

Community Detection and Correlated Attribute Cluster Analysis on Multi-Attributed Graphs

Hiroyoshi Ito[†], Takahiro Komamizu[‡], Toshiyuki Amagasa[‡], and Hiroyuki Kitagawa[‡]

[†]Graduate School of Systems and Information Engineering, University of Tsukuba

[‡]Center for Computational Sciences, University of Tsukuba

hiro.3188@kde.cs.tsukuba.ac.jp, taka-coma@acm.org, {amagasa, kitagawa}@cs.tsukuba.ac.jp

ABSTRACT

Multi-attributed graphs, in which each node is characterized by multiple types of attributes, are ubiquitous in the real world. Detection and characterization of communities of nodes could have a significant impact on various applications. Although previous studies have attempted to tackle this task, it is still challenging due to difficulties in the integration of graph structures with multiple attributes and the presence of noises in the graphs. Therefore, in this study, we have focused on clusters of attribute values and strong correlations between communities and attribute-value clusters. The graph clustering methodology adopted in the proposed study involves Community detection, Atttribute-value clustering, and deriving Relationships between communities and attribute-value clusters (CAR for short). Based on these concepts, the proposed multi-attributed graph clustering is modeled as CAR-clustering. To achieve CAR-clustering, a novel algorithm named CARNMF is developed based on non-negative matrix factorization (NMF) that can detect CAR in a cooperative manner. Results obtained from experiments using real-world datasets show that the CARNMF can detect communities and attribute-value clusters more accurately than existing comparable methods. Furthermore, clustering results obtained using the CARNMF indicate that CARNMF can successfully detect informative communities with meaningful semantic descriptions through correlations between communities and attribute-value clusters.

1 INTRODUCTION

Community detection is a task to detect densely connected sub-graphs as communities. Nodes in a community tend to share same or similar properties, such phenomenon is called *homophily effect* [11, 17], meaning that nodes having similar properties tend to link together. Because diverse applications are derived from the nature of real communities, community detection is important in graph/network analyses. Examples include node property estimations [7, 9, 24], community-wise information recommendations [10], and semantic reasoning for nodes/edges [1].

Moreover, using the attributes in a graph is advantageous to realize high-quality community detection as well as to understand the characteristics of communities. Multi-attributed graphs are reasonable models of real-world networks such as social networks, co-author networks, protein-protein interaction networks, etc. In fact, several works have proposed algorithms that employ attribute information (i.e., shared interests or functional behaviors of each community) to detect not only communities but also their semantic meanings [19, 23, 25, 26].

However, community detection and extraction of semantics in multi-attributed graphs remain challenging due to difficulties on integrating graph structures and multiple attributes of different types. Community detection and extraction of semantics involves multiple steps. First, useful information from each attribute must be extracted because certain node attributes describe different aspects. Second, all extracted information must be exploited to enhance community detection by effectively integrating heterogeneous information. Notice that the previous works [19, 23, 25, 26] do not differentiate multiple attributes, that is, they consider multiple attributes equally. Moreover, real-world graphs are often incomplete and noisy. That is, some edges or nodes may be missing or attribute values may contain incorrect values, leading to inappropriate results.

To overcome these difficulties, we propose a novel clustering scheme based on the following two assumptions:

- (1) *Relevant attribute values form clusters by attribute type.* This is based on the observation that an attribute reflects a node's interests in a network. Hence, an attribute tends to be associated to a specific group of values related to an interest. For example, in a co-author network where the nodes correspond to the authors (researchers), each author typically has specific research interests (e.g., AI, data mining, and database). Thus, attributes (e.g., paper title and conference) present biased values according to interests. Consequently, it is possible to identify clusters of attributes values (attribute-value clusters) reflecting a node's interests.
- (2) *Communities are strongly correlated with attribute-value clusters.* This is related to the previous assumption. Consider the example above. The nodes in a community share similar interests (e.g., research interests) and consequently, similar attribute-value clusters (e.g., research topics, and conferences). Conversely, if some nodes (researchers) have similar attribute values, they should share similar interests and can be grouped in the same community.

Exploiting the correlation between communities and multiple attributes should improve the quality of community detection as well as attribute-value clustering. Using the information from different sources (attributes) to alleviate the effect of noise (e.g., missing values and errors), we simultaneously implement community detection and attribute-value clustering.

Based on the aforementioned ideas, we study a novel clustering scheme for multi-attributed graphs, called CAR-clustering. CAR includes Community detection, Atttribute-value clustering, and deriving Relationships between communities and attribute-value clusters for multi-attributed graphs. Additionally, we develop a novel clustering algorithm called CARNMF, which employs a non-negative matrix factorization (NMF).

The contributions of this paper are summarized as follows:

- We propose a novel clustering scheme CAR-clustering to address two technical questions. (i) Given a multi-attributed

© 2018 Copyright held by the owner/author(s). Published in the Workshop Proceedings of the EDBT/ICDT 2018 Joint Conference (March 26, 2018, Vienna, Austria) on CEUR-WS.org (ISSN 1613-0073). Distribution of this paper is permitted under the terms of the Creative Commons license CC-by-nc-nd 4.0.

graph, how can community detection and attribute-value clustering be performed for different types of attributes in a cooperative manner? (ii) How should reasonable relationships be determined between communities and attribute-value clusters for each type of attribute?

- We develop a novel algorithm CARNMF, which achieves CAR-clustering. Specifically, a dedicated loss function is designed to perform multiple NMFs simultaneously.
- We conduct experiments using real-world datasets (DBLP computer science bibliography and arXiv physics bibliography). The accuracy of CARNMF with respect to community detection and attribute-value clustering and a comparison to other methods are examined. Relative to comparative methods, CARNMF achieves a better accuracy of up to 11% for community detection and up to 22% for attribute-value clustering. Furthermore, CARNMF detects informative communities and their rich semantic descriptions by correlating multiple types of attribute-value clusters.

2 RELATED WORK

Community detection in graphs is a current topic of interest in graph analysis and AI research. Existing works for non-attributed graphs can be categorized according to the techniques used: graph separation [9, 20], probabilistic generative model [27], and matrix factorization [12, 18, 24]. [9] defined *modularity*, which indicates how separated a community is from other nodes. More comprehensive surveys can be found in [6, 22].

Recently, several works have addressed the problem of detecting communities and their semantic descriptions on node-attributed graphs. [25] proposed *CESNA*, where communities and their attributes are simultaneously detected in an efficient manner. [23] proposed *SCI* to detect communities and their semantics using NMF. [19] proposed a probabilistic generative model called the *author-topic model* to model communities and related topics. [2] proposed *COMODO* to detect communities with shared properties using subgroup discovery techniques. Likewise, [26] proposed *LCTA*, where communities and their topics are modeled separately, and then their relationships are modeled using a probabilistic generative model. A comprehensive survey over these works can be found in [5].

