

Chronological Scientific Information Recommendation via Supervised Dynamic Topic Modeling

Zhuoren Jiang
College of Transportation Management
Dalian Maritime University
Dalian, China, 116026
jzr_dimu@hotmail.com

ABSTRACT

Scientific information recommendation is crucial to assist scholars for their researches. Citation recommendation is an important field of scientific recommendation. Traditional approaches ignore the chronological nature of the citation recommendation task. In this study, I propose the "Chronological Citation Recommendation," which assumes initial user information need could shift while they are looking for the papers in different time slices. Specifically, I employed a supervised dynamic topic model to characterize the content "time-varying" dynamics and constructed a novel heterogeneous graph that contains dynamic topic-based information, time-decay citation information and word-based information. I applied different meta-paths for different ranking hypotheses, which carried different types of information for citation recommendation in different time slices along with information need shifting. I plan to generate the final "Chronological Citation Recommendation" rankings by feature integration using Learning to Rank. "Chronological Citation Recommendation" will recommend time-series ranking lists based on initial user textual information need. Preliminary experiments on the ACM corpus show that chronological citation recommendation will significantly improve the citation recommendation performance.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval

Keywords

Information Need Shifting, Chronological Citation Recommendation, Heterogeneous Graph

1. INTRODUCTION AND MOTIVATION

As an important research field of scientific recommendation, citation recommendation can help researchers to find and read the potential candidate articles for citation. The past decade has seen the rapid development of citation recommendation in both academic and industrial research. Google Scholar and Microsoft Aca-

ademic Search have been considered as the indispensable tools for researchers. In order to recommend the high-quality papers, a number of different solutions have been proposed, classic textual information retrieval approaches, topic information of the paper content and citation relationships between papers, have proven to enhance recommendation performance effectively.

However, two important characteristics of citation recommendation are ignored in most previous studies. 1) Unlike other web resource recommendations, citation recommendation is more sensitive to the publication date: scholars read and cite recently published papers for innovative ideas or algorithms, and explore classic models and foundation theories in the papers published in the early stage. 2) A researcher's information need could implicitly shift over different time slices, for instance, given the same query, the scholar might be interested in the publications focusing on "*Latent Semantic Indexing*" or "*Probabilistic Latent Semantic Indexing*" between 1990 and 2000 and "*Word Clustering*" or "*Probability Theory*" studies published before 1990. In this study, I define this phenomenon as "**Information Need Shifting**," and it challenges most of the existing text or citation based search methods which generate a general ranking list of papers statistically similar to the initial query.

I propose the "**Chronological Citation Recommendation**" in this study, which will 1) provide a time-series candidate citation recommendation result based on the user's textual queries to solve the problem of time preferences of different research tasks; 2) emphasize the dynamic evolution of scientific content, explicitly characterizing scientific information need dynamics to deal with the impact of "Information Need Shifting." One hypothesis is to consider that the additional dynamic content evolution information could make the scholarly recommendation more accurate and practical. With chronological citation recommendation, researchers can get a set of recommendation paper lists for different time slices (given a textual query), and freely choose the cite paper of their favored time slice. I hope this study can provide researchers a new and distinctive choice for citation recommendation. This study presents two challenges. The first challenge is how to model the "Information Need Shifting." The second challenge is the manner in which to naturally and efficiently combine dynamic content information and the interlinked information.

The contribution of this paper is threefold. First, I propose a novel chronological citation recommendation problem to assist scholars better identify and access the recommendation results and better understand the discipline's development (given a textual query). Chronological citation recommendation, unlike other existing methods, is able to generate multiple lists of ranked papers for a number of time slices. Second, I characterize the information need shifting for citation recommendation by employing time-decayed

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
WSDM '15, February 2–6, 2015, Shanghai, China.
Copyright 2015 ACM 978-1-4503-3317-7/15/02 ...\$15.00.
<http://dx.doi.org/10.1145/2684822.2697036>.

citation relations between papers/topics, dynamic topic-based information, and word information. All of these indicators are addressed by a number of meta-path plus random walk based features on an innovative heterogeneous graph. Last but not least, I plan to apply learning-to-rank algorithms for graph-based feature integration to generate the final time-series recommendation ranking models. Each ranking model will be utilized to optimize the publication ranking results for a distinct time slice. To the best of my knowledge, few prior studies addressed these problems.

