., ERKAN, G., ET AL. Blind men and elephants: What do citation summaries tell us about a research article? *Journal of the American Society for Information Science and Technology (JASIST)* 59, 1 (2008), 51–62.
- [6] FRANTZI, K., ANANIADOU, S., AND TSUJII, J. The c-value/nc-value method of automatic recognition for multi-word terms. *Research and Advanced Technology for Digital Libraries (ECDL)* (2009), 520–520.
- [7] GLASER, H., JAFFRI, A., AND MILLARD, I. Managing co-reference on the semantic web. In *Linked Data on the Web (LDOW)* (2009), CEUR-WS, pp. 1–6.
- [8] HE, Q., PEI, J., KIFER, D., MITRA, P., AND GILES, L. Context-aware citation recommendation. In *International Conference on World Wide Web (WWW)* (2010), ACM, pp. 421–430.
- [9] HUANG, Z., AND QIU, Y. A multiple-perspective approach to constructing and aggregating citation semantic link network. *Future Generation Computer Systems (FGCS)* 26, 3 (2010), 400–407.
- [10] JONES, K., AND WILLETT, P. *Readings in information retrieval*. Morgan Kaufmann Pub, 1997.
- [11] KAPLAN, D., IIDA, R., AND TOKUNAGA, T. Automatic extraction of citation contexts for research paper summarization: a coreference-chain based approach. In *Workshop on text and citation analysis for scholarly digital libraries (NLPIR4DL)* (2009), Association for Computational Linguistics, pp. 88–95.
- [12] KESHAV, S. How to read a paper. *ACM SIGCOMM Computer Communication Review (CCR)* 37, 3 (2007), 83–84.
- [13] KOREN, Y. Factorization meets the neighborhood: a multi-faceted collaborative filtering model. In *International Conference on Knowledge Discovery and Data Mining (SIGKDD)* (2008), ACM, pp. 426–434.
- [14] LI, L., CHU, W., LANGFORD, J., AND SCHAPIRE, R. A contextual-bandit approach to personalized news article recommendation. In *International Conference on World Wide Web (WWW)* (2010), ACM, pp. 661–670.
- [15] LIANG, H., XU, Y., LI, Y., AND NAYAK, R. Collaborative filtering recommender systems using tag information. In *IEEE/WIC/ACM International Conference on Web Intelligence (IAT)* (2009), vol. 3, IEEE, pp. 59–62.
- [16] MOOIJ, J. M. libDAI: A free and open source C++ library for discrete approximate inference in graphical models. *Journal of Machine Learning Research (JMLR)* 11 (Aug. 2010), 2169–2173.
- [17] NAAK, A., HAGE, H., AND AIMEUR, E. A multi-criteria collaborative filtering approach for research paper recommendation in papyres. *E-Technologies: Innovation in an Open World* (2009), 25–39.
- [18] PAN, C., AND LI, W. Research paper recommendation with topic analysis. In *Computer Design and Applications (ICCD)* (2010), vol. 4, IEEE, pp. V4–264.
- [19] PEARL, J. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988.
- [20] QAZVINIAN, V., AND RADEV, D. Scientific paper summarization using citation summary networks. In *International Conference on Computational Linguistics (COLING)* (2008), Association for Computational Linguistics, pp. 689–696.
- [21] RADEV, D., MUTHUKRISHNAN, P., AND QAZVINIAN, V. The acl anthology network corpus. In *Workshop on text and citation analysis for scholarly digital libraries (NLPIR4DL)* (2009), Association for Computational Linguistics, pp. 54–61.
- [22] RASANEN, E., KIKTA, R., SORVARI, A., SALMENKAITA, J., HUHTALA, Y., MANNILA, H., TOIVONEN, H., OINONEN, K., AND MURTO, J. Location-based novelty index value and recommendation system and method, Mar. 4 2008. US Patent App. 20,080/214,210.
- [23] SHI, X., LESKOVEC, J., AND MCFARLAND, D. Citing for high impact. In *Joint Conference on Digital Libraries (JCDL)* (2010), ACM, pp. 49–58.
- [24] TOUTANOVA, K., AND MANNING, C. Enriching the knowledge sources used in a maximum entropy part-of-speech tagger. In *Proceedings of the 2000 Joint SIGDAT conference on Empirical methods in natural language processing and very large corpora: held in conjunction with the 38th Annual Meeting of the Association for Computational Linguistics-Volume 13* (2000), Association for Computational Linguistics, pp. 63–70.
- [25] TSO-SUTTER, K., MARINHO, L., AND SCHMIDT-THIEME, L. Tag-aware recommender systems by fusion of collaborative filtering algorithms. In *ACM Symposium on Applied Computing (SAC)* (2008), ACM, pp. 1995–1999.
- [26] VOORHEES, E. The trec-8 question answering track report. *National Institute of Standards and Technology Special Publication (NISTSP)* (2000), 77–82.
- [27] WAN, S., PARIS, C., MUTHUKRISHNA, M., AND DALE, R. Designing a citation-sensitive research tool: an initial study of browsing-specific information needs. In *Workshop on text and citation analysis for scholarly digital libraries (NLPIR4DL)* (2009), Association for Computational Linguistics, pp. 45–53.