Constrained-meta-path-based ranking in heterogeneous information network

Chuan Shi, Yitong Li, Philip S. Yu & Bin Wu

Knowledge and Information Systems An International Journal

An International Journal

ISSN 0219-1377 Volume 49 Number 2

Knowl Inf Syst (2016) 49:719-747 DOI 10.1007/s10115-016-0916-1





Your article is protected by copyright and all rights are held exclusively by Springer-Verlag London. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your article, please use the accepted manuscript version for posting on your own website. You may further deposit the accepted manuscript version in any repository, provided it is only made publicly available 12 months after official publication or later and provided acknowledgement is given to the original source of publication and a link is inserted to the published article on Springer's website. The link must be accompanied by the following text: "The final publication is available at link.springer.com".





REGULAR PAPER

Constrained-meta-path-based ranking in heterogeneous information network

Chuan Shi¹ · Yitong Li¹ · Philip S. Yu² · Bin Wu¹

Received: 7 July 2014 / Revised: 22 November 2015 / Accepted: 11 January 2016 / Published online: 28 January 2016 © Springer-Verlag London 2016

Abstract Recently, there is a surge of interests on heterogeneous information network analysis, where the network includes different types of objects or links. As a newly emerging network model, heterogeneous information networks have many unique features, e.g., complex structure and rich semantics. Moreover, meta path, the sequence of relations connecting two object types, is widely used to integrate different types of objects and mine the semantics information in this kind of networks. The object ranking is an important and basic function in network analysis, which has been extensively studied in homogeneous networks including the same type of objects and links. However, it is not well exploited in heterogeneous networks until now, since the characteristics of heterogeneous networks introduce new challenges for object ranking. In this paper, we study the ranking problem in heterogeneous networks and propose the HRank method to evaluate the importance of multiple types of objects and meta paths. Since the traditional meta path coarsely embodies path semantics, we propose a constrained meta path to subtly capture the refined semantics through confining constraints on objects. Based on a path-constrained random walk process, HRank can simultaneously determine the importance of objects and constrained meta paths through applying the tensor analysis. Extensive experiments on three real datasets show that HRank can effectively evaluate the importance of objects and paths together. Moreover, the constrained meta path shows its potential on mining subtle semantics by obtaining more accurate ranking results.

Keywords Heterogeneous information network · Ranking · Random walk · Tensor analysis

Chuan Shi shichuan@bupt.edu.cn

Beijing Key Lab of Intelligent Telecommunications Software and Multimedia, Beijing University of Posts and Telecommunications, Beijing, China

² University of Illinois at Chicago, Chicago, IL, USA

1 Introduction

It is an important research problem to evaluate object importance or popularity, which can be used in many data mining tasks. Many methods have been developed to evaluate object importance, such as PageRank [16], HITS [9] and SimRank [6]. In this literature, objects ranking is done in a homogeneous network in which objects or relations are the same. For example, both PageRank and HITS rank the web pages in WWW.

However, in many real network data, there are many different types of objects and relations, which can be organized as heterogeneous network. Formally, heterogeneous information networks (HIN) are the logical networks involving multiple types of objects as well as multiple types of links denoting different relations [4]. For example, the movies recommendation data include multiple types of objects: movies, actors and directors and their relations [17]. It is clear that heterogeneous information networks are ubiquitous and form a critical component of modern information infrastructure [4]. Recently, many data mining tasks have been exploited in this kind of networks, such as similarity measure [17,22], clustering [23] and classification [7], among which ranking is an important but not yet exploited task.

Figure 1a shows an HIN example in bibliographic data, and Fig. 1b illustrates its network schema which depicts object types and their relations. In this example, it contains objects from four types of objects: papers (P), authors (A), labels (L, categories of papers) and conferences (C). There are links connecting different types of objects. The link types are defined by the relations between two object types. For example, links exist between authors and papers denoting the writing or written-by relations, between conferences and papers denoting the published-in relations. In this network, several interesting, yet seldom exploited, ranking problems can be proposed.

- One may be interested in the importance of one type of objects and ask the following questions:
 - *Q. 1.1 Who are the most influential authors?*
 - Q. 1.2 Who are the most influential authors in data mining field?
- As we know, some object types have an effect on each other. For example, influential authors usually publish papers in reputable conferences. So one may pay attention to the importance of multiple types of objects simultaneously and ask the following questions: *Q. 2.1 Who are the most influential authors and which reputable conferences did those influential authors publish their papers on?*

Q. 2.2 Who are the most influential authors and which reputable conferences did those influential authors publish their papers on in data mining field?



Fig. 1 A heterogeneous information network example on bibliographic data. **a** shows heterogeneous objects and their relations, **b** shows the network schema

 Furthermore, one may wonder which factor mostly affects the importance of objects, since the importance of objects is affected by many factors. So he may ask the questions like this:

Q. 3 Who are the most influential authors and which factors make those most influential authors be most influential?

Although the ranking problem in homogeneous networks has been well studied, the above ranking problems are unique in HIN (especially Q. 2 and Q. 3), which are seldom studied until now. Since there are multiple types of objects in HIN, it is possible to analyze the importance of multiple types of objects (i.e., Q. 2) as well as affecting factors (i.e., Q. 3) together. However, the ranking analysis in HIN also faces the following research challenges.

- There are different types of objects and links in HIN. If we simply treat all objects equally and apply the random walk as PageRank does in homogeneous network, the ranking result will mix different types of objects together.
- Different types of objects and links in heterogeneous networks carry different semantic meanings. The random walk along different meta paths has different semantics, which may lead to different ranking results. Here the meta path [22] means a sequence of relations between object types. So a desirable ranking method in HIN should be pathdependent.

In this paper, we study the ranking problem in HIN and propose a ranking method, HRank, to evaluate the importance of multiple types of objects and meta paths in HIN. For *Q. 1* and *Q. 2*, a path-based random walk model is proposed to evaluate the importance of single or multiple types of objects. The different meta paths connecting two types (same or different types) of objects have different semantics and transitive probability, and thus lead to different random walk processes and ranking results. Although meta path has been widely used to capture the semantics in HIN [17,22], it coarsely depicts object relations. By employing the meta path, we can answer the *Q. 1.1* and *Q. 2.1*, but cannot answer the *Q. 1.2* and *Q. 2.2*. For example, "*Author–Paper–Author*" describes the collaboration relation among authors. However, it cannot depict the fact that Philip S. Yu and Jiawei Han have many collaborations in data mining field, but they seldom collaborate in information retrieval field. In order to overcome the shortcoming existing in meta path, we propose the *constrained meta path* concept, which can effectively describe this kind of subtle semantics. The constrained meta path, we can answer the *Q. 1.2* and *Q. 2.2*.

Moreover, in HIN, based on different paths, the objects have different ranking values. The comprehensive importance of objects should consider all kinds of factors (the factors can be embodied by constrained meta paths), which have different contributions to the importance of objects. For example, although Jiawei Han and W. Bruce Croft both are influential authors in computer science, the achievements on data mining and information retrieval fields contribute to their reputation, respectively. In order to evaluate the importance of objects and meta paths simultaneously (i.e., answer Q. 3), we further propose a co-ranking method which organizes the relation matrices of objects on different constrained meta paths as a tensor. A random walk process is designed on this tensor to co-rank the importance of objects and paths simultaneously. That is, random walkers surf in the tensor, where the stationary visiting probability of objects and meta paths is considered as the HRank score of objects and paths. In addition, in order to speed up the matrix multiplication process in HRank, we design three fast computation strategies whose effectiveness has been validated by experiments.

In all, this paper has the following contributions.

- We propose the constrained meta path concept to describe the subtle semantic relation in HIN. Compared to the original meta path, the constrained meta path can depict object relation with finer granular through setting constraint condition on meta path.
- We propose a path-based ranking method to evaluate the importance of same or different types of objects in HIN by setting constrained meta path. Extensive experiments not only validate that the objects have different importance based on different constrained meta paths, but also show that the ranking results of constrained meta paths more comply with our common sense.
- A co-ranking method is proposed to simultaneously evaluate the importance of objects and paths. The method not only can comprehensively evaluate the importance of objects by considering all constrained meta paths, but also can rank the contribution of different constrained meta paths. The experiments on two real datasets illustrate that the proposed method can accurately identify the importance of objects and the corresponding paths.

The preliminary work has been published in [13]. However, this paper substantially extends the original work in the following aspects. First, to improve the efficiency of HRank, we design three fast computation strategies of which the effectiveness is verified by the corresponding experiments. Second, we propose a new version of HRank on symmetric constrained meta paths to extend the capability, and the added experiments validate its effectiveness. Third, the added qualitative and quantitative experiments are provided to extensively validate the effectiveness and efficiency of HRank.

The rest of the paper is organized as follows. In Sect. 2, we summarize and compare the related work. In Sect. 3, we describe notations in this paper and some preliminary knowledge. In Sect. 4, we present the proposed method, and the fast computation strategies are introduced in Sect. 5. Extensive experiments are done to validate the proposed method in Sect. 6. Finally, Sect. 7 concludes this paper.

2 Related work

Ranking is an important data mining task in network analysis. Many ranking methods have been proposed. For example, PageRank [16] evaluates the importance of objects through a random walk process, HITS [9] ranks objects using the authority and hub scores, and SimRank [6] evaluates the similarity of two objects by their neighbors' similarities. The recently proposed RoleSim measures the role similarity between any two nodes from network structure [8]. These approaches only consider the same type of objects in homogeneous networks, so they cannot be applied in heterogeneous networks. To rank tweets effectively by capturing the semantics and importance of different linkages, Huang et al. [5] propose the Tri-HITS model to iteratively propagate ranking scores across heterogeneous networks, it only focuses on ranking one type of objects.

