

Mining Research Communities in Bibliographical Data ^{*}

Osmar R. Zaiane, Jiyang Chen, and Randy Goebel

University of Alberta, Canada
{zaiane, jiyang, goebel}@cs.ualberta.ca

Abstract. Extracting information from very large collections of structured, semi-structured or even unstructured data can be a considerable challenge when much of the hidden information is implicit within relationships among entities in the data. Social networks are such data collections in which relationships play a vital role in the knowledge these networks can convey. A bibliographic database is an essential tool for the research community, yet finding and making use of relationships comprised within such a social network is difficult. In this paper we introduce **DBconnect**, a prototype that exploits the social network coded within the **DBLP** database by drawing on a new random walk approach to reveal interesting knowledge about the research community and even recommend collaborations.

1 Introduction

A social network is a structure made up of nodes, representing entities from different groups, that are linked with different types of relations. Viewing and understanding social relationships between individuals or other entities is known as *Social Network Analysis* (SNA). SNA methods [26] are used to study organizational relations, analyze citation or computer mediated communications, etc. There are many applications such as studying the spread of disease, understanding the flow of communication within and between organizations, and so on. As an important field in SNA, *Community Mining* [5, 16] has received considerable attention over the last few years. A community can be defined as a group of entities that share similar properties or connect to each other via certain relations. Identifying these connections and locating entities in different communities is an important goal of community mining and can also have various applications. We are interested in the application for finding potential collaborators for researchers by discovering communities in an author-conference social network, or recommending books (or other products) for users based on the borrowing records of other members of their communities in a library system. In this paper we are focusing on the social network implicit in the DBLP database which includes information about authors, their papers and the conferences they published in. DBLP [13, 3] is an on-line resource providing bibliographic information on major computer science conference proceedings and journals¹. It is such an essential index for the community that it was included in the

^{*} This work is based on an earlier work: DBconnect: mining research community on DBLP data, in Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis, COPYRIGHT ACM, 2007, <http://portal.acm.org/citation.cfm?doid=1348549.1348558>.

¹ In December 2007, DBLP comprised more than 970,000 entries.

In SNA, the closeness of two related concepts in the network is usually measured by a relevance score, which is based on selected relationships between entities. It can be computed with various techniques, e.g., *Euclidean distance* or *Pearson correlation* [26]. Here we use the random walk approach to determine the relevance score between two entities. A random walk is a sequence of nodes in a graph such that when moving from one node n to the subsequent one in the sequence, one of n 's neighbours is selected at random but with an edge weight taken into account. The closeness of a node b with respect to a node a is the static steady-state probability that the sequence of the nodes would include b when the random walk starts in a . This probability is computed iteratively until convergence, and is used as an estimated relevance score. In this paper, we adapt a variation of this idea, which is the random walk with restart (RWR): given a graph and a starting node a , at each step, we move to a neighbour of the current node at random, proportionally to the available edge weights, or go back to the initial node a with a restart probability c . RWR has been applied to many fields, e.g. anomaly detection [23], automatic image captioning [18], etc.

In this paper, we use DBLP data to generate bipartite (author-conference) and tripartite (author-conference-topics) graph models, where topics are frequent n-grams extracted from paper titles and abstracts. Moreover, we present an iterative random walk algorithm on these models to compute the relevance score between authors to discover the communities. We take into consideration the co-authorship while designing graphical models and the algorithm. We also present our ongoing work DBconnect, which provides an interactive interface for navigating the DBLP community structure online, as well as recommendations and explanations for these recommendations.

The rest of the paper is organized as follows. We discuss related work in Section 2. Graph models and Random walk algorithms for computing the relevance score are described in Section 3. The result and the ongoing DBconnect work is reported in Section 4 before the paper is concluded in Section 5.

2 Related Work

Community Mining

The ability to find communities within large social networks could be of important use, e.g., communities in a biochemical network might correspond to functional units of the same kind [8]. Since social networks can be easily modeled as graphs, finding communities in graphs, where groups of vertices within which connections are dense, but between which connections are sparser, has recently received considerable interests. Traditional algorithms, such as the spectral bisection method [19], which is based on the eigenvectors of the graph Laplacian, and the Kernighan-Lin algorithm [11], which greedily optimizes the number of within- and between-community edges, suffer from

² <http://acm.org/sigmod/dblp/db/anthology.html>

the fact that they only bisect graphs. While a larger number of communities can be identified by repeated bisection, there is no hint of when to stop the repeated partitioning process. Another approach to find communities is hierarchical clustering based on similarity measures between objects, but it cannot handle the case where some vertices are not close enough to any of the principal communities to be clustered. In the last few years, several methods have been developed based on iterative removal of between-community edges [5, 20, 25]. Important results on researcher community mining have been revealed by analysis of a co-relation (e.g., co-authorship in a paper or co-starring in a movie) graph. Nascimento et al. [15] and Smeaton et al. [21] show co-authorship graphs for several selected conferences are small world graphs³ and calculate the average distance between pairs of authors. Similarly, the Erdős Number Project⁴ and the Oracle of Bacon⁵ compute the minimum path length between one fixed person and all other people in the graph.

Community Information System

A related contribution in the context of recommending future collaborators based on their communities is W-RECMAS, which is an academic recommendation system developed by Cazella et al. [14]. The approach is based on collaborative filtering on the user profile data of the Brazilian e-government's database and can aid scientists by identifying people in the same research field or with similar interests in order to help exchange ideas and create academic communities. However, researchers need to post and update their research interests and personal information in the database before they can be recognized and recommended by the system, which makes the approach impractical. In order to efficiently browse the DBLP bibliographical database [13], Klink et al. [12] developed a specialized tool, DBL-Browser, which provides an interactive user interface and essential functionalities such as searching and filtering to help the user navigate through the complex data of DBLP. Another project to explore information for research communities is the DBLife system⁶. It extracts information from web resources, e.g., mailing list archives, newsletters, well-known community websites or research homepages, and provides various services to exploit the generated entity relationship graph [4]. While they do not disclose the process and the means used, they provide related researchers and related topics to a given researcher. In addition to the DBLife project supported by Yahoo, Microsoft Research Asia also developed a similar project called Libra⁷, which discovers connected authors, conferences and journals etc. However, in our own experience of using the two systems, we found some incorrect instances of these related entities. Distinct from DBLife and Libra, our DBconnect focuses on finding related researchers more accurately based on a historical publication database and explicit existing relationships in the DBLP coded social network. Moreover, DBLife and Libra do not provide recommendations like DBconnect does.

³ A small-world graph is a graph in which most nodes are not neighbors of one another, but can be reached from every other by a small number of hops or steps [2].