The aforementioned works only consider single textual attributes or uniformly handle multiple attributes without any distinction. In reality, each attribute represents different aspects of the nodes. In our research, we deal with heterogeneous attributes individually. In addition to community detection, we perform clustering over attribute values for each attribute, which, in turn, can be used to improve the quality of communities detected.

Some works have investigated clustering over networks containing different types of nodes and/or edges. [3] studied community detection with characterization from multidimensional networks, which is defined as a graph consisting of a set of nodes and multiple types of edges. [4] studied subgraph detection from multi-layer graphs with edge labels. In contrast, we assume a different model where each node is characterized by multiple attributes. [21] proposed a scheme of ranking-based clustering for multi-typed heterogeneous networks, where two or more types of nodes are included. Similarly, [16] proposed an NMF-based method for such networks. These methods differ from ours in that they define a cluster consisting of all types of nodes. In other words, these methods cannot handle each attribute in a unique way. In contrast, our work deals with different attributes

individually, but solves community detection and attribute-value clustering in a unified manner.

3 PROBLEM STATEMENT

In this work, we deal with multi-attributed graphs, where each node is characterized by two or more attributes. Given such a graph, *CAR-clustering* is used to solve the following three sub-problems: community detection, attribute-value clustering, and derivation of relationships between communities and attribute-value clusters, which have been independently studied. Below, we provide the formal definitions which are necessary to define the clustering scheme.

3.1 Multi-Attributed Graph

Multi-attributed graph \mathbb{G} is defined by extending weighted graph \mathbb{G}' with several attributed graphs \mathbb{G}_t for attribute $t \in \mathbb{T}$. The following are formal definitions.

DEFINITION 1 (WEIGHTED GRAPH). *Weighted graph \mathbb{G}' is defined by a triplet, $\langle \mathbb{V}, \mathbb{E}, \mathbb{W} \rangle$, where \mathbb{V} is a set of nodes, $\mathbb{E} (\subseteq \mathbb{V} \times \mathbb{V})$ is a set of edges, and $\mathbb{W} : \mathbb{E} \rightarrow \mathbb{R}^+$ is a map of edge weights. \square*

DEFINITION 2 (ATTRIBUTED GRAPH). *Attributed graph $\mathbb{G}_t = \langle \mathbb{V} \cup \mathbb{X}_t, \mathbb{E}_t, \mathbb{W}_t \rangle$ of attribute $t \in \mathbb{T}$ is a bipartite graph consisting of set \mathbb{V} of nodes, set \mathbb{X}_t of attribute-values, a set of edges $\mathbb{E}_t \subseteq \mathbb{V} \times \mathbb{X}_t$, and $\mathbb{W}_t : \mathbb{E}_t \rightarrow \mathbb{R}^+$ is a map of edge weights. \square*

DEFINITION 3 (MULTI-ATTRIBUTED GRAPH). *Given weighted graph $\mathbb{G}' = \langle \mathbb{V}, \mathbb{E}, \mathbb{W} \rangle$ and a set of attributed graphs $\{\mathbb{G}_t\}_{t \in \mathbb{T}}$ where $\mathbb{G}_t = \langle \mathbb{V} \cup \mathbb{X}_t, \mathbb{E}_t, \mathbb{W}_t \rangle$, multi-attributed graph $\mathbb{G} = \langle \mathbb{G}', \{\mathbb{G}_t\}_{t \in \mathbb{T}} \rangle$ is a union of these graphs. \square*

3.2 CAR-clustering

Given a multi-attributed graph, information can be extracted from different perspectives. In this work, we extract *communities*, *attribute-value clusters*, and the *relationship* between them.

Community. For a multi-attributed graph, a set of nodes with the following properties is regarded as a community. (1) Nodes in a community are densely connected with each other and sparsely connected with other nodes. (2) Nodes in a community tend to share common values in distinct attributes. This study assumes that communities can overlap. That is, each node belongs to more than one community. This assumption is reasonable for real applications. Formally, given the number of communities ℓ , node $n \in \mathbb{V}$ belonging to community $c \in \mathbb{C}$ is described by probability distribution $p(n | c)$, where $|\mathbb{C}| = \ell$.

Attribute-value cluster. For attribute $t \in \mathbb{T}$ in a multi-attributed graph, similar or highly correlated attribute values can be grouped into attribute-value clusters. Herein, we assume overlapping clusters. That is, each attribute-value belongs to more than one cluster. Formally, given the number of clusters k_t of attribute $t \in \mathbb{T}$, cluster member $x \in \mathbb{X}_t$ for attribute-value cluster $s_t \in \mathbb{S}_t$ is described by probability distribution $p(x | s_t)$, where $|\mathbb{S}_t| = k_t$.

Relationship between a community and an attribute-value cluster. Nodes in a community often share common attribute-value clusters. Detecting such relationship is useful in many applications. Given community $c \in \mathbb{C}$ and attribute-value cluster $s_t \in \mathbb{S}_t$ of attribute $t \in \mathbb{T}$, the probability that c is related to s_t is described as the relationship between c and s_t . In this work, a community may be related to more than one attribute-value cluster. Formally, this is described by probability distribution $p(s_t | c)$.

CAR-clustering. CAR-clustering is formally defined by Definition 4.

DEFINITION 4 (CAR-CLUSTERING). *Given a multi-attributed graph \mathbb{G} , CAR-clustering is to perform community detection, attribute-value clustering, and detection of the relationship between the communities and the attribute-value clusters simultaneously.* \square

Solving these sub-problems simultaneously is more beneficial than evaluating each one independently because, in many cases, communities and attribute-value clusters are mutually correlated. Solving the problems simultaneously exploits this correlation, leading to improved results.

4 CARNMF – ALGORITHM FOR CAR-CLUSTERING

In this section, we propose an NMF (non-negative matrix factorization)-based algorithm, called CARNMF, for CAR-clustering. CARNMF models communities and attribute-value clusters. Additionally, we introduce an auxiliary matrix to maintain the relationship between the communities and the attribute-value clusters. A unified loss function is used to solve the different NMFs in a unified manner. It is assumed that the user gives the number ℓ of communities and the number k_t of clusters for each attribute $t \in \mathbb{T}$.

