

Decayed DivRank: Capturing Relevance, Diversity and Prestige in Information Networks

Pan Du^{*}, Jiafeng Guo, and Xue-Qi Cheng
Institute of Computing Technology, CAS
Beijing 100190, China
{dupan, guojiafeng}@software.ict.ac.cn, cxq@ict.ac.cn

ABSTRACT

Many network-based ranking approaches have been proposed to rank objects according to different criteria, including relevance, prestige and diversity. However, existing approaches either only aim at one or two of the criteria, or handle them with additional heuristics in multiple steps. Inspired by DivRank, we propose a unified ranking model, Decayed DivRank (DDRank), to meet the three criteria simultaneously. Empirical experiments on paper citation network show that DDRank can outperform existing algorithms in capturing relevance, diversity and prestige simultaneously in ranking.

Categories and Subject Descriptors

H.2.8 [Database Applications]: Data Mining

General Terms

Algorithms, Experimentation

Keywords

Diversity, Relevance, Prestige, Multi-objective Ranking

1. INTRODUCTION

Information networks are widely used to characterize the relationships between data objects. Ranking objects over information networks is a key problem in many real applications, where relevance, prestige and diversity are the three most recognized criteria. Two canonical approaches for network-based ranking are PageRank [1] and HITS [4], which mainly focus on the prestige of objects in ranking. Later, Personalized PageRank [3] was introduced to capture relevance as well as prestige. Recently, diversity in ranking has been widely studied [6] and different methods have been proposed, including GrassHopper [7], MRSP [2] and DivRank [5]. However, existing approaches either only aim at one or two of the three criteria, or handle them with additional heuristics (e.g. greedy vertex selection for diversity) in multiple steps.

In this paper, we propose a novel ranking model, Decayed DivRank (DDRank), to address relevance, prestige and diversity in a unified way. Our work is inspired by the model DivRank proposed by Mei et al. [5]. DivRank is basically a query independent ranking model, which balances prestige and diversity elegantly in ranking through a

^{*}This work was supported by the National High-tech R&D Program of China under grant No.2010AA012500.

vertex-reinforced random walk process. One may turn DivRank into a query dependent model by simply introducing a query dependent prior. However, the reinforce process in DivRank essentially changes the structure of information network. In this way, the relevance between objects and the query, which is fully depend on the primitive network structure, may not be well captured. Different from DivRank, DDRank partially preserves the local relationship between objects around the query, and enhances the competitiveness of these objects in the process. The preserving and reinforcing operations over the network structure “cooperate” to acquire relevance and diversity simultaneously during the time-variant random walk process. Experiments conducted on paper citation network demonstrate the effectiveness of the proposed DDRank in capturing relevance, prestige and diversity simultaneously.

2. DECAYED DIVRANK

Before describing the DDRank approach, we first introduce the original DivRank. The iteration process of DivRank is described as equation 1:

$$f_{i+1}^T = \alpha f_i^T (P_0 N_t) D_i^{-1} + (1 - \alpha) r^T, 0 \leq \alpha \leq 1 \quad (1)$$

where $P_0 = \beta P + (1 - \beta)I$, $0 \leq \beta \leq 1$ and $D_t(i, i) = \sum_{j=1}^n P_0(i, j) N_t(j, j)$, f^T is the ranking score vector and r^T is the prior vector about relevance. P is the primitive transition matrix acquired from the adjacent relationship of a weighted network. P_0 is the new transition matrix on which the vertex-reinforced random walk depend. I is an identity matrix to forge self-links. The self-links in P_0 help to prevent the vertices from losing the profit already acquired during the reinforcement. N_t is a diagonal matrix with each diagonal element recording the visiting times of corresponding object. It acts as the reinforcing factor during the random walk process. Matrix D_t is to re-normalize $P N_t$ into a transition matrix P_t , and to make sure the process will eventually converge.

We now describe our DDRank model, a query dependent ranking model where relevance, prestige and diversity are addressed simultaneously. To capture the relevance, we try to preserve the local structure around the query, and improve the competitiveness of these relevant objects during the DivRank process. For this purpose, we modify the DivRank algorithm by re-weighting the reinforcement on each object according to the relevance between the corresponding object and the query. In this way, we can achieve two goals: 1) If the object is more relevant to the query, it will be more competitive. 2) The competition between objects near the query is weaker than that away from it, which makes

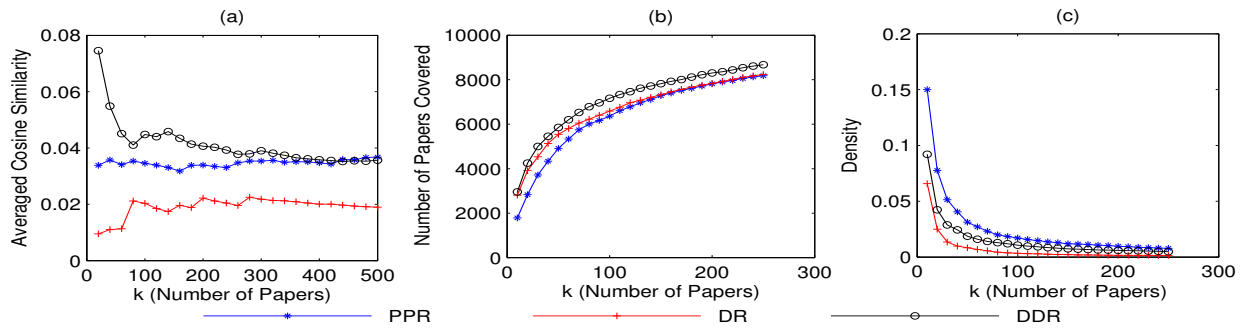


Figure 1: Performance comparison on relevance, prestige and diversity.

tradeoff between relevance and diversity. Therefore, we can balance relevance, prestige and diversity during the vertex-reinforced random walk.