⁴ <http://www.oakland.edu/~grossman/erdoshp.html>

⁵ <http://www.oracleofbacon.org/>

⁶ <http://dblfe.cs.wisc.edu/>

⁷ <http://libra.msra.cn/>

Random Walk Algorithm

As a popular metric to measure the similarity between entities, the random walk algorithm has received increasing attention after the undeniable success of the Google search engine, which applies a random walk approach to rank web pages in its search result as well as the list of visited pages to re-index [1]. Specifically, Page-Rank [17] learns ranks of web pages, which are N-dimensional vectors, by using an iterated method on the adjacency matrix of the entire web graph. In order to yield more accurate search results, Topic-Sensitive PageRank [6] pre-computes a set of biased PageRank vectors, which emphasize the effect of particular representative topic keywords to increase the importance of certain web pages. Those are used to generate query-specific importance scores. Alternatively, SimRank [9] computes a purely structural score that is independent of domain-specific information. The SimRank score is a structure similarity measure between pairs of pages in the web graph with the intuition that two pages are similar if they are related by similar pages. Unfortunately, SimRank is very expensive in computation since it needs to calculate similarities between many pairs of objects. According to the authors, a pruning technique is possible to approximate SimRank by only computing a small part of the object pairs. However, it is very hard to identify the right pairs to compute at the beginning, because the similarity between objects may only be recognized after the score between them is calculated. Similar random walk approaches have been used in other domains. For example, the Mixed Media Graph [18] applies a random walk on multimedia collection to assign keywords to the multimedia object, such as images and video clips, but a similarity function for each type of the involved media is required from domain experts. He et al. [7] propose a framework named MRBIR using a random walk on a weighted graph for images to rank related images given an image query. Sun et al. [23] detect anomaly data for datasets that can be modeled as bipartite graphs using the random walk with restart algorithm. Recent work by Tong et al. [24] proposed a fast solution for applying random walk with restart on large graphs, to save pre-computation cost and reduce query time with some cost on accuracy. While random walk algorithms such as SimRank computes links recursively on all pairs of objects, LinkClus [28] takes advantage of the power law distribution of links, and develops a hierarchical structure called SimTree to represent similarities in a multi-granularity manner. By merging computations that go through the same branches in the SimTree, LinkClus is able to avoid the high cost of pairwise similarity computations but still thoroughly explores relationships among objects without random walk.

In this paper, we apply a random walk approach on tripartite graphs to include topic information, and increase the versatility of the random walk by expanding the original graph model with virtual nodes that take the co-authorship into consideration for the DBLP data. These extensions are explained in the following section.

3 Proposed Method

Searching for relevant conferences, similar authors, and interesting topics is more important than ever before, and is considered an essential tool by many in the research

community such as finding reviewers for journal papers or inviting program committee members for conferences. However, finding relationships between authors and thematically similar publications is becoming more difficult because of the mass of information and the rapid growth of the number of scientific workers [12]. Moreover, except direct co-authorships which are explicit in the bibliographical data, relationships between nodes in this complex social network are difficult to detect by traditional methods. In order to understand relations between entities and find accurate researcher communities, we need to take into consideration not only the information of who people work with, i.e. co-authors, but also where they submit their work to, i.e., conferences, and what they work on, i.e. topics. In this section, we first present the models that incorporate these concepts, then discuss the algorithms that compute the relevance scores for these models.

Given the DBLP database $D = (C \cup A)$, where conference set $C = \{c_i | 1 \leq i \leq n\}$ and author set $A = \{a_j | 1 \leq j \leq m\}$, we can model D as an undirected bipartite graph $G = (C, A, E)$: conference nodes and author nodes are connected if the corresponding author published in the conference and there are no edges in E within the same group of nodes, i.e., author to author or conference to conference. Figure 1 (a) shows an example of the bipartite graph, representing social relationships between conference and author entities. The weights of the edges are publishing frequency of different authors in a certain conference.

3.1 Adding Topic Information

As mentioned before, the research topic is an important component to differentiate any research community. Authors that attend the same conferences might work on various topics. Therefore, topic entities should be treated separately from conference and author entities. Figure 2 shows an example of linked author, conference, and topic entities. DBLP contains table of contents of some conference proceedings. These table of contents include session titles that could be considered as topics. Unfortunately, very few conference proceedings have their table of contents included in DBLP, and in the affirmative, session titles are often absent. To extract relevant topics from DBLP we resorted to the paper titles instead. Moreover, we obtained as many paper abstracts as possible from Citeseer⁸, then extracted topics based on keyword frequency from both titles and abstracts. We found that frequent co-locations in title and abstract text constitute reliable representation of topics. We concede that other methods are possible to get effective research topics.

We now consider a publication database $D = (C \cup A \cup T)$, where topic set $T = \{t_i | 1 \leq i \leq l\}$. We naturally use a tripartite graph to model such data: author/conference nodes are related to a topic node if they have a paper on that topic, the edge weight is the topic matching frequency. We apply the random walk algorithm on a tripartite graph by adjusting the walking sequence. For example, previously the random walker turns back to author nodes when it reaches a conference node; now it will go forward to topic nodes

