

A Ranking System for Scholarly Data Analysis

Author List

SKLSDE Lab, Beihang University, Beijing, China

Beijing Advanced Innovation Center for Big Data and Brain Computing, Beijing, China

{email list}@buaa.edu.cn

Abstract—Scholarly search engine brings great convenience in querying articles and discovering knowledge. Up to now, some systems have been developed with a tremendous amount of scholarly data and providing a range of functions to query keywords. However, these systems barely support different types ranking or various entities ranking (*e.g.*, venues, affiliation and authors) when researchers tend to get better knowledge about the keywords. In the paper, we design and develop a novel ranking system for scholarly data analysis that uses graph engine to efficiently manage the heterogeneous, evolving and dynamic natures of structured scholarly data and integrates different types ranking and various entities ranking based on SARank.

I. INTRODUCTION

Academic publications accelerate the dissemination of scientific discoveries all around the world. Meanwhile, data on academic publications is so huge and keeps growing at explosive way. Building a ranking system for retrieving and analysing academic data is a critical and significant task.

Researchers prefer to retrieve influential and recent works in massive scholarly articles that will inspire their future works. In order to do this, we need to rank articles by making full use of the involved entities (*i.e.*, articles, venues and authors) based on SARank [1] and assessing the relevance of searching conditions and articles.

Generally speaking, a ranking is a function that assigns each item a numerical score. Scholarly articles involve with multi entities such as articles, authors, venues, affiliations and references which form a complex heterogeneous graph. Hence, scholarly articles ranking is essentially a problem of assessing the importance of nodes in a heterogeneous graph.

Ranking such a huge heterogeneous graph and retrieving academic information remain some challenges. Firstly, high quality structured scholarly data is difficult to manage and operate that has exceeded 126 million articles and 114 million authors. Furthermore, academic data keeps increasing at around 5.7 million per year [2]. Secondly, if we are only to rank one type of entities (*i.e.*, scholarly articles), the other type of entities such as venues and authors are closely involved. Moreover, the impact of different types of entities on the ranking of scholarly articles differ from each other. Finally, the importance of entities change with time in a complex manner. Recently published articles are more likely to have a increasing impacts in next few years and those published many years ago tend to have decreasing impacts since researchers potentially interested in latest discoveries.

Actually academic search engine has drawn significant attentions from both academic and industry. Some academic

search engines are prevalent over the past decade *e.g.*, Google Scholar [3], Microsoft Academic [2], Semantic Scholar [4], CiteSeerX [5], Aminer [6] and AceMap [7]. Their common goal is to help researchers discover academic information and present academic data. However, there still remains several shortages in these systems, especially in ranking articles. For example, when rank queried articles some of them focus on citation of articles [6], [7]. While CiteSeerX evaluate articles employing weight PageRank and rank authors by coauthorship networks [8], [9]. They fails to capture the evolving nature of entities and take the involved entities (*i.e.*, venues and authors) into consideration. Furthermore, some of them fail to support to query or rank affiliation, venue information which is also practicable to researchers [3], [7].

Contributions. To this end, we build a novel ranking system for scholarly data analysis which captures the evolving nature of entities and assembles the importance of involved entities. Our system aims to provide better rankings for researchers that realizes a range of functions including: (1) easily and efficiently search academic information by keywords; (2) rank various scholarly entities (*i.e.*, article, author, affiliation and venue) by different ranking type (*i.e.*, SARank, relevance rank, citation and year); (3) check home pages of affiliation, author, venue *etc.*; (4) find detailed information about an article.

(1) To describe the essence of articles reference, we construct a huge heterogenous academic data property graph based on graph engine that achieve great performance in querying academic data.

(2) We develop a ranking model for various entities (*i.e.*, articles, author, affiliation and venue) by evaluating the importance of the entity in academic graph. And the importance captures the temporal nature of entities.

(3) We build a prototype system for ranking and analysing scholarly data that provide better services for researchers.

Organization. The rest of our paper is organized as follows. Section 2 introduces the ranking model including author ranking, venue ranking, affiliation ranking and article ranking. We give a system overview in Section 3, which gives a brief description of our system. And the conclusion is shown in Section 4.

II. RANKING MODEL

We first introduce Time-Weighted PageRank for evaluating the importance of entities, defined as a combination of the prestige and popularity, and then introduce entity ranking

including author ranking, venue ranking, affiliation ranking and type ranking *e.g.*, relevance ranking.