2. BACKGROUND AND RELATED WORK

Scientific Information Recommendation

Scientific information recommendation occurs when a scientific publication, venue, or author is recommended to users based on the similarity between the recommended resource and user profiles or samples of in-progress text. Similar to previous studies in text information retrieval and text mining, a number of studies employed text similarity approaches such as paper/author content matching and unsupervised topic modeling [3, 10].

Scholarly or bibliographic networks are another important approach that have also been used to recommend scientific resources. For instance, Shi, Leskovec and McFarland [13] developed citation projection graphs by investigating citations among publications cited by a given paper. In this study, I employed citation importance decay, which was similar to CiteRank [15], an enhanced ranking algorithm compared to PageRank, which enabled the ranking method to estimate the traffic $T_i(\tau_{dir}, \alpha)$ to a given *paper_i*. For this method, a recent paper was more likely to be selected with a probability that was exponentially discounted according to the age of the paper, τ_{dir} .

However, to the best of my knowledge, few existing studies have investigated the chronological citation recommendation problem. In [6], I used a recursive model "dynamic topic/citation influence model" to address this problem, but the processing time of model is not satisfactory. In this study, I will propose a novel method to address this problem effectively.

Dynamic Topic Modeling

Topic dynamics has been recently investigated. Blei and Lafferty [1] proposed a dynamic topic model (DTM), which explicitly characterized the chronological nature of sequential corpora by utilizing a Markov chain of term distributions over time along with multiple LDA-based topic distributions. Based on [1], Gerrish and Blei proposed the document influence model (DIM) [5], this model respected the ordering of the documents, not only tracked how underlying theme has changed over time, but also captured how past articles exhibit varying influence on future articles. Unfortunately, both DTM and DIM may not be an appropriate solution for this task for three reasons. First, given a textual information need, DTM or DIM alone cannot discover optimized dynamic topic chains in terms of user textual information needs. Second, scientific publications are interlinked by citation, which contains important information. DTM or DIM, however, assumes documents are independent of each other, and the topic number, as other unsupervised LDA approaches, are not given. Liu et al., [10] found that the unsupervised LDA approach is not ideal for scientific publications, and scientific metadata, i.e., author assigned keywords, should be used to enhance the topic modeling performance.

In order to solve this problem, I utilized Labeled-LDA (LLDA) [12], which used corpus metadata, i.e., keywords and categories, to enhance the topic modeling process, constrained topic models by defining a one-to-one correspondence between latent topics and labels which solved the problems of topic number setting and topic

interpretability. I proposed my work based on DIM, and added supervised technology into DIM by combining LLDA.

Heterogeneous Graph Mining

The concept of meta-path was first proposed in [14], which could systematically capture the semantic relation between objects in a heterogeneous information network scenario. Different meta-path-based mining tasks were studied including similarity search, relationship prediction, user-guided clustering, and recommendation. It turned out that meta-path served as a very critical feature extraction tool for most of the mining tasks in a heterogeneous information network. In this study, I employed meta-path-based approach as the ranking features, which were extracted from an innovative, dynamic scholarly heterogeneous network. More recently, Lao and Cohen [7] used both supervised and unsupervised methods with the Random Walk with Restart (RWR) algorithm for citation, author, and venue recommendation. Similar approach was also used by [9]. It had been demonstrated that by using the heterogeneous link information in a network, mining functions, such as similarity search, ranking, clustering and classification could be significantly enhanced [14].

