

ARTICLE TYPE

Factor Analysis on Citations, Using a Combined Latent and Logistic Regression Model

Namjoon Suh¹ | Xiaoming Huo^{*1} | Eric Heim² | Timothy Van Slyke² | Lee Seversky²

¹School of Industrial and Systems Engineering, Georgia Institute of Technology, 755 Ferst Dr, Atlanta, GA, USA

²Information System Division (AFRL/RIS), the Air Force Research Laboratory, Department of the Air Force, Air Force Materiel Command, AFRL/RIK - Rome, 26 Electronic Parkway, Rome, New York, USA

Correspondence

*Xiaoming Huo, Georgia Institute of Technology, Email: huo@gatech.edu

Present Address

755 Ferst Dr, Atlanta, GA, USA.

Summary

We propose a combined model, which integrates the latent factor model and the logistic regression model, for the citation network. It is noticed that neither a latent factor model nor a logistic regression model alone is sufficient to capture the structure of the data. The proposed model has a latent (i.e., factor analysis) model to represents the main technological trends (a.k.a., factors), and adds a sparse component that captures the remaining ad-hoc dependence. Parameter estimation are carried out through the construction of a joint-likelihood function of edges and properly chosen penalty terms. The convexity of the objective function allows us to develop an efficient algorithm, while the penalty terms push towards a low-dimensional latent component and a sparse graphical structure. Simulation results show that the proposed method works well in practical situations. The proposed method has been applied to a real application, which contains a citation network of statisticians (Ji and Jin, 2016 [11]). Some interesting findings are reported.

KEYWORDS:

citation network, matrix decomposition, latent variable model, logistic regression model, convex optimization, ADMM

1 | INTRODUCTION

We study a citation network, where each node (i.e., item) can be a technical report or a publication. A node may cite another node. Associated with a pair of nodes i and j , we denote a binary random variable X_{ij} , where $1 \leq i, j \leq n$ and n is the total number of nodes. We have $X_{ij} = 1$ if and only if either node i cites node j or vice versa; otherwise $X_{ij} = 0$. For each node i , we assume that there is an associated binary vector $f_i \in \mathbb{R}^K$, such that the k th entry of f_i , $f_{ik} = 1$, if and only if node i is related to topic (i.e., factor) k , $1 \leq k \leq K$. Here K is the total number of underlying topics (i.e., factors, or trends). We assume a logistic model for X_{ij} 's: for $1 \leq i, j \leq n$,

$$\mathbb{P}(X_{ij} = 1) = \frac{e^{\alpha + f_i^T D f_j}}{1 + e^{\alpha + f_i^T D f_j}}, \quad (1)$$

⁰Abbreviations: fused factor analysis and logistic graphical models, ADMM, network model

where $\alpha \in \mathbb{R}$ is a parameter and matrix $D \in \mathbb{R}^{K \times K}$ is a diagonal matrix: $D = \text{diag}\{d_1, d_2, \dots, d_K\}$. We assume $d_{ii} > 0$ for $1 \leq i \leq K$. Another way to put (1) is

$$\mathbb{P}(X_{ij} = 1) = \frac{\exp\left\{\alpha + \sum_{k=1}^K f_{ik}f_{jk}d_k\right\}}{1 + \exp\left\{\alpha + \sum_{k=1}^K f_{ik}f_{jk}d_k\right\}}. \quad (2)$$

A justification of the above model is that when both node i and node j are related to topic k , they have a higher chance to cite one way or the other. We have assumed a common strengthen coefficient d_k ($1 \leq k \leq K$) for factor k , despite different nodes. We denote a matrix $F = \{f_1, f_2, \dots, f_n\} \in \mathbb{R}^{K \times n}$. Each column i in matrix F contains the factor loadings associated with the node i ($1 \leq i \leq n$). Given the diagonal matrix D and the factor loading matrix F , we assume that X_{ij} 's are independent; therefore we have the total conditional probability function as follows:

$$\mathbb{P}(\{X_{ij}, 1 \leq i, j \leq n\}) = \prod_{1 \leq i < j \leq n} \mathbb{P}(X_{ij}) = \prod_{1 \leq i < j \leq n} \frac{e^{X_{ij}(\alpha + f_i^T D f_j)}}{1 + e^{\alpha + f_i^T D f_j}}, \quad (3)$$

where $\mathbb{P}(X_{ij})$ is given in (2). The last equation holds because X_{ij} only takes binary (i.e., 0 or 1) values. Recall that the dot product of two matrices with same dimensionality, $A, B \in \mathbb{R}^{a \times b}$, is defined as $A \bullet B = \text{trace}(A^T B) = \sum_{i=1}^a \sum_{j=1}^b a_{ij}b_{ij}$. The above (3) can be further rewritten as

$$\mathbb{P}(\{X_{ij}, 1 \leq i, j \leq n\}) = \frac{\exp(\alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet (F^T D F))}{\prod_{1 \leq i < j \leq n} 1 + e^{\alpha + f_i^T D f_j}}, \quad (4)$$

where we assume $X_{ii} = 0$ for all i ($1 \leq i \leq n$) and $X_{ij} = X_{ji}$ for all i and j ($1 \leq i, j \leq n$), i.e., the matrix X is symmetric. The above delivers a factor analysis model. Various linear and nonlinear latent variable models have been studied extensively in the literature (e.g., [12, 15, 14, 17, 10, 13]).

