
Structure-Preserving Embedding of Multi-layer Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper investigates structure-preserving embedding for multi-layer networks
2 with community structure. We propose a novel generative tensor-based latent space
3 model (TLSM) that allows heterogeneity among vertices. It embeds vertices into
4 a low-dimensional latent space so that vertices within the same community are
5 close to each other in the ambient space, and captures layer heterogeneity through
6 a layer-effect factor matrix. With a general and flexible tensor decomposition
7 on the expected network adjacency tensor, TLSM is dedicated to preserving the
8 original vertex relations and layer-specific effects in the network embedding. An
9 efficient alternative updating scheme is developed to estimate the model parameters
10 and conduct community detection simultaneously. Theoretically, we establish the
11 asymptotic consistencies of TLSM in terms of both multi-layer network estimation
12 and community detection. The theoretical results are supported by extensive
13 numerical experiments on both synthetic and real-life multi-layer networks.

14 1 Introduction

15 Network has arisen as one of the most common structures to represent the relations among entities.
16 In many complex systems, entities can be multi-relational in that they may interact with each other
17 under various circumstances. A multi-layer network, which consists of a common vertex set across all
18 network layers representing the entities and an edge set at each layer to characterize a particular type
19 of relation among entities, is faithful to represent these relations. Examples of multi-layer networks
20 include social networks of multiple interaction channels [42, 15], biological networks of different
21 collaboration schemes [49, 31, 29] and world trading networks [1, 37] of various goods.

22 In this paper, we propose a structure-preserving embedding framework for multi-layer networks
23 via a tensor-based latent space model. Specifically, TLSM utilizes the factorization of network
24 adjacency tensor as a building block, embeds the vertices into a low dimensional latent space, and
25 captures the heterogeneity among different layers through a layer-effect factor matrix. Consequently,
26 the community structure of the multi-layer network can be detected from a network embedding
27 perspective, such that vertices within the same community are closer to one another in the ambient
28 space than those in different communities. In addition, one key feature of TLSM is that it introduces
29 a sparsity factor into the vanilla logit transformation of the network adjacency tensor, which allows
30 TLSM to model sparse multi-layer networks in a more explicit fashion and accommodate relatively
31 sparser multi-layer networks as the ones considered in literature [22]. More importantly, this sparsity
32 factor can be estimated from the network adjacency tensor directly.

33 The main contribution of this paper is three-fold. First, the proposed TLSM is flexible and general
34 in that it includes many popular network models as special cases. It also relaxes the layer-wise
35 positive semi-definite condition that has been frequently employed in literature [6, 35]. Second, a
36 joint modeling framework is constructed for TLSM, consisting of the multi-layer network likelihood

37 and a clustering type penalty, to estimate the multi-layer network and conduct community detection
 38 simultaneously. Its advantages are supported by extensive numerical experiments on both synthetic
 39 and real-life multi-layer networks. Third, the asymptotic consistencies of TLSM are established in
 40 terms of both multi-layer network estimation and community detection. Notably, the established
 41 theoretical results imply that the proposed methods can accommodate the sparsest multi-layer
 42 networks considered in literature.

43 The rest of the paper is organized as follows. The remaining of Section 1 discusses related works and
 44 introduces necessary notations. Section 2 presents the proposed TLSM and its estimation scheme with
 45 an efficient algorithm. In Section 3, we establish the asymptotic consistencies of TLSM. Extensive
 46 numerical performance of TLSM on synthetic and real-life multi-layer networks as well as ablation
 47 studies on two novel components of the proposed method are carried out in Section 4. Section 5
 48 concludes the paper. The supplementary materials contains technique proofs and necessary lemmas,
 49 additional simulation studies, detailed parameter tuning process, among others.

50 1.1 Related work

51 While there is a growing number of literature focusing on community detection in single-layer
 52 network [48, 28, 13], community detection in multi-layer network is still in its infancy. One classical
 53 approach is to detect community structure in each layer separately [4, 5], which fails to leverage
 54 the homogeneity across different layers. Another approach is to aggregate multi-layer networks
 55 into a single-layer one [41, 12, 35], which heavily relies on the assumption of homogeneous linking
 56 pattern across multiple layers. Recently, [26] proposed to aggregate the biased-adjusted version of
 57 the squared adjacency matrix in each layer to alleviate the information loss in aggregation. yet it
 58 requires the average node degree to grow at a sub-optimal order.

59 In terms of multi-layer network generative models, [34] extended the seminal stochastic block
 60 model (SBM; 19) to the multi-layer stochastic block model (MLSBM; 34), where the probability for
 61 any two vertices to form an edge in a given layer depends only on their community memberships.
 62 Clearly, MLSBM heavily relies on the assumption of homogeneous vertices within communities.
 63 The framework of MLSBM has also been incorporated in degree-corrected network estimation [36],
 64 spectral clustering [6, 35, 26], least square estimation [27] and likelihood-based approaches [45]. In
 65 addition, network response regression model [46] and tensor factorization methods [8, 22] have also
 66 been proposed to detect community structures in multi-layer networks.

67 To allow heterogeneous vertices, the latent space model [18] and random dot product graph model
 68 [3] have been extended to multi-layer networks[47, 32, 2]. In addition, graph neural network and
 69 graph convolutional networks has been extended to multi-layer network for learning the multi-layer
 70 network embedding [14, 23, 17, 39].

71 1.2 Notations

72 Throughout the paper, we use boldface calligraphic Euler scripts (\mathcal{A}) to denote tensors, boldface
 73 capital letters (\mathbf{A}) or Greece letters (α, β) to denote matrices, boldface lowercase letters (\mathbf{a}) to
 74 denote vectors, and regular letters (a) to denote scalars. For an order three tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$,
 75 $\mathcal{A}_{i,\dots} \in \mathbb{R}^{I_2 \times I_3}$, $\mathcal{A}_{\cdot,j,\dots} \in \mathbb{R}^{I_1 \times I_3}$, and $\mathcal{A}_{\dots,m} \in \mathbb{R}^{I_1 \times I_2}$ are the i -th horizontal slide, j -th lateral slide
 76 and m -th frontal slide of \mathcal{A} , respectively. Similarly, for a matrix \mathbf{A} , $\mathbf{A}_{i,\cdot}$ denotes its i -th row and $\mathbf{A}_{\cdot,j}$
 77 denotes its j -th column. For a vector \mathbf{a} , $\text{diag}(\mathbf{a})$ stands for the diagonal matrix whose diagonal is \mathbf{a} .
 78 We use $\|\cdot\|$, $\|\cdot\|_\infty$, and $\|\cdot\|_F$ to denote the l_2 -norm, l_∞ -norm of a vector, and the Frobenius norm
 79 of matrix or tensor, respectively. For any integer n , denote $[n] = \{1, 2, \dots, n\}$.