3. RESEARCH METHODS

Dynamic Topic Training

As aforementioned, "information need shifting" has a significant impact for citation recommendation, which requires characterization of the "time-varying" content dynamics. This study will model the content dynamics on topic level. Unlike most previous studies that treat topic as a static term-distribution, this model assumes topic can evolve in a historical topic trajectory. I analyze the corpus using a supervised dynamic topic model which is depicted in Figure 1.

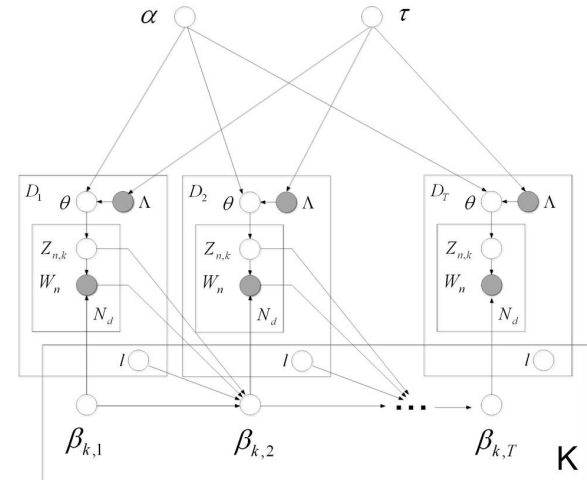


Figure 1: Graphical representation of a supervised dynamic topic model

Here, $\beta_{k,t}$ is considered as the term distribution at time slice t for topic k , $z_{n,k}$ is the indicator that the n^{th} word in document d is assigned to topic k , l is a normally distributed influence score assigned to each article in a specific time slice, N_d represents the word number of document d , θ is a multinomial distribution, represents the topic mixture proportions, α is a Dirichlet prior distribution, Λ is a binary vector of document's labels (keywords metadata), τ represents a labeling prior probability.

This model is based on document influence model (DIM) [5], DIM is a variant of dynamic topic model (DTM) [1], which has been demonstrated as a powerful tool which can explicitly model

the dynamics of the underlying topics and an exploratory tool for discovering influential articles. it introduces a Markov chain of term distributions to capture probabilities that drift over time, where the topics associated with slice t evolve from the topics associated with slice $t-1$:

$$\beta_{k,t} | \beta_{k,t-1}, (w, l, z)_{t-1,1:D} \sim \mathbb{N}(\beta_{k,t-1} + \exp(-\beta_{k,t-1} \sum_d l_{d,k} \sum_n w_{d,n} z_{d,n,k}, \sigma^2 I)).$$

It assumes that the higher the influence, the more the words of the article affect how the topic drifts. While traditional DIM uses LDA [2] to model documents of each time slice, I model them with Labeled LDA (L-LDA) [12], that's only difference between DIM and this model. By combining the supervised technology from L-LDA, I make this model as a supervised dynamic topic model, which overcomes the limitation of traditional DIM: the topics learned from the model are hard to interpret, and the setting of topic number in corpus is arbitrary. L-LDA associates individual words in a document with their most appropriate labels. In this study, the labels come from author-provided metadata such as keywords, and the number of topics in L-LDA is the number of unique labels K in the corpus. The supervised label association is straightforward, I generate the document's labels vector Λ using keyword metadata with the labeling prior probability τ , which can be obtained by observation, then restrict θ to be defined only over the topics that correspond to document's labels. This process will ensure that all the topic assignments are limited to the document's keyword labels. The detailed posterior inference can be found in [5].

By pre-training scholarly corpus using supervised document influence model, I can obtain three time-series distributions in which K is the number of topic in all time slices, T is the number of time slices, time slice T is the latest time slice, time slice 1 is the earliest, and D is the document set. $D = \{D_1, D_2, \dots, D_T\}$, where D_i ($1 \leq i \leq T$) is the candidate documents published in time slice i , and $V = \{w_1, w_2, \dots, w_{|V|}\}$ represents the vocabulary of corpus:

(1) Time-series document-topic distributions (**TDTD**)

$$\theta = \left\{ \theta_{t_1, d_1^1, 1:K}, \theta_{t_2, d_2^1, 1:K}, \dots, \theta_{t_T, d_T^1, 1:K} \right\}$$

Where, $\theta_{t_i, d_m^i, k}$ is k^{th} topic's ($1 \leq k \leq K$) proportion in d_m^{th} document ($1 \leq d_m \leq |D_i|$) published in time slice i ($1 \leq i \leq T$), represents the topic mixture proportions of a document.