Some researches have begun to pay attention to the co-ranking on multiple types of objects. For example, Zhou et al. [27] co-rank authors and their publications by coupling two random walk processes, and the co-HITS [3] incorporates the bipartite graph with the content information and the constraints of relevance. Soulier et al. [21] propose a bi-type entity ranking algorithm to rank jointly documents and authors in a bibliographic network by combining content-based and network-based features. Although these methods can rank different types of objects existing in HIN, they are restricted to bipartite graphs. Recently, MultiRank [14] determines the importance of both objects and relations simultaneously for

723

multi-relational data, and HAR [12] is proposed to determine hub and authority scores of objects and relevance scores of relations in multi-relational data for query search. These two methods focus on same type of objects with multi-relations, not multiple types of objects.

In recent years, there is a surge on the HIN analysis. Many data mining tasks have been exploited in HIN, such as similarity measure [17,22], clustering [19,23] and classification [10]. As a unique feature of HIN, the links connecting different types of objects contain semantics. So the meta path [22], connecting object types via a sequence of relations, has been widely used to capture the relation semantics. Sun et al. [22] put forward the concept of meta path to describe the rich semantic relations, and studied similarity search on symmetric meta paths. As an extension of Sun's work, Yu et al. [26] use a meta-path-based ranking model ensemble to represent semantic meanings for similarity queries. HeteSim is also proposed by Shi et al. [17] to measure the relevance scores of heterogeneous objects in HIN. PathSelClus [23] integrates meta path selection with user-guided clustering to cluster objects in networks. Kong et al. [10] develop an HCC solution to assign labels to a group of instances that are interconnected through different meta paths. In addition, Shi et al. [20] propose a semantic path-based personalized recommendation method SemRec to explore various semantics. Although meta path may convey semantic information in HIN, it is too coarse to capture the subtle semantics in some applications.

3 Preliminary

In this section, we describe notations used in this paper and present some preliminary knowledge.

A heterogeneous information network is a special type of information network with the underneath data structure as a directed graph, which contains either multiple types of objects or multiple types of links.

Definition 3.1 (*Information Network* [22]) Given a schema $S = (A, \mathcal{R})$ which consists of a set of entity types $\mathcal{A} = \{A\}$ and a set of relations $\mathcal{R} = \{R\}$, an information network is defined as a directed graph G = (V, E) with an object type mapping function $\varphi : V \to \mathcal{A}$ and a link type mapping function $\psi : E \to \mathcal{R}$. Each object $v \in V$ belongs to one particular object type $\varphi(v) \in \mathcal{A}$, and each link $e \in E$ belongs to a particular relation $\psi(e) \in \mathcal{R}$. When the types of objects $|\mathcal{A}| > 1$ or the types of relations $|\mathcal{R}| > 1$, the network is called *heterogeneous information network*; otherwise, it is a *homogeneous information network*.

In heterogeneous information networks, there are multiple object types and relation types. We use the network schema to depict the object types and the relations existing among object types. For a relation *R* existing from type *S* to type *T*, denoted as $S \xrightarrow{R} T$, *S* and *T* are the *source type* and *target type* of relation *R*, which is denoted as *R*.*S* and *R*.*T*, respectively. The inverse relation R^{-1} holds naturally for $T \xrightarrow{R^{-1}} S$. Generally, *R* is not equal to R^{-1} , unless *R* is symmetric and these two types are the same. Figure 1b shows a network schema of bibliographic information network, which describes the object types and their relations in the HIN.

Different from homogeneous networks, two objects in a heterogeneous network can be connected via different paths and these paths have different meanings. As an example shown in Fig. 1b, authors can be connected via "Author–Paper–Author" (*APA*) path, "Author–Paper–Conference–Paper–Author" (*APCPA*) path and so on. These paths are called meta paths which can be defined as follows.

Definition 3.2 (*Meta path* [22]) A meta path \mathcal{P} is a path defined on a schema $S = (\mathcal{A}, \mathcal{R})$ and is denoted in the form of $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_l} A_{l+1}$ (abbreviated as $A_1A_2 \dots A_{l+1}$), which defines a composite relation $R = R_1 \circ R_2 \circ \dots \circ R_l$ between type A_1 and A_{l+1} , where \circ denotes the composition operator on relations.

It is obvious that semantics underneath these paths are different. The *APA* path means authors collaborating on the same papers (i.e., co-author relation), while the *APCPA* path means the authors' papers publishing on the same conferences (i.e., co-conference relation). Based on different meta paths, there are different relation networks, which may result in different importance of objects. For example, the importance of authors under *APA* path has bias on the authors who write many papers having many authors, while the importance of authors under *APCPA* path emphasizes the authors who publish many papers on those productive conferences. So the importance of objects depends on the meta path in the heterogeneous networks. As an effective semantic capturing method, the meta path has been widely used in many data mining tasks in HIN, such as similarity measure [17,22], clustering [23] and classification [10]. However, meta path fails to capture some subtle semantics. Taking Fig. 1b as an example, the *APA* path cannot reveal the co-author relations in a certain research field, such as data mining and information retrieval. Although Jiawei Han has co-work many papers with Philip S. Yu in the data mining field, they never co-work in the operation system field. The *APA* path cannot subtly reflect this difference.

In order to overcome the shortcomings in meta path, we propose the concept of constrained meta path, defined as follows.

Definition 3.3 (*Constrained meta path*) A constrained meta path is a meta path based on a certain constraint which is denoted as CP = P|C. $P = (A_1A_2...A_l)$ is a meta path, while C represents the constraint on the objects in the meta path.

Note that the C can be one or multiple constraint conditions on objects. Taking Fig. 1b as an example, the constrained meta path APA|P.L = "DM" represents the co-author relations of authors in data mining field through constraining the label of papers with DM. Similarly, the constrained meta path APCPA|P.L = "DM" &&C = "CIKM" represents the co-author relations of authors in CIKM conference and the papers of authors are in data mining field. Obviously, compared to meta path, the constrained meta path conveys richer semantics by subdividing meta paths under distinct conditions. Particularly, when the length of meta path is 1 (i.e., a relation), the constrained meta path degrades to a *constrained relation*. In other words, the constrained relation confines constraint conditions on objects of the relation.

For a relation $A \xrightarrow{R} B$, we can obtain its transition probability matrix as follows.

Definition 3.4 (*Transition probability matrix*) W_{AB} is an adjacent matrix between type A and B on relation $A \xrightarrow{R} B$. U_{AB} is the normalized matrix of W_{AB} along the row vector, which is the transition probability matrix of $A \xrightarrow{R} B$.

Then we make some constraints on objects of the relation $A \xrightarrow{R} B$ (i.e., constrained relation). We can have the following definition.

Definition 3.5 (*Constrained transition probability matrix*) W_{AB} is an adjacent matrix between type A and B on relation $A \xrightarrow{R} B$. Suppose there is a constraint C on object type A. The constrained transition probability matrix U'_{AB} of constrained relation R|C is $U'_{AB} = M_C U_{AB}$, where M_C is the constraint matrix generated by the constraint condition C on object type A.

The constraint matrix M_C is usually a diagonal matrix whose dimension is the number of objects in object type A. The element in the diagonal is 1 if the corresponding object satisfies the constraint, else the element in the diagonal is 0. For example, in the path PC|C = "CIKM", M_C is a diagonal matrix of conferences, where the "CIKM" column is 1 and the others are 0. Similarly, we can confine the constraint on the object type B or both types. Note that the transition probability matrix is a special case of the constrained transition probability matrix, when we let the constraint matrix M_C be the identity matrix I.

Given a network G = (V, E) following a network schema S = (A, R), we can define the meta path-based reachable probability matrix as follows.

Definition 3.6 (*Meta path-based reachable probability matrix*) For a meta path $\mathcal{P} = (A_1A_2...A_{l+1})$, the meta path-based reachable probability matrix PM is defined as $PM_{\mathcal{P}} = U_{A_1A_2}U_{A_2A_3}...U_{A_lA_{l+1}}$. $PM_{\mathcal{P}}(i, j)$ represents the probability of object $i \in A_1$ reaching object $j \in A_{l+1}$ under the path \mathcal{P} .

Similarly, we have the following definition for constrained meta path.

Definition 3.7 (*Constrained meta path-based reachable probability matrix*) For a constrained meta path $CP = (A_1A_2 \dots A_{l+1}|C)$, the constrained meta path-based reachable probability matrix is defined as $PM_{CP} = U'_{A_1A_2}U'_{A_2A_3} \dots U'_{A_lA_{l+1}}$. $PM_{CP}(i, j)$ represents the probability of object $i \in A_1$ reaching object $j \in A_{l+1}$ under the constrained meta path P|C.

In fact, if there is no constraint on the objects of a relation $A_i \xrightarrow{R} A_{i+1}$, $U'_{A_iA_{i+1}}$ is equal to $U_{A_iA_{i+1}}$. If there is a constraint on the objects, we only consider the objects that satisfy the constraint. For simplicity, we use the reachable probability matrix and the $M_{\mathcal{P}}$ to represent the constrained meta path-based reachable probability matrix in the following section.