⁸ <http://citeseer.ist.psu.edu/>

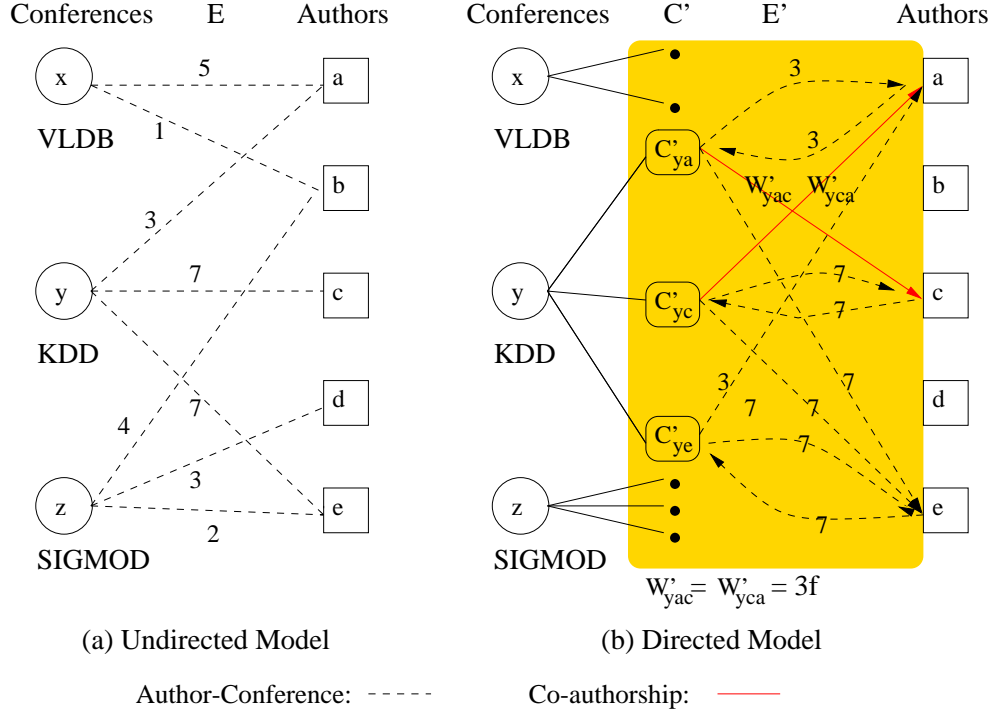


Fig. 1. Bipartite Model for Conference-Author Social Network

first, and then walk back to author nodes. By such modifications, the relevance score now contains both conference and topic influences, i.e., in a tripartite model, authors with high relevance score share similar conference experiences and paper topics with the given author.

3.2 Adding Co-author Relations

Table 1 shows the number of publications of five authors a, b, c, d, e in three conferences VLDB KDD and SIGMOD. Authors a and c have co-authored 3 papers in KDD, a and b co-authored 1 paper in VLDB and d, e co-authored 2 papers in SIGMOD. Unfortunately, the corresponding bipartite graph, which is shown in Figure 1 (a), fails to represent any co-authorships. For example, author a and c co-authored many papers at KDD, but there are no edges in the bipartite graph that can be used to represent this information: edge $e(y, a)$ and $e(y, c)$ are both used by relations between conference and author. On the other hand, author e seems more related to author c since the weights of edges connecting them to KDD are the heaviest ($w_{yc} = 7, w_{ye} = 7$). The influence of the important co-author a is neglected because the model only represents publication frequency.

To capture the co-author relations, just adding a link between a and c does not suf-

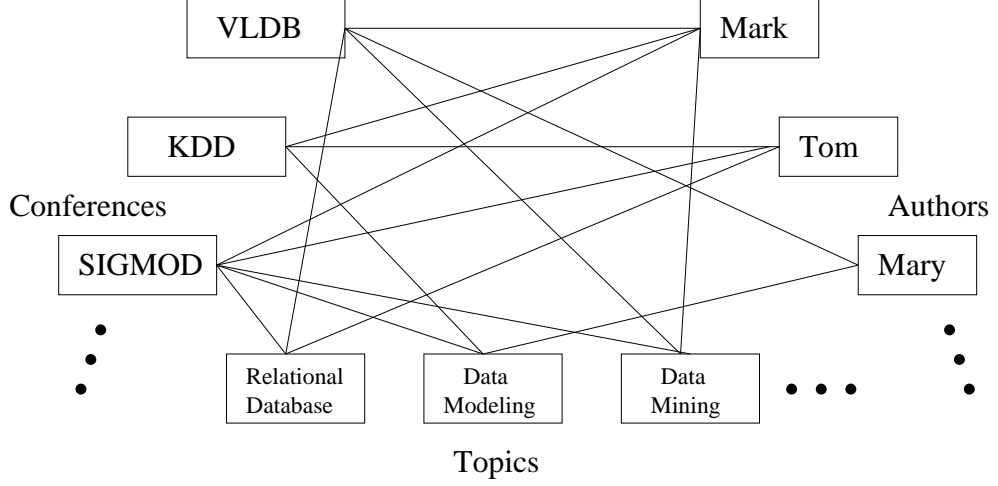


Fig. 2. Tripartite Graph Model for Author-Conference-Topic

| | Publication Records |
|---------------|---------------------|
| VLDB(x) | a(4), ab(1) |
| KDD(y) | ac(3), c(4), e(7) |
| SIGMOD(z) | b(4), d(1), de(2) |

Table 1. Author Publication Records in Conferences. For example, a, b, c, d, e are authors, $ac(3)$ means that author a and c published three papers together in a certain conference.

face, since it misses the role of KDD, where the co-authorship happens. Making the link connecting a and c to KDD directional does not work either, as from KDD there are edges to many other authors, which would make the random walk infeasible (i.e., yielding undesirable results). Moreover, adding additional nodes to represent each co-author relation is impractical when there is a huge number of such relations. For instance, adding “Papers” between Authors and Conferences to make a tri-partite graph would actually not only add a significant number of edges since many authors have multiple papers per conference series, but also, this scheme does not allow the random walk to favor co-authorship as any author or co-author gets the same probability to be visited.

Our approach is to re-structure the bipartite model by adding surrogate nodes to replace the KDD node and having them link to a and c so that the random walk calculation can be applied while the connection between related nodes remains the same. In more detail, we add a virtual level of nodes to replace the conference partition, and add direction to the edges. Figure 1 (b) shows details of node KDD as an example. We first split y into 3 nodes to represent relations between y and authors who published there (a, c and e). These author nodes then connect to their own splitted relation nodes with the original weight ($e'(a, C'_{ya}), e'(c, C'_{yc}), e'(e, C'_{ye})$). Then we connect from C' nodes to all author nodes that have published at KDD. If the author node has a co-author relation

with the author included in the C' node, the edge is weighted by co-author frequency multiplied by a parameter f (which is explained in the following), otherwise, the edge is weighted as original. We can see that the co-authorship, which is missed in the simple bipartite graph, is now represented by extra weight of edge $e'(C'_{yc}, a)$ and $e'(C'_{ya}, c)$, which shows author a is more related to c than author e through KDD due to their collaborations. The parameter f is used to control the co-author influence, usually we set $f = k$ (k is the total author number of a conference).

3.3 Random Walk on DBLP Social Network

Before presenting the random walk algorithms, we define the problems we are solving: given an author node $a \in A$, we compute a relevance score for each author $b \in A$. The result is a one-column vector containing all author scores with respect to a . We measure closeness of researchers so we can discover implicit communities in the DBLP data.