Formally, the DDRank model can be described as follows:

$$f_{t+1}^T = \alpha f_t^T (P_0 N_t^{1-r}) D_t^{-1} + (1 - \alpha) r^T, \quad 0 \leq \alpha \leq 1 \quad (2)$$

where $P_0 = \beta P + (1 - \beta)I$, $0 \leq \beta \leq 1$ and $D_t(i, i) = \sum_{j=1}^n P_0(i, j) N_t(j, j)^{1-r}$. If the network is ergodic, after a sufficiently large t , the reinforced random walk defined by equation 2 also converges to a stationary distribution π . Then this distribution is used to rank the vertices in the information network by DDRank.

From equation 2 we can find that the ratio of N_t between a couple of neighbors (denoted as $\tilde{a} \approx \left[\frac{N_t(1)}{N_t(2)} \right]^{1-r} = a^{1-r}$) is suppressed according to their relevance r because the first order derivative form of \tilde{a} on r is $\tilde{a}'_r = -\ln(a) a^{1-r}$. We have $\tilde{a}'_r > 0$ if $a < 1$, and $\tilde{a}'_r < 0$ if $a > 1$. Similarly, the first order derivative form of reinforcement (denoted as $\tilde{e} = N_t^{(1-r)}$) on r is $\tilde{e}'_r = -\ln(N_t) N_t^{1-r}$. Noting that $N_t \leq 1$, we have $\tilde{e}'_r \geq 0$. This means the more relevant objects are attached more competitiveness than the less relevant ones in DDRank. Relevance r is approximated by the similarity between the query and others in our experiment.

3. EMPIRICAL RESULTS

In this section, we evaluate the effectiveness of DDRank empirically. Due to the space limitation, we only demonstrate the experiments conducted on a paper citation network. The paper citation network is known as the ACL Anthology Network (AAN)¹, which covers 11609 papers, with unit-weighted edges from each paper to the papers it cited. Given a paper, the task aims to find the most relevant, prestigious and diverse papers from the paper-citation network.

Ideally, the contents of the top-K selected papers should be relevant to that of the given paper. We leverages the average cosine similarity of the top-K papers' contents to evaluate the relevance performance. The assumption is that the relevant papers should be similar with each other, since they are all similar with the given paper. Hence, higher average cosine similarity implies good relevance of the top-K papers. Following [5], we use "impact coverage" of the top-K papers as the prestige measure. "Impact Coverage" counts the number of unique papers citing the top-K papers. The basic assumption is that the ideal top-K papers should be cited by as many unique papers as possible. One could notice that the "impact coverage" measure is not a pure prestige measure, where there is also an implicit notion of diversity.

¹Downloadable at <http://clair.si.umich.edu/clair/anthology/>. Density measure evaluate the inverse diversity of the top-K papers. The density of a network is defined as the number of edges presenting in the network divided by the maximal possible number of edges in the network. Hence we can use the density of the subgraph $d(G_K)$ composed by the top-K papers as an inverse measure of diversity in top-K ranked papers. The assumption is that the smaller $d(G_K)$ is, the more independent the top-K papers are, thus higher diversity is contained by the top ranked papers.

We compare our method with personalized PageRank and DivRank with a simple prior, which records the cosine similarity of each paper with the given paper. The performance of relevance, prestige and diversity are demonstrated in Figure 1(a), 1(b) and 1(c) respectively. Personalized PageRank is denoted as PPR, DivRank as DR, and Decayed DivRank as DDR in the figures. From Figure 1, we can see that our approach outperforms Personalized PageRank and DivRank on relevance (as shown in Figure 1(a), the higher the better) and prestige (as shown in Figure 1(b), the higher the better). The performance on diversity (as shown in Figure 1(c), the lower the better) lies between personalized PageRank and DivRank, as can be expected.

In summary, experiments conducted on the paper citation network from AAN show that DDivRank balances the ranking criteria of relevance, diversity and prestige successfully in a unified way. We will perform more thorough theoretical analysis and empirical experiments to make the conclusion more convincing in our future work.

4. REFERENCES

- [1] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. *Comput. Netw. ISDN Syst.*, 30(1-7):107–117, 1998.
- [2] P. Du, J. Guo, J. Zhang, and X. Cheng. Manifold ranking with sink points for update summarization. In *Proceedings of the 19th ACM CIKM*, pages 1757–1760. ACM, 2010.
- [3] T. H. Haveliwala. Topic-sensitive pagerank. In *Proceedings of WWW*, pages 517–526. ACM, 2002.
- [4] J. M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, 1999.
- [5] Q. Mei, J. Guo, and D. Radev. Divrank: the interplay of prestige and diversity in information networks. In *Proceedings of the 16th ACM SIGKDD*, pages 1009–1018. ACM, 2010.
- [6] F. Radlinski, P. N. Bennett, B. Carterette, and T. Joachims. Redundancy, diversity and interdependent document relevance. *SIGIR Forum*, 43:46–52, December 2009.
- [7] X. Zhu, A. Goldberg, J. Van Gael, and D. Andrzejewski. Improving diversity in ranking using absorbing random walks. In *HLT 2007: NAACL*, pages 97–104. ACL, 2007.