

Research Article

Detecting Communities in 2-Mode Networks via Fast Nonnegative Matrix Trifactorization

Liu Yang, Wang Tao, Ji Xin-sheng, Liu Caixia, and Xu Mingyan

National Digital Switching System Engineering & Technological R&D Center, Zhengzhou 450002, China

Correspondence should be addressed to Liu Yang; liuyang198610@163.com

Received 13 August 2014; Accepted 16 October 2014

Academic Editor: Hamid R. Karimi

Copyright © 2015 Liu Yang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With the rapid development of the Internet and communication technologies, a large number of multitype relational networks widely emerge in real world applications. The bipartite network is one representative and important kind of complex networks. Detecting community structure in bipartite networks is crucial to obtain a better understanding of the network structures and functions. Traditional nonnegative matrix factorization methods usually focus on homogeneous networks, and they are subject to several problems such as slow convergence and large computation. It is challenging to effectively integrate the network information of multiple dimensions in order to discover the hidden community structure underlying heterogeneous interactions. In this work, we present a novel fast nonnegative matrix trifactorization (F-NMTF) method to cocluster the 2-mode nodes in bipartite networks. By constructing the affinity matrices of 2-mode nodes as manifold regularizations of NMTF, we manage to incorporate the intratype and intratype information of 2-mode nodes to reveal the latent community structure in bipartite networks. Moreover, we decompose the NMTF problem into two subproblems, which are involved with much less matrix multiplications and achieve faster convergence. Experimental results on synthetic and real bipartite networks show that the proposed method improves the slow convergence of NMTF and achieves high accuracy and stability on the results of community detection.

1. Introduction

Community structure is an important feature of real-world networks as it is crucial for us to study and understand the functional characteristics of the real complex systems. A community is usually thought of as a group of nodes with more interactions among its members than outside of the network. With the rapid growth of the Internet and computational technologies in the past decade, numerous community detection algorithms [1–3] have advanced swiftly from the simple clustering of homogeneous datasets (1-mode network) to heterogeneous datasets (multiple mode networks). Unlike homogeneous networks that only contain 1-mode nodes and have explicit community structure, the community structures of heterogeneous networks are usually obscure and complicated owing to the coexistence of multiple-mode interactions. Therefore, it is challenging to effectively integrate the network information of multiple dimensions in order to discover the hidden community structure underlying heterogeneous interactions.

The 2-mode network [4] is a simple and important kind of heterogeneous networks, which is composed of 2 kinds of disjoint node sets (for simplicity, we call them as subject-node set $N_S = \{S_1, \dots, S_m\}$ and object-node set $N_O = \{O_1, \dots, O_n\}$). So the links only connect two end nodes from different sets and no links between nodes of the same set (Figure 1). In many studies, 2-mode networks are called bipartite networks by most researchers. Bipartite network is also the most common complex networks in the real world. Many real-world networks may be naturally modeled as bipartite graphs, such as actor-movie, author-paper, word-document, product-consumer, and context-tag networks. Thereby, it is necessary and valuable for us to study community detection in bipartite networks.

Due to the special link patterns and multirelational nodes of bipartite networks, it is more suitable for mining the special community structure via coclustering methods. Recent works [5–7] have shown that clustering multimode datasets simultaneously, that is, coclustering, can effectively improve

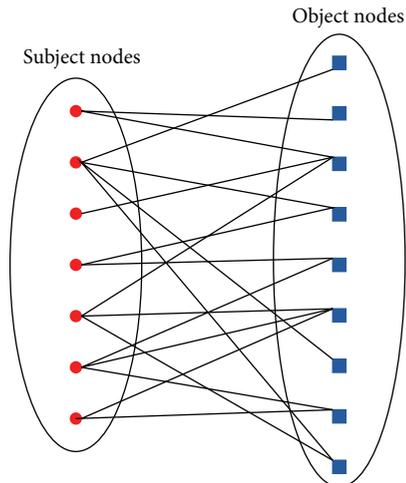


FIGURE 1: An illustrative example of the bipartite network.

the clustering performance in the sense that coclustering can make full use of the dual interdependence between heterogeneous nodes to discover certain hidden community structures. In contrast to traditional one-side clustering (on either columns or rows), coclustering treats the data matrix in a symmetric form that a partitioning of rows can induce a partitioning of columns, and vice versa. By clustering both rows and columns of a data matrix simultaneously, coclustering can effectively deal with the multimode networks and discover the structure of them.

Among different approaches of coclustering, nonnegative matrix trfactorization (NMTF) provides a superior model for co-clustering, in which the relations between row and column clusters are explicitly embodied. There are three major advantages with NMTF methods. First, as shown in (1), NMTF introduces one more factor matrix S to absorb the different scales of X , F , and G and provides increased degrees of freedom such that the low-rank matrix approximation remains accurate. Second, instead of being independent, the clustering tasks of multimode nodes in heterogeneous networks are often closely related. NMTF coclusters both the rows and columns of the original networks simultaneously by making efficient use of the duality (intratype and intertype) information between 2-mode nodes. Finally, in NMTF-based coclustering scenario, the three factor matrices jointly determine the appropriate latent community structure of bipartite networks, where F and G , respectively, cluster the subject nodes and object nodes and S is the correlation matrix reflecting the relationship between subject clusters and object clusters. Through the NMTF model, we can seamlessly integrate multiple node information to discover the underlying community structure in bipartite networks, which is highly valuable in many real world applications.

Here we adopt nonnegative matrix trfactorization as a tool to find the communities because of its powerful interpretability and applicability for coclustering in heterogeneous networks. In this work, we present a novel fast NMTF solution for community detection in bipartite networks. To summarize, the main contributions of this work include the

following. (1) We construct the affinity matrices of 2-mode nodes as manifold regularizations of NMTF, which efficiently incorporate the intertype and intratype information of subject and object spaces, to enhance the community detection in bipartite networks. (2) We present an optimization algorithm for fast nonnegative matrix trfactorization (F-NMTF), which decouple the original optimization problem into two smaller subproblems requiring much less matrix multiplications, and then cocluster the 2-mode relational nodes from the heterogeneous interactions.