Our work is motivated from a recent work *Fused Latent and Graphical (FLaG) model* (Chen et al, 2016, [6]). They assume that majority of variation of responses can be accounted by low dimensional latent vector, and remaining dependent structure of responses can be explained by sparse graphical structure. Thus, the resulting model contains a low-dimensional latent vector and a sparse conditional graph. Their key idea is to separate these two dependent structures so that they can facilitate the statistical inference. In our model, we also assume that there exist two dependent structures among citation edges in a network. A low-dimensional version of the aforementioned latent vector model is largely correct and majority of the citations among the nodes are induced by these common latent vectors f_i 's (with weight coefficients d_i 's). There is still a small remainder due to the sparse graphical component.

Though it may seem similar to Chen et al [6], we work on a different model formulation in several aspects. We summarize the differences as follows.

1. FLaG is built to analyze the Eysenck's Personality Questionnaire that consists of items designed to measure Psychoticism, Extraversion, and Neuroticism. So there are p questions that need to be answered, and each questions fall into above three categories. If there are n respondents to questions, we have n independent data generated from same distribution. In our case, the observed citation network can be thought of as one realization of a random graph.
2. In FLaG model, a collection of binary responses for each question in the questionnaire follows a joint distribution, which is a combination of the Item Response Theory (IRT) model and the Ising model. We model the citation edges among papers as random variables, whose dependent structure is characterized by the combination of the Latent Factor Analysis model and the sparse graphical model.
3. FLaG approximates the original likelihood through constructing pseudo-likelihood function by taking advantage of conditional independence among the nodes. In our model, likelihood function is directly accessible due to the conditional independence among edges given parameters.

The proposed modeling framework is also related with the analysis of decomposing a matrix into low-rank and sparse components ([1, 4, 5, 19]). Specifically, [5] studies statistical inference of a multivariate Gaussian model whose precision matrix admits the form of a low-rank matrix plus a sparse matrix. The inference and optimization of the current model are different from the aforementioned cases. We will construct a regularized-likelihood function, based on which estimator will be proposed for simultaneous model selection and parameter estimation. The objective function in the optimization problem for

the regularized estimator is convex, for which we will develop an efficient algorithm through the alternating direction method of multiplier (ADMM, [3, 8, 9]).

The rest of the paper is organized as follows. In Section 2, we will give a presentation on how to build a model, which can encode both the latent dependent structure due to the common topics and the remaining sparse ad-hoc dependent structure. In Section 3, we will discuss the assumptions in our model and penalization on likelihood function that is constructed in Section 2. Section 5 gives the detailed procedure on how to compute the estimator of the optimization problem formulated in Section 3. In Section 6, we will present simple numerical experiments on synthetic data, and application of our model on real citation network of statisticians. We finally conclude this work in Section 7 with several open questions and some directions for future research.

2 | MODEL FORMULATION

Recall the following graphical model that was established in (4), which is essentially a factor model (latent variable model):

$$\mathbb{P}(\{X_{ij}, 1 \leq i, j \leq n\}) = \frac{\exp(\alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet (F^T D F))}{\prod_{1 \leq i < j \leq n} 1 + e^{\alpha + f_i^T D f_j}},$$

where $X_{ij}, 1 \leq i, j \leq n$, are binary random variables indicating either node i cites node j or vice versa, matrix $X = \{X_{ij}\} \in \mathbb{R}^{n \times n}$ is symmetric with diagonal entries all being equal to zero, factor loading matrix $F = [f_1, f_2, \dots, f_n] \in \mathbb{R}^{K \times n}$ records the relation between nodes and the underlying topics, F^T is the transpose of F , and matrix $D \in \mathbb{R}^{K \times K}$ is diagonal with entries being the weight coefficients of factors.