4 The HRank method

Since the importance of objects is related to the meta path designated by users, we propose the path-based ranking method HRank in heterogeneous networks. In order to answer the three ranking problems proposed in Sect. 1, we design three versions of HRank, respectively.

4.1 Ranking based on symmetric constrained meta paths

In order to evaluate the importance of one type of objects (i.e., Q. I), we design the HRank-SY method based on symmetric constrained meta paths, since the constrained meta paths connecting one type of objects are usually symmetric, such as APA|P.L = "DM".

For a symmetric constrained meta path $\mathcal{P} = (A_1A_2 \dots A_l|C)$, \mathcal{P} is equal to \mathcal{P}^{-1} and A_1 and A_l are the same. Similar to PageRank [16], the importance evaluation of object A_1 (i.e., A_l) can be considered as a random walk process in which random walkers wander from type A_1 to type A_l along the path \mathcal{P} . The HRank value of object A_1 (i.e., $R(A_1|\mathcal{P})$) is the stable visiting probability of random walkers, which is defined as follows:

$$R(A_1|\mathcal{P}) = \alpha R(A_1|\mathcal{P})M_{\mathcal{P}} + (1-\alpha)E \tag{1}$$

where $M_{\mathcal{P}}$ is the constrained meta path-based reachable probability matrix as defined above. *E* is the restart probability vector for convergence. It is set equally for all objects of type A_1 , which is $1/|A_1|$. α is the decay factor, which can be set with 0.85 as the parameter experiments suggested. HRank-SY and PageRank both have the same idea that the importance of objects



is decided by the visiting probability of random surfers. Different from PageRank, the random surfers in HRank-SY should wander along the constrained meta path to visit objects.

As shown in Fig. 2, the red broken line illustrates an example of the process of calculating rank values, where the CP is APA|P.L = "DM". The concrete calculating process is as follows:

$$R(\text{Author}|\mathcal{CP}) = \alpha R(\text{Author}|\mathcal{CP})M_{\mathcal{CP}} + (1-\alpha)E$$

$$M_{\mathcal{CP}} = U_{AP}^{'}U_{PA}^{'} = U_{AP}M_{P}M_{P}U_{PA}$$
(2)

where M_P is the constraint matrix on object type P (paper).

4.2 Ranking based on asymmetric constrained meta paths

For the question Q. 2, we propose the HRank-AS method based on asymmetric constrained meta paths, since the paths connecting different types of objects are asymmetric. For an asymmetric constrained meta path $\mathcal{P} = (A_1 A_2 \dots A_l | \mathcal{C})$, \mathcal{P} is not equal to \mathcal{P}^{-1} . Note that A_1 and A_l are either of the same or different types, such as APC|P.L = "DM" and PCPLP|C = "CIKM".

Similarly, HRank-AS is also based on a random walk process that random walkers wander between A_1 and A_l along the path. The ranks of A_1 and A_l can be seen as the visiting probability of walkers, which are defined as follows:

$$R(A_{l}|\mathcal{P}^{-1}) = \alpha R(A_{1}|\mathcal{P})M_{\mathcal{P}} + (1-\alpha)E_{A_{l}}$$

$$R(A_{1}|\mathcal{P}) = \alpha R(A_{l}|\mathcal{P}^{-1})M_{\mathcal{P}^{-1}} + (1-\alpha)E_{A_{1}}$$
(3)

where $M_{\mathcal{P}}$ and $M_{\mathcal{P}^{-1}}$ are the reachable probability matrix of path \mathcal{P} and \mathcal{P}^{-1} . E_{A_1} and E_{A_l} are the restart probability of A_1 and A_l . Obviously, HRank-SY is the special case of HRank-AS. When the path \mathcal{P} is symmetric, Eq. 3 is the same with Eq. 1.

The blue broken line in Fig. 2 illustrates an example which simultaneously evaluates the importance of authors and conferences. Here the CP is APC|P.L = "DM". The concrete calculating process is as follows:

$$R(\text{Conf.}|\mathcal{CP}) = \alpha R(\text{Aut.}|\mathcal{CP})M_{\mathcal{CP}} + (1-\alpha)E_{\text{Conf.}}$$

$$R(\text{Aut.}|\mathcal{CP}) = \alpha R(\text{Conf.}|\mathcal{CP})M_{\mathcal{CP}^{-1}} + (1-\alpha)E_{\text{Aut.}}$$

$$M_{\mathcal{CP}} = U'_{AP}U'_{PC} = U_{AP}M_PM_PU_{PC}$$

$$M_{\mathcal{CP}^{-1}} = U'_{CP}U'_{PA} = U_{CP}M_PM_PU_{PA}$$
(4)

where M_P is the constraint matrix on object type P (paper).

Deringer



Fig. 3 An example of multi-relations of objects generated by multiple paths: **a** the graph representation; **b** the corresponding tensor representation

4.3 Co-ranking for objects and relations in HIN

Until now, we have created methods to rank same or different types of objects under a certain constrained meta path. However, there are many constrained meta paths in heterogeneous networks. It is an important issue to automatically determine the importance of paths [22,23], since it is usually hard for us to identify which relation is more important in real applications. To solve this problem (i.e., Q. 3), we propose the HRank-CO to co-rank the importance of objects and relations. The basic idea is based on an intuition that important objects are connected to many other objects through a number of important relation and important relations connect many important objects. So we organize the multiple relation networks with a tensor and a random walk process is designed on this tensor. The method not only can comprehensively evaluate the importance of objects by considering all constrained meta paths, but also can rank the contribution of different constrained meta paths.

In Fig. 3a, we show an example of multiple relations among objects, generated by multiple meta paths. There are three objects of type A, three objects of type B and three types of relations among them. These relations are generated by three constrained meta paths with type A as the source type and type B as the target type. To describe the multiple relations among objects, we use the representation of tensor which is a multidimensional array. We call $X = (x_{i,j,k})$ a third-order tensor, where $x_{i,j,k} \in R$, for i = 1, ..., m, j = 1, ..., l, k = 1, ..., n. $x_{i,j,k}$ represents the times that object i is related to object k through the jth constrained meta path. For example, Fig. 3b is a three-way array, where each two-dimensional slice represents an adjacency matrix for a single relation. So the data can be represented as a tensor of size $3 \times 3 \times 3$. In the multi-relational network, we define the transition probability tensor to present the transition probability among objects and relations.

Definition 4.1 (*Transition probability tensor*) In a multi-relational network, X is the tensor representing the network. F is the normalized tensor of X along the column vector. R is the normalized tensor of X along the tube vector. T is the normalized tensor of X along the row vector. F, R and T are called the transition probability tensor which can be denoted as follows:

$$f_{i,j,k} = \frac{x_{i,j,k}}{\sum_{i=1}^{m} x_{i,j,k}}$$
 $i = 1, 2, ..., m$

Deringer

$$r_{i,j,k} = \frac{x_{i,j,k}}{\sum_{j=1}^{l} x_{i,j,k}} \qquad j = 1, 2, \dots, l$$

$$t_{i,j,k} = \frac{x_{i,j,k}}{\sum_{k=1}^{n} x_{i,j,k}} \qquad k = 1, 2, \dots, n$$
(5)

 $f_{i,j,k}$ can be interpreted as the probability of object *i* (of type *A*) being the visiting object when relation *j* is used and the current object being visited is object *k* (of type *B*), $r_{i,j,k}$ represents the probability of using relation *j* given that object *k* is visited from object *i*, and $t_{i,j,k}$ can be interpreted as the probability of object *k* being visited, given that object *i* is currently the visiting object and relation *j* is used. The meaning of these three tensors can be defined formally as follows:

$$f_{i,j,k} = \operatorname{Prob}(X_t = i | Y_t = j, Z_t = k)$$

$$r_{i,j,k} = \operatorname{Prob}(Y_t = j | X_t = i, Z_t = k)$$

$$t_{i,j,k} = \operatorname{Prob}(Z_t = k | X_t = i, Y_t = j)$$
(6)

in which X_t , Z_t and Y_t are three random variables representing visiting at certain object of type A or type B and using certain relation, respectively, at the time t.

Now, we define the stationary distributions of objects and relations as follows

$$x = (x_1, x_2, \dots, x_m)^{T}$$

$$y = (y_1, y_2, \dots, y_l)^{T}$$

$$z = (z_1, z_2, \dots, z_n)^{T}$$
(7)

in which

$$x_{i} = \lim_{t \to \infty} \operatorname{Prob}(X_{t} = i)$$

$$y_{j} = \lim_{t \to \infty} \operatorname{Prob}(Y_{t} = j)$$

$$z_{k} = \lim_{t \to \infty} \operatorname{Prob}(Z_{t} = k)$$
(8)

From the above equations, we can get:

$$Prob(X_{t} = i) = \sum_{j=1}^{l} \sum_{k=1}^{n} f_{i,j,k} \times Prob(Y_{t} = j, Z_{t} = k)$$

$$Prob(Y_{t} = j) = \sum_{i=1}^{m} \sum_{k=1}^{n} r_{i,j,k} \times Prob(X_{t} = i, Z_{t} = k)$$

$$Prob(Z_{t} = k) = \sum_{i=1}^{m} \sum_{j=1}^{l} t_{i,j,k} \times Prob(X_{t} = i, Y_{t} = j)$$
(9)

where $\operatorname{Prob}(Y_t = j, Z_t = k)$ is the joint probability distribution of Y_t and Z_t , $\operatorname{Prob}(X_t = i, Z_t = k)$ is the joint probability distribution of X_t and Z_t , and $\operatorname{Prob}(X_t = i, Y_t = j)$ is the joint probability distribution of X_t and Y_t .