80 The mode-1 product between a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ and a matrix $\mathbf{U} \in \mathbb{R}^{J_1 \times I_1}$ is a tensor $\mathcal{A} \times_1 \mathbf{U} \in$
 81 $\mathbb{R}^{J_1 \times I_2 \times I_3}$ such that its (j_1, i_2, i_3) -th entry is defined as $(\mathcal{A} \times_1 \mathbf{U})_{j_1, i_2, i_3} = \sum_{i_1=1}^{I_1} \mathcal{A}_{i_1, i_2, i_3} \mathbf{U}_{j_1, i_1}$.
 82 The mode-2 or mode-3 product between \mathcal{A} and any matrix of appropriate dimension are defined
 83 similarly. The CANDECOMP/PARAFAC (CP) decomposition of \mathcal{A} has the form

$$\mathcal{A} = \sum_{r=1}^R \mathbf{a}^{(r)} \circ \mathbf{b}^{(r)} \circ \mathbf{c}^{(r)}, \quad (1)$$

84 where $\mathbf{a}^{(r)} \in \mathbb{R}^{I_1}$, $\mathbf{b}^{(r)} \in \mathbb{R}^{I_2}$, and $\mathbf{c}^{(r)} \in \mathbb{R}^{I_3}$ for $r \in [R]$, and \circ stands for the vector outer product.
 85 The CP-rank [24] of the tensor $\mathbf{a}^{(r)} \circ \mathbf{b}^{(r)} \circ \mathbf{c}^{(r)}$ is defined to be 1, for $r \in [R]$. The minimal number

86 of rank-1 tensors in the CP decomposition of \mathcal{A} is called the CP-rank of \mathcal{A} . Let $\mathcal{I} \in \{0, 1\}^{R \times R \times R}$
87 be the identity tensor such that $\mathcal{I}_{i_1, i_2, i_3} = 1$ if $i_1 = i_2 = i_3$ and 0 otherwise, and let $\mathbf{A} \in \mathbb{R}^{I_1 \times R}$,
88 $\mathbf{B} \in \mathbb{R}^{I_2 \times R}$, and $\mathbf{C} \in \mathbb{R}^{I_3 \times R}$ such that $\mathbf{A}_{\cdot, r} = \mathbf{a}^{(r)}$, $\mathbf{B}_{\cdot, r} = \mathbf{b}^{(r)}$, and $\mathbf{C}_{\cdot, r} = \mathbf{c}^{(r)}$. Equation (1)
89 then can be equivalently written as $\mathcal{A} = \mathcal{I} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$.

90 2 Structure-preserving embedding

91 In this paper, we consider multi-layer networks that can be represented as an undirected and un-
92 weighted M -layer graph $\mathcal{G} = (V, \mathcal{E})$, where $V = [n]$ consists of the common n vertices across
93 different layers, and $\mathcal{E} = \{E^{(m)}\}_{m=1}^M$ with $E^{(m)} \subset V \times V$ representing the m -th relation network
94 among vertices. A order three adjacency tensor $\mathcal{A} = (a_{i,j,m}) \in \{0, 1\}^{n \times n \times M}$ is then defined to
95 represent \mathcal{G} with entries $a_{i,j,m} = 1$ if $(i, j) \in E^{(m)}$ and 0 otherwise.

96 2.1 Tensor-based latent space model

97 To fully characterize the multi-layer network structure, we propose the following generative tensor-
98 based latent space model (TLSM). For any $i \leq j \in [n]$, and $m \in [M]$,

$$a_{i,j,m} = a_{j,i,m} \stackrel{ind.}{\sim} \text{Bernoulli}(p_{i,j,m}), \text{ with} \quad (2)$$

$$\theta_{i,j,m} = \log \left(\frac{p_{i,j,m}}{s_n - p_{i,j,m}} \right), \text{ and} \quad (3)$$

$$\Theta = \mathcal{I} \times_1 \alpha \times_2 \alpha \times_3 \beta, \quad \alpha \in \Omega_\alpha, \beta \in \Omega_\beta, \quad (4)$$

99 where \mathcal{I} is the order three R -dimensional identity tensor. Basically, (2) follows the standard routine
100 in the multi-layer network literature [34, 35, 27, 22] to model that $a_{i,j,m} = a_{j,i,m}$ are independently
101 generated from a Bernoulli distribution, for $i \leq j \in [n]$ and $m \in [M]$. Denote $\mathcal{P} = (p_{i,j,m}) \in$
102 $\mathbb{R}^{n \times n \times M}$ as the network underlying probability tensor, and then $\Theta = (\theta_{i,j,m}) \in \mathbb{R}^{n \times n \times M}$ is
103 the entry-wise transformation of \mathcal{P} by (3). We call the transformation (3) as the modified logit
104 transformation in that the constant 1 in the standard logit transformation is replaced by a sparsity
105 factor s_n , which may vanish with n and M . We further assume all entries of \mathcal{P} are of the order s_n ; that
106 is, there exists a constant $\frac{1}{2} \leq \xi < 1$ such that $(1 - \xi)s_n \leq p_{i,j,m} \leq \xi s_n$, for $i, j \in [n]$ and $m \in [M]$.
107 Thus, s_n essentially controls the overall network sparsity and the entries of Θ are ensured to locate in
108 the interval $[-\log \frac{\xi}{1-\xi}, \log \frac{\xi}{1-\xi}]$. More importantly, (4) models the CP decomposition of Θ by the
109 factor matrices $\alpha \in \mathbb{R}^{n \times R}$ and $\beta \in \mathbb{R}^{M \times R}$ with CP-rank R , which can greatly reduce the number of
110 free parameters from $n(n+1)M/2$ to $(n+M)R$. Throughout the paper, the CP-rank R is allowed
111 to diverge with n . In the CP decomposition of Θ , α is the vertex latent position matrix with each row
112 $\alpha_{i,\cdot}$ serving as the embedding of vertex i , and β captures heterogeneity across different layers. Herein,
113 we define the constraint sets for α and β as $\Omega_\alpha = \{\alpha \in \mathbb{R}^{n \times R} : \|\alpha_{i,\cdot}\| \leq \sqrt{\log \frac{\xi}{1-\xi}}, \text{ for } i \in [n]\}$
114 and $\Omega_\beta = \{\beta \in \mathbb{R}^{M \times R} : \|\beta_{\cdot, r}\| = 1, r \in [R]\}$. Note that the constraint on β is necessary for
115 model identification, and detailed discussion will be presented shortly. The constraint set $\Omega_\alpha \times \Omega_\beta$
116 is sufficient to maintain the bounded condition of Θ since a general Hölder inequality yields that
117 $|\theta_{i,j,m}| = |\mathcal{I} \times_1 \alpha_{i,\cdot}^T \times_2 \alpha_{j,\cdot}^T \times_3 \beta_{m,\cdot}| \leq \|\alpha_{i,\cdot}\| \|\alpha_{j,\cdot}\| \|\beta_{m,\cdot}\|_\infty \leq \log \frac{\xi}{1-\xi}$. To conclude this
118 paragraph, we remark that the parameter ξ is introduced for theoretical purpose and it is not treated as
119 a tuning parameter. One can choose ξ sufficiently close to 1 in empirical studies so that the restriction
120 on α will be alleviated.

121 We make several essential observations of the proposed TLSM. First and foremost, TLSM is flexible
122 and general. It includes the celebrated MLSBM [34, 43, 35, 27, 26, 36, 22] as special case. Specif-
123 ically, suppose the vertices comes from K disjoint communities, the standard MLSBM assumes
124 that the underlying network probability tensor $\mathcal{P} = \mathcal{B} \times_1 \mathbf{Z} \times_2 \mathbf{Z}$, where $\mathcal{B} \in \mathbb{R}^{K \times K \times M}$ is a
125 semi-symmetric core probability tensor with $\mathcal{B}_{k_1, k_2, m} = \mathcal{B}_{k_2, k_1, m}$ for $k_1, k_2 \in [K]$ and $m \in [M]$,
126 and $\mathbf{Z} \in \{0, 1\}^{n \times K}$ is the community membership matrix with $Z_{i,k} = 1$ if vertex i comes from the
127 k -th community and 0 otherwise. That is, the probability of any vertex pair to form an edge in a
128 particular layer depends only on their community memberships. Equivalently, under the modified
129 logit transformation (3), we have $\Theta = \tilde{\mathcal{B}} \times_1 \mathbf{Z} \times_2 \mathbf{Z}$, where $\tilde{\mathcal{B}}$ is the entry-wise transformation
130 of \mathcal{B} under (3). Taking R to be the CP-rank of $\tilde{\mathcal{B}}$, the CP-decomposition of $\tilde{\mathcal{B}}$ then has the form

131 $\tilde{\mathbf{B}} = \mathcal{I} \times_1 \mathbf{C} \times_2 \mathbf{C} \times_3 \boldsymbol{\beta}$ for some matrix $\mathbf{C} \in \mathbb{R}^{K \times R}$ and $\boldsymbol{\beta} \in \mathbb{R}^{M \times R}$ due to semi-symmetry.
 132 This leads to the CP decomposition of Θ has the form (4) with $\boldsymbol{\alpha} = \mathbf{Z}\mathbf{C}$. It is clear that MLSBM
 133 requires vertices within the same community are homogeneous and exchangeable, while TLSM
 134 allows vertices to have different embeddings even when they are in the same community.

