Community Detection in Partially Observable Social Networks

Cong Tran, Student Member, IEEE, Won-Yong Shin, Senior Member, IEEE, and Andreas Spitz

Abstract—The discovery of community structures in social networks has gained significant attention since it is a fundamental problem in understanding the networks' topology and functions. However, most social network data are collected from partially observable networks with both missing nodes and edges. In this paper, we address a new problem of detecting overlapping community structures in the context of such an incomplete network, where communities in the network are allowed to overlap since nodes belong to multiple communities at once. To solve this problem, we introduce KroMFac, a new framework that conducts community detection via regularized nonnegative matrix factorization (NMF) based on the Kronecker graph model. Specifically, from an inferred Kronecker generative parameter matrix, we first estimate the missing part of the network. As our major contribution to the proposed framework, to improve community detection accuracy, we then characterize and select influential nodes (which tend to have high degrees) by ranking, and add them to the existing graph. Finally, we uncover the community structures by solving the regularized NMF-aided optimization problem in terms of maximizing the likelihood of the underlying graph. Furthermore, adopting normalized mutual information (NMI), we empirically show superiority of our KroMFac approach over two baseline schemes by using both synthetic and real-world networks.

Index Terms—Community detection, influential node, Kronecker graph model, matrix factorization, overlapping community, partially observable social network

1 Introduction

1.1 Backgrounds

Real-world networks extracted from various biological, social, technological, and information systems usually contain inhomogeneities that reveal a high level of hierarchical and structural properties. Research on community detection, which is one of the most important tasks in network analysis, has thus become crucial in understanding the fundamental features (e.g., topology and functions) of these networks [1]. In general terms, communities can be regarded as the sets of points that are "close" to each other with respect to a predefined measure of distance or similarity. Since applications of community detection are diverse, there exist a variety of graph-theoretic approaches [2] that conduct optimization based on measures such as modularity [3] and conductance [4], whose performance depends heavily on network topology.

On the one hand, community detection algorithms for online social networks should be designed by taking into account their inherently overlapping and imprecise nature, since community memberships in social networks are allowed to overlap as nodes belong to multiple clusters at once [5]. The extraction of such overlapping communities is known to be more challenging than non-overlapping community detection due to higher complexity and higher computational demands.

In practice, on the other hand, most social network data are collected from partially observable networks with both missing nodes and edges [6], which further complicates the detection of communities. For example, due to limited resources, a person or an organization may be allowed to obtain only a subset of data within a specific geographic query region. This is further compounded due to privacy settings specified by the users that may partially or entirely hide some of their traces or friendships [7]. For example, 52.6% of Facebook users in New York City hid their friend lists in June 2011 [8]. Such types of incomplete networks constitute a severe obstacle for topology-based optimization methods in detecting the true community structures. Surprisingly, while some research exists into the recovery of edges and nodes in such incomplete networks, the problem of community detection under such conditions and its solutions have not yet been investigated.

C. Tran is with the Department of Computer Science and Engineering, Dankook University, Yongin 16890, Republic of Korea, and the Department of Computational Science and Engineering, Yonsei University, Seoul 03722, Republic of Korea.
 E-mail: congtran@ieee.org.

W.-Y. Shin is with the Department of Computational Science and Engineering, Yonsei University, Seoul 03722, Republic of Korea.
 E-mail: wy.shin@yonsei.ac.kr.

A. Spitz is with the School of Computer and Communication Sciences, cole Polytechnique Fdrale de Lausanne, Lausanne 1015, Switzerland. E-mail: andreas.spitz@epfl.ch.

1.2 Motivation and Main Contributions

In this work, we formulate a new problem of detecting overlapping community structures in the context of such an incomplete network in which some of the nodes and edges are missing. To solve the problem, we present KroMFac, a new framework that intelligently combines network completion and community detection methods into one unified framework, which is the first attempt in the literature. To this end, KroMFac first estimates the missing part of the network using a Kronecker generative parameter matrix acquired under the Kronecker graph model [9], which basically differs from the well-known link prediction task in machine learning since node labels would never be acquired by link prediction. Our important contribution to the proposed framework is based on the insight that including the entirety of recovered nodes and edges in the existing graph may be detrimental to enhancement of community detection accuracy. This is because adding more recovered nodes and edges would cause the inference errors to accumulate. To address this problem, we characterize and select influential nodes by centrality ranking, which tend to have high degrees, in the effort of limiting the accumulated errors in our model. Finally, we perform community detection using a state-ofthe-art algorithm along with the recovered graph in which influential nodes are added. In this study, we focus on nonnegative matrix factorization (NMF)-based community detection since it is beneficial in accelerating the computation speed while allowing the NMF-based method to be scalable in large-scale networks. Thus, we formulate a regularized NMF-aided optimization problem in terms of maximizing the likelihood of the underlying graph to discover the community structures. After solving the problem, we assign nodes to communities depending on the values of each entry in the factorized affiliation matrix.

Adopting normalized mutual information (NMI) as a popular information-theoretic performance metric, we empirically verify the superior performance of our proposed approach over two baselines that 1) do not infer missing nodes and edges (Baseline 1) and 2) leverage completion of the entire network (Baseline 2). In summary, our main contributions are five-fold and summarized as follows:

- design of a new framework, named KroMFac, that intelligently combines network completion and community detection in our incomplete network;
- formulation of a regularized NMF-aided joint optimization problem;
- characterization and selection of influential nodes via ranking, which play a vital role in improving community detection accuracy;
- validation of our KroMFac approach through inten-

TABLE 1: Summary of notations

Notation	Description
G	observable graph
V	set of observable nodes
E	set of observable edges
V_M	set of missing nodes
E_M	set of missing edges
N	number of observable nodes
M	number of missing nodes
G'	true graph
\mathbf{F}	affiliation matrix
C	number of communities
\mathbf{A}	adjacency matrix of the observable graph G
\mathbf{A}'	adjacency matrix of the true graph G'
i	number of recovered nodes
$R^{(i)}$	recovered graph after connecting i nodes
$\mathbf{A}_R^{(i)}$	adjacency matrix of the recovered graph $R^{(i)}$
$\mathbf{Z_1}$	matrix containing links between recovered nodes and existing nodes
$\mathbf{Z_2}$	matrix containing links between between recovered nodes
H	number of influential nodes
λ	regularization parameter
Θ	Kronecker parameter matrix
Θ_{init}	initialized Kronecker parameter matrix
K	number of Kronecker products
Cen	degree centrality
\mathcal{D}	loss function
ϵ	threshold for determining influential nodes
ψ	set of communities
δ	threshold determining communities
r	ranking vector

sive experiments based on parameter search using both synthetic and real-world datasets;

analysis and empirical validation of the computational complexity.

Our framework takes an important first step towards establishing a new line of research and towards a better understanding of jointly conducting both network recovery and community detection in partially observable networks.

1.3 Organization

The remainder of this paper is organized as follows. In Section 2, we summarize significant studies that are related to our work. In Section 3, we explain the methodology of our work, including the problem definition and the overview of our KroMFac framework. Section 4 describes implementation details of the proposed KroMFac framework. Experimental results are provided in Section 5 with comparison to two baseline approaches. Finally, we summarize the paper with some concluding remarks in Section 6.

1.4 Notations

Throughout this paper, \mathbb{R} and $\mathbb{P}(\cdot)$ indicate the field of real numbers and the probability, respectively. Unless otherwise stated, all logarithms are assumed to be to the base e. Table 1 summarizes the notations used in this paper. These notations will be formally defined in the following sections when we introduce our network model and technical details.

2 RELATED WORK

The framework that we propose in this paper is related to four broader areas of research, namely community detection in graphs, detection of overlapping communities, community detection in incomplete networks with missing edges, and network completion in social networks.

Community detection in graphs. Since research into community detection in complex networks constitutes a very active field, there are many efforts devoted to community detection in graphs. The most popular techniques include modularity optimization [3], stochastic block models [10], spectral graph-partitioning [11], clique percolation [12], clustering [13], and label propagation [14]. However, these techniques focus on graphs in which nodes can be partitioned into communities and do not address the inherent overlapping nature of community structures in many real-world networks.

Detection of overlapping communities. To cope with this contrast, a recently emerging topic covers the detection of overlapping communities by investigating the structural properties of such communities, especially in the case of social networks [15], [16]. As the computational complexity increases drastically for the recovery of overlapping communities instead of partitioned communities, the research has focused either on detecting communities based on local expansion [17] and link embedding [18] or on employing scalable techniques such as NMF [19], [20], [21] and label propagation [22], [23].

Community detection in incomplete networks. Recently, research on community detection in incomplete networks with *missing edges* has attracted wide attention due to a lack of information caused by users' privacy settings and limited resources. Most of these studies predict the missing links between nodes based on the incorporated additional information [24] or the similarity of topological structures [25], [26], and then discover communities in the underlying recovered networks. In [27], a hierarchical gamma process infinite edge partition model was presented to detect communities and recover missing edges in parallel. In contrast to edge recovery, approaches for node recovery are largely still missing.