To obtain x_i , y_j and z_k , we assume that X_t , Y_t and Z_t are all independent from each other which can be denoted as below:

$$Prob(X_t = i, Y_t = j) = Prob(X_t = i)Prob(Y_t = j)$$

$$Prob(X_t = i, Z_t = k) = Prob(X_t = i)Prob(Z_t = k)$$

$$Prob(Y_t = j, Z_t = k) = Prob(Y_t = j)Prob(Z_t = k)$$
(10)

🖄 Springer

Consequently, through combining the equations with the assumptions above, we get

$$x_{i} = \sum_{j=1}^{l} \sum_{k=1}^{n} f_{i,j,k} y_{j} z_{k}, \quad i = 1, 2, ..., m,$$

$$y_{j} = \sum_{i=1}^{m} \sum_{k=1}^{n} r_{i,j,k} x_{i} z_{k}, \quad j = 1, 2, ..., l,$$

$$z_{k} = \sum_{i=1}^{m} \sum_{j=1}^{l} t_{i,j,k} x_{i} y_{j}, \quad k = 1, 2, ..., n.$$

(11)

The equations above can be written in a tensor format:

$$x = Fyz, \quad y = Rxz, \quad z = Txy \tag{12}$$

with

$$\sum_{i=1}^{m} x_i = 1, \quad \sum_{j=1}^{l} y_j = 1, \text{ and } \sum_{k=1}^{n} z_k = 1.$$

According to the analysis above, we can design the following algorithm to co-rank the importance of objects and relations.

Algorithm 1 HRank-CO Algorithm

Input: Three tensors *F*, *T* and *R*, three initial probability distributions x_0 , y_0 and z_0 and the tolerance ϵ . **Output:** Three stationary probability distributions *x*, *y* and *z*.

Procedure: Set *t* = 1; **repeat** Compute $x_t = Fy_{t-1}z_{t-1}$; Compute $y_t = Rx_tz_{t-1}$; Compute $z_t = Tx_ty_t$; **until** $||x_t - x_{t-1}|| + ||y_t - y_{t-1}|| + ||z_t - z_{t-1}|| < \epsilon$

4.4 Discussion

First we analyze the connection of three versions of HRank. We have stated that HRank-SY is a special version of HRank-AS when the asymmetric path degrades to a symmetric path. We can also find that HRank-AS is the special version of HRank-CO. When there is only one relation in HRank-CO generated by path \mathcal{P} , T and F are the transition probability matrices between type A and B along path \mathcal{P} and \mathcal{P}^{-1} (i.e., $M_{\mathcal{P}}$ and $M_{\mathcal{P}^{-1}}$), respectively. Moreover, R and y become 1. In this condition, Eq. 12 turns into Eq. 3 without considering the restarting probability.

Then we estimate the space and time complexity of HRank. For simplicity, we assume that there are *r* relations, *n* objects for each type and *t* iterations for convergence. For HRank-CO, the space complexity is $O(rn^2)$ to store the transition probability tensor. The time complexity of HRank-CO comes from two parts: iteratively compute rank values (see Algorithm 1) and construct the transition probability tensor (see Definition 4.1). The time complexity of rank computation is $O(trn^2)$. For the *l* length path, the complexity of constructing the transition

probability tensor is $O(rln^3)$. So the whole time complexity is $O(trn^2 + rln^3)$. For HRank-AS and HRank-SY, the number of relations (i.e., r) is 1, so their space complexity are $O(n^2)$ and time complexity are $O(tn^2 + ln^3)$. For real applications, the relation matrices are very sparse, so the real time complexity is much smaller than the theoretical analysis.

5 Fast computation strategies

HRank has a high computation demand for time, and it is not affordable for online query in large-scale information networks. So the fast computation of HRank is necessary and important to improve its efficiency. According to the time complexity analysis above, we can find that the main time-consuming component of HRank lies in constructing the reachable probability matrix with the complexity $O(rln^3)$. Therefore, the main idea is to speed up the matrix multiplication process to construct the reachable probability matrix.

The main challenge of fast computation is the trade-off of accuracy and efficiency. The high accuracy may limit the efficiency improvement, while the high-efficiency improvement may lead to the loss of accuracy. In real applications, we may have different requirements under different situations. As a result, we design three memory-based fast computation strategies to satisfy various scenarios. When the network size in real applications is too large to be contained in memory, we need to design the parallelized version of HRank with parallel models (e.g., MapReduce [2] or BSP [25]). That is our future work.

5.1 Dynamic programming strategy

As we know, the matrix multiplication obeys the associative property, i.e., $(M_1 \times M_2) \times M_3$ is equal to $M_1 \times (M_2 \times M_3)$. However, the different sequences of matrix multiplication have different running time. For example, for three two-dimensional matrix multiplication $(M_1(9, 2) \times M_2(2, 9)) \times M_3(9, 2)$ (the numbers in the parenthesis represent the row and column number of matrix), it needs $9 \times 2 \times 9 + 9 \times 9 \times 2$ (i.e., 324) addition operations, while it only needs $2 \times 9 \times 2 + 9 \times 2 \times 2$ (i.e., 72) addition operations for $M_1(9, 2) \times$ $(M_2(2, 9) \times M_3(9, 2))$. So we can design a **Dyn**amic **P**rogramming strategy (DynP) to speed up matrix multiplication through changing their computation order. The basic idea of DynP is to assign small-dimensioned matrix with the high computation priority. For a meta path $\mathcal{P} = R_1 \circ R_2 \circ \cdots \circ R_l$, we can calculate the expected minimal number of addition operations by the following equation and record the computation sequence in *i*.

$$C(R_{1} \circ \dots \circ R_{l}) = \begin{cases} 0 & l = 1 \\ |R_{1}.S| \times |R_{1}.T| \times |R_{2}.T| & l = 2 \\ \arg\min_{i} \{C(R_{1} \cdots R_{i}) + C(R_{i+1} \cdots R_{l}) + \\ |R_{1}.S| \times |R_{i}.T| \times |R_{l}.T| \} & l > 2 \end{cases}$$
(13)

where the function *C* represents the number of addition operations, and $|R_i.S|$ ($|R_i.T|$) represents the row (column) number of the relation matrix R_i . The above equation can be easily solved with the $O(l^2)$ complexity through adopting the dynamic programming method. The running time can be omitted, since *l* is much smaller than the matrix dimension. The DynP does not change relevance results, so it is an information-lossless strategy.

5.2 Truncation strategy

The basic hypothesis of Truncation strategy (Trun) is that removing the probability on those less important nodes would not significantly degrade the performance. It has been proved by many researches [1,11]. Through keeping matrix sparse, the truncation strategy could greatly reduce the amount of space and time consumption. We can add a truncation step at each step of the matrix multiplication. In the truncation step, the probability value is set with 0 for those nodes whose values are smaller than a threshold ε . Although a static threshold is usually used in many methods (e.g., ref. [11]), it faces the following disadvantage: It may truncate nothing for matrix whose elements all have high probability and it may truncate most nodes for matrix whose elements all have low probability. Since we usually pay close attention to the top *k* objects in query task, the threshold ε can be set as the top *k* relevance value for each search object. The *k* is dynamically adjusted as follows.

$$k = \begin{cases} L & \text{if } L \le W \\ \lfloor (L - W)^{\beta} \rfloor + W(\beta \in [0, 1]) & \text{otherwise} \end{cases}$$

where *L* is the vector length and *W* is the number of top objects, decided by users. The *W* and β determine the truncation level. The larger *W* or β will cause the larger *k*, which means a denser matrix. It is expensive to determine the top *k* relevance value for each object, so we can estimate the value by the top *kM* value for the whole matrix (*M* is the number of objects). However, it is also time-consuming to calculate the top *kM* value for the whole matrix. It can be approximated with the sample data from the raw matrix. The sample ratio is γ . The larger γ leads to more accurate approximation with longer running time. In summary, the truncation strategy is an information-loss strategy. It can keep matrix sparse with small sacrifice on accuracy. In addition, it needs additional time to estimate the threshold ε .

5.3 Monte Carlo strategy

The basic idea of Monte Carlo method (MonC) is that, for each node u, K independent random walkers are simulated starting from u. The distribution of u is approximated by the normalized counts of the number of times the random walkers visit a node. So the reachable probability $PM_{\mathcal{P}}(a, b)$ can be approximated by the normalized count of the number of times that the walkers visit the node b from a along the path \mathcal{P} .

$$PM_{\mathcal{P}}(a, b) = \frac{\text{#times the walkers visit } b \text{ along } \mathcal{P}}{\text{#walkers } from a}$$

The number of walkers from a (i.e., K) controls the accuracy and amount of computation. The larger K will achieve more accurate estimation with more time cost. An advantage of the MonC strategy is that its running time is not affected by the dimension and sparsity of matrix. However, the high-dimensional matrix needs larger K for high accuracy. As a sampling method, the MonC is also an information-loss strategy.

6 Experiments

In this section, we do experiments to validate the effectiveness of three versions of HRank on three real datasets, respectively.