135 Second, TLSM is identifiable when both $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ have full column ranks. When both $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$
 136 have full column ranks, the Kruskal's k-ranks [25] of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ satisfy $k_{\boldsymbol{\alpha}} = k_{\boldsymbol{\beta}} = R$, then Θ has
 137 CP-rank R . Hence, $k_{\boldsymbol{\alpha}} + k_{\boldsymbol{\alpha}} + k_{\boldsymbol{\beta}} \geq 2R + 2$ as long as $R \geq 2$. By Theorem 1 of [40], the fixed
 138 column l_2 -norm constraint of $\boldsymbol{\beta}$ implies that the tensor factorization in (4) is unique up to column
 139 permutations of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ and column sign flip of $\boldsymbol{\alpha}$. It is important to remark that the community
 140 structure encoded in $\boldsymbol{\alpha}$ remains unchanged under any column permutation or sign flip.

141 Third, introducing a sparsity factor s_n via a modified logit transformation into the TLSM is non-
 142 trivial. We take a single-layer network as an example to illustrate the limitation of the standard
 143 logit transformation in handling sparse network. Suppose a vanilla logit link is used to connect
 144 the network underlying probability matrix \mathbf{P} and its transformation Θ , and the latent space model
 145 usually assumes that $\Theta = \boldsymbol{\alpha}\boldsymbol{\alpha}^T$. A sparse network requires the entries of Θ diverge to negative
 146 infinite due to the small magnitude of edge probability, which leads to unstable estimation of $\boldsymbol{\alpha}$ in
 147 numerical experiments. Moreover, this may conflict with the assumption that vertices within the same
 148 community tend to be close in the embedding space and their inner product is likely to be positive.
 149 These difficulties can be naturally circumvented when an appropriate s_n is chosen in (3).

150 2.2 Regularized likelihood

Given a network adjacency tensor \mathcal{A} and number of communities K , our goal is to estimate the
 multi-layer network embedding $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ and conduct community detection on the vertices. Throughout
 this paper, we assume the number of potential communities K is given and may diverge with n . Under
 the TLSM framework, with slight abuse of notation, we denote the average negative log-likelihood
 function of the multi-layer network \mathcal{G} is $\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathcal{A}) = \mathcal{L}(\Theta; \mathcal{A})$ with

$$\mathcal{L}(\Theta; \mathcal{A}) = \frac{1}{\varphi(n, M)} \sum_{m=1}^M \sum_{i \leq j} L(\theta_{i,j,m}; a_{i,j,m}),$$

151 where $\varphi(n, M) = \frac{1}{2}n(n+1)M$ is the number of potential edges, and $L(\theta; a) = \log\left(1 + \frac{s_n}{1-s_n+e^{-\theta}}\right) -$
 152 $a \log\left(\frac{s_n}{1-s_n+e^{-\theta}}\right)$ is a negative log-density of a Bernoulli random variable a . We now introduce a
 153 novel regularization term to detect the potential communities in \mathcal{G} ,

$$J(\boldsymbol{\alpha}) = \min_{\mathbf{Z} \in \Gamma, \mathbf{C} \in \mathbb{R}^{K \times R}} \frac{1}{n} \|\boldsymbol{\alpha} - \mathbf{Z}\mathbf{C}\|_F^2, \quad (5)$$

154 where \mathbf{C} encodes the vertex embedding centers and $\Gamma \subset \{0, 1\}^{n \times K}$ is the set of all possible
 155 community membership matrices; that is, for any $\mathbf{Z} \in \Gamma$, each row of \mathbf{Z} consists of only one 1
 156 indicating the community membership and all others entries being 0. This leads to the proposed
 157 regularized cost function,

$$\mathcal{L}_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathcal{A}) = \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathcal{A}) + \lambda_n J(\boldsymbol{\alpha}), \quad (6)$$

158 where λ_n is a positive tuning parameter that strikes the balance between network estimation and
 159 community detection in the cost function. It is clear that the embeddings of vertices with similar
 160 linking pattern will be pushed towards the same center, and thus close to each other in the ambient
 161 space, leading to the desired community structure in \mathcal{G} .

162 2.3 Projected gradient descent algorithm

163 We develop a scalable projected gradient descent (PGD) algorithm to optimize the penalized cost
 164 function (6), which is highly non-convex and can be solved only locally. PGD, which alternatively
 165 conducts gradient step and projection step, is one of the most popular and computationally fast
 166 algorithm in tackling non-convex optimization problem [7, 33, 47, 9].

167 To compute the gradients of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, we introduce the following notations. Define $\mathcal{T} \in \mathbb{R}^{n \times n \times M}$
 168 with entries $\mathcal{T}_{i,j,m} = \frac{\exp(-\theta_{i,j,m})}{1-s_n+\exp(-\theta_{i,j,m})} (p_{i,j,m} - a_{i,j,m})$, and $\mathbf{X}_{\mathcal{T}(2,3)}^{\boldsymbol{\alpha}, \boldsymbol{\beta}} \in \mathbb{R}^{n \times R}$ whose i -th row

169 consists of the diagonal elements of the slice $(\mathcal{T} \times_2 \alpha^T \times_3 \beta^T)_{i,\dots}$. That is, $\mathbf{X}_{\mathcal{T}(2,3)}^{\alpha,\beta}(i,r) =$
 170 $(\mathcal{T} \times_2 \alpha^T \times_3 \beta^T)_{i,r,r}$. Similarly, we define $\mathbf{X}_{\mathcal{T}(1,2)}^{\alpha,\alpha} \in \mathbb{R}^{R \times M}$, $\mathbf{X}_{\mathcal{T}(3)}^{\beta} \in \mathbb{R}^{n \times R}$, and $\mathbf{X}_{\mathcal{T}(1,2)} \in$
 171 $\mathbb{R}^{n \times M}$, such that $\mathbf{X}_{\mathcal{T}(1,2)}^{\alpha,\alpha}(r,m) = (\mathcal{T} \times_1 \alpha^T \times_2 \alpha^T)_{r,r,m}$, $\mathbf{X}_{\mathcal{T}(3)}^{\beta}(i,r) = (\mathcal{T} \times_3 \beta^T)_{i,i,r}$, and
 172 $\mathbf{X}_{\mathcal{T}(1,2)}(i,m) = \mathcal{T}_{i,i,m}$. Consequently, when the vertex membership matrix \mathbf{Z} and the community
 173 center matrix \mathbf{C} are fixed, we can derive the gradients of $\mathcal{L}_\lambda(\alpha, \beta; \mathcal{A})$ with respect to α and β , as

$$\frac{1}{\varphi(n, M)} (\mathbf{X}_{\mathcal{T}(2,3)}^{\alpha,\beta} + \mathbf{X}_{\mathcal{T}(3)}^{\beta} * \alpha) + 2\lambda_n(\alpha - \mathbf{Z}\mathbf{C}) \text{ and } \frac{1}{2\varphi(n, M)} ((\mathbf{X}_{\mathcal{T}(1,2)}^{\alpha,\alpha})^T + \mathbf{X}_{\mathcal{T}(1,2)}^T(\alpha * \alpha)),$$

174 respectively. Herein, $*$ denotes the Hadamard product (entry-wise product) between two matrices.

175 Let $(\tilde{\alpha}, \tilde{\beta})$ denote the solution given by one-step gradient descent, we then project $(\tilde{\alpha}, \tilde{\beta})$ onto
 176 $\Omega_\alpha \times \Omega_\beta$ in the following steps.

177 *Step 1.* Multiply the r -th column of $\tilde{\alpha}_{\cdot,r}$ by $\|\tilde{\beta}_{\cdot,r}\|^{1/2}$ for $r \in [R]$. Denote the resultant matrix as $\tilde{\alpha}'$.

