

Literature Recommendation by Researchers' Publication Analysis*

Jinpeng Chen

*College of Computer Science
Inner Mongolia University*

Hohhot, Inner Mongolia Autonomous Region, China

sdchjp1990@163.com

Zhijie Ban[†]

*College of Computer Science
Inner Mongolia University*

Hohhot, Inner Mongolia Autonomous Region, China

banzhijie@imu.edu.cn

Abstract – Scholarly paper recommendation has been an important research topic in the field of information filtering because scholars find thousands of publications that match their search queries but are largely irrelevant to their latent information needs. Many existing methods are based on the users' online behaviours to construct user interest model. However, an author's published works constitute a clean signal of his latent interests. In previous research, Sugiyama et al. recognized the important user profile and proposed the scholarly paper recommendation via user's recent research interests. However, in reality, a user's interests can be diverse and a researcher's published papers often involve his multiple research directions. On the other hand, it is often more desirable to recommend papers based not only on content relevance but also on the reputation of venues. To deal with the above mentioned limitations and problems, in this paper, a novel scholarly paper recommendation method is proposed. The main distinctive features of the proposed model include: (1) user's information needs are generated in term of scholar's multiple research directions; (2) user's interests are represented by his research content and recommendation standard; (3) research content model and recommendation standard of each scholar is generated from his published papers, citation papers and reference papers; (4) in order to distinguish the scholar's different attention on his research contents, we present the penalty factor concept. Extensive experiments are designed to evaluate the effectiveness of the proposed model by using one real data set. The results show that the proposed model significantly outperforms existing methods.

Index Terms - Literature Recommendation, Interest Profile, Recommendation Standard, Penalty Factor

I. INTRODUCTION

The explosive growth of literatures leads to academic information overload, where users find an overwhelming number of publications that match their search queries but are largely irrelevant to their latent information needs [1]. Academic paper recommendation is one of the key technologies to solve this problem. Some digital electronic libraries have developed personalized recommendation systems according to the user's subscription records. However, these recommendation systems not only require users to spend extra time and effort on describing their interests, but also it is difficult to recommend related papers accurately for users if users' description exists a deviation.

Traditional academic recommendation models were developed using an approach that is based on user historical be-

haviours such as publishing, keeping, tagging, commenting and scoring [4-6]. But these models suffer from the problem of collecting user behaviour information and the slow process of constructing personalized recommendation model. To overcome the limitation of behaviours-based approaches, citation network or co-author network based techniques have been used to calculate the correlation between documents [7-11]. However, a scholar's published papers reflect his research content and it is easy to acquire researcher's all publications without his participation. In previous research, Sugiyama et al. recognized the important user profile and proposed the scholarly paper recommendation via user's recent research interests [1-3]. A key part of their model is to enhance the profile derived directly from past works with information coming from the past works' referenced papers as well as papers that cite the works. However, in reality, user's interests can be diverse and a researcher often has multiple research directions. For example, the research interests of the expert Jie Tang include social network theories, data mining methodologies, machine learning algorithms, and semantic web technologies [12]. On the other hand, it is often more desirable to recommend literature based not only on content relevance but also on the reputation of authors and venues. In this paper, we propose to find multiple interests and recommendation standard of users, which is called literature recommendation based on a researcher's multiple research directions and recommendation standard, abbreviated as MDRS, to generate a scholar's profile.

The original contributions of the proposed MDRS to the field of academic paper recommendation can be described as follows:

1) We propose to model users' interests with multiple research interests rather than single interest under the assumption that some scholars' research directions can be diverse.

2) MDRS includes research content model and recommendation standard of each scholar, which is generated from his published papers, citation papers and citing papers. We propose to integrate data mining techniques with document vector space model to generate a multiple interests-based model to represent a researcher's interests for research content. The recommendation standard of each scholar is constructed by analysing the level of his published papers.

3) In order to distinguish the scholar's different attention on his multiple study directions, we present the penalty factor

*This work is supported with natural science foundation of Inner Mongolia Autonomous Region (2014MS0603).

[†]Corresponding author.

concept. The proposed model MDRS consists of interest distributions describing interest preferences of one researcher.

In Section 2, we discuss the related work about some academic recommendation models and related techniques. Section 3 presents the details of our proposed model. Then, extensive experiments on the proposed model and baseline models have been conducted on a real data collection in Section 4. Section 5 concludes the whole work and presents ideas for future work.