Fig. 4 The network schema of three heterogeneous datasets. **a** DBLP bibliographic dataset, **b** ACM bibliographic dataset, **c** IMDB movie dataset

6.1 Datasets

We use three heterogeneous information networks for our experiments, including DBLP dataset, ACM dataset and IMDB dataset. They are summarized as follows:

DBLP dataset [17,22] The DBLP dataset is a sub-network collected from DBLP Web site ¹ involving major conferences in two research areas: database (DB) and information retrieval (IR), which naturally form two labels. The dataset contains 9682 authors, 20 conferences (or journals) and 22,185 papers which are all labeled with one of the two research areas. The network schema is shown in Fig. 4a.

ACM dataset [17] The ACM dataset was downloaded from ACM digital library² in June 2010. The ACM dataset comes from 14 representative computer science conferences: KDD, SIGMOD, WWW, SIGIR, CIKM, SODA, STOC, SOSP, SPAA, SIGCOMM, MobiCOMM, ICML, COLT and VLDB. These conferences include 196 corresponding venue proceedings (e.g., KDD conference includes 12 proceedings, such as KDD'10 and KDD'09). The dataset has 12,499 papers, 17431 authors, 1903 terms and 1804 author affiliations. The network also includes 73 labels of these papers in ACM category (e.g., H.2 is database management). The network schema of ACM dataset is shown in Fig. 4b.

IMDB dataset [18] We crawled movie information from the Internet movie database³ to construct the network. The IMDB movie data collect 1591 movies before 2010. The related objects include movies, actors, directors and movie types, which are organized as a star schema shown in Fig. 4c. Movie information includes 5324 actors, 1591 movies, 551 directors and 112 movie types (e.g., comedy and romance).

6.2 Ranking of homogeneous objects

Since the homogeneous objects are connected by symmetric constrained meta paths, the experiments validate the effectiveness of HRank-SY on symmetric constrained meta paths.

6.2.1 Experiment study on symmetric constrained meta paths

This experiment ranks the same type of objects by designating a symmetric constrained meta path on ACM dataset. Here we rank the importance of authors through the symmetric meta

¹ http://www.informatik.uni-trier.de/~ley/db/.

² http://dl.acm.org/.

³ www.imdb.com/.

Rank	APA	APA P.L = "H.3"	APA P.L = "H.2"	PageRank	Degree
1	Jiawei Han	W. Bruce Croft	Jiawei Han	Ming Li (1522)	Jiawei Han
2	Philip Yu	ChengXiang Zhai	Christos Faloutsos	Wei Wei (2072)	Philip Yu
3	Christos Faloutsos	James Allan	Philip Yu	Jiawei Han (5385)	ChengXiang Zhai
4	Zheng Chen	Jamie Callan	Jian Pei	Tao Li (6090)	Zheng Chen
5	Wei-Ying Ma	Zheng Chen	H. Garcia-Molina	Hong-Jiang Zhang (6319)	Christos Faloutsos
6	ChengXiang Zhai	Ryen W. White	Jeffrey F. Naughton	Wei Ding (6354)	Ravi Kumar
7	W. Bruce Croft	Wei-Ying Ma	Divesh Srivastava	Jiangong Zhang (7285)	W. Bruce Croft
8	Scott Shenker	Jian-Yun Nie	Raghu Ramakrishnan	Christos Faloutsos (7895)	Wei-Ying Ma
9	H. Garcia- Molina	Gerhard Weikum	Charu C. Aggarwal	Feng Pan (8262)	Gerhard Weikum
10	Ravi Kumar	C. Lee Giles	Surajit Chaudhuri	Hongyan Liu (8440)	Divesh Srivastava

 Table 1
 Top ten authors of different methods on ACM dataset

The number in the parenthesis of the fifth column means the rank of authors in the whole ranking list returned by PageRank

path *APA*, which considers the co-author relations among authors. We also employ two constrained meta paths APA|P.L = "H.2" and APA|P.L = "H.3", where the categories of ACM *H*.2 and *H*.3 represent "database management" and "information storage/retrieval," respectively. That is, two constrained meta paths subtly consider the co-author relations in database/data mining field and information retrieval field, respectively. We employ HRank-SY to rank the importance of authors based on these three paths. As the baseline methods, we rank the importance of authors with PageRank and the degree of authors (called Degree method). We directly run PageRank on the whole ACM network by ignoring the heterogeneity of objects. Since the results of PageRank mix all types of objects, we select the author type from the ranking list as the final results.

The top ten authors of each method are given in Table 1. We can find that these ranking lists all have some common influential authors except that of PageRank. The results of PageRank include some not very well-known authors in DB/IR field, such as Ming Li and Wei Wei, although they may be very influential in other fields. We know that the PageRank values of objects are decided by their degrees to a large extent, so the rank values of affiliation objects are high due to their high degrees. It improves the rank values of author objects connecting multiple high-ranking affiliations. The bad results of PageRank show that the ranking in heterogeneous networks should consider the heterogeneity of objects. Otherwise, it cannot distinguish the effect of different types of links. Moreover, we can also observe that the results of HRank with constrained meta paths have obvious bias on the field it assigns. For example, the path APA|P.L = "H.3" reveals the important authors in information retrieval field, such as W. Bruce Croft, ChengXiang Zhai and James Allan. However, the path APA|P.L = "H.2" returns the influential authors in database and data mining field, such as Jiawei Han and Christos Faloutsos. For the meta path APA, it mingles well-known authors in these two fields. The results illustrate that the constrained meta paths are able to capture subtle semantics by deeply disclosing the most influential authors in a certain field.

6.2.2 Quantitative comparison experiments

Based on the results returned by five methods, we can obtain five candidate ranking lists of authors in ACM dataset. To evaluate the results quantitatively, we crawled data as ground truth from two well-known Web sites. The first ground truth provides the author ranks from Microsoft Academic Search.⁴ Specifically, we crawled two standard ranking lists of authors in two academic fields: DB and IR. Then we compare the difference between our candidate ranking lists and the standard ranking lists. In order to measure the quality of the ranking results, we use the *Distance* criterion proposed in [15], which is defined as follows.

$$D(R, R') = \frac{\sum_{i=1}^{n} [(n-i) \times \sum_{j=1 \wedge R'_{j} \notin \{R_{1}, \dots, R_{i}\}}^{i} 1]}{\sum_{i=1}^{\lfloor \frac{n}{2} \rfloor} [(n-i) \times i] + \sum_{i=\lfloor \frac{n}{2} \rfloor + 1}^{n} [(n-i) \times (n-i)]}$$
(14)

where R_i represents the *i*th object in ranking list R, while R'_j denotes the *j*th object in ranking list R'. And *n* is the total number of objects in the ranking lists. Note that the numerator of the formula measures the real distance between the two rankings, and the denominator of the formula is used to normalize the real distance to a number between 0 and 1. So the criterion not only measures the number of mismatches between these two lists, but also considers the position of these mismatches. The smaller *Distance* means the smaller difference (i.e., better performance).

In this experiment, we compare the five candidate ranking lists with each of the two standard ranking lists from Microsoft Academic Search, and the *Distance* results are shown in Fig. 5. We can observe an obvious phenomenon: The results obtained by the constrained meta paths have the smallest *Distance* on its corresponding field, while they have the largest *Distance* on other fields. For example, HRank with the path APA|P.L = "H.2" has the smallest *Distance* on the DB field in Fig. 5a, while it has the largest *Distance* on IR field in Fig. 5b. The reason lies in that the path APA|P.L = "H.2" focuses on the authors in the DB field. Meanwhile, these authors deviate from those in the IR field. The results further illustrate that the constrained meta path (i.e., APA) considers the co-author relationship on all fields, it achieves mediocre performances to PageRank and Degree methods. It implies that the constrained meta path in HRank indeed helps to improve the ranking performances in a specific field.

Furthermore, we quantitatively evaluate the results according to the second ground truth from ArnetMiner [24] that offers comprehensive search and mining services for academic community.⁵ Specifically, we crawl the first 200 authors as experts in DB and IR fields through searching "data mining" and "information retrieval." Since these 200 experts have no ranking order, we evaluate the accuracy of the top *k* authors of five candidate ranking lists with the F1 score. From the results shown in Fig. 6, we can observe the same phenomena. That is, the constrained meta paths always achieve the best performances on their corresponding fields, while they have the worst performances on other fields (note that the higher F1 score means the better performances). Moreover, the meta paths also have the moderate performances. The experiments on both ground truths confirm that HRank is able to improve the ranking performances in a specific field through assigning constrained meta paths.

⁴ http://academic.research.microsoft.com/.

⁵ http://arnetminer.org/.



Fig. 5 The distances between the ranking lists obtained by different methods and the standard ranking lists on different fields on ACM dataset. The ground truth is from Microsoft Academic Search. **a** DB field, **b** IR field



Fig. 6 The F1 accuracy of the ranking lists obtained by different methods on different fields on ACM dataset. The ground truth is from ArnetMiner. **a** DB field, **b** IR field

6.3 Ranking of heterogeneous objects

Then the experiments validate the effectiveness of HRank-AS on asymmetric constrained meta paths.

6.3.1 Experiment study on asymmetric constrained meta paths

The experiments are done on the DBLP dataset. We evaluate the importance of authors and conferences simultaneously based on the meta path APC, which means authors publish papers on conferences. Two constrained meta paths (APC|P.L = "DB") and APC|P.L = "IR") are also included, which means authors publish DB(IR)-field papers on conferences. Similarly, the experiments also include two baseline methods (i.e., PageRank and Degree) in above experiments with the same experimental process.