Network completion in social networks. In addition to the studies on community detection, network completion thus plays an important role in our research since it should precede the community detection process. As the most influential study, KronEM, an approach based on Kronecker graphs to solving the problem of discovering both missing nodes and edges was suggested by Kim and Leskovec [9], where the expectation maximization (EM) algorithm is applied. For cases in which only a small number of edges are missing, vertex similarity [28] was shown to be useful in recovering the original networks. Another graph recovery approach was also developed in [29] to solve the missing node identification problem when the information of con-

nections between missing nodes and observable nodes is available.

Despite these contributions, there has been no prior work in the literature that combines the contexts of community detection in incomplete social networks with the recovery of both missing nodes and edges. In the following, we therefore present such an approach that seamlessly integrates the recovery of missing parts of a network with subsequent community detection on the recovered network, while benefiting from a resulting more complete community structure.

3 METHODOLOGY

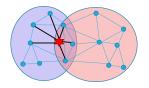
As basis for the algorithms in Section 4, we discuss network fundamentals, formalize the problem definition, and then introduce the generative graph model for network completion, our node selection strategy, and the community detection method.

3.1 Problem Definition

3.1.1 Network Model and Basic Assumptions

Let us denote the partially observable network as G =(V, E), where V and E are the set of vertices and the set of edges, respectively. The network G with N = |V| nodes can be interpreted as a subgraph taken from an underlying true social network $G' = (V \cup V_M, E \cup E_M)$, where V_M is the set of unobservable nodes and E_M is the set of unobservable edges. If we assume G' to be a scale-free network, then the degree distribution of G' can be approximated as $\mathbb{P}(k) \sim k^{-\gamma}$, where the probability $\mathbb{P}(k)$ of a node in the network is inversely proportional to its degree k raised to the power of an exponent parameter γ . While not all realworld social networks necessarily follow a power-law distribution [31], fitting a power-law model to the long tail of the distribution is usually sufficient for practical applications. Therefore, other types of networks following a heavy-tailed degree distribution can also serve as suitable input for our work. In real-world social networks, the nodes and edges of the network G' correspond to users and their relationships, respectively, with little additional information being available. In the following, we thus consider both G and G' to be undirected unweighted networks. Furthermore, we assume that the number of missing nodes $M = |V_M|$ is either known or can be approximated by standard methods for estimating the size of hidden or missing populations [32]. To detect overlapping communities, we assume that social networks follow the affiliation graph model (AGM) [15], which

^{1.} The degree distribution can be estimated via least squares approximation just by taking at most 1% of the samples using a sublinear approach as indicated in [30].





- (a) Before node deletion
- (b) After node deletion

Fig. 1: An example that illustrates the difference between the community structures before and after an influential node (star) and its incident edges (black) have been deleted from the graph. Potential communities are depicted with different colors.

states that the more communities a pair of nodes shares, the higher the probability that these two nodes are connected. The number of communities in the network is denoted by C. The AGM can be represented by a non-negative weight affiliation matrix $\mathbf{F} \in \mathbb{R}^{(N+M)\times C}$ such that each element \mathbf{F}_{uc} represents the degree of membership of a node $u \in (V \cup V_M)$ to the community c. The probability $\mathbb{P}(u,v)$ of a connection between two nodes u and v then depends on the value of \mathbf{F} and is given by $\mathbb{P}(u,v) = 1 - \exp(-\mathbf{F}_u\mathbf{F}_v^\top)$, where $\mathbf{F}_u \in \mathbb{R}^C$ and $\mathbf{F}_v \in \mathbb{R}^C$ are the row vectors that correspond to nodes u and v, respectively [19].

3.1.2 Problem Formulation

As illustrated in Fig. 1, the network structures of partially observable networks are potentially distorted significantly due to the effect of both missing nodes and edges. As a result, methods established for detecting communities may appear to perform well on the partially observable networks but are not effective in extracting the true community structures of the underlying true network. The recovery of overlapping communities in such incomplete, partially observable networks has not been investigated before in the literature. To address this task, we thus present KroMFac, a novel framework for recovering a partially observable network and then discovering the overlapping community structures of the recovered underlying graph. To this end, we first recover missing nodes and edges, which is equivalent to filling in the missing part of the binary adjacency matrix $\mathbf{A}' \in \{0,1\}^{(N+M)\times(N+M)}$ of the graph G' based on the topological information of the observable matrix A (refer to Section 3.2). This inference of missing parts of the network is not without risk since adding more recovered nodes and edges may also accumulate more errors. However, many such nodes and edges may not be very relevant to the subsequent community detection. This motivates us to propose a node selection strategy that aims to characterize and include only a small number of nodes that have a high impact on the community detection. Specifically, we need to selectively recover nodes. We start by formally defining the adjacency matrix that is acquired after the iterative addition of nodes (and their edges) according to an importance ranking strategy.

Definition 1. Let $R^{(i)}$ be the selectively recovered graph formed by connecting $i \in \{0,1,\cdots,M\}$ nodes to the existing graph G according to a predefined selection order. Based on the fact that G and $R^{(i)}$ correspond to $\mathbf{A} \in \{0,1\}^{N \times N}$ and $\mathbf{A}_R^{(i)} \in \{0,1\}^{(N+i) \times (N+i)}$, respectively, $\mathbf{A}_R^{(i)}$ can be written as the following partitioned block matrix:

 $\mathbf{A}_R^{(i)} = egin{bmatrix} \mathbf{A} & \mathbf{Z}_1 \ \mathbf{Z}_1^ op & \mathbf{Z}_2 \end{bmatrix},$

where the matrix $\mathbf{Z}_1 \in \{0,1\}^{N \times i}$ contains the links between recovered nodes and existing nodes and the matrix $\mathbf{Z}_2 \in \{0,1\}^{i \times i}$ contains the links between recovered nodes.

By definition, if we select the top i nodes, then we obtain a unique matrix $\mathbf{A}_R^{(i)}$. As special cases, it follows that $\mathbf{A}_R^{(0)} = \mathbf{A}$ and $\mathbf{A}_R^{(M)} = \mathbf{A}'$. To limit the accumulated errors to a certain level in our model, we only take into account top $H \in \{0,1,\cdots,M\}$ nodes in the ranked list, termed influential nodes (refer to Definition 2 in Section 3.3).

The next step is the detection of communities, which is equivalent to estimating the affiliation matrix $\mathbf{F} \in \mathbb{R}^{(N+i) \times C}$. For each i (i.e., the number of newly added nodes), estimation of the affiliation matrix \mathbf{F} leads to a probabilistic approximation of the matrix $\mathbf{A}_R^{(i)}$ (refer to Section 3.4). When $\hat{\mathbf{F}}$ has the highest chance to generate the graph $R^{(\hat{i})}$, we formulate a joint optimization problem as follows:

$$(\hat{\mathbf{F}}, \hat{i}) = \underset{\mathbf{F} \ge 0, i \in \{0, 1, \dots, H\}}{\arg \max} \log \mathbb{P}(\mathbf{A}_R^{(i)} | \mathbf{F}) + \lambda \log(i+1), \quad (1)$$

which corresponds to a maximum log-likelihood problem with *regularization* for a given value of i.² Here, $\log(i+1)$ indicates a regularization term, which is used to compensate the log-likelihood $\mathbb{P}(\mathbf{A}_R^{(i)}|\mathbf{F})$ reduced by increasing the size of the matrix $\mathbf{A}_R^{(i)}$ since the probability $\mathbb{P}(\mathbf{A}_R^{(i)}|\mathbf{F})$ tends to decrease with the number of elements of $\mathbf{A}_R^{(i)}$; and the parameter $\lambda>0$ determines the impact of the regularization term and needs to be properly set according to the size of observable graphs (which will be specified in Section 5.4). The role of i in (1) is especially important as it is the key factor in handling the total error of the recovered graph. The overall procedure of our approach is visualized in Fig. 2.

3.2 Generative Graph Model

For recovery of the true network structures, the missing part of the network can be inferred by investigating the connec-

2. Here, the superscript $\hat{\ }$ is used to indicate the optimal argument.

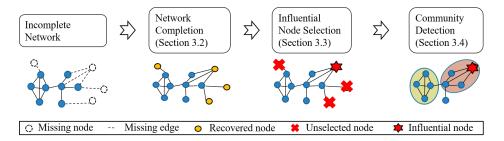


Fig. 2: The schematic overview of our **KroMFac** framework.

tivity patterns in the observable part. To this end, *generative* models for graphs have been developed. The two major generative graph models with this aim include the stochastic block model [10] and the Kronecker graph model [33]. For our research, we adopt the Kronecker graph model since it is scalable and can be used to efficiently model a probability distribution over the missing part of social networks [33]. Thus, we briefly describe the Kronecker graph model before proceeding to network completion.

The model is based on the Kronecker product of two graphs [34]. For two given adjacency matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{m' \times n'}$, the Kronecker product $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{mm' \times nn'}$ is defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \dots & a_{mn}\mathbf{B} \end{bmatrix},$$

where a_{uv} denotes the (u,v)th element of the matrix $\mathbf A$ for $u \in \{1,\cdots,m\}$ and $v \in \{1,\cdots,n\}$. The Kronecker graph model is then defined by a Kronecker generative parameter matrix $\mathbf \Theta \in [0,1]^{N_0 \times N_0}$, where $N_0 \in \mathbb N$. By Kronecker-powering the parameter $\mathbf \Theta$, we obtain increasingly larger and larger stochastic graph adjacency matrices. Since every entry of the matrix $\mathbf \Theta$ can be interpreted as a probability, the Kronecker graph model is then equivalent to a probability distribution of edges over networks.