Recall that we extend the bipartite model into a directed bipartite graph $G' = (C', A, E')$, where A has m author nodes, C' is generated base on C and has $n * m$ nodes (we assume every node in C is split into m nodes). The basic intuition of our approach is to apply random walks on the adjacency matrix of graph G' starting from a given author node. To form the adjacency matrix, we first generate a matrix for directional edges from C' to A , which is $M_{(n*m) \times m}$, then form a matrix for edges from A to C' , which is $N_{m \times (n*m)}$. In these two matrices, $M(\alpha, \beta)$ or $N(\alpha, \beta)$ indicates the weight of the directed edge from node α to node β in G' (0 means no such edge). A random walk starting from a node represented by row α in M (the same applies to N) takes the edge (α, β) based on the probability which is proportional to the edge weight over the sum of weight of all outgoing edges of α . Therefore, we normalize M and N such that every row sums up to 1. We can then construct the adjacency matrix J of G' :

$$J_{(n*m+m) \times (m+n*m)} = \begin{pmatrix} 0 & (Norm(N))^T \\ (Norm(M))^T & 0 \end{pmatrix}$$

We then transform the given author node α into a one-column vector v_α consisting of $(n * m + m)$ elements. The value of the element corresponding to author α is set to 1. We now need to compute a steady-state vector u_α , which contains relevance scores of all nodes in the graph model. The scores for the author nodes are the last m elements of the vector. The result is achieved based on the following lemma and the RWR approach.

Lemma 1 *Let c be the probability of restarting random walk from node α . Then the steady-state vector u_α satisfies the following equation:*

$$u_\alpha = (1 - c)Ju_\alpha + cv_\alpha$$

See [22] for proof of the lemma.

Algorithm 1 applies the above lemma repeatedly until u_α converges. We set c to be 0.15 and ϵ to be 0.1, which gives the best convergence rate according to [23]. The bipartite structure of the graph model is used to save the computation of applying Lemma

Algorithm 1 The Random Walk with Restart Algorithm

Input: node $\alpha \in A$, bipartite graph model G , restarting probability c , converge threshold ϵ .

Output: relevance score vector A for author nodes.

1. Construct graph model G' for co-authorship based on G . Compute the adjacency matrix J of G' .

2. Initialize $v_\alpha = 0$.
set value for α to 1: $v_\alpha(\alpha) = 1$.

3. While $(\Delta u_\alpha > \epsilon)$

$$u_\alpha = (1 - c) \left(\frac{(Norm(N))^T u_{\alpha(n*m+1:n*m+m)}}{Norm(M)^T u_{\alpha(1:n*m)}} \right) + c v_\alpha$$

4. Set vector $A = u_{\alpha(n*m+1:n*m+m)}$

5. Return A .

1 in step 3. The last m elements of the result vector $u_{\alpha(n*m+1:n*m+m)}$ contains the relevance score for all author nodes in A .

We extend algorithm 1 for the tripartite graph model $G'' = (C, A, T, E'')$. Assume we have n conferences, m authors and l topics in G' , we can represent all relations using three corresponding matrices: $U_{n \times m}$, $V_{m \times l}$ and $W_{n \times l}$. We normalize them such that every column sum up to 1: $Q(U) = col_norm(U)$, $Q(U^T) = col_norm(U^T)$. We then construct the adjacency matrices of G'' after normalization:

$$J_{CA} = \begin{pmatrix} 0 & Q(U) \\ Q(U^T) & 0 \end{pmatrix}$$

$$J_{CT} = \begin{pmatrix} 0 & Q(W) \\ Q(W^T) & 0 \end{pmatrix}$$

$$J_{AT} = \begin{pmatrix} 0 & Q(V) \\ Q(V^T) & 0 \end{pmatrix}$$

Similarly, given a node $\alpha \in C$, we want to compute a relevance score for all nodes that are in C, A, T . We apply the RWR approach following the visiting sequence until convergence, e.g., walk from author to conference, to topic, and back to author if we want to rank authors (see Algorithm 2). There are $m + n + l$ elements for all nodes in the graph model in the result relevance score vector. The value of the corresponding node, either starting author, topic or conference, is initialized to 1. After the random walk algorithm terminates, scores for conference, author and topic nodes are recorded from 1 to n , from $n + 1$ to $n + m$ and from $n + m + 1$ to $n + m + l$ in the vector, respectively.

Here we show a simple random walk on the conference-author network example we give in Table 1. The relational matrix M of the network is shown as follows.

$$M = \left(\begin{array}{c|ccccc} & a & b & c & d & e \\ \hline VLDB & 5 & 1 & 0 & 0 & 0 \\ KDD & 3 & 0 & 7 & 0 & 7 \\ SIGMOD & 0 & 4 & 0 & 3 & 2 \end{array} \right)$$

Then we build and normalize the adjacency matrix J of the graph shown in Figure 1.

$$J = \left(\begin{array}{c|cccccc} & VLDB & KDD & SIGMOD & a & b & c & d & e \\ \hline VLDB & 0 & 0 & 0 & 5 & 1 & 0 & 0 & 0 \\ KDD & 0 & 0 & 0 & 3 & 0 & 7 & 0 & 7 \\ SIGMOD & 0 & 0 & 0 & 0 & 4 & 0 & 3 & 2 \\ a & 5 & 3 & 0 & 0 & 0 & 0 & 0 & 0 \\ b & 1 & 0 & 4 & 0 & 0 & 0 & 0 & 0 \\ c & 0 & 7 & 0 & 0 & 0 & 0 & 0 & 0 \\ d & 0 & 0 & 3 & 0 & 0 & 0 & 0 & 0 \\ e & 0 & 7 & 2 & 0 & 0 & 0 & 0 & 0 \end{array} \right)$$

$$J_{normalize} = \left(\begin{array}{c|cccccc} & VLDB & KDD & SIGMOD & a & b & c & d & e \\ \hline VLDB & 0 & 0 & 0 & 0.62 & 0.2 & 0 & 0 & 0 \\ KDD & 0 & 0 & 0 & 0.38 & 0 & 1.0 & 0 & 0.77 \\ SIGMOD & 0 & 0 & 0 & 0 & 0.8 & 0 & 1.0 & 0.22 \\ a & 0.84 & 0.18 & 0 & 0 & 0 & 0 & 0 & 0 \\ b & 0.16 & 0 & 0.44 & 0 & 0 & 0 & 0 & 0 \\ c & 0 & 0.41 & 0 & 0 & 0 & 0 & 0 & 0 \\ d & 0 & 0 & 0.33 & 0 & 0 & 0 & 0 & 0 \\ e & 0 & 0.41 & 0.22 & 0 & 0 & 0 & 0 & 0 \end{array} \right)$$

A random walk on this graph moves from one node to one of its neighbours at random but the probability of picking a particular edge is proportional to the weight of the edge out of the sum of weights of all edges that connect to this node. For example, if we start from node SIGMOD, we build \mathbf{u} as the start vector:

$$\mathbf{u} = \{0, 0, 1, 0, 0, 0, 0, 0\}^T$$

After the first step of the first iteration, the random walk hits the author nodes with $b = 1 * 0.44, d = 1 * 0.33, e = 1 * 0.22$.

$$\mathbf{u} = \{0, 0, 0, 0, 0.44, 0, 0.33, 0.22\}^T$$

In the next step of the first iteration, the chance that the random walk goes back to SIGMOD is $0.44 * 0.8 + 0.33 * 1 + 0.22 * 0.22 = 0.73$. The other 0.27 goes to the other two conference nodes.

$$\mathbf{u} = \{0.09, 0.18, 0.73, 0, 0, 0, 0, 0\}^T$$

The vector will converge after a few iterations and gives a stable score to every node, which is the probability of a random walk may hit this node. However, the fact that these scores are always the same no matter where the walk begins makes the approach incapable for ranking for different given starting points. This problem can be solved by random walk with restart: in each random walk iteration, the walker goes back to the start node with a restart probability. Therefore, nodes that are closer to the starting node now have a higher chance to be visited and obtain larger ranking score.