The remainder of the paper is organized as follows. In Section 2, we define the NMTF model of community detection and demonstrate the theoretical foundations of our approach along with an illustrative example. In Section 3, we formulate the coclustering NMTF method for community detection in bipartite networks and present an optimal algorithm to achieve fast convergence. Then we test our algorithm on a variety of artificial and real bipartite networks and present the experimental results in Section 4. Finally, Section 5 concludes the paper.

2. Preliminary

2.1. Related Work. In the past years, NMF-based algorithms [8–11] for detecting network communities have gained great attention, because the data matrix factorization effectively reflects the community structure of the networks and promises a meaningful community interpretation that is independent of the network topology. In addition to a quantification of how strongly each node participates in each community, nonnegative matrix factorization (NMF) does not suffer from the drawbacks of modularity optimization methods [12], such as the resolution limit [13]. Psorakis et al. [8] proposed a Bayesian NMF model to extract overlapping modules. This method can automatically determine the number of communities, but also may mislead the factorization to return a bad solution, when some errors come out with its estimate of the community number. Cao et al. [9] used nonnegative matrix factorization with I-divergence as the cost function and introduce two approaches which are respectively applied to the directed and undirected networks. In [10], Nguyen et al. developed a novel NMF model where vertices are measured by their centrality in communities and detect overlapping communities, hubs, and outliers from the NMF framework altogether. Based on the importance of each node when forming links in each community, He et al. [11] use nonnegative matrix factorization to form a generative model and take it as an optimization problem to discover the structure of link communities.

However, most community detection methods still focus on homogeneous networks and might not work well in the bipartite networks. There are two main reasons. First, the special link patterns of bipartite networks greatly limit the effectiveness of these methods, which tend to cluster the heterogeneous nodes by constructing the node or edge similarities of them, as they did in the unipartite networks. But for bipartite networks, the similarities among one-mode nodes sometimes can only be defined by the nodes of the

other mode. That made these methods cannot keep working well in bipartite networks. Second, bipartite networks contain multiple types of nodes that are related to each other. Tackling each type independently will lose these interactions. So it is necessary for us to utilize the 2-mode nodes to gain a full understanding of the bipartite networks.

Motivated by recent progress in coclustering and matrix factorization, several novel NMF-based algorithms have been proposed to detect the underlying community structure, including GNMf [14], DRCC [15], and BNMTF [16]. These researches showed that the NMTF model is more applicable to discover the hidden structures in the bipartite networks. Compared to the classical clustering algorithms, nonnegative matrix trifactORIZATION provides a good model for coclustering, in which the relations between row and column clusters are explicitly embodied. In view of these facts [17–20], our approach aims to cocluster both the rows and columns of the bipartite networks simultaneously by making full use of the background information of 2-mode nodes. It is believed that adding intratype information as additional constraints can definitely improve the clustering performance.

In this work, we extend the NMTF model by adding the intratype information of 2-mode nodes as manifold constraints and develop an optimization solution to improve the coclustering performance of bipartite networks. In addition, due to the fact that NMTF-based methods often require the community numbers of networks to be specified beforehand, several methods [10, 21, 22] have been developed to solve this problem. Due to the simplicity and practicability of the existed method in [10], here we choose it to get the community numbers.

2.2. An Illustrative Example. Before the introduction of our NMTF-based method for community discovery, let us see an illustrative example on how NMTF works on the bipartite networks. Considering a bipartite network $G(N_S, N_O, E)$, where N_S is the subject-node set, N_O is the object-node set, and E denotes the link set of the network G , $|N_S| = m$, $|N_O| = n$. Given an asymmetric adjacent matrix $X \in \mathbb{R}^{m \times n}$ denoting the bipartite network G , NMTF methods approximates X by looking for 3 low-rank factor matrices with the form

$$X \approx FSG^T, \quad F \geq 0, G \geq 0, \quad (1)$$

where $F \in \mathbb{R}^{m \times l}$ and $G \in \mathbb{R}^{n \times r}$, respectively, cluster the subject nodes and object nodes and jointly determine the appropriate latent community structure in bipartite networks, $S \in \mathbb{R}^{l \times r}$. r is associated with the number of the subject-node clusters, and l is the number of object-node clusters. In most cases, let $l = r$ as the presetted community number. In this way, NMTF can simultaneously group the subject-node set and the object-node set into r clusters, where each community is the mixture of the 2-mode heterogeneous nodes.

Let X_{ij} represent the weight of the edge connecting subject-node i and object-node j . Generally, the original

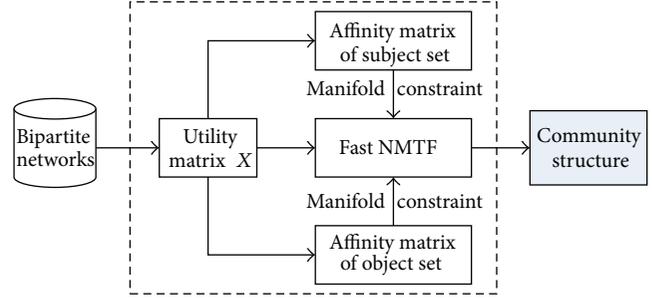


FIGURE 2: The solution framework of fast nonnegative matrix trifactORIZATION.

NMTF can be computed using the following multiplicative update rules [23]:

$$\begin{aligned} G_{ij} &\leftarrow G_{ij} \frac{(X^T FS)_{ij}}{(GG^T X^T FS)_{ij}}, \\ F_{ij} &\leftarrow F_{ij} \frac{(XGS^T)_{ij}}{(FF^T XGS^T)_{ij}}, \\ S_{ij} &\leftarrow S_{ij} \frac{(F^T XG)_{ij}}{(F^T FSG^T G)_{ij}}. \end{aligned} \quad (2)$$

The NMTF framework of community detection on bipartite networks is shown in Figure 2. Here the product FSG^T can be regarded as an approximate form of the network. Thus, the results of NMTF methods can be interpreted in which F_{ij} indicates the membership internal-strength of the subject-node i in the j th community and G_{ij} denotes the membership internal-strength of the object-node j in the i th community.