(2) Time-series topic-word distributions (**TTWD**)

$$\beta = \left\{ \beta_{t_1, w_1, 1:K}, \beta_{t_2, w_1, 1:K}, \dots, \beta_{t_T, w_{|V|}, 1:K} \right\}$$

Where, $\beta_{t_i, w_n, k}$ is the probability of term w_n ($1 \leq n \leq |V|$) for k^{th} topic ($1 \leq k \leq K$) in time slice i ($1 \leq i \leq T$), represents the relevance between term and topic.

(3) Time-series Document-topic influence distributions (**TDTL**)

$$L = \left\{ l_{t_1, d_1^1, 1:K}, l_{t_2, d_2^1, 1:K}, \dots, l_{t_{T-1}, d_{T-1}^{T-1}, 1:K} \right\}$$

Where, $l_{t_i, d_m^i, k}$ is the influence of d_m^{th} document ($1 \leq d_m \leq |D_i|$) on k^{th} topic ($1 \leq k \leq K$) published in time slice i ($1 \leq i \leq T-1$), represents the document's influence on topic.

These three time-series topic-based distributions will provide the dynamic topic-based information in chronological citation recommendation.

Heterogeneous Graph Construction

In this study, the chronological citation recommendation problem is defined on the sequential scientific publication dataset which can be formatted into a heterogeneous graph. A heterogeneous graph, namely a heterogeneous information network, is a directed graph which contains multiple types of objects and links. In a heterogeneous graph, two different types of vertices can be connected via different paths which always carry different semantic information. A meta-path represents a composite relation over two vertices.

For example, $W - T_{time} - P$ denotes a meta-path between a word and a paper connect by a historical topic that has the word and is related to the paper. [14] gives the formal definition of "Heterogeneous information network," "Network schema" and "Meta-path."

I constructed a novel heterogeneous graph (Figure 2) with three kinds of vertices: "Word," W ; "Historical Topic," T_{time} ; "Paper published in specific time slice," P_{time} . The historical topic is the topic in a specific time slice, that means even the topics share the same label, they are different topics if located in different time slices, i.e. *Language model*_(2000–2002) and *Language model*_(2006–2007) are two different topics in this study.

Figure 2: Constructed heterogeneous graph

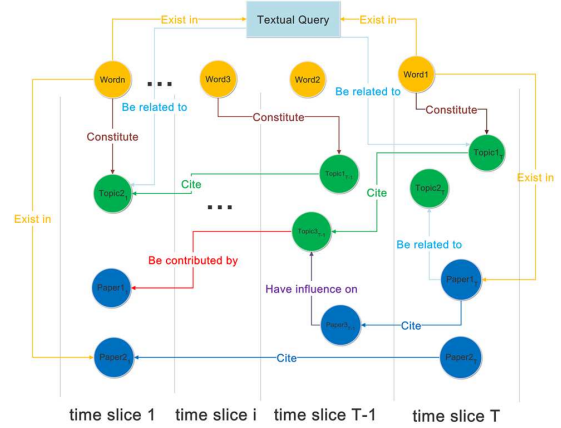


Table 1: Edges in the constructed heterogeneous graph

Edge	Description	Information Carrying
$W \xrightarrow{e} P$	Word exist in paper	Word-based information
$W \xrightarrow{cons} T_{time}$	Word constitute historical topic	Dynamic topic-based information
$T_{time} \xrightarrow{cite} T_{time}$	Historical topic cite historical topic	Time-decay topic citation information
$P_{time} \xrightarrow{cite} P_{time}$	Paper cite paper.	Time-decay paper citation information
$T_{time} \xrightarrow{cont} P_{time}$	Historical topic is contributed by paper	Time-decay paper citation information and dynamic topic-based information
$P_{time} \xrightarrow{infl} T_{time}$	Paper has influence on historical topic	Dynamic topic-based information
$P_{time} \xrightarrow{r} T_{time}$	Paper is related to historical topic	Dynamic topic-based information