4.1 Matrix representation

We represent a multi-attributed graph by two sorts of matrices: an adjacency matrix $A \in \mathbb{R}^{|\mathbb{V}| \times |\mathbb{V}|}$ and attribute matrices $X^{(t)} \in \mathbb{R}^{|\mathbb{V}| \times |\mathbb{X}_t|}$ for $t \in \mathbb{T}$. An element $A_{u,v}$ of A corresponds to an edge $e_{u,v} = (u,v) \in \mathbb{E}$. $A_{u,v} = \mathbb{W}(e_{u,v}) / \sum_{e_{i,j} \in \mathbb{E}} \mathbb{W}(e_{i,j})$, indicating the joint probability for the presence of edge $e_{u,v}$. Similarly, for $t \in \mathbb{T}$, an element $X_{u,x}^{(t)}$ in $X^{(t)}$ corresponds to an edge $e_{u,x}^{(t)} \in \mathbb{E}_t$. $X_{u,x}^{(t)} = \mathbb{W}_t(e_{u,x}^{(t)}) / \sum_{v,y \in \mathbb{E}_t} \mathbb{W}_t(e_{v,y}^{(t)})$, indicating the joint probability of the presence of edge $e_{u,x}^{(t)}$.

4.2 Loss Function

We achieve CAR-clustering in terms of several NMFs, which correspond to the aforementioned sub-problems. To achieve CAR-clustering, we introduce loss functions for the sub-problems followed by a unified loss function.

Loss function for community detection. In CARNMF, communities \mathbb{C} are denoted by a matrix $U^* \in \mathbb{R}^{|\mathbb{V}| \times \ell}$, where each row and column correspond to a node $u \in \mathbb{V}$ and a community $c \in \mathbb{C}$, respectively. A cell $U_{u,c}^*$ represents probability $p(u | c)$. In probability $p(u, v | c)$, u and v are connected through community c , and is represented by $U_{u,c}^* U_{v,c}^*$. Moreover, joint probability $p(u, v)$, or the existence of edge $e_{u,v} \in \mathbb{E}$, is expressed as $\sum_{c \in \mathbb{C}} U_{u,c}^* U_{v,c}^*$. Therefore, when U^* minimizes the following loss function, U^* is the best approximation of the edges in the graph.

$$\arg \min_{U^* \geq 0} \left\| A - U^* (U^*)^T \right\|_F^2 \quad (1)$$

$$s.t. \forall 1 \leq c \leq \ell, \|U_{*,c}^*\|_1 = 1$$

where $\|\cdot\|_F^2$ and $\|\cdot\|_1$ represents the Frobenius norm and the ℓ_1 norm, respectively.

Loss function for attribute-value clustering. In CARNMF, attribute-value clusters \mathbb{S}_t of attribute $t \in \mathbb{T}$ are represented as

a matrix $V^{(t)} \in \mathbb{R}^{|\mathbb{X}_t| \times k_t}$, where each row and column correspond to an attribute $x \in \mathbb{X}_t$ and an attribute cluster $s_t \in \mathbb{S}_t$, respectively. A cell $V_{x,s_t}^{(t)}$ represents probability $p(x | s_t)$.

To derive $V^{(t)}$ from $X^{(t)}$, we introduce a matrix $U^{(t)} \in \mathbb{R}^{|\mathbb{V}| \times k_t}$, which denotes the relationships between the nodes and attribute-value clusters with probability $p(u | s_t)$. Using both matrices $U^{(t)}$ and $V^{(t)}$, probability $p(u, x | s_t)$, which is the existence of edge $e_{u,x}^{(t)} \in \mathbb{E}_t$ in terms of attribute-value cluster s_t , is calculated as $U_{u,s_t}^{(t)} V_{x,s_t}^{(t)}$. Moreover, probability $p(u, x)$ is derived as $\sum_{s_t \in \mathbb{S}_t} U_{u,s_t}^{(t)} V_{x,s_t}^{(t)}$. Therefore, when $U^{(t)}, V^{(t)}$ minimize loss function, $U^{(t)}, V^{(t)}$ represent the best approximation of the edges in the graph.

$$\arg \min_{U^{(t)}, V^{(t)} \geq 0} \left\| X^{(t)} - U^{(t)} (V^{(t)})^T \right\|_F^2 \quad (2)$$

$$s.t. \forall 1 \leq r \leq k_t, \|V_{*,r}^{(t)}\|_1 = 1$$

Loss function for relationship detection. In CARNMF, the relationships between communities and attribute-value clusters of attribute $t \in \mathbb{T}$ are represented as a matrix $R^{(t)} \in \mathbb{R}^{\ell \times k_t}$, where each row and column corresponds to a community $c \in \mathbb{C}$ and an attribute-value cluster $s_t \in \mathbb{S}_t$, respectively. The cell contains the probability $p(s_t | c)$. We assume $R^{(t)}$ is a linear transformation that maps U^* into $U^{(t)}$, where U^* and $U^{(t)}$ derived by Equation 1 and Equation 2, respectively. Therefore, when $R^{(t)}$ minimizes the loss function, $R^{(t)}$ represents the relationships between the communities and the attribute-value clusters.

$$\arg \min_{U^{(t)}, U^*, R^{(t)} \geq 0} \left\| U^{(t)} - U^* R^{(t)} \right\|_F^2 \quad (3)$$

$$s.t. \forall 1 \leq p \leq \ell, \|U_{*,p}^*\|_1 = 1, \|R_{p,*}^{(t)}\|_1 = 1$$

Equation 3 can be regarded as an NMF that decomposes the matrix of the node-by-attribute value cluster into node-by-community and community-by-attribute value cluster matrices. In other words, Equation 3 indicates the effect of the relationship between nodes and attribute-value clusters against communities.

Unified loss function. To achieve CAR-clustering, the aforementioned three sub-problems must be solved. In this work, we attempt to solve them simultaneously by introducing a unified loss function, which is expressed as

$$L = \arg \min_{U^*, \{U^{(t)}, V^{(t)}, R^{(t)}\}_{t \in \mathbb{T}}} \left\| A - U^* (U^*)^T \right\|_F^2$$

$$+ \sum_{t \in \mathbb{T}} \left\{ \left\| X^{(t)} - U^{(t)} (V^{(t)})^T \right\|_F^2 + \lambda_t \left\| U^{(t)} - U^* R^{(t)} \right\|_F^2 \right\} \quad (4)$$

$$s.t. \forall 1 \leq r \leq k_t, \forall 1 \leq p \leq \ell, \forall t \in \mathbb{T},$$

$$\|U_{*,p}^*\|_1 = 1, \|V_{*,r}^{(t)}\|_1 = 1, \|R_{p,*}^{(t)}\|_1 = 1$$

where λ_t for attribute $t \in \mathbb{T}$ is a user-defined parameter to control the effect of attribute-value clusters for community detection. Higher λ_t yields a stronger effect of the attribute-value clusters in community detection.

4.3 Optimization

Similar to the ordinary NMF, the loss function in Equation 4 is not simultaneously convex for all variables. Hence, we consider the NMF to be a Frobenius norm optimization, where update equations are derived based on [14].