Consider a toy network with $|N_S| = 6$, $|N_O| = 8$ nodes and $|E| = 36$ edges of varying weights (Figure 3(a)). Specifically, the NMTF coclusters the 2-mode node sets to yield a comprehensive network partition solution (Figure 3(b)), where F is the subject nodes' community indicator, and G reflects the community structure of object nodes. The larger square indicates the larger value of a corresponding element in the matrix. Specially, it is worthwhile to note that the matrix S represents how subject-clusters are related to object-clusters. Each column in S reflects which subject-clusters make contribution to each object-cluster. In this case, subject-cluster 1 and 3 contribute to object-cluster 1, while subject-clusters 2 and 3 contribute to object-cluster 2.

Hence, our NMTF framework can simultaneously cocluster both the subject set and object set. By incorporating the duality information of 2-mode nodes, we can effectively capture the heterogeneous community structure in bipartite networks. From the following illustrative example, we can see that the coclustering results of NMTF intuitively agree with the real community structure of the bipartite network and directly indicate the relationship between subject clusters and object clusters.

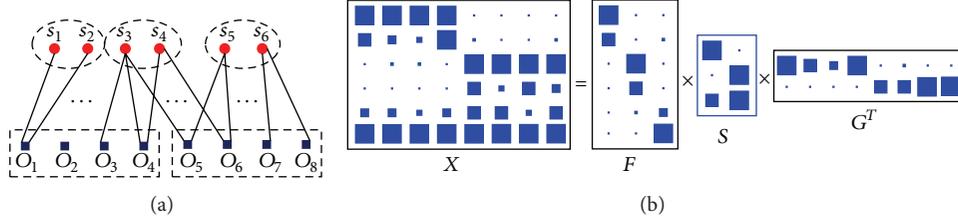


FIGURE 3: (a) A toy bipartite network with $|N_S| = 6$, $|N_O| = 8$. (b) The NMTF procedure for community detection on the toy bipartite network $X \in \mathbb{R}^{6 \times 8}$ with 3 subject-clusters and 2 object-clusters.

3. Fast NMTF Algorithm for Coclustering Community Detection

In this section, we propose a fast nonnegative matrix trifactorization (F-NMTF) method as the solution of community detection on bipartite networks. First, we respectively formulate the affinity matrices of 2-mode node sets in bipartite networks and incorporate them as two manifold regularizations of subject-set and object-set in the objective function. Then we decompose the original optimization problem into two smaller subproblems and present an optimization algorithm to develop the iterative updating rules of three factor matrices. For convenience, we present in Notations section the important notations used in this paper.

3.1. Constructing the Affinity Matrices. For the original NMTF framework (on (1)), it just considers the intertype information of 2-mode nodes. Such formulation assumes each subnetwork to be independent and fails to model the bipartite networks in a unified way. Recently, some researchers [15, 24] have found that coclustering data on manifold is well applied to bipartite networks for regularization-based clustering, because it can promote the performance of the intrinsic structure discovery in multimode networks. As a result, by constructing the two affinity matrices, the optional intratype information of 2-mode nodes is incorporated into NMTF as manifold regularizations. More importantly, we can exploit the manifold structures in both subject and object spaces to group like-minded users from different social perspectives, thus strengthening the community detection in bipartite networks.

First, we construct the affinity matrix W_F as [15, 25] whose entries correspond to subject nodes $N_S = \{S_1, \dots, S_m\}$. Generally, if nodes S_i and S_j share most connections to the nodes of the other mode, it means they are close to each other, and then their corresponding indicator vectors F_i and F_j should be close as well. It is formulated as follows:

$$\frac{1}{2} \sum_{ij} \|F_i - F_j\|^2 \cdot W_F(i, j), \quad (3)$$

where $\|\cdot\|$ is the Frobenius norm, F_i is the i th row of the subject indicator matrix F and indicates the community membership of subject-node S_i . $W_F(i, j)$ is the affinity

relationship between community indicator vectors F_i and F_j . For simplicity, we define the affinity matrix W_F as follows:

$$W_F(i, j) = \begin{cases} 1 & \text{if } S_j \in \mathcal{N}(S_i) \text{ or } S_i \in \mathcal{N}(S_j) \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $\mathcal{N}(S_i)$ denotes the k -nearest neighbor of S_i . In addition, other kinds of affinity can also be adopted, for example, heat kernel [26].

And (3) can be further rewritten as

$$\begin{aligned} & \frac{1}{2} \sum_{ij} \|F_i - F_j\|^2 \cdot W_F(i, j) \\ &= \sum_{i,j} F_i W_F(i, j) F_i^T - \sum_{i,j} F_i W_F(i, j) F_j^T \\ &= \sum_{i,j} F_i D_{ii}^F F_i^T - \sum_{i,j} F_i W_F(i, j) F_j^T \\ &= \text{tr}(F^T (D^F - W^F) F) \\ &= \text{tr}(F^T L_F F), \end{aligned} \quad (5)$$

where $D_{ii}^F = \sum_j W_F(i, j)$ is the diagonal degree matrix and $L_F = I - D_F^{-1/2} W_F D_F^{-1/2}$ is the normalized graph Laplacian of the subject-node set N_S .