Each edge has its own weight. The weight of $W_n \xrightarrow{e} P_i$ is the TF-IDF value of *word_n* in *paper_i*. The weight of $W_n \xrightarrow{cons} T_k^{t_i}$ is the probability $\beta_{t_i, w_n, k}$ from **TTWD**. The weight of $T_m^{t_i} \xrightarrow{cite} T_n^{t_j}$ is the time-decay topic citation probability $\frac{1}{out_{T_m^{t_i}}} * \varepsilon_{t_i}$, $out_{T_m^{t_i}}$ is the number of outlinks of historical topic $T_m^{t_i}$, ε_{t_i} is the time decay factor, which can make sure that the citation from historical topic in recent time slice gives more credits, $\varepsilon_{t_i} = \frac{1}{\log_2(2+T-i)}$. Similarly, the weight of $P_m^{t_i} \xrightarrow{cite} P_n^{t_j}$ is the time-decay paper citation probability $\frac{1}{out_{P_m^{t_i}}} * \varepsilon_{t_i}$. The weight of $P_m^{t_i} \xrightarrow{infl} T_k^{t_i}$ is the influence score of $l_{t_i, d_m^i, k}$ from **TDTL**. The weight of $P_m^{t_i} \xrightarrow{r} T_k^{t_i}$ is the proportion $\gamma_{t_i, d_m^i, k}$ from **TDTD**.

In order to estimate the contribution of each paper to a historical topic, I calculated the paper importance given a historical topic $T_n^{t_j}$ by using PageRank with Prior algorithm [16]. The normalized topical PageRank authority score is used for the weights of

$T_k^{t_j} \xrightarrow{\text{cont}} P_m^{t_i}$. For this step, I used classical homogeneous graph that the vertex is a paper. The citation relationship between the vertices is utilized to calculate the PageRank authority scores. Each vertex is also characterized by a topic prior vector which is $\gamma_{t_i, d_m^i, 1:K}$ from *TDTD*. The result of PageRank (with prior) is the historical topic based paper topic authority vector, i.e., for $P_m^{t_i}$, the result is a paper authority vector $\text{Authority}(P_m^{t_i} | T_{1:K}^{t_j})$, the authority score of $P_m^{t_i}$ given the historical topics $T_{1:K}^{t_j}$. Note that $T_k^{t_j}$ is contributed by $P_m^{t_i}$ ($T_k^{t_j} \xrightarrow{\text{cont}} P_m^{t_i}$) doesn't necessarily mean $P_m^{t_i}$ is related to $T_k^{t_j}$ ($P_m^{t_i} \xrightarrow{r} T_k^{t_j}$). For example, some "topic modeling" papers published in the year of 2003 can be important for historical topic *Information Retrieval*_(2006–2007).

The input of chronological citation recommendation is a piece of text to briefly summarize the information need, i.e., paper abstract or research idea description. I will use starting vertices to represent the user textual query. There are two kinds of starting vertices: W^* and T_{time}^* . W^* is the word in the query, T_{time}^* is the query related historical topics which come from the LDA inference of query [2] for each time slice. The end vertex $P_{time}^?$ is the candidate paper in a specific time slice. The output is several lists of ranked candidate papers in different time slices, which could potentially be cited given user's input.