A. Time-Weighted PageRank

PageRank is a typical method in scholarly articles ranking as we can easily make use of the reference between different articles. Due to the following problems, (1) xxx, (2)xxx. We introduce Time-Weighted PageRank(TWPageRank) by extending a time decay factor because of the impact of an article decay with time after peak time.

We present TWPageRank that evaluate the prestige of nodes in a directed graph, in which each node attached with time information. And we use *the impact weight* to describe the relative weight from the edge sources to targets. Formally, the impact weight on a directed edge (u, v) , *i.e.*, an edge u from v , is defined as:

$$w(u, v) = \begin{cases} 1 & T_u < Peak_v \\ e^{\sigma(T_u - Peak_v)} & T_u \geq Peak_v \end{cases} \quad (1)$$

where T_u is the time of node u , $Peak_v$ is the peak time of node v after which the impact weights of edges to v decay with time, and σ is a negative number controlling the decaying speed of the impacts. By default, we use years as its time granularity in Eq. (1).

Thus, the update rule of Time-weighted PageRank is

$$PR(v) = d \sum_{(u,v) \in E} \frac{w(u, v)PR(u)}{W(u)} + \frac{1-d}{n} \quad (2)$$

where $PR(u)$ and $PR(v)$ are the TWPageRank score of u and v . And E is a set of edges, $W(u) = \sum_v w(u, v)$ is the sum of the impact weights on all edges from u , n is the number of nodes and d is a damping parameter in $(0, 1)$.

B. Entity Ranking and Type Ranking

Scholarly entities(*e.g.*, affiliation, venue, author and article) ranking is a problem of assessing the importance of nodes in a heterogeneous network. The importance is a combination of *prestige* and *popularity* to capture the evolving nature of entities. The prestige of scholarly entities is derived by applying TWPageRank on the citation graph G , and each type of entity is assigned the corresponding TWPageRank score as its prestige score Prs .

To learn about the popularity of different scholarly entities, we first introduce the popularity of an article. The popularity of an article v is the sum of all its citation freshness, *i.e.*, the closeness to the current year:

$$Pop(v) = \sum_{(u,v) \in E^c} e^{\sigma(T_0 - T_u)} \quad (3)$$

Here, T_0 is the current year, T_u is the largest year among all articles, σ is the negative decaying factor in Eq. (1).

Intuitively, prestige favors those with many citations soon after the publication of articles or associated articles of venues and authors, and popularity capture the temporal nature of entities.

Affiliation Ranking. We computes the importance of affiliations by their associated articles. As the importance of an affiliation evolves with time, we treat the affiliation importance in each year individually, and its importance is the sum of importance in all individual years.

We construct an affiliation graph $G^a(V^a, E^a)$ using the citation information among affiliations, in which a node represents an affiliation in a specific year and a direct edge (s, t) means that there exists articles of affiliation s citing articles of affiliation t . Thus, the impact weights are defined as sums of impact weights from affiliation s to affiliation t , *i.e.*,

$$w_a(s, t) = \sum_{u \in C(s), v \in C(t)} w(u, v) \quad (4)$$

Here, $C(s)$ and $C(t)$ are the sets of articles of affiliation s and affiliation t , and $w(u, v)$ is the impact weight of articles u and v .

The prestige of an affiliation in a specific year (Prs_a) is computed using the impact weights Eq. (4) and the update rule in Eq. (2). The popularity of an affiliation in a specific year (Pop_a) is defined as the average popularity of its articles that is computed using Eq. (3). Thus, the *affiliation importance score* (Imp_a) is defined as a combination of its prestige and popularity:

$$Imp_a(v) = Prs_a(v)^\lambda Pop_a(v)^{1-\lambda} \quad (5)$$

Here, $\lambda \in [0, 1]$ is the importance factor, indicates the weighting about prestige and popularity.

Venue Ranking. We computes the importance of venue using their associated articles which is similar with affiliation ranking. We treat the venue importance in each year, and construct a venue graph $G^v(V^v, E^v)$ using citation information among venues. And then we combine the prestige of a venue (Prs_v) and the popularity of venue (Pop_v) as the *venue importance score* (Imp_v).

Author Ranking. We evaluate the importance of each author, and compute the average importance of the authors of an article as its *author importance score*. However, it is obvious that the author citation graph is too large to handle. Hence, we evaluate the prestige, popularity of the author by using the average prestige, popularity of all her/his published articles, respectively. Then, the author importance score (Imp_{aut}) is the combination of its prestige and popularity, similar to affiliation ranking.