178 *Step 2.* Regularize each row of α as $\alpha_{i,\cdot} = \tilde{\alpha}'_{i,\cdot} \cdot \min\{\sqrt{\log \frac{\xi}{1-\xi}}, \|\tilde{\alpha}'_{i,\cdot}\|\} / \|\tilde{\alpha}'_{i,\cdot}\|$, for $i \in [n]$.

179 *Step 3.* Normalize the columns of β as $\beta_{\cdot,r} = \tilde{\beta}_{\cdot,r} / \|\tilde{\beta}_{\cdot,r}\|$, for $r \in [R]$.

180 Next, when (α, β) are given, we apply a $(1 + \delta)$ -approximation K-means algorithm on $\tilde{\alpha}$ to update
 181 the vertex community membership matrix \mathbf{Z} and community center matrix \mathbf{C} .

182 The above steps will be alternatively conducted until convergence or reaching the maximum number
 183 of iterations. We further summarized the developed alternative updated scheme in Algorithm 1 in
 184 Appendix A of the supplementary materials

185 Several remarks on the algorithm are in order. First, Algorithm 1 can only be guaranteed to converge
 186 to a stationary point but not any local minimizer. We hence employ a transformed higher order
 187 orthogonal iteration (HOOI) algorithm for warm initialization in all the numerical experiments in
 188 Section 4 and 5. Specifically, given a user-specific value τ , we define $\tilde{\Theta}$ to mimic the magnitude
 189 of Θ such that $\tilde{\Theta}_{i,j,m} = -\tau$ if $a_{i,j,m} = 0$ and $\tilde{\Theta}_{i,j,m} = \tau$ otherwise. A standard HOOI algorithm
 190 [11] is applied to $\tilde{\Theta}$ to obtain $\alpha^{(0)}$ and $\beta^{(0)}$. We set $\tau = 100$ in all the numerical experiments.
 191 Second, the sparsity factor s_n is an intrinsic quantity of the multi-layer network data, and it should be
 192 estimated from the network directly. Note that the minimal and maximal probabilities for any vertex
 193 pair to form an edge in any layer are $p_{\min} = (1 - \xi)s_n$ and $p_{\max} = \xi s_n$, respectively. Interestingly,
 194 $p_{\min} + p_{\max} = s_n$, which does not depend on ξ any more. Therefore, we propose to estimate s_n as

$$\hat{s}_n = \min_{i \in [n]} \frac{1}{nM} \sum_{m=1}^M \sum_{j=1}^n a_{i,j,m} + \max_{i \in [n]} \frac{1}{nM} \sum_{m=1}^M \sum_{j=1}^n a_{i,j,m}, \quad (7)$$

195 which is the sum of the minimal and maximal frequencies of a vertex to form edges with all other
 196 vertices in all layers. Third, to optimally choose λ_n , we extend the network cross-validation by
 197 edge sampling scheme in [30] to multi-layer networks. The detailed tuning procedure is relegated to
 198 Appendix B in the supplementary materials.

199 3 Asymptotic theory

200 3.1 Consistency in estimating Θ^*

201 Let $\Omega = \{\Theta = \mathcal{I} \times_1 \alpha \times_2 \alpha \times_3 \beta : \alpha \in \Omega_\alpha, \beta \in \Omega_\beta\}$ be the parameter space of the problem and
 202 $\Theta^* = \mathcal{I} \times_1 \alpha^* \times_2 \alpha^* \times_3 \beta^*$ be the true underlying transformed network probability tensor. Denote
 203 $KL(\Theta^* || \Theta) = \varphi^{-1}(n, M) \sum_{m=1}^M \sum_{i \leq j} E(L(\theta_{i,j,m}; a_{i,j,m}) - L(\theta_{i,j,m}^*; a_{i,j,m}))$ be the averaged
 204 Kullback–Leibler divergence of the network generation distributions parametrized by Θ^* and Θ , for
 205 any $\Theta \in \Omega$. The following large deviation inequality is derived to quantify the behavior of $\mathcal{L}_\lambda(\Theta; \mathcal{A})$
 206 for any Θ in the neighborhood of Θ^* defined by $KL(\Theta^* || \Theta)$.

207 **Proposition 1.** *Suppose $\lambda_n J(\alpha^*) \leq \epsilon_n$, and $(n + M)R\varphi^{-1}(n, M)\epsilon_n^{-1} \log(\epsilon_n^{-1/2}) \leq c_1$ for some
 208 constant c_1 . Then with probability at least $1 - 2 \exp\left(-\frac{\varphi(n, M)\epsilon_n}{156 \frac{\xi}{1-\xi} + 28 \log 2}\right)$, we have*

$$\mathcal{L}_\lambda(\Theta^*; \mathcal{A}) \leq \inf_{\{\Theta \in \Omega | KL(\Theta^* || \Theta) \geq 4\epsilon_n\}} \mathcal{L}_\lambda(\Theta; \mathcal{A}) - \epsilon_n.$$

209 Proposition 1 basically states that any estimators with sufficiently small objective value should
 210 be close enough to Θ^* in terms of $KL(\Theta^*||\Theta)$. We next study the asymptotic behavior of these
 211 estimators more precisely. Let $(\hat{\alpha}, \hat{\beta}) \in \Omega_\alpha \times \Omega_\beta$ be any estimator of (α^*, β^*) such that

$$\mathcal{L}_\lambda(\hat{\alpha}, \hat{\beta}; \mathcal{A}) \leq \mathcal{L}_\lambda(\alpha^*, \beta^*; \mathcal{A}) + \epsilon_n, \quad (8)$$

212 and denote $\hat{\Theta} = \mathcal{I} \times_1 \hat{\alpha} \times_2 \hat{\alpha} \times_3 \hat{\beta}$. we have the following theorem.

Theorem 1. *Under the condition of Proposition 1, if $(\hat{\alpha}, \hat{\beta})$ satisfies (8), then with probability at least $1 - 2 \exp\left(-\frac{\varphi(n, M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28 \log 2}\right)$, we have*

$$\frac{1}{n\sqrt{M}} \|\hat{\Theta} - \Theta^*\|_F \leq \frac{4\sqrt{2}\sqrt{\epsilon_n}}{(1-\xi)\sqrt{\xi s_n}}.$$

213 The condition that $\lambda_n J(\Theta^*) \leq \epsilon_n$ in Proposition 1 is mild. It implies that the true em-
 214 beddings of vertices within the same community are close to one another. We remark that
 215 $\lambda_n J(\Theta^*)$ exactly equals to zero under the MLSBM discussed in Section 2.2. The condition that
 216 $(n+M)R\varphi^{-1}(n, M)\epsilon_n^{-1} \log(\epsilon_n^{-1/2})$ vanishes with n is also mild. When $R = O(1)$, we can take any
 217 ϵ_n such that $\epsilon_n \gg \frac{\log n}{n \min\{n, M\}}$. Consequently, to ensure $\hat{\Theta}$ converges to Θ^* , Theorem 1 implies the
 218 smallest sparsity factor one can take is $s_n \gg \epsilon_n \gg \frac{\log n}{n \min\{n, M\}}$, which means that the average degree
 219 of a vertex in any particular layer can be as small as ns_n . We remark that a common assumption
 220 $M = O(n)$ that appears in literature, such as [27] and [22], is not necessary in our theory. If we
 221 further assume $M = O(n)$, we find that the average degree of a vertex in any layer under the
 222 proposed TLSM set up can be smaller than that in [27] by a factor $(M \log n)^{-1/2}$ and in [22] by a
 223 factor $(\log n)^{-3}$, showing that our theoretical result accommodates sparser multi-layer networks.