I will use meta-path from starting vertices to candidate paper vertices in the heterogeneous graph as ranking function. In order to quantify the ranking score of candidate vertices relevant to the starting vertices following the meta-path, I propose the simulated random walk based measure to relevance score: $r(v_i^{(1)}, v_j^{(l+1)}) = \sum_{t=v_i^{(1)} \rightsquigarrow v_j^{(l+1)}} RW(t)$. Where t is a tour from $v_i^{(1)}$ to $v_j^{(l+1)}$ following the meta-path, and $RW(t)$ is the simulated random walk probability of the tour t . Suppose $t = (v_{i1}^{(1)}, v_{i2}^{(2)}, \dots, v_{il+1}^{(l+1)})$, the random walk probability is then $RW(t) = \prod_j w(v_{ij}^{(j)}, v_{i,j+1}^{(j+1)})p(v_{ij}^{(j)})$, $w(v_{ij}^{(j)}, v_{i,j+1}^{(j+1)})$ is the weight of edge $v_{ij}^{(j)} \rightarrow v_{i,j+1}^{(j+1)}$, $p(v_{ij}^{(j)})$ is the vertex prior probability.

Two or more parallel meta-paths leading to the same candidate vertices will be considered. The ranking score of candidate cited paper $P_{time}^?$ is the linear combination of the ranking scores for all sub-meta-paths. The meta-paths preliminarily investigated are listed in Table 2.

Modeling Information Need Shifting via Learning to Rank

In this study, I propose various meta-paths to address different ranking hypotheses where each meta-path carries one or more types of ranking information. It's vital to combine all the information from the meta-paths. Learning-to-rank [8] optimizes ranking performance using machine learning techniques. In the learning-to-rank framework, documents are usually represented by ranking feature vectors. I will use learning-to-rank to combine different meta-path based ranking features, while avoiding manual parameter tuning. Each candidate paper is assigned different ranking scores via various meta-paths. As mentioned in section 1, user information need can vary in different time slices (information need shifting). I will train multi-ranking models, $\{\Phi_1, \Phi_2 \dots \Phi_T\}$ (total T time slices), for different time slices, and I assume the importance of different ranking features can be significantly different for different time slices.

I plan to train the ranking models by using a list-wise algorithm (AdaRank [17], Coordinate Ascent [11] or LambdaMART [4]) on the training data. I will use a scientific candidate publication set and a set of testing paper for training and evaluation. The abstracts of testing papers will be used as the experiment query and the (author-provided) reference lists will be employed as the judgment data. The cite time of different cited papers indicates the im-

portance of each cited paper (for nDCG relevance score). N-fold cross-validation will be used in this study.

4. PRELIMINARY EXPERIMENT

Data

I extracted 45,372 publications with 107,691 citation relations from 6,818 journals and 4,455 conference proceedings or workshops on computer science between 1951 and 2011 (from the *ACM digital library*). I used the title and the abstract information to represent the content of the paper. All the selected papers' contents had more than 80 words after removing stop words and stemming each term to its root. The vocabulary size was 6,517. Most frequent 5% and less than 15 times words were removed. Author-provided keywords were used as topic labels. The initial number of labels was 1,082, regardless of time slice, and the total historical topic number was 8,072. To verify and compare different recommendation algorithms, I split the corpus into eight time slices: pre-1990, 1990-1999, 2000-2002, 2003-2005, 2006-2007, 2008-2009, the year 2010, the year 2011. I generated 2,957,476 word exist-in paper relations, 107,691 paper cite paper relations, 2,242,547 historical topic cite historical topic relations.

For citation recommendation evaluation, I used a test collection with 274 papers. The selected papers met the following conditions: (1) the selected papers were excluded from the 45,372 publication candidate citation collection; (2) each selected paper had more than 15 citations from the candidate citation collection, and (3) each paper's abstract had at least 150 words. The abstract was used as a working context to represent information need, and I recommended citations from the candidate citation collection.

Ranking performance comparison

In the experiment, I applied proposed meta-paths for chronological citation recommendation.