Considering the Karush-Kuhn-Tucker (KKT) first-order conditions applied to our problem, we derive:

$$U^* \geq \mathbf{0}, U^{(t)} \geq \mathbf{0}, V^{(t)} \geq \mathbf{0}, R^{(t)} \geq \mathbf{0} \quad (5)$$

$$\nabla_{U^*} L \geq \mathbf{0}, \nabla_{U^{(t)}} L \geq \mathbf{0}, \nabla_{V^{(t)}} L \geq \mathbf{0}, \nabla_{R^{(t)}} L \geq \mathbf{0} \quad (6)$$

$$\begin{aligned} U^* \odot \nabla_{U^*} L &= \mathbf{0}, U^{(t)} \odot \nabla_{U^{(t)}} L = \mathbf{0}, \\ V^{(t)} \odot \nabla_{V^{(t)}} L &= \mathbf{0}, R^{(t)} \odot \nabla_{R^{(t)}} L = \mathbf{0} \end{aligned} \quad (7)$$

where \odot is the element-wise product. From the Karush-Kuhn-Tucker (KKT) conditions, we derive derivatives corresponding to the variables:

$$\begin{aligned} \nabla_{U^*} L &= -2A^T R U^* + 2U^* (U^*)^T U^* \\ &\quad + \sum_{t \in \mathbb{T}} \lambda_t (-U^{(t)} (R^{(t)})^T + U^* R^{(t)} (R^{(t)})^T) \end{aligned} \quad (8)$$

$$\begin{aligned} \nabla_{U^{(t)}} L &= -X^{(t)} V^{(t)} + U^{(t)} (V^{(t)})^T V^{(t)} \\ &\quad + \lambda_t (U^{(t)} - U^* R^{(t)}) \end{aligned} \quad (9)$$

$$\nabla_{V^{(t)}} L = - (X^{(t)})^T U^{(t)} + (V^{(t)})^T (U^{(t)})^T U^{(t)} \quad (10)$$

$$\nabla_{R^{(t)}} L = - (U^*)^T U^{(t)} + (U^*)^T U^* R^{(t)} \quad (11)$$

By substituting the corresponding gradients in Equation 4, we derive the following update rules:

$$U^* \leftarrow U^* \odot \frac{A^T U^* + \sum_{t \in \mathbb{T}} \lambda_t U^{(t)} (R^{(t)})^T}{2U^* (U^*)^T U^* + \sum_{t \in \mathbb{T}} U^* R^{(t)} (R^{(t)})^T} \quad (12)$$

$$U^{(t)} \leftarrow U^{(t)} \odot \frac{X^{(t)} V^{(t)} + \lambda_t U^* R^{(t)}}{U^{(t)} (V^{(t)})^T V^{(t)} + \lambda_t U^{(t)}} \quad (13)$$

$$V^{(t)} \leftarrow V^{(t)} \odot \frac{(X^{(t)})^T U^{(t)}}{(V^{(t)})^T (U^{(t)})^T U^{(t)}} \quad (14)$$

$$R^{(t)} \leftarrow R^{(t)} \odot \frac{(U^*)^T U^{(t)}}{(U^*)^T U^* R^{(t)}} \quad (15)$$

The aforementioned update rules monotonically decrease Equation 4. However, these variables may violate the probability definition (i.e., their sum does not equal one). To satisfy constraints, $\|U^*_{:,p}\|_1 = 1$, $\|V^*_{:,r}\|_1 = 1$ and $\|R^*_{p,:}\|_1 = 1$, the variables are normalized immediately after updating. The normalization is expressed as

$$U^* \leftarrow U^* (Q^*)^{-1} \quad (16) \quad U^{(t)} \leftarrow U^{(t)} Q^{(t)} \quad (18)$$

$$V^{(t)} \leftarrow V^{(t)} (Q^{(t)})^{-1} \quad (17) \quad R^{(t)} \leftarrow R^{(t)} (Q^R)^{-1} \quad (19)$$

where $Q^* = \text{Diagonalize}(U^*)$, $Q^{(t)} = \text{Diagonalize}(V^{(t)})$, and $Q^R = \text{Diagonalize}(R^{(t)})$.

$$\text{Diagonalize}(Z \in \mathbb{R}^{a \times b}) = \text{Diag} \left(\sum_{i=1}^a Z_{i,1}, \dots, \sum_{i=1}^a Z_{i,b} \right) \quad (20)$$

$\text{Diag}(\cdot)$ provides a diagonal matrix where the diagonals are the input sequence.

Algorithm 1 shows the optimization algorithm based on the aforementioned update rules. Matrix normalization is applied after the updates. Without normalization, each matrix may have significantly different values, leading to inconsistent results. Algorithm 1 describes the order of update rules and normalizations.

Algorithm 1 Optimization Algorithm

Input: $A, \{X^{(t)}\}_{t \in \mathbb{T}}, \{\lambda_t\}_{t \in \mathbb{T}}, \delta$

Output: $U^*, \{U^{(t)}, V^{(t)}, R^{(t)}\}_{t \in \mathbb{T}}$

```

1:  $U^*, \{U^{(t)}, V^{(t)}, R^{(t)}\}_{t \in \mathbb{T}} \leftarrow$  random non-negative init
2:  $\epsilon' \leftarrow \text{maxFloat}, \epsilon \leftarrow \frac{\epsilon'}{2}$ 
3: while  $\text{abs}(\epsilon' - \epsilon) \geq \delta$  do
4:    $U^* \leftarrow U^* \odot \frac{A^T U^* + \sum_{t \in \mathbb{T}} \lambda_t U^{(t)} (R^{(t)})^T}{2U^* (U^*)^T U^* + \sum_{t \in \mathbb{T}} U^* R^{(t)} (R^{(t)})^T}$ 
5:    $U^* \leftarrow U^* (Q^*)^{-1}$ 
6:   for  $t \in \mathbb{T}$  do
7:      $U^{(t)} \leftarrow U^{(t)} \odot \frac{X^{(t)} V^{(t)} + \lambda_t U^* R^{(t)}}{U^{(t)} (V^{(t)})^T V^{(t)} + \lambda_t U^{(t)}}$ 
8:      $V^{(t)} \leftarrow V^{(t)} \odot \frac{(X^{(t)})^T U^{(t)}}{(V^{(t)})^T (U^{(t)})^T U^{(t)}}$ 
9:      $U^{(t)} \leftarrow U^{(t)} Q^{(t)}$ 
10:     $V^{(t)} \leftarrow V^{(t)} (Q^{(t)})^{-1}$ 
11:     $R^{(t)} \leftarrow R^{(t)} \odot \frac{(U^*)^T U^{(t)}}{(U^*)^T U^* R^{(t)}}$ 
12:     $R^{(t)} \leftarrow R^{(t)} (Q^R)^{-1}$ 
13:   end for
14:    $\epsilon' \leftarrow \epsilon$ 
15:    $\epsilon \leftarrow L(U^*, \{U^{(t)}, V^{(t)}, R^{(t)}\}_{t \in \mathbb{T}})$ 
16: end while