224 3.2 Consistency in community detection

225 We now turn to establish the consistency of community detection in multi-layer network
 226 \mathcal{G} . Let $\psi^* : [n] \rightarrow [K]$ be the true community assignment function such that $\psi^* =$
 227 $\arg \min_{\psi} \min_{C_1, \dots, C_K} \sum_{i=1}^n \|\alpha_i^* - C_{\psi_i}\|^2$, and then the community detection error of any esti-
 228 mated community assignment function $\hat{\psi}$ can be evaluated by the minimum scaled Hamming distance
 229 between $\hat{\psi}$ and ψ^* under permutations, which is defined as

$$\text{err}(\psi^*, \hat{\psi}) = \min_{\pi \in S_K} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{\psi_i^* \neq \pi(\hat{\psi}_i)\}, \quad (9)$$

230 where $\mathbf{1}\{\cdot\}$ is the indicator function and S_K is the symmetric group of degree K . Such a scaled
 231 or unscaled Hamming distance has become a popular metric in quantifying the performance of
 232 community detection [21, 22].

233 Denote $N_k^* = \{i : \psi_i^* = k\}$ be the k -th true underlying community whose cardinality is n_k . Let
 234 $\mathbf{C}^* \in \mathbb{R}^{K \times R}$ be the true underlying community centers of the network embedding with $\mathbf{C}_{k^*}^* =$
 235 $\frac{1}{n_k} \sum_{\psi_i^* = k} \alpha_i^*$, and let $\mathbf{B}^* = \mathcal{I} \times_1 \mathbf{C}^* \times_2 \mathbf{C}^* \times_3 \beta^*$. The following assumptions are made to ensure
 236 that communities within the multi-layer networks are asymptotically identifiable.

237 **Assumption A.** *Assume the difference between any two distinct horizontal slides of \mathbf{B}^* satisfies that*

$$\min_{k, k' \in [K], k \neq k'} \frac{1}{\sqrt{KM}} \|\mathbf{B}_{k, \dots}^* - \mathbf{B}_{k', \dots}^*\|_F \geq \gamma_n,$$

238 where $\gamma_n > 0$ may vanish with n .

Assumption B. *Assume the tuning parameter λ_n satisfies that*

$$\lambda_n \epsilon_n s_n^{-2} (\log s_n^{-1})^{-1} \geq c_2,$$

239 for an absolute constant c_2 that does not depend on any model parameter.

Assumption C. *Denote $n_{\min} = \min_{k \in [K]} n_k$ as the minimal community size. Assume*

$$\frac{\gamma_n n_{\min} \sqrt{K}}{n} \geq c_\xi \sqrt{\frac{\epsilon_n}{s_n}},$$

240 where $c_\xi = \frac{4\sqrt{2}}{(1-\xi)\sqrt{\xi}} + c_3 \sqrt{\frac{(1+\delta) \min\{M, R\}}{M}}$ and c_3 is a constant that depends on ξ only.

241 Assumption A is the minimal community separation requirement, and similar assumption has been
 242 employed in [27] with a constant γ_n . Together with the condition $\lambda_n J(\boldsymbol{\alpha}^*) \leq \epsilon_n$ in Proposition 1,
 243 Assumption B gives a feasible interval for λ_n . Assumption C allows for unbalanced communities
 244 with vanishing n_{\min}/n if the network is not too sparse. Note that c_ξ can be further bounded by
 245 $\frac{4\sqrt{2}}{(1-\xi)\sqrt{\xi}} + c_3\sqrt{1+\delta}$, and the first term of c_ξ will dominate the second term if $R = o(M)$.

Theorem 2. *Suppose all the assumptions in Theorem 1 as well as Assumptions A, B and C are satisfied, it holds true that*

$$\text{err}(\psi^*, \hat{\psi}) \leq \frac{c_\xi^2 n \epsilon_n}{n_{\min} K \gamma_n^2 s_n},$$

246 *with probability at least $1 - \frac{1}{n^2} - 2 \exp\left(-\frac{\varphi(n, M) \epsilon_n}{156 \frac{\xi}{1-\xi} + 28 \log 2}\right)$.*

247 Theorem 2 assures that the community structure in a multi-layer network can be consistently recovered
 248 by the proposed TLSM. As a theoretical example, we consider a sparse case with $s_n = \frac{(\log n)^{1+\tau_1}}{n \min\{n, M\}}$,
 249 where $0 < \tau_1 < 1$, $n_{\max} = O(n_{\min})$, $\frac{1}{\sqrt{n}} \|\boldsymbol{\alpha}^* - \mathbf{Z}^* \mathbf{C}^*\|_F \leq (\log n)^{-3/2}$, and both γ_n , R and K
 250 are of constant orders. With $\lambda_n = \frac{(\log n)^{2+2\tau_1}}{n \min\{n, M\}}$, Theorems 1 and 2 imply that $\epsilon_n = \frac{(\log n)^{1+\tau_2}}{n \min\{n, M\}}$ with
 251 $0 < \tau_2 < \tau_1$ and $\text{err}(\psi^*, \hat{\psi}) = o_p(1)$.

252 4 Numerical experiments

253 In this section, we evaluate the numerical performance of the proposed TLSM in a variety of synthetic
 254 as well as real-life multi-layer networks, compare it against four competitors in literature, including
 255 the mean adjacency spectral embeddings (MASE; 16), least square estimation (LSE; 27), Tucker
 256 decomposition with HOSVD initialization (HOSVD-Tucker; 22), and spectral kernel (SPECK; 35),
 257 and conduct some ablation studies. The implementations of LSE and SPECK are available at the
 258 authors' personal websites, HOSVD-Tucker is implemented in the routine "tucker" of the Python
 259 package "tensorly", and TLSM and MASE are implemented in Python by ourselves.

260 4.1 Synthetic networks

261 The multi-layer network $\mathcal{A} = (a_{i,j,m}) \in \{0, 1\}^{n \times n \times M}$ is generated as follows. First, we randomly
 262 select $K = 4$ elements uniformly from $\{2.5 * (b_1, b_2, \dots, b_R) : b_r \in \{-1, 1\}, r \in [R]\}$ as community
 263 centers, which are denoted as \mathbf{c}_k , $k \in [K]$. Second, the latent space embedding of vertex i is
 264 generated as $\boldsymbol{\alpha}_i = \mathbf{c}_{\psi_i} + \mathbf{e}_i$ with $\mathbf{e}_i \sim N(\mathbf{0}_R, 1.5 * I_R)$, and $\psi_i \in [K]$ are independently drawn
 265 from the multinomial distribution $\text{Multi}(1; \frac{1}{K} \mathbf{1}_K)$. Third, we generate $\boldsymbol{\beta} = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_M]^T$ with
 266 $\boldsymbol{\beta}_{m,r}$ being independent standard normal random variables, for $m \in [M]$ and $r \in [R]$. We then
 267 rescale the column norms of $\boldsymbol{\beta}$ to be 1 for model identifiability. Finally, we generate \mathcal{A} according
 268 to the proposed TLSM with $s_n = 0.1$. For the sake of fair comparisons, the embedding dimension
 269 R is set as K in all scenarios. We aim to illustrate the community detection performance of
 270 all methods as the number of vertices and number of layers increase. To this end, we consider
 271 $(n, M) \in \{200, 400, 600, 800\} \times \{5, 10, 15, 20\}$. The averaged hamming errors and their standard
 272 errors over 50 independent experiments of all methods are reported in Table 1.

273 It is evident that TLSM consistently outperforms its competitors, and the performances of LSE
 274 and HOSVD-Tucker are better than those of MASE and SPECK. This is expected since TLSM,
 275 LSE and HOSVD-Tucker work on the multi-layer network adjacency tensor directly, while MASE
 276 and SPECK are matrix aggregation methods that suffer from information loss. Furthermore, as the
 277 number of vertices and number of layers increase, the community detection errors of all methods
 278 decrease rapidly. Notably, TLSM and LSE converge faster than the other methods, and attain stable
 279 performance even for relatively small n and M . Additional simulation studies for various network
 280 sparsity and unbalanced community sizes are relegated to Appendix C in the supplementary materials.

