

---

# Citation Recover with Weighted Citation Network Links and Contextual Semantic Information

---

Jonathan C. Barker   Kartik Goyal   Zhengzhong Liu  
Language Technologies Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
{jcbarker, kartikgo, liu}@cs.cmu.edu

## Abstract

In this work, we address the problem of missing citation prediction. Current approaches to this problem includes making prediction through the citation network structure and identifying them through semantic similarity. In this work, we attempt to incorporate both these information to help make better prediction. We spot the phenomenon that not all cited papers are equally important to the citing paper. Thus we adopt a supervised random walk method to learn the appropriate weights on these links using textual semantic information and cross-document topic models. We evaluate our weighted network by predicting missing citations for a paper, given the paper's current citations.

## 1 Introduction

Identifying relevant literature from electronic collection is becoming difficult with the rapid growth of the collection. This makes it very challenging for researchers to conduct efficient and comprehensive literature review. Thus automation of such process become necessary. We consider the problem as a traditional network link prediction problem. Existing citation networks only provide the links between citing and cited papers. Such networks can provide us binary information about whether a link exists between two papers, but fail to encode information about the relative importance and relevance. As a result, it is difficult to draw profound conclusions from these networks about the relationships between different network nodes. We address this problem by assigning weights to network links according to their importance and relevance to the cited paper in order to mine deep information. Tasks like automatic paper recommendation and missing citation prediction can be made possible. Authors can use it to detect any papers they might not have cited in their initial draft, while journal reviewers can be aided in detecting any relevant works which are not cited by the submitted papers.

In practice, various techniques like random walks [12, 13] have been employed to find node relevance in networks. These techniques generally work with a uniformly weighted network and utilize information encoded in graph structure. Our belief is that these techniques can be used more effectively if the network graph were to contain weighted links. The weights on the links would represent the similarity between citing and cited papers according to their semantic correlation, citing patterns, shared authors and temporal proximity. Using these weights should provide more meaningful results for tasks like missing citation prediction when the existing network based techniques are applied. The major challenge in this approach is figuring out how various sources of information can be combined with the network structure. In doing so we must determine the relative importance of these sources/features with respect to the relationship score of two papers. When have finally synthesized this information to create our weighted network we aim to show that we can use it to produce more meaningful results for tasks such as missing citation prediction and recovery.

## 2 Related Work

Nallapati et al. [8] use topic models created with Latent Dirichlet Allocation(LDA) [2] to predict the existence of a citation link between two papers. Apart from generating words, they also generate the reference links from the topic model, thus creating a topical dependency between cited and citing papers.

Backstorm and Leskovec [1] perform supervised random walks on networks to learn the weights on the links between community members. They use this weighted network to predict future links between the members of the social network. The weights are learnt by minimizing an error function on the probability of predicting incorrect and correct future links.

LDA can be used for determining semantic correlations between two documents. Celikyilmaz et.al studied the similarity measures based upon LDA on the QA systems. These similarity measures used both, topic proportions and importance of similar words for estimating similarity.

## 3 Method

**Problem formulation:** A major challenge of this work is how to learn the relevant weights of features of according, we aim to learn a function that can properly assign weights based on the dataset. Here we adopt a method similar to [1] which optimize a function that can assign higher random walk probability to real missing nodes versus other nodes. Consider we are given a uni-directed graph  $G(V, E)$ (citation links only come from one paper to another paper.), a node  $s$  representing the paper we are interested in. Node  $s$  links to several nodes in the network. The task is to identify the other potentially connectable links  $D = d_1, \dots, d_m$  and avoid linking to the  $L = l_1, \dots, l_n$  where there should not be a link. For each edge  $e_{uv} \in E$  we denote the set of pairwise features that associated with node  $(u, v)$  as  $\psi_{uv}$ . By creating multiple pairs of  $d$  and  $l$ , our objective is to maximize the possibility  $p_d$  of recovering nodes in  $D$  and minimize the possibility  $p_l$  of linking to nodes in  $L$ . Here we use a random walk method to calculate these possibilities. In order to direct the random walk to our objective using the features, we assign each edge  $e \in E$  a weight using a parameterized function  $f_w(\psi)$  which combine the feature  $\psi$  with the parameters  $w$ . So the  $w$  is exactly the parameter vector that we are going to learn. Similar to [1], our target is to minimize the function:

$$F(w) = L(w) + h(p_l - p_d) \quad (1)$$

with respect to  $w$  with certain regularization function  $L$  on  $w$  and the loss function  $h$ . By taking the partial derivative of  $F(w)$  we have:

$$\frac{\partial F(w)}{\partial w} = \frac{\partial L(w)}{\partial w} + \sum_{l,d} \frac{\partial(p_l - p_d)}{\partial w} \quad (2)$$

The relation between the possibility of linking to node  $p_u$  can be related with  $w$  using the transition matrix  $Q$ . The basic transition matrix  $Q^{(0)}$  is given by:

$$Q_{uv} = \begin{cases} \frac{f_w(\psi_{uv})}{\sum_t f_w(\psi_{ut})} & (u, v) \in E \\ 0 & \text{otherwise} \end{cases}$$

In our problem we are interested a more general  $Q$  which consider the preference vector, which is given by:

$$Q = (1 - \alpha)Q^{(0)} + \alpha \mathbf{1}(u \in P) \quad (3)$$

where  $P$  is the preference vector.