I used test paper's original citations as ground truth with the mention times as citation importance. I also split each test paper's citations into eight time slices according cited paper's publication date. Meanwhile, I employed the following baseline feature groups. 1. **PageRank**: Apply classical PageRank algorithm in paper citation network (query independent). 2. **BM25**: BM25 model for test query. 3. **Language Model with Dirichlet prior smoothing (LM-D)**: Language Model with Dirichlet smoothing for test query. 4. **Language Model with Jelinek-Mercer smoothing (LM-JM)**: Language Model with Jelinek-Mercer smoothing for test query. For each baseline feature, I would first generate the whole ranked paper list, then split the list into eight time slices.

As Table 3 shows, experiment result is promising, there are four meta-paths that outperform all the baseline algorithms with averaged MAP indicator, and all of them are significantly better (with t-test $p < 0.005$ for averaged MAP in all evaluated time slices). Meanwhile, there are also four meta-paths that outperform all the baseline algorithms with averaged nDCG indicator, three of them (Meta-Path 1, Meta-Path 4 and Meta-Path 5) are significantly better (with t-test $p < 0.005$ for averaged nDCG in all evaluated time slices).

5. ISSUES FOR DISCUSSION

My study is still preliminary, there are several research issues need to be discussed:

(1) The quality of experiment dataset. After splitting into different time slices, the citations of test corpus is sparse in some time slices, which may cause "bias" for evaluation (this problem is more critical in more recent slices), is there any better solution to solve this problem? (2) I plan to extract more relations, adjust the weight on edges, implement more meta-paths and use them as ranking fea-

Table 2: The meta-path for chronological citation recommendation

No.				Meta-Path	Publication Ranking Hypothesis
1				$W^* \xrightarrow{e} P_{time} \xrightarrow{cite} P_{time}^?$	The candidate paper is important, if the candidate paper is cited by the papers have the important words in the query
2				$T_{time}^* \xrightarrow{cite} T_{time} \xrightarrow{cont} P_{time}^?$	The candidate paper is important, if the candidate paper contributed the historical topics which are cited by the related historical topics of query
3				$W^* \xrightarrow{e} P_{time}^? \xrightarrow{r} T_{time}^*$	The candidate paper is important, if the important word in the query is important to the candidate paper, and the candidate paper is also very related to the related historical topics of query
4				$W^* \xrightarrow{e} P_{time} \xrightarrow{cite} P_{time}^? \xleftarrow{cont} T_{time}^*$	The candidate paper is important, if the candidate paper is cited by the papers which not only have important word in the query but also contribute the related historical topics of query
5				$W^* \xrightarrow{e} P_{time}^? \xleftarrow{cite} P_{time} \xrightarrow{r} T_{time}^*$	The candidate paper is important, if the candidate paper has important word in the query and is cited by the papers which are related to the related historical topics of query
6				$ \begin{array}{c} T_{time}^* \xrightarrow{cite} T_{time} \xleftarrow{r} P_{time}^? \xrightarrow{infl} T_{time} \xleftarrow{cite} T_{time}^* \\ T_{time}^* \xrightarrow{cite} T_{time} \xrightarrow{cont} \nearrow \quad \quad \quad \nwarrow W^* \end{array} $	The candidate paper is important, if the important word in query also is important to the candidate paper or the candidate paper is related to or has influence on or contributes the historical topics which are cited by the related historical topics of query

■ word information
■ time-decay citation information
■ time-decay topic citation information
■ dynamic topic-based information

Table 3: The Average MAP and nDCG Performance

					MAP10	MAP30	MAP50
				PageRank	0.012500	0.021829	0.028700
				BM25	0.012500	0.015200	0.016114
				LM-D	0.013329	0.015200	0.016129
				LM-JM	0.011229	0.014286	0.015143
				Meta-Path 1	0.029057*	0.045514*	0.055686*
				Meta-Path 2	0.010914	0.019286	0.024671
				Meta-Path 3	0.007071	0.010214	0.011471
				Meta-Path 4	0.026100	0.043857	0.052757
				Meta-Path 5	0.028871	0.038186	0.044429
				Meta-Path 6	0.026886	0.030471	0.032271
				nDCG@10		nDCG@30	nDCG@50
				PageRank	0.022500	0.025900	0.032429
				BM25	0.015129	0.024557	0.025529
				LM-D	0.020229	0.025571	0.027886
				LM-JM	0.019043	0.022557	0.023814
				Meta-Path 1	0.029371	0.037571	0.044329
				Meta-Path 2	0.020243	0.023743	0.031857
				Meta-Path 3	0.024557	0.027471	0.030157
				Meta-Path 4	0.033643	0.054529*	0.062286
				Meta-Path 5	0.044486*	0.053257	0.064800*
				Meta-Path 6	0.040643	0.047643	0.051043