For network completion, we use the KronEM algorithm [9], which is built upon the Kronecker generative graph model and is the current state-of-the-art algorithm in the literature. Based on an observable network G, KronEM estimates the parameter matrix Θ used to generate the full network Θ^K representing the Kth Kronecker power of Θ , where K is a positive integer such that $N_0^{K-1} < N + M \le N_0^K$. Let $(\mathbf{A}_R^{(M)}, \sigma)$ denote a permutation matrix, where σ indicates a permutation of the set $\{1, \cdots, N+M\}$ and $\sigma(u)$ is the index of node u in the graph $R^{(M)}$ after permutation. The first N elements of σ map the nodes in G while the remaining M elements map the nodes in the missing part.

Then, the likelihood $\mathbb{P}(\mathbf{A}_R^{(M)}, \sigma | \mathbf{\Theta})$ can be expressed as

$$\mathbb{P}(\mathbf{A}_R^{(M)}, \sigma | \mathbf{\Theta}) = \prod_{a_{uv} = 1} [\mathbf{\Theta}^K]_{\sigma(u)\sigma(v)} \prod_{a_{uv} = 0} (1 - [\mathbf{\Theta}^K]_{\sigma(u)\sigma(v)}),$$

where a_{uv} denotes the (u,v)th element of the matrix $\mathbf{A}_R^{(M)}$, and $[\boldsymbol{\Theta}^K]_{\sigma(u)\sigma(v)}$ denotes the $(\sigma(u),\sigma(v))$ th element of the matrix $\boldsymbol{\Theta}^K$. As the matrix $\boldsymbol{\Theta}^K$ is a probabilistic representation of $\mathbf{A}_R^{(M)}$, we also obtain the missing parts \mathbf{Z}_1 and \mathbf{Z}_2 by assigning the value of every entry in \mathbf{Z}_1 and \mathbf{Z}_2 to be zero or one according to a series of *Bernoulli coin-tosses* with the mapped entries in $\boldsymbol{\Theta}^K$ as the probabilities. The detailed steps will be discussed in Section 4.

3.3 Influential Node Selection by Ranking

Network completion can be seen as a statistical learning process, as we predict the value of entries in the missing part of the network by leveraging information in the observed part. Thus, after obtaining the missing part via inference, we note that using the whole recovered nodes may lead to an inaccurate detection of communities due to two types of errors. One type stems from the prediction model, while the other stems from random errors that occur during the Bernoulli series used to project a probabilistic value to zero or one. While the prediction error can be reducible, the random error is irreducible. For this reason, the more recovered nodes we include, the higher the sum of errors. On the other hand, using just a very small portion of missing nodes is unlikely to provide correct community structures, since there is insufficient information available to the community detection model.

Since our eventual goal is to recover the true community structures, it is intuitive to add only nodes that are useful in the community detection process. To assess the usefulness of nodes to a social network, we rely on the concept of centrality and adopt the *degree centrality*, which measures the number of immediate neighbors of a node. Given an undirected graph, the degree centrality of node u, denoted by Cen(u), is defined as the number of connections of a node (i.e., the number of incident edges of a node), and is

^{3.} The parameter N_0 is typically set to two to model the structure of social networks [33], but it can also be set to any integer so that there is no limit in the network size.

computed as

$$\operatorname{Cen}(u) = \sum_{v=1}^{N+M} a_{uv}.$$
 (2)

To select a subset of important nodes to recover, we rank the inferred nodes by first calculating their centrality measures and then sorting them in order of descending centrality. In the following, we formally define the concept of *influential nodes*.

Definition 2. Let Cen(u) denote a centrality measure of node $u \in (V \cup V_M)$. Then, u is defined as an influential node if $u \notin V$ and $Cen(u) \geq \epsilon$, where $\epsilon > 0$ is a predefined threshold, V is the set of observable nodes, and V_M is the set of missing nodes. Here, $H \in \{0, 1, \cdots, M\}$ denotes the cardinality of the set of influential nodes.

The threshold ϵ signifies when the amount of acquired information outweighs the incurred errors from recovering parts of the true network, which means that recovering more than H nodes can be harmful to the community detection process.

3.4 Community Detection via Regularized NMF

Matrix factorization-based approaches are commonly used tools in the detection of overlapping communities in social networks [19], [20], [21]. The benefit of these approaches lies in their scalability since many efficient techniques for solving NMF problems have been developed [35]. Additionally, the NMF-based approaches aim at detecting community structures in the whole given network in a deterministic manner, which often requires less effort to find the optimize parameter setting than those based on the label propagation. In this subsection, we describe how community detection can be transformed into a regularized NMF-aided optimization problems and elaborate on detailed steps for solving the combined problem of graph inference and community detection.

From the previous steps, recall that we have an observation matrix \mathbf{A} , a list of influential nodes, and their corresponding ranking by the degree centrality measure. Furthermore, all recovered matrices $\mathbf{A}_R^{(i)}$ for $i \in \{0,1,\cdots,H\}$ are available. Due to the fact that $\mathbb{P}(u,v)=1-\exp(-\mathbf{F}_u\mathbf{F}_v^\top)$ according to the AGM, the likelihood $\mathbb{P}(\mathbf{A}_R^{(i)}|\mathbf{F})$ in (1) can be rewritten as

$$\mathbb{P}(\mathbf{A}_R^{(i)}|\mathbf{F}) = \prod_{a_{uv}^{(i)}=1} (1 - \exp(-\mathbf{F}_u\mathbf{F}_v^\top)) \prod_{a_{uv}^{(i)}=0} (\exp(-\mathbf{F}_u\mathbf{F}_v^\top)),$$

where $a_{uv}^{(i)}$ denotes the (u,v)th element of the matrix $\mathbf{A}_R^{(i)}$. Thus, we have

$$\log(\mathbb{P}(\mathbf{A}_R^{(i)}|\mathbf{F})) = \sum_{a_{uv}^{(i)}=1} \log(1 - \exp(-\mathbf{F}_u \mathbf{F}_v^\top)) - \sum_{a_{uv}^{(i)}=0} (\mathbf{F}_u \mathbf{F}_v^\top).$$

Suppose that $f(\mathbf{X}) = 1 - \exp(-(\mathbf{X}))$ for a matrix \mathbf{X} . Then, we obtain a matrix $f(\mathbf{F}\mathbf{F}^\top)$ that probabilistically approximates the adjacency matrix $\mathbf{A}_R^{(i)}$. To estimate the difference between two matrices $\mathbf{A}_R^{(i)}$ and $f(\mathbf{F}\mathbf{F}^\top)$, instead of using the Euclidean distance metric, we utilize the negative log-likelihood in (3) as a loss function \mathcal{D} , which indicates that

$$\mathcal{D}(\mathbf{A}_{R}^{(i)}, f(\mathbf{F}\mathbf{F}^{\top})) = -\log(\mathbb{P}(\mathbf{A}_{R}^{(i)}|\mathbf{F})). \tag{4}$$

As a result, the optimization problem in (1) can then be cast into a regularized NMF formulation as

$$(\hat{\mathbf{F}}, \hat{i}) = \operatorname*{arg\,min}_{\mathbf{F} \geq 0, i \in \{0, 1, \cdots, H\}} \mathcal{D}(\mathbf{A}_R^{(i)}, f(\mathbf{F}\mathbf{F}^\top)) - \lambda \log(i+1), \tag{5}$$

where the objective function in (5) is referred to as the *regularized loss* in our setup.

4 PROPOSED KROMFAC FRAMEWORK

In this section, to provide a complete solution to the problem of community detection in a partially observable graph, we present KroMFac, a novel framework that consists of the following three major phases: 1) network completion, 2) node ranking and selection, and 3) community detection. The overall procedure is described in Algorithm 1. The observable graph G, the number of missing nodes, M, and the number of communities, C, are the key input parameters of the algorithm. The dimension of the parameter matrix Θ is given by $N_0 \times N_0$, and we initialize Θ as a randomly generated matrix $\Theta_{\text{init}} \in [0,1]^{N_0 \times N_0}$. Further parameters serve as control parameters. In particular, ϵ plays a central role in determining the set of influential nodes and are introduced in detail in Section 3.3, and λ controls the impact of regularization, which can be quantified via an empirical study. The parameter $\delta > 0$ serves as a threshold to decide to which communities each node belongs (i.e., the degree of membership of nodes), and can be estimated for a given network [19]. Finally, the parameter η_{detect} is an arbitrarily small positive constant used as stopping criterion during community detection (i.e., convergence criteria). As the output of Algorithm 1, we define ψ as the set of detected communities.