```

4.4 Complexity Analysis

Here, we analyze the computational complexity of the proposed algorithm. The equations in our algorithm have the following complexities:

- Updating U^* (Eqs. 12 and 16) needs $O(|\mathbb{E}| \ell + |\mathbb{V}| \ell^2 \sum_t k_t)$.
- Updating $U^{(t)}$ (Eqs. 13 and 18) and $V^{(t)}$ (Eqs. 14 and 18) needs $O((|\mathbb{V}| + |\mathbb{X}_t|) k_t^2 + |\mathbb{E}_t| k_t)$.
- Updating $R^{(t)}$ (Eqs. 15 and 19) needs $O(|\mathbb{V}| (\ell k_t + \ell^2))$.

In summary, the time complexity of our algorithm is follows, where $iter$ is the number of outer iterations (lines 3–16 in our algorithm).

$$O \left(iter \sum_t \left(|\mathbb{V}| (\ell^2 k_t + k_t^2) + |\mathbb{X}_t| k_t^2 + |\mathbb{E}| \ell + |\mathbb{E}_t| k_t \right) \right) \quad (21)$$

5 EXPERIMENTAL EVALUATIONS

To demonstrate the applicability and effectiveness of CARNMF, we conducted a set of experiments using real-world datasets. Specifically, the performance of the proposed scheme was compared to simple baseline and state-of-the-art methods.

The experiments were performed on a PC with an Intel Core i7 (3.3 GHz) CPU with 16 GB RAM running Ubuntu14.04. CARNMF was implemented by Python 2.7.6 with Numpy 1.9.0.

5.1 Datasets

We used two datasets: DBLP and arXiv.

- *DBLP*: Digital Bibliography Project¹ is a bibliographic database in the computer science area. DBLP contains publication information, such as authors and conferences. We used a part of the dataset by extracting conferences similar to [8]. We extracted four research areas: data mining, databases, machine learning, and information retrieval, and five major conferences for each area. Consequently,

¹<http://dblp.uni-trier.de/>

Table 1: Selected conferences on four research areas.

DB	DM	ML	IR
SIGMOD, VLDB PODS, EDBT ICDT	KDD, ICDM PKDD, SDM PAKDD	NIPS, ICML ECML, UAI COLT	SIGIR, ECIR JCDL, ECDL TREC

Table 2: Selected journals on four research areas.

math-ph
Communications in Mathematical Physics Reviews in Mathematical Physics Letters in Mathematical Physics Journal of Mathematical Physics
nucl-th
Annual Review of Nuclear and Particle Science Progress in Particle and Nuclear Physics Atomic Data and Nuclear Data Tables Journal of Nuclear Materials
astro-ph
Astronomical Journal Annual Review of Astronomy and Astrophysics New Astronomy Reviews Space Science Review
cond-mat
Nature Nanotechnology Smart Materials and Structures Semi Conductor Science Journal of Materials Science

10,491 papers in 20 conferences (shown in Table 1) were selected.

- *arXiv*: arXiv² is a repository of electronic preprints in various scientific fields. Similar to above, we chose four research areas: mathematical physics (math-ph), nuclear (nucl-th), astrophysics (astro-ph), and materials (part of cond-mat), and four major journals for each area. Consequently, 12,497 papers in 16 journals (shown in Table 2) were selected.

Multi-attributed graphs were constructed from the datasets as follows: The nodes correspond to the authors. If two authors co-author a paper, we placed a weighted edge between the authors. The weighting denotes the number of co-authored papers. Each author has attributes *term*, *paper*, and *conference/journal*, which are defined below:

- *term*: Each term is regarded as a node. An edge is generated between an author and a term if the author uses the term in the titles of at least one paper. The edge weight denotes the term frequency for each author. As a preprocessing, we applied stop-word elimination and stemming.
- *paper*: Each paper is regarded as a node. An edge is generated if the author publishes the paper. The edge weight is always 1.0 because each paper can only be published once.
- *conference/journal*: Each conference or journal corresponds with a node. An edge is created between an author and a conference/journal if the author publishes at least one paper at the conference/journal. The edge weight is the total number of publications at the conference/journal.

5.2 Results of CAR-clustering

Figure 1 shows examples of the detected communities and their associated attribute-value clusters in DBLP. The number of communities and the number of term clusters were each 50, whereas the number of conference clusters and the number paper clusters were each 4. The red, blue and gray rectangles correspond to communities, term clusters, and conference clusters, respectively.

²<https://arxiv.org>

Each rectangle shows the top contributing nodes in the community/cluster, and the edge weights show the strength of the relationship between the community and the corresponding cluster. We chose famous researchers in different research domains (i.e., *Jiawei Han*, *Michael Stonebraker*, and *Michael I. Jordan*).

Figure 1(a) show the community and the correlated attribute-value clusters of *Jiawei Han*, who is a leading researcher in data mining and database areas. The results show that (1) he collaborates with Chinese researchers, (2) he publishes many papers related to data mining and database conferences (i.e., *KDD*, *ICDM*, *SDM*, *PAKDD*, and *VLDB*), and (3) his researches are highly correlated with topics in data mining, such as *clustering* and *classification on large graph*.

Similarly, Figure 1(b) shows the result for *Michael Stonebraker*, a renowned database researcher. His community is strongly related to conferences in databases (*SIGMOD*, *VLDB*, *PODS*, *EDBT*, and *ICDT*). Topics such as *view management*, *distance metric*, and *query evaluation* are detected. Figure 1(c) shows the result for *Michael I. Jordan*, an expert in machine learning research. This community is strongly related to the conferences of machine learning, (*NIPS*, *ICML*, *UAI*, *COLT*, and *ECML*) and the topics like *learn network*, *expert model*, and *prediction*.

The detected communities and the associated attribute-value clusters seem to be reasonable.