Likewise, we also define the affinity matrix W_G whose entries correspond to object nodes $N_O = \{O_1, \dots, O_n\}$ as follows:

$$W_G(i, j) = \begin{cases} 1 & \text{if } O_j \in \mathcal{N}(O_i) \text{ or } O_i \in \mathcal{N}(O_j) \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

where $\mathcal{N}(O_i)$ denotes the k -nearest neighbor of O_i . And (3) for object-node set can be further rewritten as

$$\frac{1}{2} \sum_{ij} \|G_i - G_j\|^2 \cdot W_G(i, j) = \text{tr}(G^T L_G G), \quad (7)$$

where G_i is the i th column of the object indicator matrix G and indicates the community membership of object-node O_i . $D_{ii}^G = \sum_j W_G(i, j)$ is the diagonal degree matrix, and $L_G = I - D_G^{-1/2} W_G D_G^{-1/2}$ is the normalized graph Laplacian of the object-node set N_O .

Here we construct two affinity matrices based on the intratype information of 2-mode nodes and introduce them as manifold constraints to explore the hidden community structures of bipartite networks. In the following, we impose these two constraints on NMTF to achieve additional flexibility and incorporated the intratype information to enhance the orthogonality and accuracy on matrix factorization.

3.2. Objective Function of F-NMTF. Applying the two manifold regularizations in (1), the objective of our F-NMTF approach is transformed to minimize

$$J = \|X - FSG^T\|^2 + \alpha \operatorname{tr}(G^T L_G G) + \beta \operatorname{tr}(F^T L_F F) \quad (8)$$

s.t. $F \in \mathbb{R}^{m \times r} \geq 0$, $G \in \mathbb{R}^{n \times r} \geq 0$,

where $\alpha \geq 0$ and $\beta \geq 0$ are regularization parameters to balance the reconstruction error of F-NMTF in the first term and manifold regularizations in the second and third terms. Adding regulations to NMTF is a common strategy, since it not only improves interpretability, but also enhances numerical stability of the estimation by making the NMTF optimization less underconstrained.

Since F and G are constrained to be cluster indicator matrices, it is difficult to solve (8) in general. Hence we simplify this problem by using the following proposition.

Proposition 1. Let a symmetric matrix $A = D^{-1/2}WD^{-1/2}$; the term $\min \operatorname{tr}[C^T(I - A)C]$ is equivalent to the following optimization subproblem:

$$\min_{Q^T Q = I} \operatorname{tr} \|C - BQ\|^2. \quad (9)$$

Proof. $\min \operatorname{tr}[C^T(I - A)C]$ is equivalent to $\max \operatorname{tr}[C^T A C]$ that is further equivalent to $\min \operatorname{tr} \|CC^T - A\|^2$. By definition, the low-rank approximation of A is given by $A = BQ(BQ)^T$; then the objective term becomes $\min_{Q^T Q = I} \operatorname{tr} \|CC^T - BQ(BQ)^T\|^2$. C approximating BQ is equivalent to CC^T approximating $BQ(BQ)^T$. Hence, $\min \operatorname{tr}[C^T(I - A)C]$ can be reasonably transformed to (9), thus completing the proof of Proposition 1. \square

The two manifold regularizations terms in (8) can be rewritten as

$$\begin{aligned} \operatorname{tr}(F^T L_F F) &= F^T (I - D_F^{-1/2} W_F D_F^{-1/2}) F, \\ \operatorname{tr}(G^T L_G G) &= G^T (I - D_G^{-1/2} W_G D_G^{-1/2}) G. \end{aligned} \quad (10)$$

Then, applying Proposition 1 in (8), the objective of our F-NMTF approach is transformed to minimize

$$J = \|X - FSG^T\|^2 + \alpha \|G - B_G Q_G\| + \beta \|F - B_F Q_F\| \quad (11)$$

s.t. $F \in \mathbb{R}^{m \times r} \geq 0$, $G \in \mathbb{R}^{n \times r} \geq 0$,

where B_F and B_G are computed from L_F and L_G following the procedures described in Proposition 1.

This section presents a general framework of NMTF, which are developed to introduce the affinity matrices as the dual manifold regularization. Hence, our F-NMTF method successfully incorporates the intratype and intertype information of the bipartite network to coclustering multitype relational networks.

3.3. Optimization Iteration. Due to the fact that NMTF algorithms always have high computation complexity, it is essential and valuable to introduce fast iterative rules for iteration optimization. In this subsection, as a step toward accelerating convergence of NMTF, we apply a fast iterative algorithm to alternatively optimize the objective, with computing one factor matrix while fixing the other variables.

Theorem 2. Given a general optimization problem:

$$\min \|B - AQ\|^2, \quad \text{s.t. } Q^T Q = I, \quad (12)$$

when A and B are fixed, the optimum Q is given by $Q = UV^T$, where $H = A^T B$ and the singular value decomposition (SVD) of H is given by $H = U\Lambda V^T$.

Proof. When B is fixed, $\min \|B - AQ\|^2$ is equivalent to $\max_{Q^T Q = I} \operatorname{tr}(Q^T H)$. Let $\operatorname{tr}(Q^T H) = \operatorname{tr}(QU\Lambda V^T) = \operatorname{tr}(\Lambda V^T Q U) = \operatorname{tr}(\Lambda Z) = \sum_i \lambda_{ii} z_{ii}$, where $Z = V^T Q U$, λ_{ii} , and z_{ii} are the (i, i) th entry of Λ and Z , respectively.

Note that Z is orthonormal; that is, $Z^T Z = I$, thus $z_{ii} \leq 1$. λ_{ii} is the i th singular value of H , $\lambda_{ii} \geq 0$. Therefore, $\operatorname{tr}(Q^T H) = \sum_i \lambda_{ii} z_{ii} \leq \sum_i \lambda_{ii}$. That is to say, $\operatorname{tr}(Q^T H)$ reaches its maximum when $Z = I$. Therefore, the solution to $\max_{Q^T Q = I} \operatorname{tr}(Q^T H)$ is $Q = UZ^T V^T = UV^T$. Theorem 2 is proved. \square

According to Theorem 2, the optimization problem of F-NMTF can be decoupled into two subproblems with much smaller sizes, and the decoupled subproblems would involve much less matrix multiplications. In this way our approach can be computationally efficient and scale well to large-scale real world networks.