The top ten authors and conferences returned by these five methods are given in Tables 2 and 3, respectively. As shown in Table 2, the ranking results of these methods on authors all are reasonable; however, the constrained meta paths can find the most influential authors in a certain field. For example, the top three authors of APC|P.L = "DB" are Surajit Chaudhuri, Hector Garcia-Molina and H. V. Jagadish, and all of them are very influential

Rank	APC	$\begin{array}{l} APC P.L = \\ "DB" \end{array}$	$\begin{array}{l} APC P.L = \\ ``IR'' \end{array}$	PageRank	Degree
1	Gerhard Weikum	Surajit Chaudhuri	W. Bruce Croft	W. Bruce Croft (23)	Philip S. Yu
2	Katsumi Tanaka	H. Garcia-Molina	Bert R. Boyce	Gerhard Weikum (24)	Gerhard Weikum
3	Philip S. Yu	H. V. Jagadish	Carol L. Barry	Philip S. Yu (25)	Divesh Srivastava
4	H. Garcia-Molina	Jeffrey F. Naughton	James Allan	Jiawei Han (26)	Jiawei Han
5	W. Bruce Croft	Michael Stonebraker	ChengXiang Zhai	H. Garcia-Molina (27)	H. Garcia-Molina
6	Jiawei Han	Divesh Srivastava	Mark Sanderson	Divesh Srivastava (28)	W. Bruce Croft
7	Divesh Srivastava	Gerhard Weikum	Maarten de Rijke	Surajit Chaudhuri (29)	Surajit Chaudhuri
8	Hans-Peter Kriegel	Jiawei Han	Katsumi Tanaka	H. V. Jagadish (30)	H. V. Jagadish
9	Divyakant Agrawal	Christos Faloutsos	Iadh Ounis	Jeffrey F. Naughton (31)	Jeffrey F. Naughton
10	Jeffrey Xu Yu	Philip S. Yu	Joemon M. Jose	Rakesh Agrawal (32)	Rakesh Agrawal

Table 2 Top ten authors of different methods on DBLP dataset

The number in the parenthesis of the fifth column means the rank of authors in the whole ranking list returned by PageRank

researchers in the database field. The top three authors of APC|P.L = "IR" are W. Bruce Croft, Bert R. Boyce and Carol L. Barry, and they all have the high academic reputation in the information retrieval field. Similarly, as we can see in Table 3, HRank with constrained meta paths (i.e., APC|P.L = "DB" and APC|P.L = "IR") can clearly find the important conferences in DB and IR fields, while other methods mingle these conferences. For example, the most important conferences in the DB field are ICDE, VLDB and SIGMOD, while the most important conferences in the IR field are SIGIR, WWW and CIKM. Observing Tables 2 and 3, we can also find the mutual effect of authors and conferences. For example, W. Bruce Croft published many papers in SIGIR and CIKM, while Surajit Chaudhuri has many papers in SIGMOD, ICDE and VLDB.

6.3.2 Quantitative comparison experiments

To verify the effectiveness of these methods, we use the above *Distance* criterion to calculate the difference between their results and standard ranking lists crawled from Microsoft Academic Search. Figure 7 shows the differences of author ranking lists. We can observe the same phenomenon with above quantitative experiments again. That is, HRank with constrained meta paths achieves the best performances on their corresponding field. Meanwhile, they have the worst performances on other fields. In addition, compared to that of PageRank and Degree, the mediocre performances of HRank with meta path *APC* further demonstrate the importance of constrained meta path to capture the subtle semantics contained in heterogeneous networks. Similarly, we further evaluate the F1 accuracy of these methods according

Rank	APC	APC P.L = "DB"	APC P.L = "IR"	PageRank	Degree
1	CIKM	ICDE	SIGIR	ICDE (3)	ICDE
2	ICDE	VLDB	WWW	SIGIR (4)	SIGIR
3	WWW	SIGMOD	CIKM	VLDB (5)	VLDB
4	VLDB	PODS	JASIST	CIKM (6)	SIGMOD
5	SIGMOD	DASFAA	WISE	SIGMOD (7)	CIKM
6	SIGIR	EDBT	ECIR	JASIST (8)	JASIST
7	DASFAA	ICDT	APWeb	WWW (9)	WWW
8	JASIST	MDM	WSDM	DASFAA (10)	PODS
9	WISE	WebDB	JCIS	PODS (11)	DASFAA
10	EDBT	SSTD	IJKM	JCIS (12)	EDBT

Table 3 Top ten conferences of different methods on DBLP dataset

The number in the parenthesis of the fifth column means the rank of conferences in the whole ranking list returned by PageRank



Fig. 7 The distances between the candidate author ranking lists and the standard ranking lists on different fields on DBLP dataset. The ground truth is from Microsoft Academic Search. **a** DB field, **b** IR field

to the ground truth crawled from ArnetMiner. The results are shown in Fig. 8. Once again the results reveal the same findings that HRank can more accurately discover the authors ranking in a special field with the help of constrained meta path.

6.3.3 Experiments on meta path with multiple constraints

Furthermore, we validate the effectiveness of meta path with multiple constraints. In the above experiments, we employ the constraint on the label of papers in HRank with the meta path *APC*. Here we add one more constraint on conference. Specifically, in contrast to the constrained meta path APC|P.L = "DB", we employ the paths APC|P.L = "DB" &&C = "VLDB", APC|P.L = "DB" &&C = "SIGIR", and APC|P.L = "DB" &&C = "CIKM", which mean authors publish DB field papers on specified conferences (e.g., VLDB, SIGIR, and CIKM). Similarly, we add the same conference constraints on the path APC|P.L = "IR". Same with the above experiments, we calculate the rank accuracy of HRank with these constrained meta paths and the results are shown in Fig. 9.



Fig. 8 The F1 accuracy of the ranking lists obtained by different methods on different fields on DBLP dataset. The ground truth is from ArnetMiner. **a** DB field, **b** IR field



Fig. 9 The rank accuracy of HRank with different constrained meta paths on DBLP dataset. **a** DB field, **b** IR field

We know that HRank with the path APC|P.L = "DB" (APC|P.L = "IR") can reveal the influence of authors in the DB (IR) field. As ground truth, this ranking is based on the aggregation of many conferences related to the DB field. The added conference constraint in HRank further reveals the influence of authors in the specific conference of the field. So we can use the closeness to the ground truth to reveal the importance of a conference to that field. That is, if the ranking from a specific conference is quite closer to the ground truth rank, that can imply the conference is a dominating conference in that field. From Fig. 9a, we can find that the VLDB conference constraint (the blue curve) achieves the closest performances to the ground truth ranking, while the performances of the SIGIR conference constraint (the black curve) deviate most. So we can infer that the VLDB is more important than SIGIR in the DB field and the CIKM has the middle importance. Similarly, from Fig. 9b, we can infer that the SIGIR is more important than VLDB in the IR field. These findings comply with our common sense. As we know, although the VLDB and SIGIR both are the top conferences in computer science, they are very important only in their research fields. For example, the VLDB is important in the DB field, while it is not so important in the IR field. The middle importance of the CIKM conference stems from the fact that it is a comprehensive conference including papers from both DB and IR fields. In addition, we can find that the SIGIR curve almost overlaps with the ground truth over the IR field, while the VLDB curve still has a gap with the ground truth over the DB field. We think the reason is that SIGIR is the main conference in the IR field, while in the DB field, there are also other important conferences, such as SIGMOD and ICDE. In all, the experiments show that HRank with constrained meta path can not only effectively find the influential authors in each research field on a specified conference but also indirectly reveal the importance of conferences in the fields. It also implies that HRank can achieve accurate and subtle ranking results by flexibly setting the combination of constraints.

6.4 Co-ranking of objects and paths

6.4.1 Experiment study on co-ranking on symmetric constrained meta paths

In this experiment, we will validate the effectiveness of HRank-CO to rank objects and symmetric constrained meta paths simultaneously. The experiment is done on ACM dataset. First we construct a (2, 1)th order tensor X based on 73 constrained meta paths (i.e., $APA|P.L = L_j$, j = 1...73). When the *i*th and the *k*th authors co-publish a paper together, of which the label is the *j*th label (i.e., ACM categories), we add one to the entries $x_{i,j,k}$ and $x_{k,j,i}$ of X. In this case, X is symmetric with respect to the index *j*. By considering all the publications, $x_{i,j,k}$ (or $x_{k,j,i}$) refers to the number of collaborations by the *i*th and the *k*th author under the *j*th paper label. In addition, we do not consider any self-collaboration, i.e., $x_{i,j,i} = 0$ for all $1 \le i \le 17431$ and $1 \le j \le 73$. The size of X is $17431 \times 73 \times 17431$ where there are 91520 nonzero entries in X. The percentage of nonzero entries is 4.126×10^{-4} %. In this dataset, we will evaluate the importance of authors through the co-author relations; meanwhile, we will analyze the importance of paths (i.e., which paths have the most contributions to the importance of authors).