281 4.2 Real-life networks

282 We also apply the proposed TLSM method to analyze three real-life multi-layer networks, including
 283 a social network in the department of Computer Science at Aarhus University (AUUS) [38], a yeast

Table 1: The averaged hamming errors of various methods with their standard errors in Scenario I. The best performer in each case is bold-faced.

n	M	TLSM	LSE	MASE	HOSVD-Tucker	SPECK
200	5	0.1180 (0.0147)	0.1405(0.0118)	0.5086(0.0136)	0.1623(0.0126)	0.4254(0.0138)
	10	0.0585 (0.0046)	0.0751(0.0050)	0.4949(0.0131)	0.1148(0.0106)	0.2996(0.0141)
	15	0.0551 (0.0067)	0.0593(0.0045)	0.4910(0.0176)	0.1040(0.0115)	0.2505(0.0142)
	20	0.0510 (0.0037)	0.0588(0.0043)	0.4977(0.0161)	0.1023(0.0110)	0.1942(0.0156)
400	5	0.0653 (0.0066)	0.1019(0.0087)	0.3845(0.0193)	0.1220(0.0106)	0.3766(0.0195)
	10	0.0608 (0.0063)	0.0636(0.0037)	0.3859(0.0160)	0.1012(0.0092)	0.2244(0.0191)
	15	0.0511 (0.0031)	0.0595(0.0036)	0.3844(0.0221)	0.0787(0.0051)	0.1490(0.0123)
	20	0.0536 (0.0047)	0.0551(0.0036)	0.3985(0.0185)	0.0795(0.0063)	0.1409(0.0131)
600	5	0.0607 (0.0029)	0.0909(0.0040)	0.3665(0.0186)	0.1221(0.0108)	0.3038(0.0193)
	10	0.0567 (0.0029)	0.0688(0.0031)	0.3726(0.0179)	0.1003(0.0081)	0.1651(0.0127)
	15	0.0558 (0.0027)	0.0630(0.0030)	0.3803(0.0167)	0.0918(0.0076)	0.1231(0.0076)
	20	0.0548 (0.0028)	0.0586(0.0029)	0.3814(0.0185)	0.0883(0.0078)	0.1150(0.0088)
800	5	0.0556 (0.0056)	0.0768(0.0055)	0.3012(0.0194)	0.1003(0.0103)	0.2733(0.0171)
	10	0.0560 (0.0063)	0.0583(0.0034)	0.3004(0.0177)	0.0788(0.0065)	0.1424(0.0127)
	15	0.0498 (0.0030)	0.0539(0.0033)	0.3179(0.0195)	0.0812(0.0068)	0.1146(0.0098)
	20	0.0485 (0.0031)	0.0516(0.0032)	0.3184(0.0218)	0.0803(0.0075)	0.0979(0.0078)

284 Saccharomyces cerevisiae gene co-expression (YSCGC) network [44], and a worldwide agriculture
 285 trading network (WAT) [10]. Specifically, we conduct community detection on the first two networks
 286 whose vertex community memberships are available, and carry out a link prediction task on the third
 287 network whose vertex community memberships are unavailable.

288 The AUCS dataset is publicly available at <http://multilayer.it.uu.se/datasets.html>, and
 289 it is a $61 \times 61 \times 5$ multi-layer network that records pairwise relationships of 5 types among 61
 290 persons in AUCS, including current working relationships, repeated leisure activities, regularly eating
 291 lunch together, co-authorship of a publication, and friendship on Facebook. Since 54 persons in
 292 the dataset come from 7 research groups and the other 7 persons do not belong to any group, the
 293 dataset consists of 8 communities corresponding to 7 research groups and an outlier community.
 294 Applying TLSM and its competitors to the dataset, the number of misclassified vertices by TLSM,
 295 LSE, MASE, HOSVD-Tucker and SPECK, are 8, 21, 19, 23, 18, respectively. Clearly, TLSM
 296 significantly outperforms its competitors by at least reducing 16.39% of community detection error.

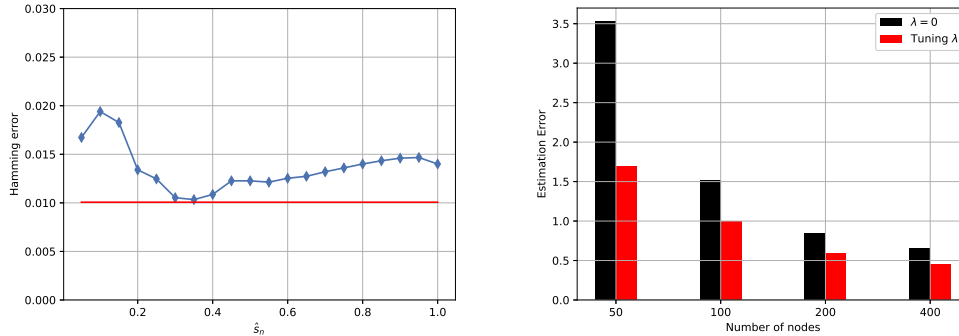
297 The YSCGC dataset is publicly available at [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC156590/)
 298 [PMC156590/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC156590/), and contains 205 genes of 4 functional categories, including protein metabolism
 299 and modification, carbohydrate metabolism and catabolism, nucleobase, nucleoside, nucleotide
 300 and nucleic acid metabolism, as well as transportation. We regard these four functional category
 301 labels as the community memberships of the genes. Further, the gene expression responses are
 302 measured by 20 systematic perturbations with varying genetic and environmental conditions in
 303 4 replicated hybridizations. We thus constructed a gene co-expression network $\mathcal{A} = (a_{i,j,m}) \in$
 304 $\mathbb{R}^{205 \times 205 \times 4}$ based on the similarities of their expressions, where each layer represents one replicated
 305 hybridization. Specifically, the similarity between genes i and j in the m -th replication is measured
 306 by $w_{i,j,m} = \exp(-\|\mathbf{x}_i^{(m)} - \mathbf{x}_j^{(m)}\|)$, where $\mathbf{x}_i^{(m)} \in \mathbb{R}^{20}$ contains the expression levels of 20
 307 perturbations in the m -th replicated hybridization for $i \in [205]$ and $m \in [4]$. The binary value $a_{i,j,m}$
 308 is obtained by thresholding $w_{i,j,m}$ with the thresholding value being the 60% quantile of all elements
 309 in $\{w_{i,j,m} : i \leq j \in [205], m \in [4]\}$. Applying TLSM and its competitors to this dataset, the number
 310 of misclassified vertices by TLSM, LSE, MASE, HOSVD-Tucker and SPECK, are 6, 9, 12, 48, 13,
 311 respectively. TLSM again outperforms its competitors in this YSCGC dataset.

312 The WAT dataset is publicly available at <http://www.fao.org>, and includes 364 agriculture
 313 product trading relationships among 214 countries in 2010. To process the data, we extract 130 major
 314 countries whose average degrees are greater than 9 from the 32 densest connected agriculture product
 315 trading relations, leading to a $130 \times 130 \times 32$ multi-layer network. Investigating the eigen-structure
 316 of the mode-1 matricization of the network adjacency tensor, we identify an elbow point [20] at the
 317 7th largest eigen-value, suggesting there are 6 potential communities among the countries, and thus
 318 we set $K = 6$. The corresponding eigen-value plot is attached in Appendix D of the supplementary
 319 materials. We then randomly selected 80% of the entries of the adjacency tensor as the training set,
 320 and conduct link prediction on the remaining 20% of the entries. Specifically, we employ TLSM

321 and the adaptations of its competitors to estimate the network expected tensor \mathcal{P} and generate
 322 estimations for the missing entries by independent Bernoulli random variables accordingly. The
 323 averaged link prediction accuracy of TLSM, LSE, MASE, HOSVD-Tucker and SPECK over 50
 324 independent replications are 79.60%, 76.66%, 75.96%, 77.78% and 79.08%, respectively, where the
 325 link prediction accuracy is defined as the percentile of the correctly predicted entries. Clearly, all 5
 326 methods are comparative in terms of link prediction, while TLSM still deliver highest averaged link
 327 prediction accuracy.