Article Ranking. If we are only to rank scholarly articles, the other type of entities such as venues and authors are closely involved. Hence, we assemble the importance of article, venue and author to produce the final scholarly articles ranking, illustrated in Fig. 1. Venue ranking and author ranking have presented in previous paragraphs. Next, introduce how to compute the importance of article using citation information.

We first construct a citation graph $G^c(V^c, E^c)$ using citation information among articles. The prestige of articles (Prs_c) is derived by using TWPageRank in citation graph G^c . The popularity of an article (Pop_c) is the sum of all its freshness which has described in Eq. (3). The importance of citation component

(Imp_c) is a combination of its prestige and popularity in the citation graph by applying Eq. (5).

Thus, the static ranking score of an article v is aggregated as follows:

$$S(v) = \alpha Imp_v(v) + \beta Imp_{aut}(v) + (1 - \alpha - \beta) Imp_c(v) \quad (6)$$

Here parameter α , β and $1 - \alpha - \beta$ regularize the contributions of the venue, author and citation information.

Relevance Ranking. We have introduced affiliation, venue, author, article ranking in the former sections. These rankings are query independent and aim to give a static ranking based on scholarly data only. However, it is vital to evaluate the similarity between the short query and sentence (*i.e.*, title of the paper) when retrieve articles by keywords.

Hence, we employ distributional semantic approach to represent words. Similarities between the term vectors indicates the corresponding semantic similarities [10]. The final ranking score of an article v that relates to semantics of article title, defined as follow:

$$F(v) = \theta S(v) + (1 - \theta) \sum_j idf(Q_j) \frac{Q_j \cdot T}{\|Q_j\| \|T\|} \quad (7)$$

Here, $S(v)$ is a static ranking score of article, $Q(j)$ is the j th of query keywords vector which have removed stop words, $idf(Q_j)$ is inverse document frequency measures how much information the word $Q(j)$ provides, T suggests we aggregate the title word vectors to their centroid and θ means the relevance factor.

III. SYSTEM OVERVIEW

As state in former section, we build a ranking system for scholarly data analysis. Fig. 1 shows the framework of our system, which consists of three main components, *Storage*, *Query Engine* and *Visualizer* respectively. And we demonstrate some scenarios that can be supported by our system.

A. System Framework

1) *Storage*: The component of storage consists of property graph and transaction. Heterogeneous academic data is stored as property graph model which possesses nodes and relationships, and we use transaction to ensure the predictability of relationship-based queries.

We apply a popular graph database Neo4j [11] to manage and operate massive scholarly data based on the following reasons. First, scholarly data is naturally structured data connected by the reference relationship. At the same time, graph database has great performance implementation handling connections. Second, Neo4j is provided with highly performant read and write scalability thanks to native graph storage and processing. Thus, it works well in managing scholarly data and searching subgraphs.

In order to manage effectively and query efficiently with Neo4j, we design a graph schema mainly based on two principles. (1) Nodes for things and relationships for structure. (2) We take into consideration of the query ability of the graph schema and adopt specific trade-offs. Thus, we model

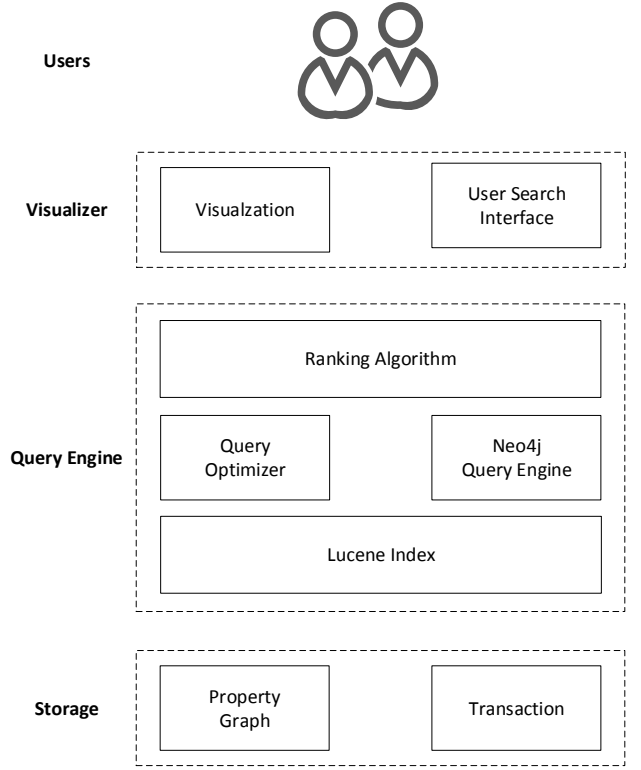


Figure 1. system framework

academic data as a huge heterogeneous graph, shown in Fig. 2, which contains more than one billion nodes and over two billion relationships.