We assume that all communities initially have no members. To find the true community structures of the incomplete input graph, we first fully recover the graph via the function GraphRecv (refer to Section 4.1). By analyzing the recovered graph via the function NodeSelect, we then select H influential nodes and determine the ranking vector $r \in \mathbb{N}^H$ that represents the indices of ranked influential nodes and plays a crucial role in community detection accuracy (refer to Section 4.2). In this step, specifically, we solve the joint optimization problem described in (5) through exhaustive search over i by sequentially connecting influential

Algorithm 1: KroMFac

```
Input: G, M, C, N_0, \Theta_{\text{init}}, \epsilon, \delta, \eta_{\text{detect}}, \lambda
      Output: \psi
 1 Initialization: \mathcal{D}_{\min} \leftarrow \infty; \psi[c] \leftarrow \{\emptyset\} for
         c \in \{1, \cdots, C\}
 2 function KroMFac
                \mathbf{A} \leftarrow \text{Adjacency matrix of } G
               \begin{aligned} \mathbf{A}_R^{(M)} &\leftarrow \mathsf{GraphRecv}(\mathbf{A}, N_0, \mathbf{\Theta}_{\mathsf{init}}, M) \\ (H, r) &\leftarrow \mathsf{NodeSelect}(\mathbf{A}_R^{(M)}, M, \epsilon) \end{aligned}
 4
 5
                for i from 1 to H do
 6
                        R^{(i)} \leftarrow \text{Connect nodes } \{r[1], \cdots, r[i]\} \text{ to } G
 7
                        \mathbf{A}_R^{(i)} \leftarrow \text{Adjacency matrix of } R^{(i)}
 8
                        (\mathcal{D}, \mathbf{F}) \leftarrow \mathsf{CommunDet}(\mathbf{A}_{R}^{(i)}, C, \eta_{\mathsf{detect}})
 9
                        \mathcal{D} \leftarrow \mathcal{D} - \lambda \log(i+1)
10
                        if \mathcal{D}_{\min} < \mathcal{D} then
11
                                 \hat{\mathbf{F}} \leftarrow \mathbf{F}
12
                                 \begin{aligned} \hat{i} &\leftarrow i \\ \mathcal{D}_{\min} &\leftarrow \mathcal{D} \end{aligned} 
13
14
               for u from 1 to N+M do
15
                        for c from 1 to C do
16
                                \begin{array}{l} \text{if } \hat{\mathbf{F}}_{uc} \geq \delta \text{ then} \\ \quad \big \lfloor \ \psi[c] \leftarrow \psi[c] \cup \{u\} \end{array}
17
18
               return \psi
19
```

nodes to the existing observable graph based on the order in r. For each i, we acquire a corresponded graph $R^{(i)}$ and its adjacency matrix $\mathbf{A}_{R}^{(i)}$ (see Definition 1). By using the function CommunDet implemented via the state-of-theart NMF-based community detection method, we are then capable of obtaining the loss function \mathcal{D} associated with affiliation matrix **F** given the input matrix $\mathbf{A}_{R}^{(i)}$. As a result, we obtain \mathcal{D}_{\min} as the smallest value of $\mathcal{D} - \lambda \log(i+1)$ (i.e., the regularized loss), which in turn provides us with the optimal $\hat{\mathbf{F}}$ and \hat{i} . Every entry $\hat{\mathbf{F}}_{uc}$ in the optimal affiliation matrix $\hat{\mathbf{F}}$ then describes the likelihood that the node $u \in (V \cup V_M)$ belongs to community $c \in \{1, \dots, C\}$. Therefore, it is possible to recover the community structures ψ by assigning node u to community c if the corresponding entry $\hat{\mathbf{F}}_{uc}$ is greater than or equal to the threshold δ . In the following, we elaborate on each major phase of the KroMFac framework.

4.1 Network Completion

The first step of our framework is the inference of missing parts of the graph from priors on the observable matrix $\bf A$ and the number of missing nodes, M, by using the function GraphRecv implemented via the KronEM algorithm. First, from an initialized $\Theta_{\rm init}$, the E-step samples the missing parts $\bf Z_1$ and $\bf Z_2$, and the permutation σ . In the M-step, a stochastic gradient descent process subsequently optimizes

the parameter matrix Θ given the samples obtained in the E-step. The EM iteration alternates between performing the E-step and M-step according to the following expressions, respectively:

E-step:

$$\left(\mathbf{Z}_1^{(t)}, \mathbf{Z}_2^{(t)}, \sigma^{(t)}\right) \sim \mathbb{P}(\mathbf{Z}_1, \mathbf{Z}_2, \sigma | \mathbf{A}, \mathbf{\Theta}^{(t)}),$$

M-step:

$$\boldsymbol{\Theta}^{(t+1)} = \mathop{\arg\max}_{\boldsymbol{\Theta} \in (0,1)^{N_0}} \mathbb{E}[\mathbb{P}(\mathbf{Z}_1^{(t)}, \mathbf{Z}_2^{(t)}, \boldsymbol{\sigma}^{(t)}, \mathbf{A} | \boldsymbol{\Theta})],$$

where the superscript (t) denotes the iteration index. Then, we generate the stochastic adjacency matrix $\mathbf{\Theta}^{K}$. To create the fully recovered matrix $\mathbf{A}_{R}^{(M)}$, for the first N rows and columns of $\mathbf{A}_R^{(M)}$, we replicate the entries of matrix \mathbf{A} in the upper left of matrix $\mathbf{A}_R^{(M)}$. To infer the missing part (i.e., the last M rows and columns of $\mathbf{A}_R^{(M)}$), we consecutively run the Bernoulli trials with the probability $\Theta^K_{\sigma(u)\sigma(v)}$ and then map the value of the missing entry in row u and column v to one if a success occurs and zero otherwise. Since the adjacency matrix $\mathbf{A}_R^{(M)}$ is symmetric, we only need to repeat this process $MN + \frac{M^2}{2}$ times. An example of this network completion phase is illustrated in Fig. 3 when $N = 6, M = 2, N_0 = 2, K = 3$, and $\sigma(u) = u$ for $u \in (V \cup V_M)$. In this example, since the last two rows and columns of Θ^3 correspond to two recovered nodes, we execute Bernoulli trials on each non-zero entry in this part to obtain the recovered matrix $\mathbf{A}_{R}^{(M)}$. If the number of missing edges can be estimated (see [9] for details), then the termination of this final step can be accelerated for sparse graphs since the Bernoulli trials can be terminated once the number of entries with a value that is equal to one exceeds the predicted number of edges.

4.2 Node Ranking and Selection

Based on the output of the GraphRecv algorithm (i.e., the recovered matrix $\mathbf{A}_R^{(M)}$), the function NodeSelect ranks missing nodes and then selects influential nodes from the set of ranked candidates. In our work, we focus on the well-known centrality measure for ranking nodes, namely degree centrality $\mathrm{Cen}(u)$ for node $u \in (V \cup V_M)$. After computing the degree centrality of missing nodes as in (2), we introduce the ranking vector $r \in \mathbb{N}^M$ to record their centrality ranking. Since we aim to select nodes whose centrality measures are greater than a given threshold ϵ (refer to Section 3.3), only $H \leq M$ most influential nodes

^{4.} Conceptually, numerous centrality measures (e.g., eigenvector centrality and Katz centrality) can also be applied to obtain such a ranking, but it turns out that they do not lead to better performance since the number of immediate neighbors of a node (rather than the number of higher-order neighbors) is important in our node selection even if it is not shown in this paper.

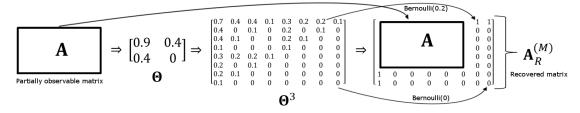


Fig. 3: An illustration of the recovery phase in our KroMFac framework. Here, parameters are set to the following values: $N=6, M=2, N_0=2, K=3$, and $\sigma(u)=u$ for $u\in (V\cup V_M)$.

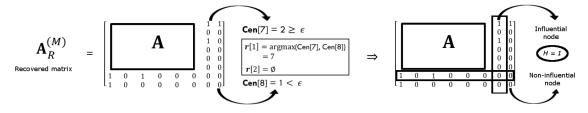


Fig. 4: An illustration of the node ranking and selection phase. Here, parameters are set to the following values: $\epsilon=2$ and M=2.

are associated with r. For example, r[u] represents the index of the node ranked at the uth position in the list. The function NodeSelect returns the number of influential nodes, H, and the ranking vector $r \in \mathbb{N}^H$. Fig. 4 shows a simple illustration of this node ranking and selection phase when $\epsilon=2$ and M=2. In this figure, due to the fact that the last two rows and columns in the input matrix $\mathbf{A}_R^{(M)}$ correspond to two recovered nodes, we calculate the degree centrality of the two nodes and select the seventh placed node as an influential node since its centrality is greater than or equal to ϵ .

4.3 Community Detection

To solve the problem of community detection, we adopt the state-of-the-art NMF-based detection algorithm, named BIGCLAM [19], via as the function CommunDet, which can however be replaced by other community detection methods that return an affiliation matrix. For given i, we solve (5) using a block coordinate gradient ascent algorithm [35]. The optimization process is terminated when the change in each iteration, denoted by $\Delta \mathcal{D} > 0$, is less than an arbitrarily small threshold $\eta_{\text{detect}} > 0$. The algorithm returns the loss function \mathcal{D} and the corresponding matrix \mathbf{F} .