Figure 10 shows the stationary probability distributions of authors and paths. It is obvious that some authors and paths have higher stationary probability, which implies these authors and paths are more important than others. Table 4 shows the top ten authors (left) and paths (right) based on their HRank values. We can find that the top ten authors all are influential researchers in the DM/IR fields, which conforms to our common senses. Similarly, the most important paths are related to DM/IR fields, such as APA|P.L = "H.3" (Information Storage and Retrieval) and APA|P.L = "H.2" (Database Management). Although the conferences in ACM dataset are from multiple fields, such as DM/DB (e.g., KDD, SIGMOD) and computation theory (e.g., SODA, STOC), there are more papers from the DM/DB fields, which makes the authors and paths in the DM/DB fields rank higher. We can also find that



Fig. 10 The stationary probability distributions of authors and constrained meta paths. a Authors, b paths

Rank	Authors	Constrained meta paths
1	Jiawei Han	H.3 (Information Storage and Retrieval)
2	Philip Yu	H.2 (Database Management)
3	Christos Faloutsos	C.2 (Computer-Communication Networks)
4	Ravi Kumar	I.2 (Artificial Intelligence)
5	Wei-Ying Ma	F.2 (Analysis of Algorithms and Problem Complexity)
6	Zheng Chen	D.4 (Operating Systems)
7	Hector Garcia-Molina	H.4 (Information Systems Applications)
8	Hans-Peter Kriegel	G.2 (Discrete Mathematics)
9	Gerhard Weikum	I.5 (Pattern Recognition)
10	D. R. Karger	H.5 (Information Interfaces and Presentation)

Table 4 Top 10 authors and constrained meta paths (note that only the constraint (L_j) of the paths $(APA|P.L = L_j, j = 1...73)$ is shown in the third column of the table)

the influence of authors and paths can be promoted by each other. The reputation of Jiawei Han and Philip Yu comes from their productive papers in the influential fields (e.g., H.3 and H.2). In order to observe this point more clearly, we show the number of co-authors of the top ten authors based on the top ten paths in Table 5. We can observe that there are more collaborations for top authors in the influential fields. For example, although Zheng Chen (rank 6) has more number of co-authors than Jiawei Han (rank 1), the collaborations of Jiawei Han focus on ranked higher fields (i.e., H.3 and H.2), so Jiawei Han has higher HRank score. Similarly, the top paths contain many collaborations of influential authors.

6.4.2 Experiment study on co-ranking on asymmetric constrained meta paths

The experiments on the Movie dataset aim to show the effectiveness of HRank-CO to rank heterogeneous objects and asymmetric constrained meta paths simultaneously. In this case, we construct a third-order tensor X based on the constrained meta paths AMD|M.T. That is, the tensor represents the actor–director collaboration relations on different types of movies. When the *i*th actor and the *k*th director cooperate in a movie of the *j*th type, we add one to the entries $x_{i,j,k}$ of X. By considering all the cooperations, $x_{i,j,k}$ refers to the number of collaborations by the *i*th actor and the *k*th director under the *j*th type of movie. The size of X is $5324 \times 112 \times 551$, and there are 36529 nonzero entries in X. The percentage of nonzero entries is 7.827×10^{-4} %.

Table 6 shows the top ten actors, directors and constrained meta paths (i.e., movie type). We observe the mutual enhancements of the importance of objects and meta paths again. Basically, the results comply with our common senses. The top ten actors are well known, such as Eddie Murphy and Harrison Ford. Similarly, these directors are also famous in filmdom due to their works. These movie types obtained are the most popular movie subjects as well. In addition, we can observe the mutual effect of objects and paths one more time. As we know, Eddie Murphy and Drew Barrymore (rank 1, 4 in actors) are famous comedy and drama (rank 1, 2 in paths) actors. Harrison Ford and Bruce Willis (rank 2,3 in actors) are popular thrill and action (rank 3,4 in paths) actors. These higher-ranked directors also prefer those popular movie subjects. Furthermore, we also compare these results with the

II	
j, j	
= <i>L</i>	
- T -	
P	
\mathbf{A}	
H P	
Ś	
ths	
pai	
e	
ft	
ō	
(in	5
2	
int	
tra	
ns	
S	
Je	
y ti	
Ę.	
it o	
tha	
ē	
noi	
s	
ath	
ğ	
eta	
Ē	
eq	
in	
itr:	
SUC	
3	
en	
p t	
5	
the	
at	
2	
ers	
oth	
hc	
vit	
ē	
rat	
pq	
lla	
3	_
IS	e)
thc	tab
aut	g
G	ft
p t	0
tol	MO.
he	st r
ut ti	fir:
tha	Je j
er	1 tł
nb	1 ir
IUL	WI
le I	ho
Th	IS S
5	3) i
le	5
q	÷
_66	

		(
Ranked Author/CP	1 (H.3)	2 (H.2)	3 (C.2)	4 (I.2)	5 (F.2)	6 (D.4)	7 (H.4)	8 (G.2)	9 (I.5)	10 (H.5)
1 (Jiawei Han)	51	176	0	0	0	0	6	2	2	0
2 (Philip Yu)	51	94	0	0	6	0	3	0	13	0
3 (C. Faloutsos)	17	107	0	5	6	0	3	4	2	0
4 (Ravi Kumar)	73	27	0	3	13	0	18	5	0	0
5 (Wei-Ying Ma)	132	26	0	6	0	0	2	0	30	10
6 (Zheng Chen)	172	6	0	6	0	0	22	0	38	6
7 (H. Garcia-Molina)	23	65	ю	0	0	0	1	0	0	4
8 (H. Kriegel)	19	28	5	0	0	0	9	0	7	4
9 (G. Weikum)	82	14	0	4	0	0	8	0	4	0
10 (D. R. Karger)	11	5	13	0	L	4	1	7	0	7

Author's personal copy

Rank	Actor	Director	Conditional meta path
1	Eddie Murphy	Tim Burton	Comedy
2	Harrison Ford	Zack Snyder	Drama
3	Bruce Willis	Marc Forster	Thriller
4	Drew Barrymore	David Fincher	Action
5	Nicole Kidman	Michael Bay	Adventure
6	Nicolas Cage	Ridley Scott	Romance
7	Hugh Jackman	Richard Donner	Crime
8	Robert De Niro	Steven Spielberg	Sci-Fi
9	Brad Pitt	Robert Zemeckis	Animation
10	Christopher Walken	Stephen Sommers	Fantasy

Table 6 Top 10 actors, directors and meta paths on IMDB dataset (note that only the constraint (T_j) of the paths $(AMD|M.T = T_i, j = 1...1591)$ is shown in the fourth column)

recommended results from the IMDB Web site.⁶ Although only a subset of movies in IMDB is included in our experiments, the 80% of the top 10 actors in our results are included in the set of the top 250 greatest movie actors in all time recommended by IMDB,⁷ and the 50% of the top 10 directors in our results are included in the set of the top 50 favorite directors recommended by IMDB.⁸ Moreover, most of movie types recommended by our method have high ranks in the popular types summarized by IMDB.⁹

6.5 Fast computation experiments

Based on the DBLP dataset, we select two meta paths with varying length (l): $(APA)^l$ and $(APCPA)^l$, where *l* means times of path repetition ranging from 1 to 5. We record the running time of matrix multiplication based on these paths with different fast computation strategies. The direct matrix multiplication is baseline. Meanwhile, we calculate the differences of results obtained by baseline and fast computation strategies (i.e., *F* norm of differences of two matrices). These differences can be considered as the accuracy measure of fast computation strategies (the smaller the better). We set the parameters in the Trun strategy as follows: *W* is 200, β is 0.5, and γ is 0.02. The number of walkers (i.e., *K*) in the MonC strategy is 500. All experiments are done on machines with Intel Xeon 8-Core CPUs of 2.13 GHz and 64 GB RAM.

Figure 11 shows the running time and accuracy of three strategies on different paths. From Fig. 11a, b, we can find that the DynP is an effective strategy to speed up matrix multiplication on both paths, while the Trun and MonC strategies only speed up matrix multiplication on the path $(APCPA)^l$. During the matrix multiplication along $(APA)^l$, the matrix is always sparse, so the baseline itself is very fast. In this condition, the Trun and MonC strategies do not work. For the path $(APCPA)^l$, the multiplication matrix becomes dense due to the low dimension of *C* (# of conferences is 20), so its running time increases greatly. In this condition, the Trun and MonC are also effective strategies to speed up matrix multiplication.

⁶ http://www.imdb.com/.

⁷ http://www.imdb.com/list/ls050720698/.

⁸ http://www.imdb.com/list/ls050131440/.

⁹ http://www.imdb.com/list/ls050782187/?view=detail&sort=listorian:asc.



Fig. 11 Running time and accuracy of matrix multiplication based on different fast computation strategies and paths. **a** Running time on $(APA)^l$, **b** running time on $(APCPA)^l$, **c** accuracy on $(APA)^l$, **d** accuracy on $(APCPA)^l$

Then, we observe their accuracy from Fig. 11c, d where the *y*-axis shows the difference on accuracy from the baseline. As an information-lossless strategy, the DynP's results are the same with the baseline. The MonC strategy has the lowest accuracy.

According to the analysis above, these strategies are suitable for different paths and scenarios. For very sparse matrix [e.g., $(APA)^l$], HRank can be computed fast without applying any fast computation strategies. In this condition, only the DynP strategy can speed up HRank without loss in accuracy. For those dense and high-dimensional matrices [e.g., $(APCPA)^l$] which have huge computation overhead, the Trun, MonC and DynP strategies can effectively improve the HRank's efficiency. To sum up, the DynP can effectively accelerate matrix multiplication without loss on accuracy, while the Trun and MonC strategies also help to speed up the multiplication of dense matrices.