The graph schema contains seven type of entities: author, affiliation, field of study, paper, venue(journal and conferences *e.g.*, KDD, ICDE), conference instance (*e.g.*, KDD 2019, ICDE 2019), years. In order to query efficiently, we introduce an additional node PAA to represent the paper-author-affiliation n-ary relationships. Intuitively, a paper get published in a journal/conference by the author means the edges among paper, author and venue node.

By employing ranking model as stated in section 2, we derive affiliation, author, venue and article ranking score using incremental computation [1]. And those score is described as a property in the graph schema. In fact, we can apply any algorithms to rank scholarly entities in the graph schema.

2) *Query Engine*: Query engine is the main component that is responsible for handling scholarly data. It consists of lucene index, query optimizer, Neo4j query engine and ranking algorithm. Lucene index is especially important for retrieving articles from Neo4j. Obviously, we index article's title after using stop words, and index the distributed vector representation of words.

Lucene Index.

why we need index. index is important for xxx what we indexed?

1. title after using stop word in desk.

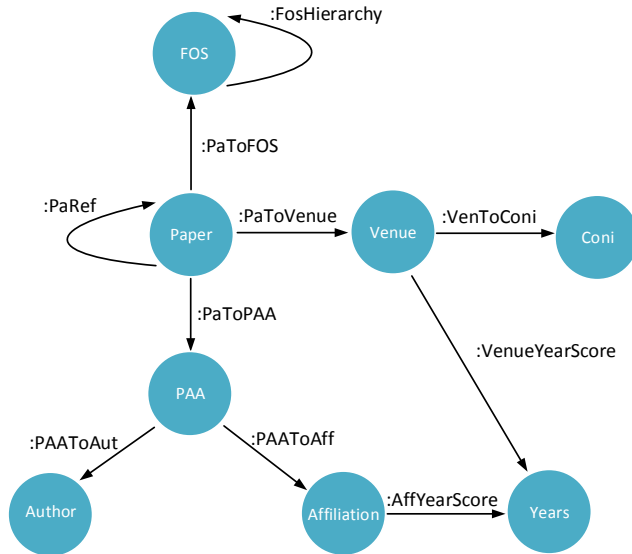


Figure 2. Neo4j schema design

2. vector, word representation

Query Optimizer and

Neo4j Query Engine

Ranking Algorithm

relevance, ranking results.

3) Visualizer: visualization

use interface

search affiliation, author, venue, articles

visualization for author, articles?

B. System Demonstration

demonstration graph

IV. CONCLUSION

conclusion

REFERENCES

- [1] S. Ma, C. Gong, R. Hu, D. Luo, C. Hu, and J. Huai, "Query independent scholarly article ranking," in *ICDE 2018: 34th IEEE International Conference on Data Engineering*, 2018.
- [2] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B.-j. P. Hsu, and K. Wang, "An overview of microsoft academic service (mas) and applications," in *Proceedings of the 24th international conference on world wide web*. ACM, 2015, pp. 243–246.
- [3] "Google scholar," <https://scholar.google.com/>.
- [4] "Semantic scholar," <https://www.semanticscholar.org/>.
- [5] H. Li, I. G. Councill, W.-C. Lee, and C. L. Giles, "Citeseerx: an architecture and web service design for an academic document search engine," in *Proceedings of the 15th international conference on World Wide Web*, 2006, pp. 883–884.
- [6] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su, "Arnetminer: extraction and mining of academic social networks," in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2008, pp. 990–998.
- [7] Z. Tan, C. Liu, Y. Mao, Y. Guo, J. Shen, and X. Wang, "Acemap: A novel approach towards displaying relationship among academic literatures," in *WWW '16 Companion Proceedings of the 25th International Conference Companion on World Wide Web*, 2016, pp. 437–442.

- [8] Y. Sun and C. L. Giles, "Popularity weighted ranking for academic digital libraries," in *ECIR'07 Proceedings of the 29th European conference on IR research*, 2007, pp. 605–612.
- [9] D. FIALA, "From citeseer to citeseer x: Author rankings based on coauthorship networks," *Journal of Theoretical & Applied Information Technology*, vol. 58, no. 1, 2013.
- [10] G. S. Corrado, J. Dean, K. Chen, and T. Mikolov, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [11] "Neo4j," <https://neo4j.com/>.