5.3 Accuracy Comparison

The proposed scheme is compared to a baseline method as well as state-of-the-art methods to quantitatively evaluate the performance of community detection and attribute-value clustering. The comparison methods include:

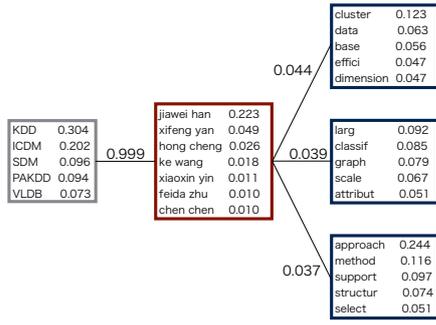
- **NMF** [15]: Baseline approaches that apply NMF for binary relationships between graph components, including author-term (A-T), author-paper (A-P), author-conference (A-C), term-paper (T-P), and term-conference (T-C)³.
- **LCTA** [26]: A probabilistic generative model for communities, topics of textual attributes, and their relationships.
- **SCI** [23]: An NMF based method for detecting communities as well as their semantic descriptions via node’s attribute values.
- **HINMF** [16]: A model that clusters objects and attributes simultaneously and takes the consensus among the binary NMFs. This work is the most similar to our proposal.

Note that, LCTA and SCI deal with a single concatenated feature of multiple attributes. Therefore, we prepare concatenated feature consisting of *term*, *document* and *conference/journal*, and apply these approaches on the feature.

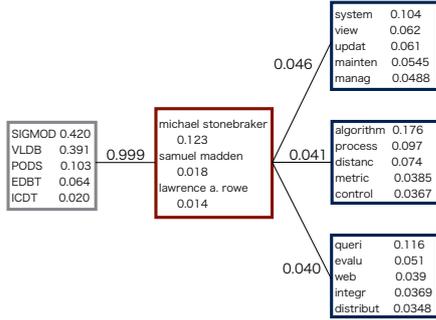
To evaluate the qualities of these methods, we compared the *accuracy* [26] w.r.t. community and attribute-value clustering w.r.t. *paper* and *conference/journal*. We designed a ground truth to measure the *accuracy*. To derive the ground truth, each author is labeled based on research areas of their papers, in other words, if the author mostly published papers for the specific area, the author is labeled as that area. Similarly, the labels for *conference/journal* and *paper* were manually given by referring to the conference categories.

DEFINITION 5 (ACCURACY). Given a set \mathbb{S} of elements, for each element $n \in \mathbb{S}$, the true label and the cluster label generated by a method are denoted by s_n and r_n , respectively. Then, the accuracy

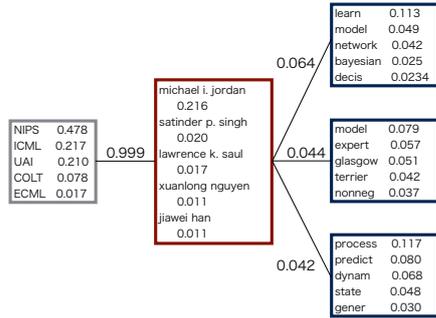
³Because NMF assumes the co-occurrences of binary relationships, paper-conference (one-to-one relationship) is excluded.



(a) Communities of “jiawei han”.



(b) Communities of “michael Stonebraker”.



(c) Communities of “michael i. jordan”.

Figure 1: Example communities with attribute-value clusters. The red, blue and gray rectangles correspond to communities, term clusters, and conference clusters, respectively.

is defined as:

$$Accuracy = \frac{\sum_{n \in \mathbb{S}} \delta(s_n, \text{map}(r_n))}{|\mathbb{S}|}$$

where $|\cdot|$ is the cardinality of a set; $\delta(x, y)$ is a delta function which returns 1 if $x = y$, otherwise 0; and $\text{map}(r_n)$ is a mapping function that maps r_n to the equivalent label in the dataset. The best mapping can be found by Kuhn-Munkres algorithm [13]. \square

Table 4 summarizes the evaluation results. The number of communities and the number of attribute-value clusters for each attribute are each four. Each cell shows the mean and the standard deviation of the accuracies for 20 trials. N/A denotes that the method does not support the category. Values in bold indicate a significant improvement using the Student-t test, where $p < 0.05$.

CARNMF achieved the best performance for community detection (author) and attribute-value clustering (paper and conference/journal) with significant gaps for DBLP dataset (respectively 11%, 22% and 7%) and for arXiv dataset (respectively 3%, 4% and 2%) relative to the comparative methods. In particular, CARNMF has an improved clustering quality compared to NMF by taking the relationships between communities and attribute-value clusters into account.

Table 5 summarizes evaluation of effects from taking multiple attributes into account. The table showcases results where different combinations of attributes are used, e.g., “(T-C)” means term and conference attributes were used. This result shows that the proposed method works the best when taking as many attributes as possible. As expected, the basic tendency is that as the number of attributes increases, the accuracy increases.

Table 3 lists the detected topics from DBLP using CARNMF when the number of topics is set to four. Our method successfully detects the four major research topics. Specifically, Topic 1 containing *retriev, inform, search, queri* and *web* seems to correspond to information retrieval, and Topic 2 containing *mine, pattern, cluster, graph* and *frequent* correspond to data mining. Topic3 contains words “*learn, network, kernel, bayesian, reinforc*”, which are typical words of machine learning. Topic4 is a topic of database containing words “*query, databas, optim, xml, manag*”, which are popular topics on database researches.

Table 3: Detected topics via CARNMF from DBLP.

Topic 1 (Information Retrieval)	Topic 2 (Data Mining)	Topic 3 (Machine Learning)	Topic 4 (Database)
retriev 0.068	mine 0.076	learn 0.063	queri 0.046
trec 0.043	pattern 0.038	model 0.036	data 0.043
inform 0.032	data 0.037	network 0.018	databas 0.043
model 0.028	cluster 0.026	algorithm 0.016	optim 0.017
document 0.027	graph 0.023	infer 0.015	xml 0.015
track 0.024	base 0.021	kernel 0.015	manag 0.015
search 0.021	frequent 0.019	bayesian 0.014	effici 0.014
queri 0.021	databas 0.019	process 0.013	base 0.013
text 0.019	effici 0.017	reinforc 0.012	system 0.012
web 0.017	larg 0.016	decis 0.012	object 0.012

5.4 Insights on Parameters

This section discusses the effect of parameter λ_t for each attribute. The larger the λ_t value, the greater the influence of the attribute-value cluster for $t \in \mathbb{T}$ is on the community. Therefore, optimal parameter setting should result in better results. Figures 2 shows the behavior of the accuracy with different values with respect to different attributes. For each evaluation, λ_s ($s \neq t$) of the other attributes were fixed. In most cases, the accuracy shows a convex form and the peak is around 10^{-2} . More importantly, the accuracy is insensitive to the setting, making tuning easier.

5.5 Convergence Analysis

In this section, we experimentally provide convergence analysis to optimize the proposed loss function in Equation 4. Figures 3(a) and (b) show the convergence curve of the loss function for DBLP and arXiv, respectively. In addition, the accuracy of each iteration is plotted. The black line shows the value of the loss function. The red, green, and blue lines show the accuracy of community detection and attribute-value clustering for author, paper, and conference/journal, respectively. As the number of iterations increases, the loss function decreases while the accuracy improves.