The random walk score is given by the stationary vector  $p$  where:

$$p^T = p^T Q \quad (4)$$

Eq. 4 and the transition matrix  $Q$  give us the connection between  $w$  and  $p$ . Then we can use the gradient descent method to optimize  $F(w)$ . Detailed derivation can be found in [1].

**Feature selection:** Pairwise features between a cited paper and a citing paper can be gathered from citation patterns, semantic correlation and the network metadata:

1. Number of citances for the cited paper. ("citance" refers to a citing sentence in a paper)

2. Co-citance score. This is defined as the number of times another paper is cited in the same citance as the target cited paper.
3. The number of common papers which both citing and cited papers cite.
4. Topical similarity between the citing and cited paper, calculated using topic proportions from Latent Dirichlet Allocation(LDA).
5. Number of common authors.
6. The difference between time of publication, measured in years.

Features 3, 5 and 6 are meta-data related to the citation network. To determine 1 and 2 we used regular expressions to extract citances and match them with the cited papers. Mallet toolkit [7] was used for running LDA, which uses gibbs sampling for parameter inference, on the raw text of all the papers. The topic proportion vectors were created and were used to calculate cosine similarity between documents.

## 4 Experiment

### 4.1 Experiment setup

We are currently conducting the experiments on the ACL Anthology Network (AAN) [6, 11] which contains the citation network of ACL paper collections and the relevant metadata from year 1965 to 2011 with over 15,000 papers. We choose to use papers from year 2010 data as the training set with papers before 2010 as the candidate set. Similarly, papers from year 2011 data (about 760 papers) are used as the testing set and all other papers are treated as the candidate set (about 14,600 papers). To generate the  $D$  and  $L$  set, we randomly remove links from each testing paper with 25% possibility. The removed links will be treated as the  $D$  set and all other links will non-existing links will be treated as  $L$  set.

As pointed out by [8], normal metrics such as precision and recall are not suitable for this task because we do not expect a paper to cite all relevant papers. We adopt their metric RKL (short for Rank Last), which represent the rank of the last citation we recovered. On the other hand, we are also interested when we can get the first correct prediction, so we introduce another metric RKF (short for Rank First).

### 4.2 General considerations on model:

In the experiments, there are several decisions that we are going to evaluate. (A) Selection of preference vector; (B) Selection of restart possibility  $\alpha$ ; (C) Selection of loss function  $h$ ; (D) Selection of regularization function  $L$  ( $L1$  or  $L2$  norm). By midterm we conduct some simple experiments and the rest will be needed to further evaluated.

**(A) Basic model and preference vector.** The purpose of our method is to mine possible citation papers. Beside from using weights to direct the random walk. The selection of preference vector also plays an important role. Random Walk with Restart (RWR) have proved to be useful in looking for similar nodes in neighbors [1, 13]. However, in a paper the author would not only cite papers that have links with each other. We make an hypothesis that author of a paper should be more likely to cite papers in multiple topic groups corresponding to the topics discussed by the paper. Thus it make sense to set the preference vector based on the topical similarity with the target document. To verify the hypothesis, we conduct some simple experiments. We consider different the combination of random walk and LDA model using Cosine similarity over topic distribution similarity. In addition, we consider different handling over the dangling nodes. In "Strongly preference" mode, we patches all dangling nodes adding transitions following our preference vector but in "Weakly preference" mode, we patches them with a uniform transition towards all other nodes. We use the WebGraph[3] package to do the experiments. Different combinations give us 5 different experiments. The preference vector  $P$  is stochastic, we assign to the target node  $s$  a preference  $\beta$  ( $0 < \beta < 1$ ). Then we assign the rest  $1 - \beta$  weights to all the other nodes in proportion to their similarity with  $s$ . In this experiment,  $\beta$  is set as 0.8 and  $\alpha$  is 0.2 .