Now we alternatively optimize the four variables of the objective function (11). First, fixing F and G , by setting the derivative of (11) with respect to S as 0, we obtain

$$S = (F^T F)^{-1} F^T X G (G^T G)^{-1}. \quad (13)$$

Second, by fixing S , F , and G , we can decouple (11) into two following subproblems:

$$\min_{Q_F^T Q_F = I} \|F - B_F Q_F\|^2, \quad \min_{Q_G^T Q_G = I} \|G - B_G Q_G\|^2. \quad (14)$$

Applying Theorem 2 to (11), $Q_F = U_F V_F^T$ where U_F and V_F are obtained by SVD on $B_F^T F$; $Q_G = U_G V_G^T$ where U_G and V_G are obtained by SVD on $B_G^T G$.

Then, we fix F , S , and Q_G to update G , and (11) is decoupled to the following simple problems:

$$\min \|X_{\cdot i} - Y_G G_{\cdot i}^T\|^2 + \alpha \|G_{\cdot i} - (Z_G)_{\cdot i}\|^2, \quad (15)$$

where $Y_G = FS$; $(Z_G)_{\cdot i}$ denotes the i th row of $Z_G = B_G Q_G$.

```

Input: matrix  $X \in \mathbb{R}^{m \times n}$ ,  $r$ ;
Output: Community detection results;
(1) Initialize the factor matrices  $F_1 \geq 0, G_1 \geq 0$ ;
(2) Calculate  $L_G, L_F, B_G$  and  $B_F$  with Proposition 1;
(3) % Obtain Community indicator matrix  $G, F$ %
(4) repeat
(5)   Compute  $S$  by (13);
(6)   Compute  $Q_G = U_G V_G^T$ ; //  $U_G$  and  $V_G$  are obtained with SVD on  $B_G^T G$ ;
(7)   Compute  $Q_F = U_F V_F^T$ ; //  $U_F$  and  $V_F$  are obtained with SVD on  $B_F^T F$ ;
(8)   Update  $G$  by (16);
(9)   Update  $F$  by (18);
(10) until Converges;
(11) % Inferring community labels from  $G, F$ %
(12)  $C_i \leftarrow \phi, \forall i = 1, 2, \dots, r$ ;
(13) for  $i \in S \cup O$  do
(14)   add subject node  $S_i$  to  $C_j$  when  $F_{ij} = 1$ ;
(15)   add object node  $O_j$  to  $C_i$  when  $G_{ij} = 1$ ;
(16) end

```

ALGORITHM 1: Community detection using F-NMTF.

Due to orthogonality and sparsity, we emphasize that in each row (column) of G (F) only one nonzero element can be set to 1, which clearly indicate the community the corresponding node belongs to. Suppose G_i corresponds to i th community of the object-node set; thus the solution G is obtained by

$$G_{ij} = \begin{cases} 1 & j = \arg \min_k (\|X_{\cdot i} - (Y_G)_{\cdot k}\|^2 - 2\alpha (Z_G)_{ik}) \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

When fixing G, S , and Q_F , let $Y_F = S G^T$, $Z_F = B_F Q_F$, and (11) is decoupled to the following problems:

$$\min \|X_{\cdot j} - F_{\cdot j} Y_F\|^2 + \beta \|F_{\cdot j} - (Z_F)_{\cdot j}\|^2. \quad (17)$$

Similarly, we obtain F as

$$F_{ji} = \begin{cases} 1 & i = \arg \min_k (\|X_{\cdot j} - (Y_F)_{\cdot k}\|^2 - 2\beta (Z_F)_{jk}) \\ 0 & \text{otherwise,} \end{cases} \quad (18)$$

where $F_{\cdot j}$ corresponds to j th community of the subject-node set. Repeat this procedure until convergence. The convergence and correctness of this alternating iterative procedure have been proofed in [15, 23]. We skip it due to space.

After iterations, we can infer the community membership of heterogeneous nodes based on the F-NMTF results. For simplicity, the community indices of subject/object nodes are determined by taking the maximum of each column/row in F/G . The detailed procedure is illustrated as shown in Algorithm 1.

In our algorithm, G and F are sparse matrices, and the computation of them only involves vector norm enumeration without matrix multiplication; thus it is more computationally efficient. Moreover, instead of minimizing each matrix factor optimally with time-consuming multiplications of

large matrices, F-NMTF decouples the original optimization problem into two smaller subproblems requiring much less matrix multiplications and coclusters the 2-mode relational nodes in bipartite networks. Compared to the other representative NMF solutions, it effectively optimizes the matrix factorization with faster convergence and lower computational complexity.

4. Experimental Results

In this section, the experiments use a series of computer-generated synthetic networks and some real networks to validate the algorithms' performance. For all the bipartite networks, we compare the experimental results with other 4 well-known algorithms of community detection: Kmeans [27], GNMF [14], DRCC [15], and BNMTF [16]. All the experiments are performed on an Intel Core2 Duo 2.0 GHz PC with 1 GB RAM, running on Windows XP.

For NMF-based methods, including GNMF, DRCC, and our F-NMTF methods, the number of nearest neighbor k is set by the grid $\{1, 2, \dots, 10\}$ as in [25], and the regularization parameters α and β are set to 0.1. In addition, we obtain the community numbers r from the method as suggested in [10], which has been shown to well predict the number of network communities. In our experiments, we repeat each method with 50 times on all the bipartite networks and compute the average results.