In Sect. 4.4, we have pointed out that the time complexity of the rank computation in HRank-CO is $O(trn^2)$. However, for real applications, the relation matrices are very sparse, so the real time complexity is linear to the number of links. This point is confirmed by the following experiment. We create three different scales of tensors (size $n \times r \times n$). We record the running time for rank computation on different link densities. The results are shown in Fig. 12. It is clear that the running time slowly and near linearly increases with the increment of link density. Moreover, the longer running time is needed for larger-scale tensor.



Тор К

6.6 Parameter study

There is a parameter α in HRank, which determines the restarting probability. In this section, we will observe the effect of different parameter settings on the performances of HRank. Based on the ACM dataset and the constrained meta path APA|P.L = "H.2", we run HRank with different α for 20 times and record the ranking accuracy (i.e., *Distance*) of HRank. The results are shown in Fig. 13. Generally, we can find that, with the increment of α , the performances of HRank rise up first and then drop down. Moreover, HRank achieves the best performances when α is 0.8 or 0.9. So we set α to be 0.85 in above experiments. In all, the parameter α in HRank complies with the same rules with the parameter α in PageRank.

6.7 Convergence experiments

In Fig. 14, we show the convergence of HRank on the previous experiments. The results illustrate that the three versions of HRank all quickly converge after no more than 20 iterations. In addition, we can also observe that HRank has different convergence speed in these three conditions. For symmetric meta paths, the HRank-SY almost converges after 6 iterations (see Fig. 14a). However, HRank-CO for co-ranking converges after 16 iterations (see Fig. 14c). We think it is reasonable, since it is more difficult to converge for more objects in the HRank-CO.



Fig. 14 The difference between two successive calculated probability vectors against iterations based on the three versions of HRank. a HRank-SY, b HRank-AS, c HRank-CO

7 Conclusions

In this paper, we first study the ranking problem in heterogeneous information network and propose the HRank framework, which is a path- based random walk method. In this framework, we introduce the constrained meta path concept to capture the more subtle and refined semantics contained in HIN. In addition, we further put forward a method to co-rank the paths and objects, since the paths make an effect on the importance of objects. Extensive experiments validate the effectiveness and efficiency of HRank on three real datasets.

There are several interesting works that are worth doing in the future. Although we have designed three memory-based fast computation strategies, the network size in many real applications is too large to be contained in memory. We can design the parallelized version of HRank with parallel models (e.g., MapReduce [2] or BSP [25]) for large-scale networks. In addition, although we have validated the effectiveness of HRank on two real datasets, we can further employ HRank on more applications to exploit its application scope.

Acknowledgments This work is supported in part by National Key Basic Research and Department (973) Program of China (No. 2013CB329606), the National Natural Science Foundation of China (No. 71231002, 61375058), the CCF-Tencent Open Fund, the Co-construction Project of Beijing Municipal Commission of Education and US NSF through Grants III-1526499.

References

- 1. Chakrabarti S (2007) Dynamic personalized pagerank in entity-relation graphs. In: WWW, pp 571-580
- Dean J, Ghemawat S (2008) Mapreduce: simplified data processing on large clusters. Commun ACM 51(1):107–113
- Deng H, Lyu MR, King I (2009) A generalized co-hits algorithm and its application to bipartite graphs. In: KDD, pp 239–248
- Han J (2009) Mining heterogeneous information networks by exploring the power of links. In: Gama J, Costa VS, Jorge AM, Brazdil PB (eds) Discovery Science. Springer, Berlin, Heidelberg, pp 13–30
- Huang H, Zubiaga A, Ji H, et al (2012) Tweet ranking based on heterogeneous networks. In: Proceedings of the 24th international conference on computational linguistics, COLING 2012, pp 1239–1256
- 6. Jeh G, Widom J (2002) Simrank: a measure of structural-context similarity. In: KDD, pp 538–543
- Ji M, Sun Y, Danilevsky M, Han J, Gao J (2010) Graph regularized transductive classification on heterogeneous information networks. In: Balcázar JL, Bonchi F, Gionis A, Sebag M (eds) Machine learning and knowledge discovery in databases. Springer, Berlin, Heidelberg, pp 570–586
- 8. Jin R, Lee VE, Hong H (2011) Axiomatic ranking of network role similarity. In: KDD, pp 922–930
- 9. Kleinberg JM (1999) Authoritative sources in a hyperlinked environment. JACM 46(5):604–632
- Kong X, Yu PS, Ding Y, Wild DJ (2012) Meta path-based collective classification in heterogeneous information networks. In: CIKM, pp 1567–1571

- Lao N, Cohen W (2010) Fast query execution for retrieval models based on path constrained random walks. In: KDD, pp 881–888
- Li X, Ng MK, Ye Y (2012) Har: Hub, authority and relevance scores in multi-relational data for query search. In: SDM, pp 141–152
- Li Y, Shi C, Yu P, Chen Q (2014) Hrank: a path based ranking method in heterogeneous information network. In: Li F, Li G, Hwang Sw, Yao B, Zhang Z (eds) Web-age information management. Lecture Notes in Computer Science, vol 8485. Springer International Publishing, pp 553–565
- Ng MK, Li X, Ye Y (2011) Multirank: Co-ranking for objects and relations in multi-relational data. In: KDD, pp 1217–1225
- 15. Nie Z, Zhang Y, Wen J, Ma W (2005) Object-level ranking: bringing order to web objects. In: WWW, pp 422–433
- 16. Page L, Brin S, Motwani R, Winograd T (1998) The pagerank citation ranking: bringing order to the web. Technical report, Stanford University Database Group
- 17. Shi C, Kong X, Yu PS, Xie S, Wu B (2012a) Relevance search in heterogeneous networks. In: EDBT, pp 180–191
- Shi C, Zhou C, Kong X, Yu PS, Liu G, Wang B (2012b) Heterecom: a semantic-based recommendation system in heterogeneous networks. In: KDD, pp 1552–1555
- 19. Shi C, Wang R, Li Y, Yu PS, Wu B (2014) Ranking-based clustering on general heterogeneous information networks by network projection. In: CIKM, pp 699–708
- Shi C, Zhang Z, Luo P, Yu PS, Yue Y, Wu B (2015) Semantic path based personalized recommendation on weighted heterogeneous information networks. In: CIKM, pp 453–462
- Soulier L, Jabeur LB, Tamine L, Bahsoun W (2013) On ranking relevant entities in heterogeneous networks using a language-based model. J Am Soc Inf Sci Technol 64(3):500–515
- 22. Sun Y, Han J, Yan X, Yu P, Wu T (2011) Pathsim: meta path-based top-k similarity search in heterogeneous information networks. In: VLDB, pp 992–1003
- Sun Y, Norick B, Han J, Yan X, Yu PS, Yu X (2012) Integrating meta path selection with user-guided object clustering in heterogeneous information networks. In: KDD, pp 1348–1356
- 24. Tang J, Zhang J, Yao L, Li J, Zhang L, Su Z (2008) Arnetminer: extraction and mining of academic social networks. In: KDD, pp 990–998
- 25. Valiant LG (1990) A bridging model for parallel computation. Commun ACM 33(8):103-111
- Yu X, Sun Y, Norick B, Mao T, Han J (2012) User guided entity similarity search using meta-path selection in heterogeneous information networks. In: CIKM, pp 2025–2029
- Zhou D, Orshanskiy SA, Zha H, Giles CL (2007) Co-ranking authors and documents in a heterogeneous network. In: ICDM, pp 739–744



Chuan Shi received the B.S. degree from the Jilin University in 2001, the M.S. degree from the Wuhan University in 2004 and Ph.D. degree from the ICT of Chinese Academic of Sciences in 2007. He joined the Beijing University of Posts and Telecommunications as a lecturer in 2007 and is a Professor and Deputy Director of Beijing Key Lab of Intelligent Telecommunications Software and Multimedia at present. His research interests are in data mining, machine learning and evolutionary computing. He has published more than 40 papers in refereed journals and conferences.



Yitong Li received the B.S. degree from the Beijing University of Posts and Telecommunications in 2014. She is currently a master student in BUPT. Her research interests are in data mining and machine learning.



Philip S. Yu is a Distinguished Professor in Computer Science at the University of Illinois at Chicago and also holds the Wexler Chair in Information Technology. Dr. Yu spent most of his career at IBM, where he was manager of the Software Tools and Techniques group at the Watson Research Center. His research interest is on big data, including data mining, data stream, database and privacy. He has published more than 920 papers in refereed journals and conferences. He holds or has applied for more than 250 US patents. Dr. Yu is a Fellow of the ACM and the IEEE. He is the Editor-in-Chief of ACM Transactions on Knowledge Discovery from Data. He is on the steering committee of the IEEE Conference on Data Mining and ACM Conference on Information and Knowledge Management and was a member of the IEEE Data Engineering steering committee. He was the Editor-in-Chief of IEEE Transactions on Knowledge and Data Engineering (2001-2004). Dr. Yu received the B.S. Degree in E.E. from National Taiwan University, the M.S. and Ph.D. degrees in E.E. from Stanford University and the M.B.A. degree from New York University.



Bin Wu received the B.S. degree from the Beijing University of Posts and Telecommunications in 1991, the M.S. and Ph.D. degrees from the ICT of Chinese Academic of Sciences in 1998 and 2002, respectively. He joined the Beijing University of Posts and Telecommunications as a lecturer in 2002 and is a professor at present. His research interests are in data mining, complex network and cloud computing. He has published more than 100 papers in refereed journals and conferences.