4.4 Analysis of Computational Complexity

In this subsection, we analyze the computational complexity of the KroMFac framework. Since our framework consists of three major phases including network recovery, influential node selection, and community detection, we elaborate on the complexity analysis of each phase. To reduce the complexity, we take advantage of the property that real-world social networks usually have a sparse and low-rank matrix

structure [36]. The network completion phase is based on the KronEM algorithm, which was shown in [9] to have the complexity of $\mathcal{O}(|E|\log|E|)$, where |E| is the number of edges in the partially observable network G. In the node selection phase, the computation of degree centrality dominates the complexity, which is given by $\mathcal{O}(|E|)$ in sparse graphs. In the community detection phase, an almost linear complexity in N can be achieved via the approach in [19], where N denotes the number of observable nodes. Here, while the community detection process is repeated H times, one can see that $H \ll M$ (see Section 5.4.1), where H and M denote the numbers of influential nodes and missing nodes, respectively. From the fact that M is smaller than N, we can deduce that the complexity of this phase is bounded by $\mathcal{O}(N)$. Hence, the total computational complexity of KroMFac is finally given by $\mathcal{O}(|E|\log|E|)$.

5 EXPERIMENTAL EVALUATION

In this section, we first describe both synthetic and real-world datasets. We also present two baseline schemes for community detection as a comparison. By adopting the NMI as a popular information-theoretic performance metric, we then present the performance of our community detection framework and compare it against the two baseline schemes.

5.1 Datasets

To evaluate the community detection performance of our approach, we rely on datasets for which ground-truth communities are explicitly labeled. In the following, both synthetic and real-world datasets across various domains are taken into account.

TABLE 2: LFR parameters of 18 synthetic graphs. Here, M and k denote 10^6 and 10^3 , respectively. NN: the number of nodes, AD: average degree, MD: maximum degree, MinC: minimum community size, MaxC: maximum community size, DE: degree exponent, CSE: community size exponent, MP: mixing parameter, ON: the number of overlapping nodes, CMN: the number of communities per node.

Graphs	NN	AD	MD	MinC	MaxC	DE	CSE	MP	ON	CMN
Graph 1	10k	10	50	10	50	3	1	0.2	500	30
Graph 2	10k	10	50	10	50	3	1	0.25	500	30
Graph 3	10k	10	50	10	50	3	1	0.28	500	30
Graph 4	10k	10	50	10	50	3	1	0.2	500	3
Graph 5	10k	10	50	10	50	3	1	0.25	500	3
Graph 6	10k	10	50	10	50	3	1	0.28	500	3
Graph 7	100k	10	50	10	500	3	1	0.2	5000	30
Graph 8	100k	10	50	10	500	3	1	0.25	5000	30
Graph 9	100k	10	50	10	500	3	1	0.28	5000	30
Graph 10	100k	10	50	10	500	3	1	0.2	5000	3
Graph 11	100k	10	50	10	500	3	1	0.25	5000	3
Graph 12	100k	10	50	10	500	3	1	0.28	5000	3
Graph 13	1M	10	50	10	5000	3	1	0.2	50000	30
Graph 14	1M	10	50	10	5000	3	1	0.25	50000	30
Graph 15	1M	10	50	10	5000	3	1	0.28	50000	30
Graph 16	1M	10	50	10	5000	3	1	0.2	50000	3
Graph 17	1M	10	50	10	5000	3	1	0.25	50000	3
Graph 18	1M	10	50	10	5000	3	1	0.28	50000	3

5.1.1 Synthetic Datasets

We construct synthetic graphs via the extended Lancichinetti-Fortunato-Radicchi (LFR) benchmark [37], which is built upon a generative model that creates nodes along with prior known community labels. The benchmark is capable of generating graphs that replicate important features of real social networks, such as the power-law degree distribution and overlapping communities. To create an LFR graph, ten parameters need to be specified, which are summarized in Table 2. While parameters such as the number of nodes, average degree, maximum degree, maximum and minimum community size, and degree exponent are rather straightforward to understand, we explain the remaining four parameters:

- The community size exponent refers to the exponent parameter from a power-law approximation of the distribution of the number of nodes in communities.
- The mixing parameter, denoted by μ , controls the proportion of random edges to total edges. For example, if $\mu=0.3$, then the LFR benchmark produces a graph such that approximately 70% of edges link to nodes within the same community, while the remaining 30% connect to nodes in other randomly selected communities. This parameter is sensitive to the performance of community detection. In general, if μ is closer to one, then the community structures become weaker. On the other hand, when μ is closer to zero, one can expect high detection performance since community structures can be easily identified.
- The number of overlapping nodes refers to the number of nodes in the graph that belong to more than one community.

TABLE 3: Statistics of the six real-world datasets. Here, M and k denote 10^6 and 10^3 , respectively

Dataset	NN	NE	NC	ACS	CMN
Amazon	0.34M	0.93M	49k	99.86	14.83
DBLP	0.43M	1.3M	2.5k	429.79	2.57
Youtube	1.1M	3.0M	30k	9.75	0.26
Facebook	4k	88k	193	22.93	1.14
Twitter	81k	2.4M	4k	33.50	1.65
Orkut	3M	117M	6.3M	34.86	95.93

 The number of communities per node indicates the average number of communities to which each of the overlapping nodes belongs.

To cover various domains of network applications, we generate 18 LFR graphs with differing parameter settings according to [5] as specified in Table 2, where we show the various values chosen for representative parameters such as NN, MP, ON, and CMN.

5.1.2 Real-World Datasets

To validate the applicability of our approach, six real-world datasets are also used for evaluation. More specifically, from the available SNAP datasets [38] that have ground-truth communities, we use the Amazon product co-purchasing network [39], the collaboration network of DBLP [40], the Youtube video-sharing social network [41], and the three friendship social networks of Facebook [42], Twitter [42], and Orkut [39]. The statistics of these datasets are summarized in Table 3, and the basic characteristics of each network are described in the following:

- *The number of nodes* (*NN*): In the Amazon network, nodes represent products. In DBLP, nodes represent authors. In the Youtube, Facebook, Twitter, and Orkut networks, nodes represent users.
- The number of edges (NE): In the Amazon network, edges connect products that are commonly purchased together. In DBLP, two authors are connected by an edge if they have co-authored a paper. In the Youtube, Facebook, Twitter, and Orkut networks, edges represent friendships between users.
- The number of communities (NC): In the Amazon network, each product category corresponds to a ground-truth community. In DBLP, the publication venues are used as ground-truth communities. In the Youtube and Orkut networks, user-created groups are used as ground-truth communities. In the Facebook and Twitter networks, circles of users are used as ground-truth communities.
- Average community size (ACS): The average number of nodes within communities.
- *Community memberships per node (CMN):* The average number of communities that a node belongs to.

5.2 Baseline Approaches

Due to the fact that community discovery in partially observable networks with both missing nodes and edges has never been studied in the literature, there is no state-of-theart method that works appropriately under our network model. For this reason, we present two types of baseline schemes by taking into account some special cases of our KroMFac framework and its variants.

5.2.1 Baseline 1 (Community Detection)

As a naïve approach, the first baseline scheme (Baseline 1) aims to directly discover community structures based on an observable network via the NMF-aided detection method without recovering any nodes and edges. To this end, Baseline 1 solves an optimization problem such that the matrix **F** is found given an adjacency matrix **A** of the incomplete network. The problem formulation is thus given by $\hat{\mathbf{F}} = \arg \max_{\mathbf{F} > 0} \mathbb{P}(\mathbf{A}|\mathbf{F})$, which corresponds to a special case with no regularization term where i is set to zero in our joint optimization problem (1). Similarly as in the methodology in Section 3.4, the optimal F can be easily acquired via the function CommunDet by replacing the input matrix $\mathbf{A}_{R}^{(i)}$ by **A**. As CommunDet results in $\hat{\mathbf{F}}$ providing fuzzy information on the community memberships, we apply the hard-decision process (refer to lines 15-18 in Algorithm 1) to find the community structure ψ .

5.2.2 Baseline 2 (Network Completion + Community Detection)

In addition to Baseline 1, to highlight the importance of node ranking and selection, we also present the second baseline scheme (Baseline 2) that performs community detection along with a full graph recovery. To this end, given an adjacency matrix $\mathbf{A}_{R}^{(M)}$ that is inferred by the network completion phase, Baseline 2 solves an optimization problem such that the matrix ${\bf F}$ is found. The problem formulation is thus given by $\hat{\mathbf{F}} = \arg\max_{\mathbf{F} \geq 0} \mathbb{P}(\mathbf{A}_R^{(M)}|\mathbf{F})$, which corresponds to another special case with no regularization term, where iis set to M in (1). Note that node ranking is not necessary since all recovered nodes are inserted into the graph. To solve this problem, we first follow the steps similar to those in Section 3.2 to recover the matrix $\mathbf{A}_{R}^{(M)}$. After the network completion phase via the function GraphRecv, the optimal **F** can be acquired via CommunDet by assuming that i = M. Then, the hard-decision process in Algorithm 1 is performed to produce the final result ψ .