In the following tests, different measures are introduced to evaluate the partition quality of the classical algorithms for community detection in bipartite networks. Since the structures of synthetic networks are all known, we adopt two standard measures widely used for clustering: normalized mutual information (NMI) [28] and execution time to quantify the partition quality of the community detection methods. For the real networks with unknown structure, we use bipartite modularity [29] and execution time for the collective validation. In particular, bipartite modularity

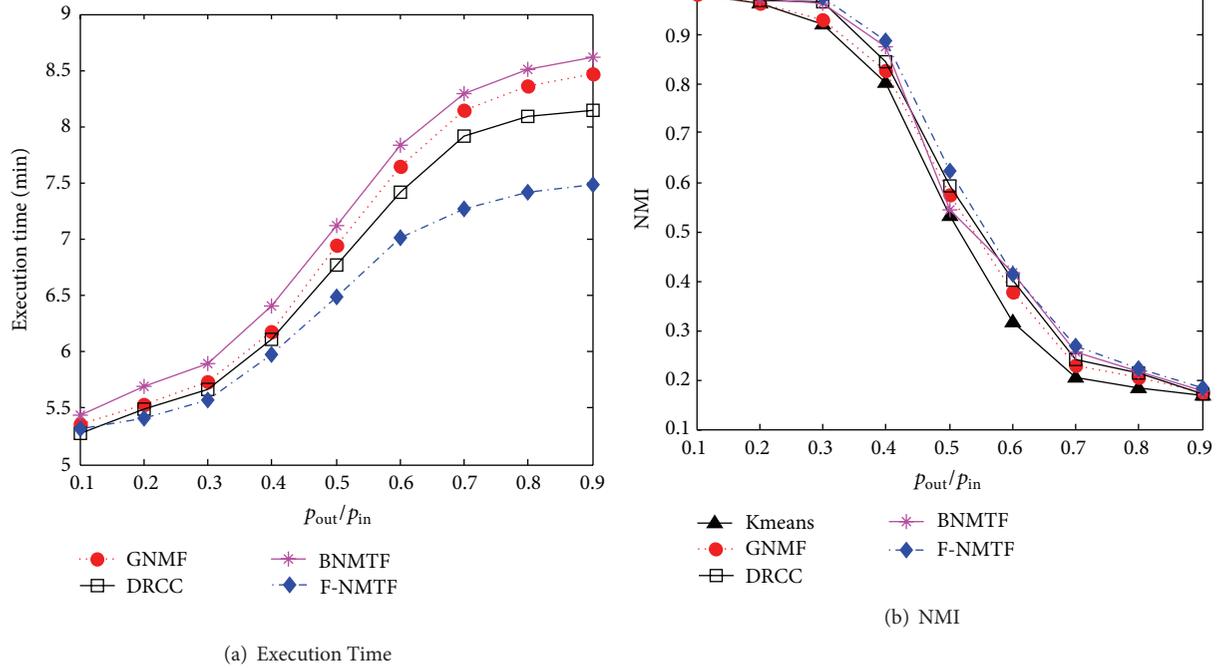


FIGURE 4: The average execution time and NMI of the community detection methods on the model bipartite networks with different p_{out}/p_{in} parameters ranging from 0.1 to 0.9.

is extended from Newman–Girvan modularity [12] and regarded as an effective standard to measure the community structure of partition results.

(1) *Synthetic Networks.* Here we have applied the model bipartite networks [29] to generate 9 groups of benchmark datasets with known structures. Each group contains 10 networks that are generated with the same parameters. In this experiment, we randomly choose 5 networks of each group to obtain the average values for comparison. Here p_{in} is fixed at the value of 0.9 while p_{out} is varied by tuning p_{out}/p_{in} from 0.1 to 0.9 with steps of 0.1. p_{in} and p_{out} , respectively, denote the intracommunity and intercommunity link probability. Generally, the higher the intracommunity link probability of the network is, the stronger community structure can be detected. Conversely, as p_{out}/p_{in} becomes bigger, the community structure of the model bipartite network should become more and more obscure. We set community number $|N_C| = 80$, subject-node number $|N_S| = 600$, and object-node number $|N_O| = 1000$.

Due to the fact that Kmeans costs much less time for computation and the fluctuation of p_{out}/p_{in} has limited effect on it, we mainly display the execution time of NMF-based methods in Figure 4(a), where p_{out}/p_{in} of the model bipartite networks ranges from 0.1 to 0.9. As p_{out}/p_{in} increases, the community structure of the model bipartite networks become weaker, and all the algorithms suffer varying degrees of performance degradation, and their execution time rises rapidly. Because iterative methods are usually necessary for NMF solution, they need to update matrix factors by multiplying each entry in each iteration round until convergence. Meanwhile, we can see that F-NMTF converges faster than

other NMF methods, and the gap of execution time between them becomes greater, because F-NMTF is decoupled into two subproblems with much smaller sizes and optimizes the factorization of each matrix. Even the community structure of the synthetic bipartite networks tends to be weaker and weaker; our algorithm also shows much better robustness than the other three algorithms.

Figure 4(b) shows the average NMI values of different algorithms against changes of p_{out}/p_{in} . When $p_{out}/p_{in} \leq 0.4$, the NMI scores of all the algorithms performance exceed 0.8. Specifically, three NMF methods have the similar well performance under the strong community structures, and only Kmeans is slightly inferior to other methods. But when $p_{out}/p_{in} \geq 0.5$, the community cohesion of model bipartite networks degrades along with p_{out}/p_{in} gradual increase, and the performance of F-NMTF still remains to have relatively higher NMI scores and keeps the stability and accuracy of community detection, rather than tending to quickly decrease like other algorithms. Our method is also superior in terms of stability as well as approximation accuracy, which means that it can achieve small recovery errors for various source networks. Specifically, the average NMI of F-NMTF is up to 5.83% better than that returned by GNMF, and 3.13% better than DRCC. In summary, the performances of F-NMTF are highly competitive to those of other methods on bipartite networks.

(2) *Real Networks.* Real networks are always more irregular and various than synthetic networks and have more complex community structures. Here we choose 5 popular real bipartite networks in different sizes: Southwomen [30], Scotland [31], Irvine forum [32], MovieLens [33], and Cond-mat [34].

TABLE 1: Basic properties of the real bipartite networks adopted in our experimentation.