328 4.3 Ablation studies

329 In this subsection, we carry out some ablation studies on two novel components of the proposed
 330 method, namely the sparsity factor s_n and the community-inducing regularizer $J(\alpha)$. To study the
 331 effectiveness of s_n , we generate a $300 \times 300 \times 5$ multi-layer network with 3 communities and the
 332 true network sparsity $s_n = 0.3$. The blue curve in the left panel of Figure 1 shows the average
 333 Hamming error of 50 independent replications given by the proposed method when employing
 334 $\hat{s}_n \in \{0.05i : i \in [20]\}$ in the optimization algorithm, and the red line indicates the averaged
 335 Hamming error of the proposed method with \hat{s}_n estimated via the proposed data-adapted estimation
 336 scheme. It is clear that the Hamming error at $s_n = 1$ is much larger than that when s_n is close
 337 to 0.3, showing the advantages of the modified logit transformation by s_n over the standard logit
 338 transformation when the network indeed reveals sparse pattern. Moreover, we observe that the red
 339 line is even lower than the minimum Hamming error in the blue curve. This further confirms the
 effectiveness of the proposed data-adapted estimation scheme for estimating s_n .



340 Figure 1: Ablation studies on s_n (left) and community-inducing regularizer (right).

341 To study the effectiveness of the community-inducing regularizer in the proposed objective function,
 342 we generate an $n \times n \times 5$ multi-layer network with 2 communities, for $n \in \{50, 100, 200, 400\}$. In
 343 the right panel of Figure 1, the black pillars indicate the network estimation error $\frac{1}{n\sqrt{5}} \|\hat{\Theta} - \Theta^*\|_F$
 344 given by the proposed method with $\lambda_n = 0$ which corresponds to the absence of $J(\alpha)$, while the
 345 red ones indicate the counterparts given by the proposed method with λ_n is selected by network
 346 cross-validation. There is a clear improvement when the community-inducing regularizer is enforced
 347 in all scenarios, particularly for small n . This showcases the helpfulness of the community-inducing
 348 regularizer in detecting network community structure.

349 5 Conclusions

350 In this paper, we propose a novel tensor-based latent space model for community detection in
 351 multi-layer networks. The model embeds vertices into a low-dimensional latent space and views
 352 the community structure from an network embedding perspective, so that heterogeneous structures
 353 in different network layers can be properly integrated. The proposed model is formulated as a
 354 regularization framework, which conducts multi-layer network estimation and community detection
 355 simultaneously. The advantages of the proposed method are supported by extensive numerical
 356 experiments and theoretical results. Particularly, the asymptotic consistencies of the proposed method
 357 are established in terms of both multi-layer network estimation and community detection, even for
 358 relatively sparse networks.

References

- [1] Luiz GA Alves, Giuseppe Mangioni, Isabella Cingolani, Francisco Aparecido Rodrigues, Pietro Panzarasa, and Yamir Moreno. The nested structural organization of the worldwide trade multi-layer network. *Scientific reports*, 9(1):1–14, 2019.
- [2] Jesús Arroyo, Avanti Athreya, Joshua Cape, Guodong Chen, Carey E Priebe, and Joshua T Vogelstein. Inference for multiple heterogeneous networks with a common invariant subspace. *Journal of Machine Learning Research*, 22(142):1–49, 2021.
- [3] Avanti Athreya, Donniell E Fishkind, Minh Tang, Carey E Priebe, Youngser Park, Joshua T Vogelstein, Keith Levin, Vince Lyzinski, and Yichen Qin. Statistical inference on random dot product graphs: a survey. *The Journal of Machine Learning Research*, 18(1):8393–8484, 2017.
- [4] Matteo Barigozzi, Giorgio Fagiolo, and Giuseppe Mangioni. Identifying the community structure of the international-trade multi-network. *Physica A: statistical mechanics and its applications*, 390(11):2051–2066, 2011.
- [5] Michele Berlingerio, Fabio Pinelli, and Francesco Calabrese. Abacus: frequent pattern mining-based community discovery in multidimensional networks. *Data Mining and Knowledge Discovery*, 27(3):294–320, 2013.
- [6] Sharmodeep Bhattacharyya and Shirshendu Chatterjee. Spectral clustering for multiple sparse networks: I. *arXiv preprint arXiv:1805.10594*, 2018.
- [7] Han Chen, Garvesh Raskutti, and Ming Yuan. Non-convex projected gradient descent for generalized low-rank tensor regression. *Journal of Machine Learning Research*, 20:1–37, 2019.
- [8] Zitai Chen, Chuan Chen, Zibin Zheng, and Yi Zhu. Tensor decomposition for multilayer networks clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3371–3378, 2019.
- [9] Eric C Chi, Brian R Gaines, Will Wei Sun, Hua Zhou, and Jian Yang. Provable convex co-clustering of tensors. *Journal of Machine Learning Research*, 21(214):1–58, 2020.
- [10] Manlio De Domenico, Vincenzo Nicosia, Alexandre Arenas, and Vito Latora. Structural reducibility of multilayer networks. *Nature communications*, 6(1):1–9, 2015.
- [11] Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors. *SIAM journal on Matrix Analysis and Applications*, 21(4):1324–1342, 2000.
- [12] Xiaowen Dong, Pascal Frossard, Pierre Vandergheynst, and Nikolai Nefedov. Clustering with multi-layer graphs: A spectral perspective. *IEEE Transactions on Signal Processing*, 60(11):5820–5831, 2012.
- [13] Junxian Geng, Anirban Bhattacharya, and Debdeep Pati. Probabilistic community detection with unknown number of communities. *Journal of the American Statistical Association*, 114(526):893–905, 2019.
- [14] Mahsa Ghorbani, Mahdich Soleymani Baghshah, and Hamid R Rabiee. Mgcn: semi-supervised classification in multi-layer graphs with graph convolutional networks. In *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 208–211, 2019.
- [15] Derek Greene and Pádraig Cunningham. Producing a unified graph representation from multiple social network views. In *Proceedings of the 5th annual ACM web science conference*, pages 118–121, 2013.
- [16] Qiuyi Han, Kevin Xu, and Edoardo Airoldi. Consistent estimation of dynamic and multi-layer block models. In *International Conference on Machine Learning*, pages 1511–1520. PMLR, 2015.
- [17] Xin He, Qiong Liu, and You Yang. Mv-gnn: Multi-view graph neural network for compression artifacts reduction. *IEEE Transactions on Image Processing*, 29:6829–6840, 2020.