Additionally, we consider two state-of-the-art algorithms for label propagation-based community detection, dubbed COPRA [22] and SLPA [23], as alternatives of BIGCLAM. COPRA enables each node to update its community assignment coefficients by allowing a node to have multiple labels,

while SLPA allows nodes to exchange labels according to pairwise interaction rules. In this case, we first recover the matrix $\mathbf{A}_R^{(M)}$ via the function GraphRecv and then obtain ψ via either COPRA or SLPA.

5.3 Performance Metric

To assess the performance of our KroMFac framework and two baseline schemes, we need to quantify the degree of agreement between the ground-truth communities and the detected communities. In particular, given a set of true labels and the set of labels assigned by the resulting community detection, we need to find the similarity between them.

While there are various ways to estimate the similarity, the NMI is one of the most widely used evaluation measures for community detection problems [19], [20], and is formally defined as in the following.

Definition 3 (NMI [43]). Assume that the community assignments are x_i and y_i , where x_i and y_i indicate the labels of vertex i in the true community $\mathcal X$ and the predicted community $\mathcal Y$, respectively. When the labels x and y are the values of two random variables X and Y, following a joint distribution $\mathbb P(x,y) = \mathbb P(X = x,Y = y)$, the NMI between $\mathcal X$ and $\mathcal Y$ is given by $\mathrm{NMI}(\mathcal X,\mathcal Y) = 1 - \frac{1}{2} \left(\frac{H(X|Y)}{H(X)} + \frac{H(Y|X)}{H(Y)} \right)$, where $H(X) = -\sum_x \mathbb P(x) \log \mathbb P(x)$ is the Shannon entropy of X and $H(X|Y) = -\sum_{x,y} \mathbb P(x,y) \log \mathbb P(x|y)$ is the conditional entropy of X given Y.

5.4 Experimental Results

To create partially observable networks *G* from the underlying true graphs G', we adopt two graph sampling strategies from [44]. The first strategy, called random node (RN) sampling, selects nodes uniformly at random to create a sample graph. The second one, forest fire (FF) sampling, starts by picking a seed node uniformly at random and adding it to a sample graph (referred to as burning). Then, FF sampling burns a fraction of the outgoing links with nodes attached to them. This process is recursively repeated for each neighbor that is burned until no new node is selected to be burned. Both sampling strategies are known not to be biased towards high degree nodes, while FF sampling is capable of preserving the degree distribution of the true graph [44]. We create partially observable networks consisting of 70% nodes from the original synthetic and real-world datasets mentioned in Sections 5.1.1 and 5.1.2, respectively, by performing two aforementioned graph sampling strategies. To

5. Instead of exploiting the number of communities as side information, COPRA and SLPA take advantage of the number of communities per node and the minimum and maximum sizes of communities, respectively, to recover the community structures.

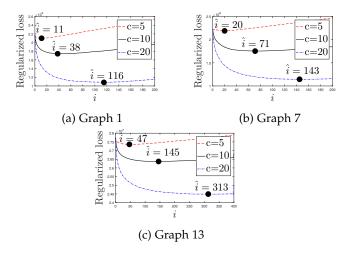


Fig. 5: Regularized loss over the number of added nodes, i, according to different values of c. Here, black circles depict the points at which the minimum losses are attained.

be consistent in evaluating the performance of community detection in the incomplete graphs, we also perform node deletion such that the ground-truth community structures ψ do not contain nodes removed from the original graphs.

Our empirical study is basically designed to answer the following four key research questions.

- *Q*1. To what extent does the parameter λ in regularization affect the performance?
- *Q*2. How close is the optimal \hat{i} that is the solution to (5) to the value of i that maximizes the NMI?
- Q3. How does the performance change when different community detection methods are adopted?
- *Q4*. How much does our **KroMFac** method enhance the performance over the baseline schemes?

5.4.1 Sensitivity Analysis (Q1 & Q2)

For the sensitivity analysis, we consider only the results where RN sampling is applied to Graphs 1, 7, and 13 of our synthetic graphs (each with different network sizes) since other cases follow similar trends. First, we analyze the sensitivity of the parameter λ , which determines the impact of regularization in our optimization problem and plays a crucial role in determining performance of our KroMFac framework. In Fig. 5, the regularized loss over the number of added influential nodes, *i*, is illustrated according to various values of c > 0, where $\lambda = cN$. We find that setting c to a small value (e.g., c = 5) results in a value of \hat{i} that is too low to compensate the loss function properly with the regularization term. In contrast, setting it to a high value (e.g., c = 20) results in a value of \hat{i} that falls out of the acceptable range of i, i.e., i > H, due to over-regularization. We empirically verify that setting $\lambda = 10N$ manifests satisfactory performance in terms of regularized loss and NMI in the following experiments. We demonstrate that our setting

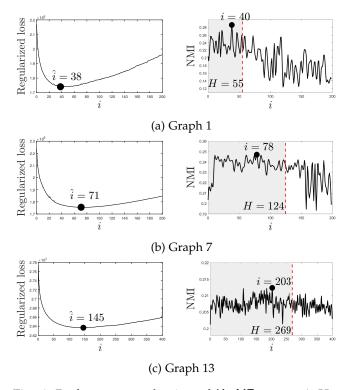


Fig. 6: Performance evaluation of KroMFac over *i*. Here, black circles depict the points (left) at which the minimum losses are attained and the points (right) at which the maximum NMIs are achieved.

of λ is robust to the overall network attributes, including the network size. Furthermore, we investigate how close the optimal \hat{i} that is the solution to (5) is to the value of i that maximizes the NMI. In Figs. 6a–6c (left), we illustrate the regularized losses (i.e., the objective function in (5)) over i, where the minimum losses in Graphs 1, 7, and 13 are attained at $\hat{i} = 38$, $\hat{i} = 71$, and $\hat{i} = 145$, respectively. In Figs. 6a–6c (right), we plot the NMI over i, where the maximum NMIs in Graphs 1, 7, and 13 are achieved at i = 40, i = 78, and i = 203, respectively. From these three graphs, we observe that adding more nodes and edges to the existing graph increases the NMI scores up to a certain number of nodes (e.g., i = 40 in Fig. 6a (right)), but drops if more nodes are added due to a higher accumulated inference error, which verifies our assertion made in Section 3. The fact that the NMIs attained by the minimum losses are close to the maximum NMIs in Figs. 6a-6c is an indication that the solutions to (5) also ensure satisfactory performance with regard to the NMI. Additionally, we empirically determine the threshold ϵ , which plays a crucial role in specifying the number of influential nodes, H. When we set $\epsilon = \frac{k_{\max}}{2}$ as in [45], the resulting value of H in Graphs 1, 7, and 13 are 55, 124, and 269, respectively, and the NMIs are computed by searching over $i \in \{1, \dots, H\}$ (see the shaded areas in Figs. 6a–6c), where $k_{\rm max} > 0$ is the maximum degree of nodes in a recovered graph. Since the threshold setting

TABLE 4: NMI of Baseline 2 according to different community detection algorithms

Name	Baseline 2-BIGCLAM (X)	Baseline 2-COPRA (Y)	Baseline 2-SLPA (Z)	Improvement rate (%)		
	baseiine 2-biGCLAM (A)	baseline 2-COFKA (1)	baseline 2-SLFA (Z)	$\frac{X-Y}{Y} \times 100$	$\frac{X-Z}{Z} \times 100$	
Graph 1	0.1882	0.0000	0.0012	-	-	
Graph 2	0.1710	0.0000	0.0020	-	-	
Graph 3	0.1354	0.0000	0.0022	-	-	
Graph 4	0.4014	0.3024	0.3121	32.74	28.61	
Graph 5	0.3245	0.2201	0.2031	47.43	59.77	
Graph 6	0.2648	0.1704	0.1425	55.40	85.83	
DBLP	0.1399	0.1112	0.1097	20.51	21.59	

 $\epsilon = \frac{k_{\mathrm{max}}}{2}$ leads to a reduction in computational complexity without loss of performance, we adopt it in our experiments in the following. This result also suggests that adding only a small number of nodes and their associated edges to the existing graph is sufficient to remarkably enhance the NMI performance.

5.4.2 Comparison With Other Community Detection Algorithms (Q3)

To see how the performance varies according to different community detection algorithms, we also evaluate the performance of COPRA [22] and SLPA [23] as specified in Section 5.2.2. Table 4 presents the NMI of Baseline 2 for fully recovered graphs, where RN sampling is applied to the DBLP network as well as Graphs 16, since the LFR parameters such as MP and CMN in synthetic datasets significantly affect the performance and the DBLP network is one of the datasets having the value of CMN comparable to that of Graphs 4–6 (refer to Table 3). The results in Table 4 demonstrate the superiority of BIGCLAM over COPRA and SLPA. It is worthwhile to note that COPRA and SLPA almost entirely fail to recover the overlapping community structures for more difficult situations that have a higher value of CMN (e.g., Graphs 1–3).