Network	Subject node	Object node	Edge	Weighted
Southwomen	14	18	89	No
Scotland	108	136	356	Yes
Irvine forum	522	899	33720	No
MovieLens	943	1682	100000	Yes
Cond-mat	16726	22015	58595	No

These networks’ basic properties (about nodes, edges) are presented in Table 1.

Southwomen shows the participation of 18 women in 14 social events over a nine-month period. The dataset was collected in the Southern United States of America in the 1930s. There is an edge for every woman who participates in an event. Irvine forum contains user posts to forums. The users are students at the University of California, Irvine. An edge represents a forum message. MovieLens consists of 100000 user-movie ratings from <https://www.movielens.org/>. An edge between a user and a movie represents a rating of the movie by the user. Cond-mat contains authorship links between authors and publications in the arXiv condensed matter section from 1995 to 1999. An edge represents an authorship connecting an author and a paper. Scotland dataset contains the corporate interlocks in Scotland in the beginning of the twentieth century. It lists the (136) multiple directors of the 108 largest joint stock companies in Scotland in 1904-1905, including 64 nonfinancial firms, 8 banks, 14 insurance companies, and 22 investment and property companies.

The average execution times found by different algorithms are shown in Table 2(a). We can see that Kmeans costs much less time than NMF-based algorithms, as it does not need the matrix factorization iterations. For all the real bipartite networks, F-NMTF effectively accelerates the convergence speed of nonnegative matrix factorization and converges in less iterations and CPU seconds than other NMF methods. Because the network scales are quite different, the corresponding performances of F-NMTF are different, too. For the larger networks, F-NMTF has a greater competitive advantage than other methods. Our method is only slower than Kmeans, which, however, has much worse clustering performance.

Table 2(b) shows the bipartite modularity values found by different algorithms. The methods using manifold constraints, including GNMF, DRCC, and F-NMTF, generally achieve better clustering results, which are consistent with the widely accepted hypothesis that clustering of both intratype and intertype information can help clustering of bipartite networks. F-NMTF method attains the maximum modularity community structure for most test cases, which means that our method has better partition quality and achieves accuracy community structure on the real bipartite networks. More importantly, our method does not suffer from the problems of modularity optimization methods and clusters both the rows and columns of the networks simultaneously by making full use of the duality information between 2-mode nodes,

TABLE 2: Community detection results: (a) execution time; (b) bipartite modularity. Each entry shows the average values of the detected community structure for the 5 community detection methods on the real bipartite network.

(a)					
Network	Kmeans	GNMF	DRCC	BNMTF	F-NMTF
Southwomen (ms)	32.3	69.2	66.9	76.9	59.1
Scotland (ms)	419.5	896.0	872.5	964.3	836.9
Irvine forum (min)	3.29	8.77	8.35	9.14	7.86
MovieLens (min)	11.68	22.63	21.90	24.85	19.40
Cond-mat (min)	10.83	17.85	17.39	18.26	16.18
(b)					
Network	Kmeans	GNMF	DRCC	BNMTF	F-NMTF
Southwomen	0.293	0.336	0.345	0.329	0.351
Scotland	0.417	0.513	0.525	0.492	0.528
Irvine forum	0.469	0.522	0.531	0.516	0.543
MovieLens	0.523	0.618	0.623	0.607	0.655
Cond-mat	0.596	0.715	0.733	0.719	0.772

which can greatly enhance the performance of clustering algorithms. Therefore, compared to other 4 methods, we can conclude that F-NMTF has competitive clustering performance in terms of both accuracy and partition quality against popular community detection methods and with much faster computational speed.

5. Conclusions

In this work, we proposed a novel fast nonnegative matrix trifactorization method for community detection in bipartite networks. Based on the idea of nonnegative matrix trifactorization, we introduce the intratype information of 2-mode nodes into NMTF via the dual manifold regularizations, thus helping to extract meaningful communities of the bipartite networks. Meanwhile, we decouple the NMTF problem into two subproblems with much smaller sizes, and the decoupled subproblems involve much less matrix multiplications, which make our approach of particular use for real world large-scale data.

Different from the traditional NMF-based methods, our work is an instructive attempt to cocluster multitype nodes in bipartite networks. In practice, our coclustering NMTF framework jointly takes intertype and intratype information of 2-mode nodes into considerations, thus making the partitioning results more reasonable and effective and detecting communities with high accuracy and quality. Experimental results on the synthetic and real-world datasets show that our algorithm is a competitive method to explore community structures in bipartite networks.

Notations

- X : Adjacent matrix of a bipartite network
- N_S : Subject-node set of a bipartite network
- N_O : Object-node set of a bipartite network

m : Number of subject-nodes in X
 n : Number of object-nodes in X
 r : Prior-known community number
 F : Indicator matrix of subject-node partition
 G : Indicator matrix of object-node partition
 F_i : i th row of F
 G_i : i th column of G
 W_F : Affinity matrix of subject-node set
 W_G : Affinity matrix of object-node set.

Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

Acknowledgments

This work was supported by the National Science and Technology Major Project of the Ministry of Science and Technology of China under Grant no. 2012ZX03006002 and the National High Technology Development 863 Program of China under Grant no. 2011AA010604.