Algorithm 2 Random Walk Algorithm for Tripartite Model

Input: node α , tripartite graph model G'' , restarting probability c , converge threshold ϵ .

Output: relevance score vector \mathbf{c} , \mathbf{a} and \mathbf{t} for author, conference and topic nodes.

1. Compute the adjacency matrices J_{CA} , J_{CT} and J_{AT} of G'' .

2. Initialize $\mathbf{v}_\alpha = 0$, set element for α to 1: $\mathbf{v}_\alpha(\alpha) = 1$.

3. While ($\Delta \mathbf{u}_\alpha > \epsilon$)

$$\mathbf{u}_{\alpha(n+1:n+m)} = (Q(U^T) * \mathbf{u}_{\alpha(1:n)})$$

$$\mathbf{u}_{\alpha(n+m+1:n+m+l)} = (Q(V^T) * \mathbf{u}_{\alpha(n+1:n+m)})$$

$$\mathbf{u}_{\alpha(1:n)} = (Q(W) * \mathbf{u}_{\alpha(n+m+1:n+m+l)})$$

$$\mathbf{u}_\alpha = (1 - c)\mathbf{u}_\alpha + c\mathbf{v}_\alpha$$

4. Set vector $\mathbf{c} = \mathbf{u}_{\alpha(1:n)}$, $\mathbf{a} = \mathbf{u}_{\alpha(n+1:n+m)}$,

$$\mathbf{t} = \mathbf{u}_{\alpha(n+m+1:n+m+l)}.$$

6. Return \mathbf{c} , \mathbf{a} , \mathbf{t} .

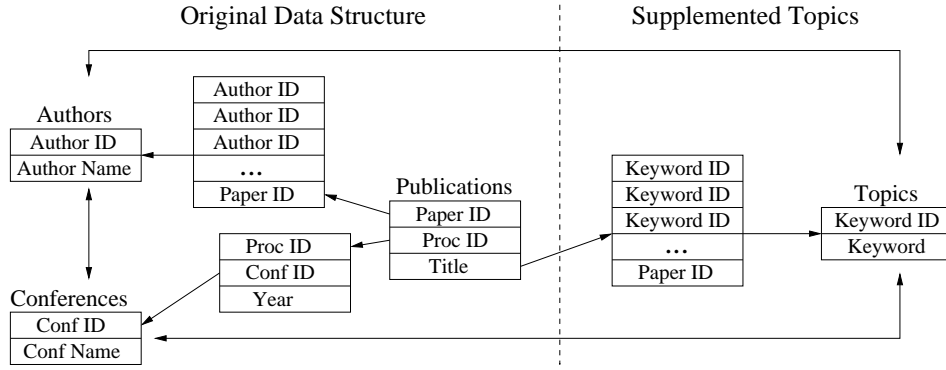


Fig. 3. Our Data Structure extracted from DBLP and Citeseer

4 Exploring DBLP Communities

In the academic world, since a researcher could usually belong to multiple related communities, e.g., Database and AI, it is unnecessary and improper to classify this researcher into any specific arbitrary communities. Therefore, in our experiment, we focus on investigating the closeness of researchers, i.e., we are interested at *how* and *why* two people are in the same community, instead of *which* community they are in.

4.1 DBLP Database

We downloaded the publication data for conferences from the DBLP website⁹ in July 2007. Any publication after that date is not included in our experimental data. Moreover, we kept only conference proceedings and removed all journals and other publications. These were minimal compared to the conference publications. The data structure

⁹ <http://www.informatik.uni-trier.de/~ley/db/>

is shown in Figure 3. We extracted topics based on keyword frequency from paper titles in DBLP data and abstracts from Citeseer¹⁰, which provides abstracts of about 10% of all papers. First we manually selected a list of stopwords to remove frequently used but non-topic-related words, e.g., “Towards”, “Understanding”, “Approach”, etc. Then we counted frequency of every co-located pairs of stemmed words and selected the top 1000 most frequent bi-grams as topics. Additionally, we manually added several tri-grams, e.g. World Wide Web, Support Vector Machine, etc., since we observe both bilateral bi-grams to be frequent (e.g. World Wide and Wide Web). We chose to use bi-grams because they can distinguish most of the research topics, e.g, Relational Database, Web Service and Neural Network, while single keywords fail to separate different topics, e.g. “Network” can be part of “Social Network” or “Network Security”.

Since the publication database is huge (it contains more than 300,000 authors, about 3,000 conferences and the selected 1,000 N-gram topics), the entire adjacency matrix becomes too big to make the random walk efficient. However, we can compute the result by first performing graph partitioning on the model and only running the random walk on the part where the given author is. This approach can only achieve an approximate result, since some weakly connected communities are separated, but it is much faster since we end-up computing with much smaller matrices. In this paper, we used the METIS algorithm [10] to partition the large graph into ten subgraphs of about the same size. Note that the proposed approach is independent of the selected partitioning method.

4.2 The DBconnect System

After the author-conference-topic data extraction from the DBLP database, we generate lists of people with high relevance scores with respect to different given researchers. Our ongoing project *DBconnect*, which is a navigational system to investigate the community connections and relationships, is built to explore the result lists of our random walk approach on the academic social network. An online demo for DBConnect is available at¹¹. Figure 4 shows a screenshot of the author interface of our *DBconnect* system. There are eight lists displayed for a given author in the current version. Clicking on any of the hyper-linked names will generate a page with respect to that selected entity. We explain details of each of the lists below.