Table 4: Quality evaluations of community detection and attribute clustering.

	DBLP dataset			arXiv dataset		
	Author	Paper	Conference	Author	Paper	Journal
NMF(A-T)	64.02 ± 5.73	N/A	N/A	32.75 ± 2.30	N/A	N/A
NMF(A-P)	43.12 ± 5.17	44.58 ± 5.89	N/A	34.19 ± 3.55	33.53 ± 1.93	N/A
NMF(A-C)	75.35 ± 6.85	N/A	87.60 ± 1.73	69.05 ± 7.30	N/A	67.81 ± 4.09
NMF(T-P)	N/A	50.02 ± 7.93	N/A	N/A	28.22 ± 0.90	N/A
NMF(T-C)	N/A	N/A	69.88 ± 6.68	N/A	N/A	69.06 ± 1.36
LCTA	48.90 ± 7.57	26.13 ± 4.36	68.50 ± 12.46	40.18 ± 4.31	33.88 ± 1.82	61.56 ± 9.94
SCI	54.78 ± 8.79	22.31 ± 1.48	58.20 ± 7.40	38.68 ± 3.56	35.36 ± 0.91	42.81 ± 3.58
HINMF	68.90 ± 9.08	56.46 ± 3.08	90.10 ± 12.63	65.41 ± 5.38	33.24 ± 2.43	61.25 ± 7.81
CARNMF	86.34 ± 2.39	78.19 ± 9.87	97.20 ± 5.21	72.65 ± 8.11	42.23 ± 3.18	71.25 ± 2.53

Table 5: Quality evaluations of community detection and attribute clustering, changing attributes to be used.

	DBLP dataset			arXiv dataset		
	Author	Paper	Conference	Author	Paper	Conference
CARNMF (T)	43.29 ± 3.05	N/A	N/A	43.19 ± 6.44	N/A	N/A
CARNMF (P)	46.84 ± 4.15	41.69 ± 5.97	N/A	33.23 ± 2.63	33.49 ± 1.68	N/A
CARNMF (C)	85.11 ± 2.57	N/A	92.40 ± 8.50	67.41 ± 6.85	N/A	69.69 ± 6.64
CARNMF (T-P)	43.28 ± 4.44	41.33 ± 3.88	N/A	35.32 ± 2.04	35.79 ± 2.28	N/A
CARNMF (T-C)	83.67 ± 6.54	N/A	95.20 ± 1.99	68.93 ± 7.40	N/A	66.88 ± 4.88
CARNMF (P-C)	86.41 ± 1.93	73.30 ± 11.90	94.10 ± 3.42	69.71 ± 8.85	40.15 ± 3.13	68.75 ± 5.59
CARNMF (T-P-C)	86.34 ± 2.39	78.19 ± 9.87	97.20 ± 5.21	72.65 ± 8.11	42.23 ± 3.18	71.25 ± 2.53

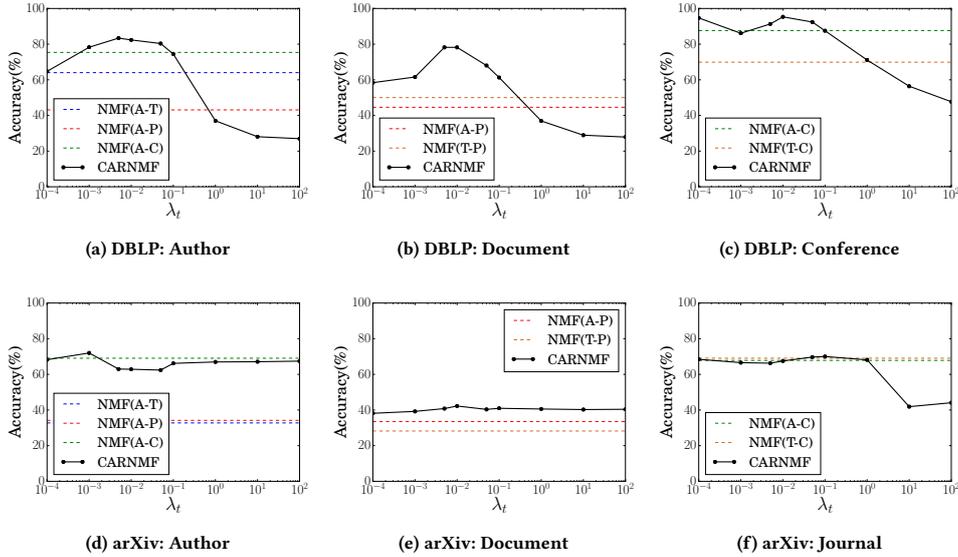


Figure 2: Accuracy for different λ_t values.

5.6 Efficiency Analysis

This section analyzes computational efficiency in terms of the numbers of communities and attribute clusters. When the numbers are fixed to four as experiments above, the running times of CARNMF on the DBLP (arXiv) dataset are 1.186 ± 0.253 s (0.682 ± 0.138 s). When changing the numbers of communities and term clusters to 50, while those of paper and conference remain four, the running times increases to 7.471 ± 0.563 s (DBLP) and 6.526 ± 0.172 s (arXiv). These values are still reasonable for various applications.

Moreover, we examine the running time of our method by changing the number of nodes in an input graph. Theoretically, as discussed in Section 4.4, the computational complexity is dependent on the number of vertices, that of edges, and that of distinct values of each attribute. As most of real-world graphs are modeled as scale-free networks, edges in a graph are very sparse, therefore, we examine the sensitivity of processing time on the proposed method in terms of the number of nodes. In this experiment, we selected all of the papers on DBLP, and construct the multi-attributed graph as same manner as described in Section 5.1. We set the number of communities and clusters are four.

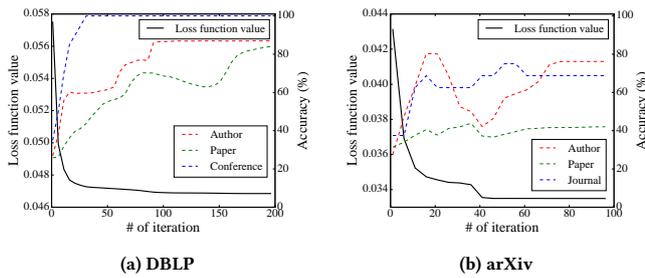


Figure 3: Convergence analysis of the proposed algorithm to optimize a loss function and the corresponding accuracy curve.