References

- [1] M. Rosvall and C. T. Bergstrom, "Maps of random walks on complex networks reveal community structure," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 105, no. 4, pp. 1118–1123, 2008.
- [2] A. Lancichinetti, F. Radicchi, J. J. Ramasco, and S. Fortunato, "Finding statistically significant communities in networks," *PLoS ONE*, vol. 6, no. 4, Article ID e18961, 2011.
- [3] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2008, no. 10, Article ID P10008, 2008.
- [4] M. Latapy, C. Magnien, and N. D. Vecchio, "Basic notions for the analysis of large two-mode networks," *Social Networks*, vol. 30, no. 1, pp. 31–48, 2008.
- [5] I. S. Dhillon, "Co-clustering documents and words using bipartite spectral graph partitioning," in *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '01)*, pp. 269–274, 2001.
- [6] I. S. Dhillon, S. Mallela, and D. S. Modha, "Information-theoretic co-clustering," in *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '03)*, pp. 89–98, August 2003.
- [7] V. Sindhwani, J. Hu, and A. Mojsilovic, "Regularized co-clustering with dual supervision," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS '08)*, vol. 21, pp. 1505–1512, 2008.
- [8] I. Psorakis, S. Roberts, M. Ebdon, and B. Sheldon, "Overlapping community detection using Bayesian non-negative matrix factorization," *Physical Review E*, vol. 83, no. 6, Article ID 066114, 2011.
- [9] X. Cao, X. Wang, D. Jin, Y. Cao, and D. He, "Identifying overlapping communities as well as hubs and outliers via nonnegative matrix factorization," *Scientific Reports*, vol. 3, article 2993, 2013.
- [10] N. P. Nguyen, T. N. Dinh, S. Tokala, and M. T. Thai, "Overlapping communities in dynamic networks: their detection and mobile applications," in *Proceedings of the 17th Annual International Conference on Mobile Computing and Networking*, pp. 85–95, ACM, September 2011.
- [11] D. He, D. Jin, C. Baquero, and D. Liu, "Link community detection using generative model and nonnegative matrix factorization," *PLoS ONE*, vol. 9, no. 1, Article ID e86899, 2014.
- [12] M. E. J. Newman, "Modularity and community structure in networks," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 103, no. 23, pp. 8577–8582, 2006.
- [13] S. Fortunato and M. Barthélemy, "Resolution limit in community detection," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 104, no. 1, pp. 36–41, 2007.
- [14] D. Cai, X. He, J. Han, and T. S. Huang, "Graph regularized non-negative matrix factorization for data representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 8, pp. 1548–1560, 2011.
- [15] Q. Gu and J. Zhou, "Co-clustering on manifolds," in *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 359–368, 2009.
- [16] Y. Zhang and D.-Y. Yeung, "Overlapping community detection via bounded nonnegative matrix tri-factorization," in *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '12)*, pp. 606–614, Beijing, China, August 2012.
- [17] H. Zhang, X. Liu, J. Wang, and H. R. Karimi, "Robust H_∞ sliding mode control with pole placement for a fluid power electrohydraulic actuator (EHA) system," *The International Journal of Advanced Manufacturing Technology*, vol. 73, no. 5–8, pp. 1095–1104, 2014.
- [18] H. Zhang, Y. Shi, and J. Wang, "Observer-based tracking controller design for networked predictive control systems with uncertain Markov delays," *International Journal of Control*, vol. 86, no. 10, pp. 1824–1836, 2013.
- [19] H. Zhang, Y. Shi, and J. Wang, "On energy-to-peak filtering for nonuniformly sampled nonlinear systems: a markovian jump system approach," *IEEE Transactions on Fuzzy Systems*, vol. 22, no. 1, pp. 212–222, 2014.
- [20] H. Zhang and J. Wang, "Combined feedback-feedforward tracking control for networked control systems with probabilistic delays," *Journal of the Franklin Institute*, vol. 351, no. 6, pp. 3477–3489, 2014.
- [21] Z.-Y. Zhang, Y. Wang, and Y.-Y. Ahn, "Overlapping community detection in complex networks using symmetric binary matrix factorization," *Physical Review E: Statistical, Nonlinear, and Soft Matter Physics*, vol. 87, no. 6, Article ID 062803, 2013.
- [22] S. Mankad and G. Michailidis, "Structural and functional discovery in dynamic networks with non-negative matrix factorization," *Physical Review E*, vol. 88, no. 4, Article ID 042812, 2013.
- [23] C. Ding, T. Li, W. Peng, and H. Park, "Orthogonal nonnegative matrix tri-factorizations for clustering," in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '06)*, pp. 126–135, August 2006.
- [24] P. Li, J. Bu, C. Chen, Z. He, and D. Cai, "Relational multimani-fold coclustering," *IEEE Transactions on Cybernetics*, vol. 43, no. 6, pp. 1871–1881, 2013.
- [25] F. Shang, L. C. Jiao, and F. Wang, "Graph dual regularization non-negative matrix factorization for co-clustering," *Pattern Recognition*, vol. 45, no. 6, pp. 2237–2250, 2012.

- [26] J. J.-Y. Wang, H. Bensmail, and X. Gao, "Multiple graph regularized nonnegative matrix factorization," *Pattern Recognition*, vol. 46, no. 10, pp. 2840–2847, 2013.
- [27] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: analysis and implementation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 881–892, 2002.
- [28] L. Danon, A. Díaz-Guilera, J. Duch, and A. Arenas, "Comparing community structure identification," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 9, Article ID P09008, 2005.
- [29] R. Guimerà, M. Sales-Pardo, and L. A. N. Amaral, "Module identification in bipartite and directed networks," *Physical Review E: Statistical, Nonlinear, and Soft Matter Physics*, vol. 76, no. 3, Article ID 036102, 2007.
- [30] A. Davis, B. B. Gardner, and M. R. Gardner, *Deep South: A Social Anthropological Study of Caste and Class*, The University of Chicago Press, Chicago, Ill, USA, 1941.
- [31] J. Scott and M. Hughes, *The Anatomy of Scottish Capital: Scottish Companies and Scottish Capital*, Croom Helm, London, UK, 1980.
- [32] T. Opsahl, "Triadic closure in two-mode networks: redefining the global and local clustering coefficients," *Social Networks*, vol. 35, no. 2, pp. 159–167, 2013.
- [33] MovieLens network dataset—KONECT, 2014, http://konect.uni-koblenz.de/networks/movielens-100k_rating.
- [34] "Arxiv cond-mat network dataset-KONECT," April 2014, <http://konect.uni-koblenz.de/networks/opsahl-collaboration>.



Hindawi

Submit your manuscripts at
<http://www.hindawi.com>