– Academic Information

Academic statistics for the given author are shown in this list, which contain three components: conference contribution, earliest publication year and average publication per year are extracted from DBLP; the H-index [27] is calculated based on information retrieved from Google Scholar¹²; approximate citation numbers are retrieved from Citeseer¹³. The query terms for Google Scholar and Citeseer are automatically generated based on the author names. Users can submit an alternative

¹⁰ <http://citeseer.ist.psu.edu/>

¹¹ <http://kingman.cs.ualberta.ca/research/demos/content/dbconnect/>

¹² <http://scholar.google.com/>

¹³ <http://citeseer.ist.psu.edu/>

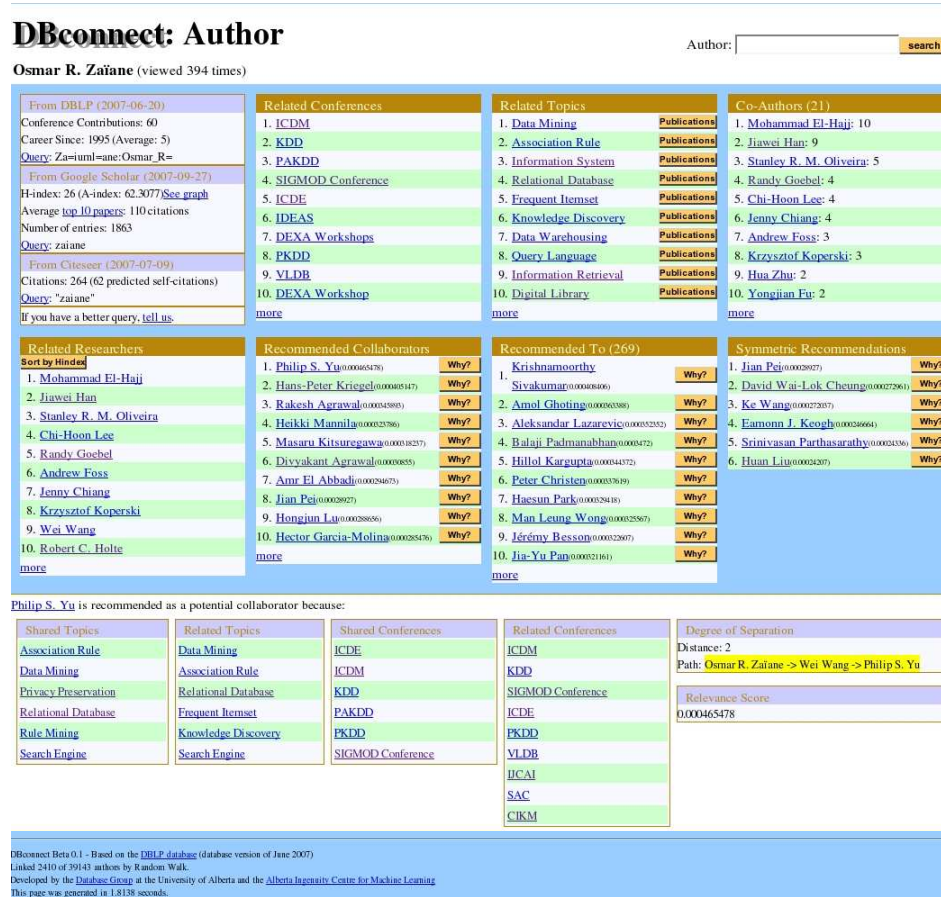


Fig. 4. DBconnect Interface Screenshot for an author

query which gives a more accurate result from the search engines. We also provide a visualization of the H-index. One can click the “See graph” link beside the H-index numbers. Figure 5 shows an example of H-index visualization.

– Related Conferences

This list is generated by the random walk, which starts from the given author, on an author-conference-topic model and is ordered by their relevance score, in descending order. These are not necessarily the conferences where the given researcher published but the conferences related to the topics and authors that are also related to the reference researcher. Clicking on the conference name leads to a new page with topics and authors related to the chosen conference.

– Related Topics

This list is ordered by the relevance scores from a random walk on the tripartite model. Clicking on the button “Publications” after each topic provides the papers that the given author has published on that topic, i.e. the papers of the given author

that contains the N-gram keywords in their titles or abstracts. Similarly, these are not necessarily the topics that the given author has worked on, but the topics most related to his topics, attended conferences and colleagues.

- *Co-authors*

The co-author list reports the number of publications that different researchers co-authored with the given person.

- *Related Researchers*

This list is based on the bipartite graph model with only conference and author entities, i.e. we apply our extended bipartite graph model to emphasize the co-authorship. The result implies that the given author is related to the same conferences and via the same co-authors as these listed researchers. In most cases, most related researchers to the given author are co-authors and co-authors of co-authors.

- *Recommended Collaborators*

This list is based on the tripartite graph author-conference-topic. Since co-authors are treated as “observed collaborators”, their names are not shown here. The result implies that the given author shares similar topics and conference experiences with these listed researchers, hence the recommendation. The relevance score calculated by our random walk is displayed following the names. Clicking on the “why” button brings the detailed information of the relationship between the two authors. For example, in Figure 4, relations between Philip Yu and Osmar Zaïane are described by the topics and conferences they share, and the degree of separation in the co-authorship chain ($A \rightarrow B$ means A and B are co-authors). Here, the “Share Topics” table lists the topics that these two authors both have publications on and the “Related Topics” table shows the topics that appear in the Related Topics lists of both authors. Similarly, the “Shared Conferences” table displays the conferences that the two authors have attended and the “Related Conferences” table shows the conferences that can be found in the Related Conferences lists of both authors.

- *Recommended To*

The recommendation is not symmetric, i.e., author A may be recommended as a possible future collaborator to author B but not vice versa. This phenomenon is due to the unbalanced influence of different authors in the social network. For example, Jiawei Han has a significant influence with his 196 conference publications, 84 co-authors and H-index 63. He has been recommended as collaborator for 6201 authors, but apparently only a few of them is recommended as collaborators to him. The Recommended To list shows the authors that have the given author in their recommendation list, ordered by the relevance score.

- *Symmetric Recommendations*

This list shows the authors that have been recommended to the given author and have the given author on their recommendation list.