II. RELATED WORK

Academic recommendation is one research branch of recommendation systems, which help users to find relevant papers in the tens of millions papers. It is mainly used in the digital libraries, papers-shared sites and academic social networking sites.

Recommendation systems can obtain user information needs from user's profiles. User's online behaviours and context information often indicates the user's interests. Wang et al. [4] proposed a method that models user historical behaviour, which built preference model for each user through collecting the operations on scientific papers of online users and carrying on the detailed analysis. The personalized recommendation model was constructed based on content-based filtering model and statistical language model. He et al. [11] given context-aware citation recommendation, and built a context-aware citation recommendation prototype. This system was capable of recommending the bibliography to a manuscript and providing a ranked set of citations to a specific citation placeholder. Wang et al. [13] proposed the merits of traditional collaborative filtering and probabilistic topic modeling. It provided an interpretable latent structure for users and items, and could form recommendations about both existing and newly published articles.

Some studies are based on the citation and co-author networks. Tang et al. [14-15] present a unified topic model to simultaneously model topical aspects of different objects in the academic network. Bolelli et al. [16] present an information theoretic approach towards measuring the significance of individual words based on the underlying link structure of the document collection. They generated a non-uniform weight distribution of the feature space which was used to augment the original corpus-based document similarities. Liang et al. [17] improved recommendation performance by defining Local Relation Strength and Global Relation Strength. Collaborative filtering has proven to be valuable for recommending items in many different domains. Mcnee et al. [7] explored the use of collaborative filtering to recommend research papers, using the citation web between papers to create the ratings matrix. Specifically, they tested the ability of collaborative filtering to recommend citations that would be suitable additional references for a target research paper.

Sugiyama et al. considered the academic achievement of a scholar could embody the scholar interests [1-3]. They proposed a generic model towards recommending scholarly papers relevant to a researcher's interests by capturing their re-

search interests through their past publications. The advantages of this approach include two sides. Firstly, it is easy to collect the user's publications in the digital library. Secondly, one scholar's scientific papers reflect his research contents. But their user model considered one scholar's all papers as one interests. In fact, the research direction of one scholar isn't unique and many scholars include several research interests. Thus, we study multiple interest user model based on the user's publications.

III. PROPOSED METHOD

A. Model Construction for Published Papers

Science papers can reflect a scholar's important research achievements, which also point out his studying methods and academic fields. A paper generally cites only a small proportion of papers. At the same time, the paper is cited by other some papers. These citation and references papers are relevant with the paper, each of which also emphasizes on the research topic from one side. Furthermore, the title, the abstract, the key words and the main body in the paper reflects its research content and its topic. Thus, the model of a paper can be constructed by its content combining with its citation and references papers. We choose the VSM (vector space model) as statistical language model. Based on the article [1], the model of paper i can be expressed as follows:

$$PM(i) = f_i + \sum_{j=1}^n si \cap f_i, f_j^r \cdot f_j^r + \sum_{k=1}^m si \cap f_i, f_k^c \cdot f_k^c. \quad (1)$$

Where f_i , f_j^r and f_k^c respectively are the VSM model of the paper i , its reference paper j and citation paper k . The VSM model is constructed by its title, abstract, the key words and the main body. n denotes the number of reference papers and m denotes the number of citation papers. $si \cap f_i, f_j^r$ is the cosine similarity value between the paper i and the reference paper j , and $si \cap f_i, f_k^c$ is the cosine similarity value between the paper i and the citation paper j .

B. Profile Construction for Researchers

The published papers of a researcher always include many aspects for his study, such as research field, research profile, studying approaches and studying thought. We can know one scholar's research interests from his publications. Many senior scholars have multiple research directions. Even if a scholar has only a research direction, he has a number of research interests in terms of different studying aspects and methods. On the other hand, when a scholar is interested in some papers, he is concerned about the papers' reputations in addition to research content. It is impossible for an expert who is focus on the publications in the junior journals when he finds some relevant papers. In order to solve the mentioned problems, our researcher model is composed of two parts. The profile of one researcher μ is defined as follows:

$$Profile_{\mu} = (M_{\mu}, S_{\mu}). \quad (2)$$

Where M_μ and S_μ are the interest model and the recommendation standard of researcher μ , respectively. We give the definition of the two parts in the following.