5.4.3 Comparison With Baseline Schemes (Q4)

Finally, the NMI of KroMFac and two baseline schemes for synthetic and real-world graphs are shown in Table 5, where both RN and FF sampling strategies are applied and BIGCLAM is adopted for community detection. The NMI performance of BIGCLAM for complete networks without deleting nodes and edges is also shown to provide an upper bound of our approach. We find that our KroMFac framework noticeably outperforms the baselines for all synthetic and real-world datasets with substantial improvement rates of up to 22.85% and 63.47% over Baselines 1 and 2, respectively. Interestingly, our findings reveal that the NMI performance of KroMFac is not likely to be influenced by the network size; for example, similar NMIs are obtained from Graphs 1, 7, and 13, whose network attributes are identical, except for the network size. From the results, it is also clear that the performance of Baseline 2 is almost comparable to or even worse than that of Baseline 1, which shows that

TABLE 5: NMI of KroMFac and two baseline schemes

		Complete	KroMFac	Baseline 1	Baseline 2	Improveme	ent rate (%)
	Name	graph	(X)	(Y)	(Z)	$\frac{X-Y}{V} \times 100$	$\frac{X-Z}{Z} \times 100$
	Graph 1 (RN)		0.2601	0.2008	0.1882	22.78	27.64
	Graph 1 (FF)	0.5125	0.3409	0.3226	0.2439	5.37	28.46
	Graph 2 (RN)		0.1917	0.1746	0.1710	8.88	10.79
	Graph 2 (FF)	0.4234	0.2667	0.2510	0.2054	5.90	23.00
	Graph 3 (RN)		0.1667	0.1410	0.1354	15.44	18.76
	Graph 3 (FF)	0.3886	0.2566	0.2234	0.1907	12.94	25.70
	Graph 4 (RN)		0.4989	0.4764	0.4014	4.52	19.55
	Graph 4 (FF)	0.7825	0.5182	0.4782	0.4320	7.71	16.62
	Graph 5 (RN)		0.3182	0.3675	0.3245	5.33	16.40
	Graph 5 (FF)	0.7214	0.3968	0.3517	0.3412	11.36	14.00
	Graph 6 (RN)		0.3558	0.3042	0.2648	14.51	25.59
		0.6821	0.3338	0.3042	0.28339	22.85	31.63
	Graph 6 (FF)					5.70	
	Graph 7 (RN)	0.4618	0.2266	0.2033	0.1931		28.48
	Graph 7 (FF)		0.3467	0.3269	0.2479	10.28	14.78
os.	Graph 8 (RN)	0.3654	0.1971	0.1800	0.1537	8.67	22.01
Synthetic datasets	Graph 8 (FF)		0.2630	0.2304	0.2153	12.39	18.12
ata	Graph 9 (RN)	0.3577	0.1709	0.1473	0.1398	13.79	18.20
5	Graph 9 (FF)		0.2704	0.2569	0.2036	4.97	24.72
het	Graph 10 (RN)	0.8125	0.4794	0.4483	0.3806	6.49	20.62
Jut.	Graph 10 (FF)	0.7548	0.5214	0.4848	0.4320	7.01	17.14
æ,	Graph 11 (RN)		0.3645	0.3472	0.2727	4.75	25.19
	Graph 11 (FF)		0.3989	0.3675	0.3104	7.88	22.18
	Graph 12 (RN)	0.6972	0.3231	0.2867	0.2560	11.26	20.76
	Graph 12 (FF)	0.07.2	0.3753	0.3128	0.2195	16.65	41.51
	Graph 13 (RN)	0.4213	0.2103	0.2001	0.1864	4.86	11.36
	Graph 13 (FF)	0.1210	0.3001	0.2871	0.1914	4.33	36.23
	Graph 14 (RN)	0.3125	0.1831	0.1732	0.1517	5.42	17.14
	Graph 14 (FF)	0.3123	0.2142	0.1855	0.1614	13.40	24.67
	Graph 15 (RN)	0.2641	0.1649	0.1373	0.1338	16.74	18.86
	Graph 15 (FF)	0.2041	0.2139	0.1941	0.1234	9.26	42.29
	Graph 16 (RN)	0.6815	0.3994	0.3416	0.3246	14.48	18.74
	Graph 16 (FF)	0.0013	0.4664	0.4021	0.3632	13.79	22.13
	Graph 17 (RN)	0.5122	0.3445	0.2972	0.2727	13.73	20.85
	Graph 17 (FF)	0.5122	0.3235	0.3024	0.2617	6.51	19.11
	Graph 18 (RN)	0.4241	0.2971	0.2347	0.2160	21.01	27.29
	Graph 18 (FF)	0.4241	0.3209	0.2944	0.1747	8.27	45.57
atasets	Amazon (RN)	0.2401	0.1448	0.1327	0.1219	9.14	18.78
	Amazon (FF)	0.3481	0.1962	0.1837	0.1727	6.79	13.59
	DBLP (RN)	0.2240	0.1667	0.1436	0.1399	16.11	19.15
	DBLP (FF)	0.3249	0.2866	0.2605	0.2583	10.02	10.97
	Youtube (RN)		0.1263	0.1179	0.1058	7.12	19.34
1 4	Youtube (FF)	0.3042	0.1556	0.1514	0.1431	2.80	8.73
Real-world datasets	Facebook (RN)		0.0741	0.0712	0.0545	3.91	26.45
	Facebook (FF)	0.1425	0.0912	0.0878	0.0745	3.71	18.31
Real	Twitter (RN)		0.0725	0.0701	0.0521	3.34	28.16
=	Twitter (FF)	0.1237	0.0815	0.0685	0.0584	15.97	28.36
	Orkut (RN)		0.0824	0.0684	0.0301	16.99	63.47
	Orkut (FF)	0.1249	0.1075	0.0845	0.0488	21.40	54.61
	JIKUI (II)		0.1075	0.0045	0.0400	21.10	54.01

including the entirety of recovered nodes and edges is not beneficial. Furthermore, we observe that the performance of both KroMFac and two baselines is higher when the FF sampling strategy is applied due to the fact that FF sampling tends to preserve the degree distribution of the true graph.

5.5 Empirical Evaluation of Complexity

We empirically show the average runtime complexity via experiments using synthetic graphs sampled by FF for different numbers of sampled nodes with $N+M=2^k$ and $k\in\{13,\cdots,19\}$. Then, 30% of nodes and their associated edges are deleted by FF sampling to create partially observable networks. Parameters of the LFR benchmark essentially follow those of Graph 1, where MaxC and ON are set proportionally to the number of nodes (refer to Table 2). Parameters of the KroMFac framework follow the

6. The runtime complexity under the LFR parameters of other synthetic graphs shows similar trends even if it is not shown in this paper. We note that real-world datasets are not usable for the evaluation of complexity since it is hardly feasible to scale up/down real-world graphs while preserving their structural properties as in the case where synthetic graphs are adopted.

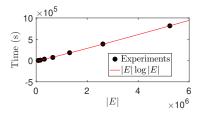


Fig. 7: The computational complexity of the KroMFac framework.

same settings as in Section 5.4.1. In Fig. 7, we illustrate the plot of the runtime complexity in seconds versus |E|, where each point is the average of experimental results obtained by executing the KroMFac process 10 times. An asymptotic curve $|E|\log |E|$ is also shown in the figure, showing a trend that is consistent with our experimental results. Moreover, we note that since the computation based on a large-size matrix is expensive in terms of memory consumption, it is necessary to adopt cost-effective techniques (e.g., coordinate format) to store and compute sparse matrices.

6 CONCLUDING REMARKS

In this paper, we introduced the problem of discovering overlapping community structures in the context of partially observable networks with both missing nodes and edges. To solve this problem, we developed a novel framework, termed KroMFac, that seamlessly incorporates network completion into community recovery. Specifically, we performed community detection via regularized NMF based on the Kronecker graph model. In particular, motivated by the insight that adding a proper number of missing nodes and edges to the existing graph would be of significant importance in improving community detection accuracy, we presented how to characterize and select influential nodes via centrality ranking. By adopting the NMI as a performance metric, we validated our proposed KroMFac framework through experiments on both synthetic and realworld datasets. Based on parameter search, we showed that our approach outperforms two baselines by a large margin on synthetic and real-world networks. Additionally, we analytically examined the computational complexity of our framework.

Potential avenues of future research in this area are the inclusion of deep generative graph models for community detection to reduce the inference error even further.

ACKNOWLEDGMENTS

This work was supported by the Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (2017R1D1A1A09000835). Won-Yong Shin is the corresponding author.