- 407 [18] Peter D Hoff, Adrian E Raftery, and Mark S Handcock. Latent space approaches to social
408 network analysis. *Journal of the American Statistical Association*, 97(460):1090–1098, 2002.
- 409 [19] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels:
410 First steps. *Social networks*, 5(2):109–137, 1983.
- 411 [20] Pengsheng Ji and Jiashun Jin. Coauthorship and citation networks for statisticians. *The Annals
412 of Applied Statistics*, 10(4):1779–1812, 2016.
- 413 [21] Jiashun Jin. Fast community detection by score. *Ann. Statist.*, 43(1):57–89, 02 2015.
- 414 [22] Bing-Yi Jing, Ting Li, Zhongyuan Lyu, and Dong Xia. Community detection on mixture
415 multilayer networks via regularized tensor decomposition. *The Annals of Statistics*, 49(6):3181–
416 3205, 2021.
- 417 [23] Muhammad Raza Khan and Joshua E Blumenstock. Multi-gcn: Graph convolutional networks
418 for multi-view networks, with applications to global poverty. In *Proceedings of the AAAI
419 Conference on Artificial Intelligence*, volume 33, pages 606–613, 2019.
- 420 [24] Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. *SIAM Review*,
421 51:455–500, 2009.
- 422 [25] Joseph B Kruskal. Three-way arrays: rank and uniqueness of trilinear decompositions, with
423 application to arithmetic complexity and statistics. *Linear algebra and its applications*, 18(2):95–
424 138, 1977.
- 425 [26] Jing Lei. Tail bounds for matrix quadratic forms and bias adjusted spectral clustering in
426 multi-layer stochastic block models. *arXiv preprint arXiv:2003.08222*, 2020.
- 427 [27] Jing Lei, Kehui Chen, and Brian Lynch. Consistent community detection in multi-layer network
428 data. *Biometrika*, 107(1):61–73, 2020.
- 429 [28] Jing Lei and Alessandro Rinaldo. Consistency of spectral clustering in stochastic block models.
430 *The Annals of Statistics*, 43(1):215–237, 2015.
- 431 [29] Dong Li, Zhisong Pan, Guyu Hu, Graham Anderson, and Shan He. Active module identification
432 from multilayer weighted gene co-expression networks: a continuous optimization approach.
433 *IEEE/ACM transactions on computational biology and bioinformatics*, 2020.
- 434 [30] Tianxi Li, Elizaveta Levina, and Ji Zhu. Network cross-validation by edge sampling. *Biometrika*,
435 107(2):257–276, 2020.
- 436 [31] Xueming Liu, Enrico Maiorino, Arda Halu, Kimberly Glass, Rashmi B Prasad, Joseph Loscalzo,
437 Jianxi Gao, and Amitabh Sharma. Robustness and lethality in multilayer biological molecular
438 networks. *Nature communications*, 11(1):1–12, 2020.
- 439 [32] Zhongyuan Lyu, Dong Xia, and Yuan Zhang. Latent space model for higher-order networks
440 and generalized tensor decomposition. *arXiv preprint arXiv:2106.16042*, 2021.
- 441 [33] Zhuang Ma, Zongming Ma, and Hongsong Yuan. Universal latent space model fitting for large
442 networks with edge covariates. *Journal of Machine Learning Research*, 21(4):1–67, 2020.
- 443 [34] Subhadeep Paul and Yuguo Chen. Consistent community detection in multi-relational
444 data through restricted multi-layer stochastic blockmodel. *Electronic Journal of Statistics*,
445 10(2):3807–3870, 2016.
- 446 [35] Subhadeep Paul and Yuguo Chen. Spectral and matrix factorization methods for consistent
447 community detection in multi-layer networks. *Ann. Statist.*, 48(1):230–250, 02 2020.
- 448 [36] Subhadeep Paul and Yuguo Chen. Null models and community detection in multi-layer networks.
449 *Sankhya A*, pages 1–55, 2021.
- 450 [37] Zhuo-Ming Ren, An Zeng, and Yi-Cheng Zhang. Bridging nestedness and economic complexity
451 in multilayer world trade networks. *Humanities and Social Sciences Communications*, 7(1):1–8,
452 2020.

- 453 [38] Luca Rossi and Matteo Magnani. Towards effective visual analytics on multiplex and multilayer
454 networks. *Chaos, Solitons & Fractals*, 72:68–76, 2015.
- 455 [39] Uday Shankar Shanthamallu, Jayaraman J Thiagarajan, Huan Song, and Andreas Spanias.
456 Gramme: Semisupervised learning using multilayered graph attention models. *IEEE transac-
457 tions on neural networks and learning systems*, 31(10):3977–3988, 2019.
- 458 [40] Nicholas D Sidiropoulos and Rasmus Bro. On the uniqueness of multilinear decomposition of
459 n-way arrays. *Journal of Chemometrics: A Journal of the Chemometrics Society*, 14(3):229–239,
460 2000.
- 461 [41] Wei Tang, Zhengdong Lu, and Inderjit S Dhillon. Clustering with multiple graphs. In *2009
462 Ninth IEEE International Conference on Data Mining*, pages 1016–1021. IEEE, 2009.
- 463 [42] Edwin JCG Van Den Oord and Ronan Van Rossem. Differences in first graders’ school
464 adjustment: The role of classroom characteristics and social structure of the group. *Journal of
465 School Psychology*, 40(5):371–394, 2002.
- 466 [43] James D Wilson, John Palowitch, Shankar Bhamidi, and Andrew B Nobel. Community
467 extraction in multilayer networks with heterogeneous community structure. *The Journal of
468 Machine Learning Research*, 18(1):5458–5506, 2017.
- 469 [44] Ka Yee Yeung, Mario Medvedovic, and Roger E Bumgarner. Clustering gene-expression data
470 with repeated measurements. *Genome biology*, 4(5):1–17, 2003.
- 471 [45] Yubai Yuan and Annie Qu. Community detection with dependent connectivity. *The Annals of
472 Statistics*, 49(4):2378–2428, 2021.
- 473 [46] Jingfei Zhang, Will Wei Sun, and Lexin Li. Network response regression for modeling popula-
474 tion of networks with covariates. *arXiv preprint arXiv:1810.03192*, 2018.
- 475 [47] Xuefei Zhang, Songkai Xue, and Ji Zhu. A flexible latent space model for multilayer networks.
476 In *International Conference on Machine Learning*, pages 11288–11297. PMLR, 2020.
- 477 [48] Yunpeng Zhao, Elizaveta Levina, and Ji Zhu. Consistency of community detection in networks
478 under degree-corrected stochastic block models. *The Annals of Statistics*, 40(4):2266–2292,
479 2012.
- 480 [49] Wei Zheng, Dingjie Wang, and Xiufen Zou. Control of multilayer biological networks and
481 applied to target identification of complex diseases. *BMC bioinformatics*, 20(1):1–12, 2019.

482 Checklist

- 483 1. For all authors...
- 484 (a) Do the main claims made in the abstract and introduction accurately reflect the pa-
485 per’s contributions and scope? [Yes] See the abstract and the third paragraph of the
486 introduction.
- 487 (b) Did you describe the limitations of your work? [Yes] The optimization algorithm can
488 only be guaranteed to converge to a stationary point.
- 489 (c) Did you discuss any potential negative societal impacts of your work? [No] There
490 should be no negative societal impacts.
- 491 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
492 them? [Yes]
- 493 2. If you are including theoretical results...
- 494 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 3.
- 495 (b) Did you include complete proofs of all theoretical results? [Yes] All technical proofs
496 are provided in Appendix E of the supplementary materials.
- 497 3. If you ran experiments...

- 498 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
499 imental results (either in the supplemental material or as a URL)? [Yes] The URLs
500 for data are included in Section 4.2, and codes with instructions are included in the
501 supplementary materials.
- 502 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
503 were chosen)? [Yes] See Section 2.3 and Appendix B in the supplementary materials.
- 504 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
505 ments multiple times)? [Yes] We show the standard errors in Table 1 and 95% confident
506 intervals of additional simulation studies in Appendix C in the supplementary materials.
- 507 (d) Did you include the total amount of compute and the type of resources used (e.g., type
508 of GPUs, internal cluster, or cloud provider)? [No]
- 509 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 510 (a) If your work uses existing assets, did you cite the creators? [Yes] We used publicly
511 available datasets and cite the creators.
- 512 (b) Did you mention the license of the assets? [Yes] All datasets we used are publicly
513 available.
- 514 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 515 (d) Did you discuss whether and how consent was obtained from people whose data you're
516 using/curating? [No]
- 517 (e) Did you discuss whether the data you are using/curating contains personally identifiable
518 information or offensive content? [No] All data we used do not contains personally
519 identifiable information or offensive content.
- 520 5. If you used crowdsourcing or conducted research with human subjects...
- 521 (a) Did you include the full text of instructions given to participants and screenshots, if
522 applicable? [N/A]
- 523 (b) Did you describe any potential participant risks, with links to Institutional Review
524 Board (IRB) approvals, if applicable? [N/A]
- 525 (c) Did you include the estimated hourly wage paid to participants and the total amount
526 spent on participant compensation? [N/A]