Interest Model

All published papers of a researcher in the paper [1] were used to model one same interest model, which represented the researcher's interest. Under this condition, it is difficult to describe multiple research directions of one scholar because the dimensions of the VSM vector could rapid increase, which can unintentionally descend cosine similarity between the single interest profile and relevant papers. In reality, the research directions of some researchers can be diverse. A senior researcher's published papers might be divided several disjoint clusters by his several study directions, and each cluster can reflect a research direction. Therefore, we propose a multiple interest model to display the multiple research directions of a researcher. In other words, we cluster the published paper set of a researcher to several subsets by clustering, and each subset represents one interest of the researcher. We use k-means clustering method because it has the advantages of simple implementation and fast speed.

Suppose PS_μ represents the published paper set of researcher μ , defined as:

$$PS_\mu = \{P_\mu^i \mid 1 \leq i \leq n\}. \quad (3)$$

Where n is the total number of papers which have been published by researcher μ , and P_μ^i is a paper in PS_μ . We build the paper model set by (1), defined as:

$$PMS_\mu = \{PM(P_\mu^i) \mid P_\mu^i \in PS_\mu\}. \quad (4)$$

We make use of the cosine similarity to measure the correlation between two vectors in PMS_μ . Then, we split PMS_μ to several disjoint subsets by the k-means clustering algorithm, which can be expressed as follows:

$$\begin{aligned} PMS_\mu &= \bigcup_{i=1}^k C_\mu^i \\ C_\mu^i &\neq \emptyset, 1 \leq i \leq k \\ C_\mu^i \cap C_\mu^j &= \emptyset, 1 \leq i, j \leq k \end{aligned} \quad (5)$$

Where k is the number of clusters. After dividing PMS_μ , we can describe the scholar's multiple interest model as following:

$$M_\mu = \{I_\mu^i \mid 1 \leq i \leq k\}. \quad (6)$$

Where I_μ^i represents an interest model, denoted as following:

$$I_\mu^i = \sum_{PM(P_\mu^j) \in C_\mu^i} PM(P_\mu^j). \quad (7)$$

Equation (7) can express the researcher's several research directions exactly.

Recommendation Standard

We believe that each researcher in a literature recommendation system should be built a recommendation standard to filter some candidate papers that is lower than the standard. Different researchers should be established different standards to match their different minimum requirements for the level of recommended papers. Generally, a scholar's academic level and research capacity can be reflected by the impact factor of the venues which he participated. Thus, we give the definition in the following:

$$S_\mu = \frac{1}{N} \sum_{i=1}^N g(p_\mu^i). \quad (8)$$

Where N is the number of the published papers of research μ , p_μ^i is his one published paper, and $g(p_\mu^i)$ represents the impact factor of the venue where p_μ^i is published.

C. Penalty Factor of Interests

It is natural to realize that not all publications attract equal attention to a researcher. In above constructing process, we ignored an issue that researchers might pay different attention to their different study directions. In other words, we should not treat each interest equally. So, we propose the penalty factor to distinguish researchers' attention level for their different interests. We believe that the interest which the most recent published paper locates in is always the main interest for researchers. For researcher μ , suppose C_μ^m represents his main interest cluster, the penalty factor can be defined as follows:

$$PF_\mu^i = \begin{cases} 1, & C_\mu^i = C_\mu^m \\ \frac{1}{1 + e^{-\lambda m}}, & C_\mu^i \neq C_\mu^m \end{cases} \quad (9)$$

Where $\lambda \in [0, 1]$ denotes the penalty coefficient, and m is the number of papers in C_μ^i , which means that researchers more emphasize on those clusters which have more papers.

D. Recommendation Process

The process of paper recommendation for a researcher is divided into several steps:

- 1) The system generates the multiple interest model and computes the recommendation standard value for the researcher.
- 2) The system recommends k papers for each interest model by the cosine similarity between the model and candidate papers.
- 3) The system punishes each paper with the penalty factor of corresponding interest group.

4) The system resorts all recommend papers by their punished cosine similarity values, and generates the final recommendation results with the first k elements.

IV. EXPERIMENTS

A. Experimental Data

The experimental data set that we used is SPRD (Scholarly Paper Recommendation Data Set) which released by National University of Singapore in 2010, and we got it from [21]. The data set takes the full 597 papers published in ACL main conference from 2000 to 2006 as the candidate papers to recommend. Besides, the data set also contains 15 junior and 13 senior researchers' published papers and reference and citation papers of these papers. What's more, each researcher's interesting list for the candidate papers is also published with the data set. All papers are published in the form of VSM in SPRD. More details of the data set can be seen in Table I.