The different experiments results are shown in table 1. The general trend shows that "weakly preference" performs better than "strongly preference". On the other hand, we can see that simple LDA

Table 1: Experiment on Basic Model and Preference Vector

Method	Average RKF	Average RKL
RWR (weakly preference)	1084.33	4523.92
RWR (strongly preference)	1692.63	6123.96
LDA similarity prediction	392.57	2460.04
RWR with LDA preference (weakly preference)	<b>151.04</b>	<b>809.23</b>
RWR with LDA preference (strongly preference)	162.83	850.92

similarity outperforms the basic RWR, this is consistent with our hypothesis that author would prefer to cite papers in multiple topical groups. However, the combination of LDA and RWR gives the best performance, which is consistent with our hypothesis. It is also a positive sign that our method of combining semantic and network information is effective.

**(B) Selection of restart possibility.**

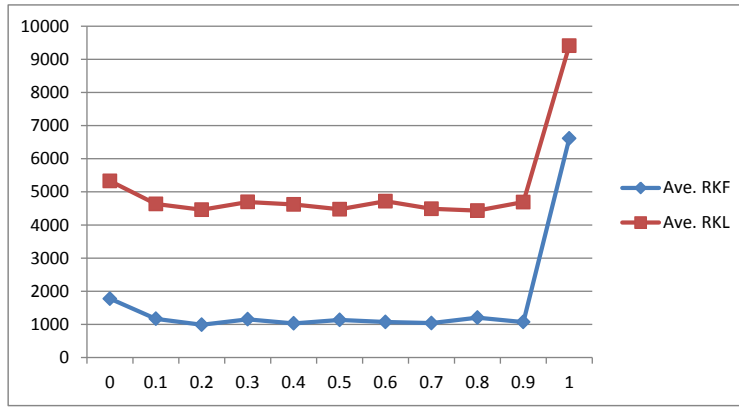


Figure 1: Behavior of different alpha value in basic random walk

**(C) Selection of loss function.**

**(D) Selection of regularization function.**

### 4.3 Observations and other considerations

We observe that self-citation also have impact on the performance

### 4.4 Evaluation result on data set

## 5 Conclusion

## References

- [1] Lars Backstrom and Jure Leskovec. Supervised random walks: predicting and recommending links in social networks. In *Proceedings of the fourth ACM international conference on Web search and data mining*, WSDM '11, pages 635–644, New York, NY, USA, 2011. ACM.

- [2] David M. Blei, Andrew Ng, and Michael Jordan. Latent dirichlet allocation. *JMLR*, 3:993–1022, 2003.
- [3] Paolo Boldi and Sebastiano Vigna. The WebGraph Framework I : Compression Techniques. *Proceedings of the 13th International World Wide Web Conference*, pages 595–602, 2004.
- [4] Asli Celikyilmaz, Dilek Hakkani-Tur, and Gokhan Tur. Lda based similarity modeling for question answering. In *Proceedings of the NAACL HLT 2010 Workshop on Semantic Search*, SS '10, pages 1–9, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [5] Aaron Clauset, Cristopher Moore, and M E J Newman. Hierarchical structure and the prediction of missing links in networks. *Nature*, 453(7191):98–101, May 2008.
- [6] Bryan Gibson Pradeep Muthukrishnan Dragomir R. Radev, Mark Thomas Joseph. A Bibliometric and Network Analysis of the field of Computational Linguistics. *Journal of the American Society for Information Science and Technology*, 2009.
- [7] Andrew Kachites McCallum. Mallet: A machine learning for language toolkit, 2002.
- [8] Ramesh Nallapati, Amr Ahmed, Eric P. Xing, and William Cohen. Joint latent topic models for text and citations. In *KDD*, pages 542–550, 2008.
- [9] Ramesh M. Nallapati, Amr Ahmed, Eric P. Xing, and William W. Cohen. Joint latent topic models for text and citations. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '08*, page 542, New York, New York, USA, 2008. ACM Press.
- [10] Alexandrin Popescul and LH Ungar. Statistical relational learning for link prediction. In *In Workshop on Learning Statistical Models from Relational Data at the International Joint Conference on Artificial Intelligence (2003)*, 2003.
- [11] Dragomir R. Radev, Pradeep Muthukrishnan, and Vahed Qazvinian. The ACL anthology network corpus. In *Proceedings, ACL Workshop on Natural Language Processing and Information Retrieval for Digital Libraries*, Singapore, 2009.
- [12] Purnamrita Sarkar and Andrew W. Moore. Fast dynamic reranking in large graphs. In *Proceedings of the 18th international conference on World Wide Web - WWW '09*, page 31, New York, New York, USA, 2009. ACM Press.
- [13] Hanghang Tong, Christos Faloutsos, and JY Pan. Fast random walk with restart and its applications. In *Sixth International Conference on Data Mining (ICDM'06)*, pages 613–622. Ieee, December 2006.
- [14] Limin Yao, David Mimno, and Andrew McCallum. Efficient methods for topic model inference on streaming document collections. In *KDD*, 2009.

## 6 Appendix