* The best result

tures to apply "Learning to Rank" algorithm to achieve better performance, which model will be suitable for this problem? (3) Is there a better evaluation method to address the chronological citation recommendation problem? (4) Can I explore more sophisticated graph mining methods for chronological citation recommendation?

6. ACKNOWLEDGMENTS

This work is sponsored by the National Natural Science Foundation of China (No. 61202232 and No. 71271034) and China Postdoctoral Science Foundation(2014M551063).

7. REFERENCES

- [1] David M Blei and John D Lafferty. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, pages 113–120. ACM, 2006.
- [2] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- [3] Peter Brusilovsky, Oliviero Stock, and Carlo Strapparava. *Adaptive hypermedia and adaptive Web-based systems*. Springer, 2000.
- [4] Christopher JC Burges. From ranknet to lambdarank to lambdamart: An overview. *Learning*, 11:23–581, 2010.
- [5] Sean Gerrish and David M Blei. A language-based approach to measuring scholarly impact. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 375–382, 2010.
- [6] Zhuoren Jiang, Xiaozhong Liu, and Liangcai Gao. Dynamic topic/citation influence modeling for chronological citation recommendation. In *Proceedings of the 5th International Workshop on Web-scale Knowledge Representation Retrieval & Reasoning*, pages 15–18. ACM, 2014.
- [7] Ni Lao and William W Cohen. Relational retrieval using a combination of path-constrained random walks. *Machine learning*, 81(1):53–67, 2010.
- [8] Tie-Yan Liu. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331, 2009.
- [9] Xiaozhong Liu, Yingying Yu, Chun Guo, Yizhou Sun, and Liangcai Gao. Full-text based context-rich heterogeneous network mining approach for citation recommendation full-text based context-rich heterogeneous network mining approach for citation recommendation full-text based context-rich heterogeneous network mining approach for citation recommendation. In *Joint Conference on Digital Libraries*, 2014.
- [10] Xiaozhong Liu, Jinsong Zhang, and Chun Guo. Full-text citation analysis: A new method to enhance scholarly networks. *Journal of the American Society for Information Science and Technology*, 64(9):1852–1863, 2013.
- [11] Donald Metzler and W Bruce Croft. Linear feature-based models for information retrieval. *Information Retrieval*, 10(3):257–274, 2007.
- [12] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D Manning. Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1*, pages 248–256. Association for Computational Linguistics, 2009.
- [13] Xiaolin Shi, Jure Leskovec, and Daniel A McFarland. Citing for high impact. In *Proceedings of the 10th annual joint conference on Digital libraries*, pages 49–58. ACM, 2010.
- [14] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu. PathSim: Meta path-based top-k similarity search in heterogeneous information networks. In *Proc. 2011 Int. Conf. Very Large Data Bases (VLDB'11)*, Seattle, WA, 2011.
- [15] Dylan Walker, Huafeng Xie, Koon-Kiu Yan, and Sergei Maslov. Ranking scientific publications using a model of network traffic. *Journal of Statistical Mechanics: Theory and Experiment*, 2007(06):P06010, 2007.
- [16] Scott White and Padhraic Smyth. Algorithms for estimating relative importance in networks. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 266–275. ACM, 2003.
- [17] Jun Xu and Hang Li. Adarank: a boosting algorithm for information retrieval. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 391–398. ACM, 2007.