TABLE I
SOME STATISTICS ABOUT THE DATASET

| Descriptions | Details |
|--------------------------------------|----------------|
| Number | 13 |
| Average number of published papers | 9.5 |
| Average number of interesting papers | 38.7 |
| Average number of citation papers | 10.5(max. 199) |
| Average number of reference papers | 19.4(max. 79) |

The impact factors of the venues involved in our experiments are derived from [22].

B. Experimental Measures

In our experiments, we adopted nDCG [18-19] and MRR [20] measures deriving from IR field to evaluate our proposed model instead of conventional precision and recall rate, because nDCG not only concerns the accuracy but also concerns the position of each item in a recommendation list, and MRR concerns the position where the first relevant item emerge, since many users are more likely to pay attention to some higher positions.

Normalized Discount Cumulative Gain (nDCG)

nDCG is suited to Top-N (in our experiments, N=5, 10) recommendation well, since nDCG could give more weight to highly ranked relevant papers. We use nDCG@5 and nDCG@10 for evaluation respectively. Equation (10) shows the details:

$$nDCG = \frac{1}{i \text{ deal } DCG} \sum_{j=1}^R \frac{2^{r(j)} - 1}{\log(1 + j)}. \quad (10)$$

Where $r(j)$ denotes the score of item j . In our experiments, we only concern whether a candidate paper is relevant with a researcher, so the score is equals to 1 if it is a relevant paper or is equals to 0 if not. The log function, apparently, will give a higher weight to a higher position. $i \text{ deal } DCG$ denotes the ideal evaluation value that a recommendation list can reach.

Mean Reciprocal Rank (MRR)

Most users have a potential habit that they notice higher items more when they browse a recommendation list, so, unlike nDCG, MRR only focuses on the position of the first relevant item in a recommendation list.

$$MRR = \frac{1}{N} \sum_{i=1}^N \frac{1}{r_i}. \quad (11)$$

Where N represents the number of researchers, and r_i represents the position of first relevant item in the recommendation list of researcher i .

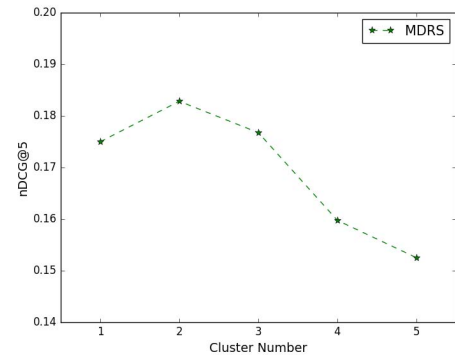
C. Experimental Results

We totally did three experiments. Experiment A is designed to find the best cluster number. Experiment B is in order to find the best penalty coefficient. Experiment C compare our model (MDRS) with baseline model (BM) which was proposed in the paper [1].

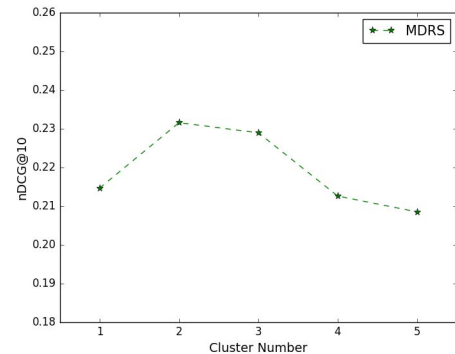
In the process of the experiments, we found that the impact factor of ACL conference is higher than the recommendation standard of each senior researcher, which means the system did not filter any candidate papers for each researcher.

Experiment A

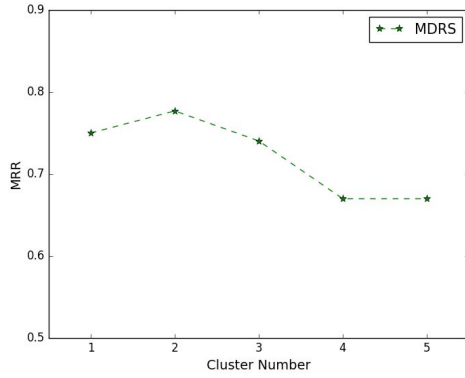
We limit CN (cluster number) from 1 to 5 because we believe that a scholar have low probability with the number of his study directions is more than 5. It should be noted that we constructed MDRS without using penalty factor since our purpose is to identify the best cluster number for scholars. Fig. 1 shows the details:



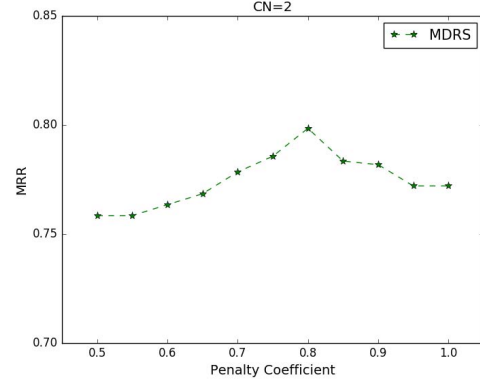
(a). nDCG@5



(b). nDCG@10



(c). MRR
Fig. 1 Experimental Accuracy vs. Cluster Number



(c). MRR
Fig. 2 Experimental Accuracy vs. Penalty Coefficient

We can see that MDRS reaches the best accuracy for nDCG@5, nDCG@10 and MRR measures when cluster number is 2. However, with the continued increase of the clustering number, the recommendation accuracy of the three measures descends rapidly.

Experiment B

We did experiment B to identify the best PC (penalty coefficient) through fixing the cluster number to 2, which is identified by experiment A. We vary the penalty coefficient from 0.5 to 1.0.

As shown in Fig.2, the recommendation accuracies are first increased and then decreased along with the increasing of the penalty coefficient for nDCG@5, nDCG@10 and MRR. Obviously, when penalty coefficient is equals to 0.8, our model reaches the best accuracy for the three measures.

Experiment C

Through experiments A and B, we have already determined the optimal value of cluster number and penalty coefficient. They are 2 and 0.8, respectively. In the third experiments, we did a comparison of MDRS and BM by fixing cluster number and penalty coefficient to these values.

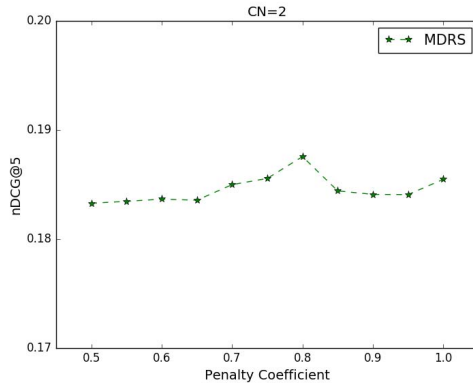
TABLE II
A COMPARISON OF MDRS AND BM UNDER CN=0.2 AND PC=0.8

| | MDRS | BM |
|---------|-------|-------|
| nDCG@5 | 0.188 | 0.174 |
| nDCG@10 | 0.262 | 0.214 |
| MRR | 0.798 | 0.75 |

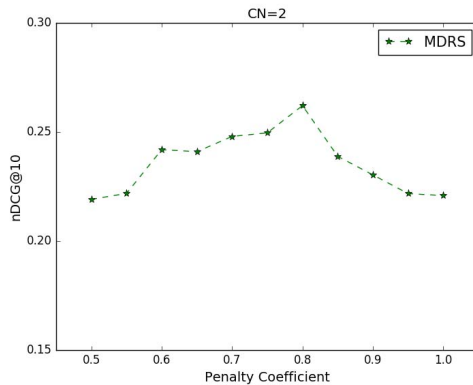
Table 2 shows that the recommendation accuracies of BM are 0.174, 0.214 and 0.75 for nDCG@5, nDCG@10 and MRR, respectively. Obviously, we can see that the recommendation accuracies of MDRS (under the optimal values) are higher than BM. They are 0.188, 0.262 and 0.798 for nDCG@5, nDCG@10 and MRR, respectively, which increase 1.4%, 4.8% and 2.3% respectively than BM.

V. CONCLUSION

In this paper, we have introduced a new kind of literature recommendation method, the so-called MDRS (literature recommendation based on a researcher's multiple research directions and recommendation standard). We have also distinguished a researcher's different attention on his several research interests, by proposing the penalty factor of a scholar's research directions. We adopted nDCG and MRR measures to evaluate our proposed model. From our well-designed experiments, we can conclude that the MDRS has better resolution to scholars who have several research directions. Especially, we observe the highest recommendation accuracy for nDCG and MRR when cluster number is equals to 2 and penalty co-



(a). nDCG@5



(b). nDCG@10

efficient is equals to 0.8, which obviously outperforms the baseline model.