Note that while there is some overlap between the list of related researchers and the list of recommended collaborators, there is a fundamental difference and the difference by no means implies that collaboration with the missing related researchers is discouraged. They are simply two different communities in the network even though they overlap. The list of related researchers is obtained from relationships derived from co-authorships and conferences by a RWR on an extended bipartite graph with co-

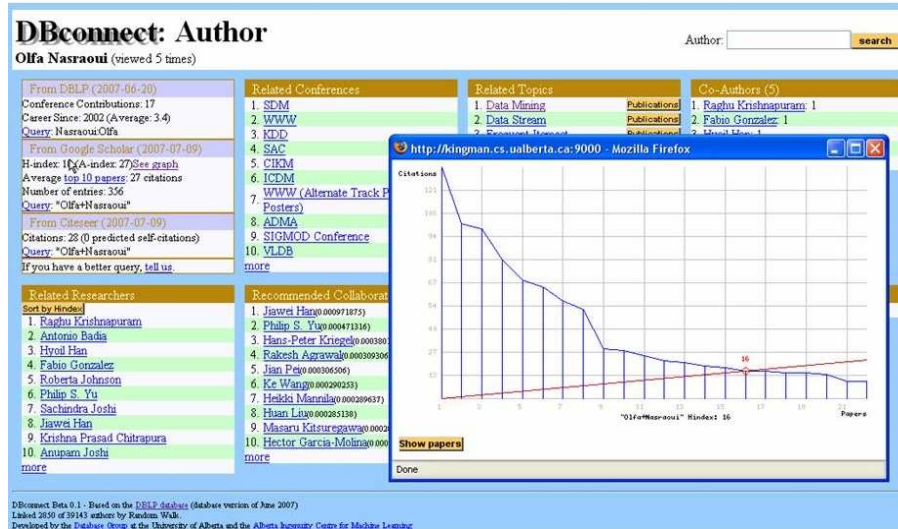


Fig. 5. DBconnect Interface Screenshot for H-Index Visualization

authorship relations. The result is a quasi-obvious list due to the closeness from co-authors. This list could create a sort of trust in the system given the clear closeness of this community. The list of recommended collaborators could be perceived as a more distant community and thus as an interesting discovery. It is obtained without co-authorship but with relations from topics. We use a RWR on a tripartite graph authors/conferences/topics. The explanation on the why collaborators are recommended (i.e. common conferences and topics, and degree of separation) establishes more trust in the recommendation. A systematic validation of these lists is difficult but the cases we manually substantiated were satisfactory and convincing.

Clicking on any conference name shows a conference page. Figure 6 illustrates an example when the entity “ICDM” is selected. Conferences have their own related conferences, authors and topics. Note that the topics here mean the most frequent topics used within titles and abstracts of papers published in the given conference.

Clicking on the topics leads to a new page with conferences, authors and topics related to the chosen topic. Note again that this relationship to topics comes from paper titles and abstracts. Figure 7 shows an example when the topic “Data Mining” is selected.

5 Conclusions and Future Work

In this paper, we extend a bipartite graph model to incorporate co-authorship, and propose a random walk approach to find related conferences, authors, and topics for a given entity. The main idea is to use a random walk with restarts on the bipartite or tripartite model of DBLP data to measure the closeness between any two entities. The

| DBconnect: Conference | | | Other Conferences |
|---|--|--------------------------------------|-----------------------------------|
| IEEE International Conference on Data Mining (6 events) (viewed 32 times) | | | |
| Related Researchers | Related Topics | Related Conferences | |
| Sort by Index | | 1. KDD | |
| 1. Philip S. Yu | 1. Data Mining | 2. ICML | |
| 2. Qiang Yang | 2. Association Rule | 3. VLDB | |
| 3. Haixun Wang | 3. Decision Tree | 4. ICDE | |
| 4. Zheng Chen | 4. Data Stream | 5. SIGMOD Conference | |
| 5. Jiawei Han | 5. Database System | 6. PAKDD | |
| 6. Hans-Peter Kriegel | 6. Time Series | 7. SIGIR | |
| 7. Wei Wang | 7. Information Retrieval | 8. CIKM | |
| 8. Sheng Ma | 8. Web Service | 9. PKDD | |
| 9. Eamonn J. Keogh | 9. Web Search | 10. INFOCOM | |
| 10. Srinivasan Parthasarathy | 10. Knowledge Discovery | 11. UAI | |
| 11. George Karypis | 11. Clustering Algorithm | 12. AAAI/IAAI | |
| 12. Wei Fan | 12. Frequent Itemset | 13. DEXA | |
| 13. Ming-Syan Chen | 13. Information System | 14. DaWaK | |
| 14. Ke Wang | 14. Bayesian Network | 15. IJCAI | |
| 15. Osmar R. Zaiane | 15. Text Classification | 16. PODS | |
| 16. Charles X. Ling | 16. Learning Algorithm | 17. ECML | |
| 17. Benyu Zhang | 17. Machine Learning | 18. SODA | |
| 18. Tao Li | 18. Digital Library | 19. ICDCS | |
| 19. Chris H. Q. Ding | 19. Relational Database | 20. DASFAA | |
| 20. Heikki Mannila | 20. Neural Network | all fewer | |
| all fewer | all fewer | | |

DBconnect Beta 0.1 - Based on the [DBLP database](#) (database version of June 2007)
 Developed by the [Database Group](#) at the University of Alberta and the [Alberta Ingenuity Centre for Machine Learning](#)

Fig. 6. DBconnect Interface Screenshot for conference ICDM

result, the relevance score, can be used to understand the relationship between entities and discover the community structure of the corresponding data. We basically use the relevance score to rank entities based on importance given a relationship.

We also present our ongoing work DBconnect, which can help explore the relational structure and discover implicit knowledge within the DBLP data collection. Not all of the more than 360,000 authors are indexed in DBconnect at the time of printing as the random walks are time consuming. A queue of authors is continuously processed in parallel and authors can be prioritized in the queue by request.

The work we presented in this paper is still preliminary. We have implemented a prototype¹⁴. However, more work is needed to verify the value of the approach. The lists of related conferences, topics and researchers to a given author are interesting and can be used to help understand the entity closeness and research communities. While the output of DBconnect is satisfactory and the manual substantiation confirms acceptable and suitable lists (as opposed to lists provided by DBLife), some systematic evaluation is still desired. However, validation of the random walk is difficult and we are considering

¹⁴ <http://kingman.cs.ualberta.ca/research/demos/content/dbconnect/>

| DBconnect: Topic | | Other Topics |
|--|--------------------------------------|--|
| Data Mining (viewed 49 times) | | |
| Related Researchers | Related Conferences | Related Topics |
| Sort By Index | | |
| 1. Jiawei Han | 1. KDD | 1. Database System |
| 2. Heikki Mannila | 2. SIGMOD Conference | 2. Association Rule |
| 3. Rakesh Agrawal | 3. VLDB | 3. Information System |
| 4. Christos Faloutsos | 4. ICDE | 4. Web Service |
| 5. Raymond T. Ng | 5. ICDM | 5. Decision Tree |
| 6. Marek Wojciechowski | 6. PKDD | 6. Management System |
| 7. Osmar R. Zaiane | 7. PAKDD | 7. Relational Database |
| 8. Hongjun Lu | 8. CIKM | 8. Information Retrieval |
| 9. Maciej Zakrzewicz | 9. IJCAI | 9. Knowledge Discovery |
| 10. Philip S. Yu | 10. SAC | 10. Query Language |
| 11. Masaru Kitsuregawa | 11. DEXA | 11. Data Modeling |
| 12. Carlo Zaniolo | 12. ICML | 12. Data Stream |
| 13. Haixun Wang | 13. AAAI | 13. Machine Learning |
| 14. Krzysztof Koperski | 14. DaWaK | 14. Data Structure |
| 15. Hans-Peter Kriegel | 15. PODS | 15. Data Warehousing |
| 16. Usama M. Fayyad | 16. ISMIS | 16. Neural Network |
| 17. Huan Liu | 17. SDM | 17. Query Optimization |
| 18. Mohammed Javeed Zaki | 18. INFOCOM | 18. Time Series |
| 19. Richard R. Muntz | 19. ECML | 19. Semantic Web |
| 20. Mika Klemettinen | 20. SIGIR | 20. Lower Bound |
| more | more | more |