REFERENCES

- [1] S. Fortunato, "Community detection in graphs," *Phys. Rep.*, vol. 486, no. 3, pp. 75–174, Feb. 2010.
- [2] J. Leskovec, K. J. Lang, and M. Mahoney, "Empirical comparison of algorithms for network community detection," in *Proc. 19th Int. Conf. World Wide Web (WWW)*, Rayleigh, NC, Apr. 2010, pp. 631– 640.
- [3] M. E. Newman, "Modularity and community structure in networks," Proc. of the National Acad. Sci., vol. 103, no. 23, pp. 8577– 8582, Apr. 2006.
- [4] I. X. Leung, P. Hui, P. Lio, and J. Crowcroft, "Towards real-time community detection in large networks," *Phys. Rev. E*, vol. 79, no. 6, p. 066107, Jun. 2009.
- [5] J. Xie, S. Kelley, and B. K. Szymanski, "Overlapping community detection in networks: The state-of-the-art and comparative study," *ACM Comp. Sur.*, vol. 45, no. 4, p. 43, Aug. 2013.
- [6] G. Kossinets, "Effects of missing data in social networks," Soc. Netw., vol. 28, no. 3, pp. 247–268, Jul. 2006.
- [7] A. Acquisti, L. Brandimarte, and G. Loewenstein, "Privacy and human behavior in the age of information," *Science*, vol. 347, no. 6221, pp. 509–514, Jan. 2015.
- [8] R. Dey, Z. Jelveh, and K. Ross, "Facebook users have become much more private: A large-scale study," in Proc. IEEE Int. Conf. Pervasive Comput. Commun. Worksh. (PERCOM), Lugano, Switzerland, 2012, pp. 346–352.
- [9] M. Kim and J. Leskovec, "The network completion problem: Inferring missing nodes and edges in networks," in *Proc. 11th SIAM Int. Conf. Data Mining (SDM)*, Mesa, AZ, Apr. 2011, pp. 47–58.
- [10] E. M. Airoldi, D. M. Blei, S. E. Fienberg, and E. P. Xing, "Mixed membership stochastic blockmodels," J. Mach. Learn. Research, vol. 9, pp. 1981–2014, Sep. 2008.
- [11] M. E. Newman, "Spectral methods for community detection and graph partitioning," Phys. Rev. E, vol. 88, no. 4, p. 042822, Oct. 2013.
- [12] N. Du, B. Wu, X. Pei, B. Wang, and L. Xu, "Community detection in large-scale social networks," in *Proc. 9th WebKDD and 1st SNA-KDD Work. Web mining and Soc. Netw. analysis*, San Jose, CA, Aug. 2007, pp. 16–25.
- [13] J. Chen and Y. Saad, "Dense subgraph extraction with application to community detection," *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 7, pp. 1216–1230, Jul. 2012.
- [14] U. N. Raghavan, R. Albert, and S. Kumara, "Near linear time algorithm to detect community structures in large-scale networks," *Phys. Rev. E*, vol. 76, no. 3, p. 036106, Sep. 2007.
- [15] J. Yang and J. Leskovec, "Structure and overlaps of communities in networks," in *Proc. 6th Work. Soc. Netw. Mining and Analysis* (SNA-KDD), Beijing, China, Aug. 2012, pp. 661–703.
- [16] S. Papadopoulos, Y. Kompatsiaris, A. Vakali, and P. Spyridonos, "Community detection in social media," *Data Mining and Knowl. Disc.*, vol. 24, no. 3, pp. 515–554, Jun. 2012.
- [17] J. J. Whang, D. F. Gleich, and I. S. Dhillon, "Overlapping community detection using neighborhood-inflated seed expansion," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 5, pp. 1272–1284, May. 2016.
- [18] J. Kim, S. Lim, J.-G. Lee, and B. S. Lee, "LinkBlackHole*: Robust overlapping community detection using link embedding," *IEEE Trans. Knowl. Data Eng.*, to appear.
- [19] J. Yang and J. Leskovec, "Overlapping community detection at scale: A nonnegative matrix factorization approach," in *Proc. 6th* ACM Int. Conf. Web. Search and Data Mining (WSDM), Rome, Italy, Feb. 2013, pp. 587–596.
- [20] I. Psorakis, S. Roberts, M. Ebden, and B. Sheldon, "Overlapping community detection using Bayesian non-negative matrix factorization," Phys. Rev. E, vol. 83, no. 6, p. 066114, Feb. 2011.

- [21] Y. Zhang and D.-Y. Yeung, "Overlapping community detection via bounded nonnegative matrix tri-factorization," in *Proc. 18th Int. Conf. Knowl. Disc. and Data Mining (SIGKDD)*, Beijing, China, Aug. 2012, pp. 606–614.
- [22] S. Gregory, "Finding overlapping communities in networks by label propagation," New J. Phys., vol. 12, no. 10, p. 103018, Oct. 2010.
- [23] J. Xie, B. K. Szymanski, and X. Liu, "SLPA: Uncovering overlapping communities in social networks via a speaker-listener interaction dynamic process," in *Proc. 11th Int. Conf. Data Mining Worksh. (ICDMW)*, Atlantic City, NJ, Nov. 2011, pp. 344–349.
- [24] J. Yang, J. McAuley, and J. Leskovec, "Community detection in networks with node attributes," in *Proc. IEEE 13th Int. Conf. Data Mining (ICDM)*, Dallas, TX, Dec. 2013, pp. 1151–1156.
- [25] B. Yan and S. Gregory, "Finding missing edges and communities in incomplete networks," J. Phys. A: Math. and Theo., vol. 44, no. 49, p. 495102, Nov. 2011.
- [26] B. Ya and S. Gregory, "Detecting community structure in networks using edge prediction methods," *J. Stat. Mech.: Theory and Exp.*, vol. 2012, no. 09, p. P09008, Sep. 2012.
- [27] M. Zhou, "Infinite edge partition models for overlapping community detection and link prediction," in *Proc. 18th Int. Conf. Artifi. Intel. and Stat. (AISTATS)*, San Diego, CA, Feb. 2015, pp. 1135–1143.
- [28] H.-H. Chen, L. Gou, X. Zhang, and C. L. Giles, "Capturing missing edges in social networks using vertex similarity," in *Proc. 6th Int. Conf. Knowl. Cap.*, Alberta, Canada, Jun. 2011, pp. 195–196.
- [29] R. Eyal, A. Rosenfeld, S. Sina, and S. Kraus, "Predicting and identifying missing node information in social networks," ACM Trans. Knowl. Disc. Data, vol. 8, no. 3, p. 14, 2014.
- [30] T. Eden, S. Jain, A. Pinar, D. Ron, and C. Seshadhri, "Provable and practical approximations for the degree distribution using sublinear graph samples," preprint. [Online]. Available: https://arxiv.org/abs/1710.08607.
- [31] A. Clauset, C. R. Shalizi, and M. E. J. Newman, "Power-Law Distributions in Empirical Data," SIAM Rev., vol. 51, no. 4, pp. 661–703, Feb. 2009.
- [32] T. H. McCormick and M. J. Salganik, "How many people you know?: Efficiently esimating personal network size," *J. Am. Stat. Assoc.*, vol. 105, no. 489, pp. 59–70, Sep. 2010.
- [33] J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, and Z. Ghahramani, "Kronecker graphs: An approach to modeling networks," J. Mach. Learn. Research, vol. 11, pp. 985–1042, Feb. 2010.
- [34] P. M. Weichsel, "The kronecker product of graphs," *J. Am. Math. Soc.*, vol. 13, no. 1, pp. 47–52, 1962.
- [35] C.-J. Hsieh and I. S. Dhillon, "Fast coordinate descent methods with variable selection for non-negative matrix factorization," in Proc. 17th Int. Conf. Knowl. Disc. and Data Mining (SIGKDD), San Diego, CA, Aug. 2011, pp. 1064–1072.
- [36] E. Richard, P.-A. Savalle, and N. Vayatis, "Estimation of simultaneously sparse and low rank matrices," in *Proc. 29th Int. Conf. Machine Learning (ICML'12)*, Edinburgh, Scotland, Jun.-Jul. 2012, pp. 51–58.
- [37] A. Lancichinetti and S. Fortunato, "Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities," *Phys. Rev. E*, vol. 80, no. 1, p. 016118, Apr. 2009.
- [38] J. Leskovec and A. Krevl, "SNAP datasets: Stanford large network dataset collection," Jun. 2014. [Online]. Available: http://snap.stanford.edu/data.
- [39] J. Yang and J. Leskovec, "Defining and evaluating network communities based on ground-truth," *Knowl. and Inf. Sys.*, vol. 42, no. 1, pp. 181–213, 2015.
- [40] L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan, "Group formation in large social networks: membership, growth, and

- evolution," in *Proc. 12th ACM Int. Conf. Knowl. Disc. and Data Mining (SIGKDD)*, Philadelphia, PA, Aug. 2006, pp. 44–54.
- [41] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhat-tacharjee, "Measurement and analysis of online social networks," in *Proc. 5th ACM/Usenix Internet Measurement Conf. (IMC'07)*, San Diego, CA, Oct. 2007, pp. 29–42.
- [42] J. Leskovec and J. J. Mcauley, "Learning to discover social circles in ego networks," in *Proc. 25th Int. Conf. Neural Inf. Process. Syst.* (NIPS'12), Lake Tahoe, NV, Dec. 2012, pp. 539–547.
- [43] A. Lancichinetti, S. Fortunato, and J. Kertész, "Detecting the overlapping and hierarchical community structure in complex networks," New J. Phys., vol. 11, no. 3, p. 033015, Mar. 2009.
- [44] J. Leskovec and C. Faloutsos, "Sampling from large graphs," in Proc. 12th Int. Conf. Knowl. Disc. and Data Mining (SIGKDD), Philadelphia, PA, 2006, pp. 631–636.
- [45] Y. Yuan, G. Wang, and Y. Sun, "FISH: A Novel Peer-to-Peer overlay network based on Hyper-deBruijn." in Proc. 11th Int. Conf. Web-Age Inf. Man., Jiuzhaigou, China, Jun. 2010, pp. 47–61.