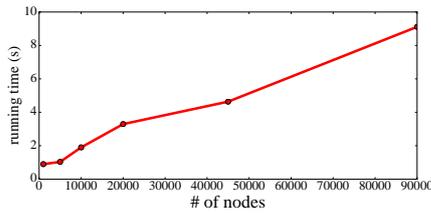


Figure 4: Time complexity of CARNMF w.r.t. the number of input nodes

Figure 4 shows that the time complexity of our method is almost linear to the number of nodes. From the figure, we ensure that the time complexity of our method is linear to the numbers of nodes and edges (as shown on Equation 21). Therefore, when the input graph is sparse, our method is highly efficient.

6 CONCLUSION

In this paper we have proposed CAR-clustering, which includes community detection, attribute-value clustering, and extraction of their relationships, for clustering over multi-attributed graphs. We have also proposed a novel algorithm CARNMF based on NMF. CARNMF employs a unified loss function to simultaneously solve different NMFs. This approach is better than the state-of-the-art methods in that it can exploit the correlation between communities and attribute-value clusters for enhancing the quality of the result. Our experiments have demonstrated that CARNMF successfully achieves CAR-clustering. CARNMF has detected reasonable communities with meaningful semantic descriptions via the relationship between communities and attribute-value clusters for real-world datasets. These results are useful for many applications such as node property estimations [7, 9, 24], community-wise information recommendations [10], and semantic reasoning for nodes/edges [1]. Additionally, CARNMF has achieved higher accuracy than comparative methods, including a baseline and the state-of-the-art methods. Our future work includes several directions. First, we will extend the proposed method for chronological analysis over temporal multi-attributed graphs. Second, we plan to automate the parameter tuning (e.g., the numbers of communities/clusters, λ_t , etc.).

ACKNOWLEDGMENT

This research was partly supported by Japan Agency for Medical Research and Development (AMED).

REFERENCES

- [1] Edoardo M Airolidi, David M Blei, Stephen E Fienberg, and Eric P Xing. 2008. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research* 9, Sep (2008), 1981–2014.
- [2] Martin Atzmueller, Stephan Doerfel, and Folke Mitzlaff. 2016. Description-oriented community detection using exhaustive subgroup discovery. *Information Sciences* 329 (2016), 965–984.
- [3] Michele Berlingerio, Michele Coscia, and Fosca Giannotti. 2011. Finding and characterizing communities in multidimensional networks. In *Advances in Social Networks Analysis and Mining (ASONAM), 2011 International Conference on*. IEEE, 490–494.
- [4] Brigitte Boden, Stephan Günnemann, Holger Hoffmann, and Thomas Seidl. 2012. Mining coherent subgraphs in multi-layer graphs with edge labels. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 1258–1266.
- [5] Cecile Bothorel, Juan David Cruz, Matteo Magnani, and Barbora Mícenková. 2015. Clustering attributed graphs: models, measures and methods. *Network Science* 3, 3 (2015), 408–444.
- [6] Santo Fortunato. 2010. Community detection in graphs. *Physics reports* 486, 3 (2010), 75–174.
- [7] Mario Frank, Andreas P Streich, David Basin, and Joachim M Buhmann. 2012. Multi-assignment clustering for boolean data. *Journal of Machine Learning Research* 13, Feb (2012), 459–489.
- [8] Jing Gao, Wei Fan, Yizhou Sun, and Jiawei Han. 2009. Heterogeneous source consensus learning via decision propagation and negotiation. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 339–348.
- [9] Michelle Girvan and Mark EJ Newman. 2002. Community structure in social and biological networks. *Proceedings of the national academy of sciences* 99, 12 (2002), 7821–7826.
- [10] Junzo Kamahara, Tomofumi Asakawa, Shinji Shimojo, and Hideo Miyahara. 2005. A community-based recommendation system to reveal unexpected interests. In *Multimedia Modelling Conference, 2005. MMM 2005. Proceedings of the 11th International*. IEEE, 433–438.
- [11] Denise B Kandel. 1978. Homophily, selection, and socialization in adolescent friendships. *American journal of Sociology* 84, 2 (1978), 427–436.
- [12] Da Kuang, Chris Ding, and Haesun Park. 2012. Symmetric nonnegative matrix factorization for graph clustering. In *Proceedings of the 2012 SIAM International Conference on Data Mining*. SIAM, 106–117.
- [13] Harold W Kuhn. 1955. The Hungarian method for the assignment problem. *Naval research logistics quarterly* 2, 1-2 (1955), 83–97.
- [14] Daniel D Lee and H Sebastian Seung. 1999. Learning the parts of objects by non-negative matrix factorization. *Nature* 401, 6755, 788–791.
- [15] Daniel D Lee and H Sebastian Seung. 2001. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*. 556–562.
- [16] Jialu Liu and Jiawei Han. 2013. HINMF: A Matrix Factorization Method for Clustering in Heterogeneous Information Networks. In *Proceedings of the international joint conference on artificial intelligence workshop*.
- [17] Peter V Marsden. 1988. Homogeneity in confiding relations. *Social networks* 10, 1 (1988), 57–76.
- [18] Ioannis Psorakis, Stephen Roberts, Mark Ebdon, and Ben Sheldon. 2011. Overlapping community detection using bayesian non-negative matrix factorization. *Physical Review E* 83, 6 (2011), 066114.
- [19] Michal Rosen-Zvi, Thomas Griffiths, Mark Steyvers, and Padhraic Smyth. 2004. The author-topic model for authors and documents. In *Proceedings of the 20th conference on Uncertainty in artificial intelligence*. AUAI Press, 487–494.
- [20] Jianbo Shi and Jitendra Malik. 2000. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence* 22, 8 (2000), 888–905.
- [21] Yizhou Sun, Yintao Yu, and Jiawei Han. 2009. Ranking-based clustering of heterogeneous information networks with star network schema. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 797–806.
- [22] Lei Tang and Huan Liu. 2010. Community detection and mining in social media. *Synthesis Lectures on Data Mining and Knowledge Discovery* 2, 1 (2010), 1–137.
- [23] Xiao Wang, Di Jin, Xiaochun Cao, Liang Yang, and Weixiong Zhang. 2016. Semantic community identification in large attribute networks. In *AAAI*. 265–271.
- [24] Jaewon Yang and Jure Leskovec. 2013. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 587–596.
- [25] Jaewon Yang, Julian McAuley, and Jure Leskovec. 2013. Community detection in networks with node attributes. In *2013 IEEE 13th International Conference on Data Mining*. IEEE, 1151–1156.
- [26] Zhijun Yin, Liangliang Cao, Quanquan Gu, and Jiawei Han. 2012. Latent community topic analysis: integration of community discovery with topic modeling. *ACM Transactions on Intelligent Systems and Technology (TIST)* 3, 4 (2012), 63.
- [27] Haizheng Zhang, Baojun Qiu, C Lee Giles, Henry C Foley, and John Yen. 2007. An LDA-based Community Structure Discovery Approach for Large-Scale Social Networks. *ISI* 200 (2007).