DBconnect Beta 0.1 - Based on the [DBLP database](#) (database version of June 2007)
 Developed by the [Database Group](#) at the University of Alberta and the [Alberta Ingenuity Centre for Machine Learning](#)

Fig. 7. DBconnect Interface Screenshot for topic Data Mining

devising methods to confirm the accuracy of the relevance score and the generated lists. Moreover, it is hard to extract correct topics for researchers since the only available information is the title of the paper, which usually does not suffice to describe the content. Some titles are even unconventionally unrelated to the content of the paper only to attract attention or are metaphoric. We are considering implementing a hierarchy of topics to group similar topics and ease the browsing of the long list of related topics in computer science. We also plan to address the issue of acronyms in titles that are currently discarded. For example HMM for Hidden Markov Model is currently eliminated due to infrequency while relevant as a topic. In addition, the matrix multiplications in the random walk process make it expensive to compute. Improving the efficiency of the random walk without jeopardizing its effectiveness is necessary since the computations for relevance score estimation need to be redone continuously as the the DBLP database never ceases to grow.

6 Acknowledgments

Our work is supported by the Canadian Natural Sciences and Engineering Research Council (NSERC), the Alberta Ingenuity Centre for Machine Learning (AICML), and the Alberta Informatics Circle of Research Excellence (iCORE).

References

1. Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. In *Seventh International World Wide Web Conference*, pages 107–117, Brisbane, Australia, 1998.
2. Mark Buchanan. *Nexus: Small worlds and the groundbreaking theory of networks*, Norton, W. W. Company, Inc, 2003.
3. DBLP (Digital Bibliography & Library Project) Bibliography database. <http://www.informatik.uni-trier.de/~ley/db/>.
4. AnHai Doan, Raghu Ramakrishnan, Fei Chen, Pedro DeRose, Yoonkyong Lee, Robert McCann, Mayssam Sayyadian, and Warren Shen. Community information management. *IEEE Data Engineering Bulletin, Special Issue on Probabilistic Databases*, 29(1), 2006.
5. Michelle Girvan and M. E. J. Newman. Community structure in social and biological networks. In *Proceedings of the National Academy of Science USA*, 99:8271–8276, 2002.
6. Taher H. Haveliwala. Topic-sensitive pagerank. In *WWW: Proceedings of the 11th international conference on World Wide Web*, pages 517–526, 2002.
7. Jingrui He, Mingjing Li, Hong-Jiang Zhang, Hanghang Tong, and Changshui Zhang. Manifold-ranking based image retrieval. In *MULTIMEDIA: Proceedings of the 12th annual ACM international conference on Multimedia*, pages 9–16, 2004.
8. P. Holme, M. Huss, and H. Jeong. Subnetwork hierarchies of biochemical pathways. *Bioinformatics*, 19:532–538, 2003.
9. Glen Jeh and Jennifer Widom. Simrank: a measure of structural-context similarity. In *KDD*, 2002.
10. George Karypis and Vipin Kumar. Multilevel k-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, 1998.
11. B. W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *Bell System Technical Journal*, 49:291–307, 1970.
12. Stefan Klink, Patrick Reuther, Alexander Weber, Bernd Walter, and Michael Ley. Analysing social networks within bibliographical data. In *DEXA*, pages 234–243, 2006.
13. Michael Ley. The DBLP computer science bibliography: Evolution, research issues, perspectives. In *SPIRE*, pages 1–10, 2002.
14. Silvio César Cazella and Luis Otávio Campos Alvares. An architecture based on multi-agent system and data mining for recommending research papers and researchers. In *Proc. of the 18th International Conference on Software Engineering and Knowledge Engineering (SEKE)*, pages 67–72, 2006.
15. Mario A. Nascimento, Jörg Sander, and Jeffrey Pound. Analysis of sigmod’s co-authorship graph. *SIGMOD Record*, 32(2):57–58, 2003.
16. M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.
17. Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. In *Technical report, Stanford University Database Group*, 1998.

18. Jia-Yu Pan, Hyung-Jeong Yang, Christos Faloutsos, and Pinar Duygulu. Automatic multimedia cross-modal correlation discovery. In *KDD*, pages 653–658, 2004.
19. A. Pothén, H. Simon, and K. P. Liou. Partitioning sparse matrices with eigenvectors of graphs. *SIAM J. Matrix Anal. Appl.*, 11:430–452, 1990.
20. Filippo Radicchi, Claudio Castellano, Federico Cecconi, Vittorio Loreto, and Domenico Parisi. Defining and identifying communities in networks. *PROC.NATL.ACAD.SCI.USA*, 101:2658, 2004.
21. A. F. Smeaton, G. Keogh, C. Gurrin, K. McDonald, and T. Soderling. Analysis of papers from twenty-five years of sigir conferences: What have we been doing for the last quarter of a century. *SIGIR Forum*, 36(2):39–43, 2002.
22. Gilbert Strang. Introduction to linear algebra, Wellesley-Cambridge Press, 3 Edition, 1998.
23. Jimeng Sun, Huiming Qu, Deepayan Chakrabarti, and Christos Faloutsos. Neighborhood formation and anomaly detection in bipartite graphs. In *ICDM*, pages 418–425, 2005.
24. Hanghang Tong, Christos Faloutsos, and Jia-Yu Pan. Fast random walk with restart and its applications. In *ICDM*, pages 613–622, 2006.
25. Joshua R. Tyler, Dennis M. Wilkinson, and Bernardo A. Huberman. Email as spectroscopy: automated discovery of community structure within organizations. *Communities and technologies*, pages 81–96, 2003.
26. S. Wasserman and K. Faust. Social network analysis: Methods and applications, Cambridge University Press, 1994.
27. Michael C. Wendl. H-index: however ranked, citations need context. *Nature*, 449(403), 2007.
28. Xiaoxin Yin, Jiawei Han, and Philip S. Yu. Linkclus: efficient clustering via heterogeneous semantic links. In *VLDB*, pages 427–438, 2006.