REFERENCES

- [1] Sugiyama K, Kan M Y, "Scholarly paper recommendation via user's recent research interests[C]," *In Proceedings of the 10th annual joint conference on Digital libraries. ACM*, 2010, pp.29-38.
- [2] Sugiyama K, Kan M Y, "Serendipitous recommendation for scholarly papers considering relations among researchers[C]," *In Proceedings of the 2011 Joint International Conference on Digital Libraries, JCDL 2011*, Ottawa, ON, Canada, pp.307-310, June 13-17, 2011.
- [3] Sugiyama K, Kan M Y, "Exploiting potential citation papers in scholarly paper recommendation[C]," *The 13th ACM/IEEE Joint Conference on Digital Libraries (JCDL 2013)*. 2013, pp.153-162.
- [4] Wang Y, Liu J, Dong X L, et al, "Personalized Paper Recommendation Based on User Historical Behavior[J]," *Communications in Computer & Information Science*, 2012, pp.1-12.
- [5] Cui T, Tang X, Zeng Q, "Usage of tagging for research paper recommendation[C]" *Advanced Computer Theory and Engineering (ICACTE)*, 2010 3rd International Conference on. *IEEE*, 2010:V2-439-V2-442.
- [6] Choochaiwattana W, "Usage of tagging for research paper recommendation[C]," *Advanced Computer Theory and Engineering (ICACTE)*, 2010 3rd International Conference on. *IEEE*, 2010:V2-439-V2-442.
- [7] Mcnee S M, Albert I, Cosley D, et al, "On the recommending of citations for research papers[C]," *Proceedings of the 2002 ACM conference on Computer supported cooperative work. ACM*, 2002, pp.116-125.
- [8] Gori M, Pucci A, "Research paper recommender systems: A random-walk based approach[C]," *Web Intelligence, 2006. WI 2006. IEEE/WIC/ACM International Conference on. IEEE*, 2006, pp.778-781.
- [9] Chen C H, Mayanglambam S D, Hsu F Y, et al, "Novelty Paper Recommendation Using Citation Authority Diffusion[C]," *Technologies and Applications of Artificial Intelligence (TAAI)*, 2011 International Conference on. *IEEE*, 2011, pp.126-131.
- [10]Xue H, Guo J, Lan Y, et al, "Personalized paper recommendation in online social scholar system[C]," *Advances in Social Networks Analysis and Mining (ASONAM)*, 2014 IEEE/ACM International Conference on. *IEEE*, 2014, pp. 612-619.
- [11]Qi He, Jian Pei, Daniel Kifer, et al, "Context-aware citation recommendation," *In Proceedings of the 19th international conference on World wide web*, 2010, pp.421-430.
- [12]<http://keg.cs.tsinghua.edu.cn/jietang/>
- [13]Chong Wang , David M Blei, "Collaborative topic modeling for recommending scientific articles," *In Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2011, pp.448-456.
- [14]J. Tang, J. Zhang, and R. Jin, et al, "Topic level expertise search over heterogeneous networks[J]," *Machine Learning Journal*, vol. 82, pp. 211-237, 2011.
- [15]Tang, Jie, Zhang, Jing, Yao, Limin, et al, "ArnetMiner: extraction and mining of academic social networks[C]," *International Conference on World Wide Web, WWW 2008*, Beijing, China, April. 2008, pp.990-998.
- [16]Bolelli L, Ertekin S, Giles C L, "Clustering Scientific Literature Using Sparse Citation Graph Analysis[C]," *Knowledge Discovery in Databases: Pkdd 2006, European Conference on Principles and Practice of Knowledge Discovery in Databases, Berlin, Germany, September 18-22, 2006, Proceedings*. 2006, pp.30-41.
- [17]Liang Y, Li Q, Qian T, "Finding Relevant Papers Based on Citation Relations[C]," *Web-Age Information Management - International Conference, WAIM 2011, Wuhan, China, September 14-16, 2011. Proceedings*. 2011, pp.403-414.
- [18]R. Torres, S. M. Mcnee, M. Abel, J. A. Konstan, and J. Riedl, "Enhancing Digital Libraries with TechLens," *In Proc. of the 4th ACM/IEEE Joint Conference on Digital Libraries (JCDL 2004)*, 2004, pp.228-236.
- [19]K. Järvelin and J. Kekäläinen, "IR Evaluation Methods for Retrieving Highly Relevant Documents," *In Proc. of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2000)*, 2000, pp 41-48.
- [20]E. M. Voorhees, "The TREC-8 Question Answering Track Report," *In Proc. Of the 8th Text Retrieval Conference (TREC-8)*, 1999, pp.77-82.
- [21]<http://www.comp.nus.edu.sg/~sugiyama/SchPaperRecData.html>.
- [22]<http://citeseer.ist.psu.edu/>