The above specifies a latent model (or equivalently a factor model). We now describe a graphical model as follows. The graphical model will complement the latent model by characterizing links that are not interpretable via common factors. For the aforementioned binary random variable $X_{ij}, 1 \leq i, j \leq n$, we define

$$\mathbb{P}(X_{ij} = 1) = \frac{e^{\alpha' + S_{ij}}}{1 + e^{\alpha' + S_{ij}}}, \quad (5)$$

where $S_{ij} \in \mathbb{R}$, for $1 \leq i, j \leq n$, denotes the relation between nodes i and j . Note that the matrix S is introduced to capture the ad-hoc links in the graph. If we have $S_{ij} \leq 0$, then it is less likely to have a citational relationship between nodes i and j . On the other hand, if $S_{ij} > 0$, then it is more likely to have a citation link between nodes i and j . Here parameter $\alpha' \in \mathbb{R}$ plays the same role as parameter α does in model (1). Denote the matrix $S = \{S_{ij}, 1 \leq i, j \leq n\} \in \mathbb{R}^{n \times n}$. Assume that given the matrix S , the binary random variables X_{ij} 's are independent; consequently, we have the total conditional probability function as follows:

$$\begin{aligned} \mathbb{P}(\{X_{ij}, 1 \leq i, j \leq n\}) &= \prod_{1 \leq i < j \leq n} \mathbb{P}(X_{ij}) \\ &= \prod_{1 \leq i < j \leq n} \frac{e^{X_{ij}(\alpha' + S_{ij})}}{1 + e^{\alpha' + S_{ij}}} \\ &= \frac{\exp(\alpha' \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet S)}{\prod_{1 \leq i < j \leq n} 1 + e^{\alpha' + S_{ij}}}. \end{aligned} \quad (6)$$

Recall that we have assumed that $X_{ii} = 0$ for all i ($1 \leq i \leq n$) and $X_{ij} = X_{ji}$ for all i and j ($1 \leq i, j \leq n$), i.e., the matrix X is symmetric. In the combined model, we integrate (4) and (6) to render the joint conditional probability function as follows:

$$\begin{aligned} \mathbb{P}(X \mid \alpha, F, D, S) &= \prod_{1 \leq i < j \leq n} \frac{e^{X_{ij}(\alpha + S_{ij} + f_i^T D f_j)}}{1 + e^{\alpha + S_{ij} + f_i^T D f_j}} \\ &= \frac{\exp\left(\alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet (F^T D F) + \frac{1}{2} X \bullet S\right)}{\prod_{1 \leq i < j \leq n} (1 + e^{\alpha + f_i^T D f_j + S_{ij}})}. \end{aligned} \quad (7)$$

3 | ESTIMATION

Note that in the model (7), the log-likelihood function has the form as follows:

$$\begin{aligned} \mathbb{L}(\alpha, F, D, S; X) = & \alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet (F^T D F) + \frac{1}{2} X \bullet S \\ & - \sum_{1 \leq i < j \leq n} \log \left(1 + e^{\alpha + f_i^T D f_j + S_{ij}} \right). \end{aligned} \quad (8)$$

If we consider maximizing the above log-likelihood function, we will encounter several technical issues that are described below.

1. We would like the matrix $S \in \mathbb{R}^{n \times n}$ to have as many zero entries as possible; i.e., matrix S is *sparse*.
2. There is an identifiability issue with the formation $F^T D F$. More specifically, let $P \in \mathbb{R}^{K \times K}$ be a signed permutation matrix, then we have $P^T P = I_n$, where $I_n \in \mathbb{R}^{K \times K}$ is the identity matrix. Notice that matrix $F' = P F$ is also a factor loading matrix, and matrix $D' = P D P^T$ is still a diagonal matrix; we have

$$F^T D F = F^T P^T P D P^T P F = (F')^T D' F',$$

i.e., the choice of F and D is not unique.

3. We would like the number of nonzeros in each column of F to be small, reflecting that each node is associated with a small number of underlying topics.
4. Overall, the rank of matrix $F^T D F$ cannot be larger than $\min\{n, K\}$. With the application that we have in mind, in this paper, we assume that K is much smaller than n .
5. To ensure the separation of matrices $\alpha \mathbb{1} \mathbb{1}^T$ and L , we assume that the eigen-vector of L is centered, that is,

$$J L J = L \quad \text{where} \quad J = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T,$$

where $\mathbb{1}$ denotes a n -dimensional vector whose entries are all 1s. This condition uniquely identifies F up to a common orthogonal transformation of its columns.

Directly maximizing the objective function in (8) is not going to be an easy task. Following the approaches that were mentioned in Introduction, we propose to relax $F^T D F$ to L , where L is a low rank matrix. Consequently, the log-likelihood function in (8) can be rewritten as

$$\begin{aligned} \mathbb{L}_n(\alpha, L, S; X) = & \alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet L + \frac{1}{2} X \bullet S \\ & - \sum_{1 \leq i < j \leq n} \log \left(1 + e^{\alpha + L_{ij} + S_{ij}} \right). \end{aligned} \quad (9)$$

We propose a penalized likelihood estimation approach as follows:

$$(\hat{\alpha}, \hat{L}, \hat{S}) = \arg \min_{\alpha, L, S} \left\{ -\frac{1}{n} \mathbb{L}_n(\alpha, L, S; X) + \gamma \|S\|_1 + \delta \|L\|_* \right\}, \quad (10)$$

where $\gamma > 0$ and $\delta > 0$ are algorithmic parameters whose values will be discussed later, the L_1 norm of matrix S is defined as $\|S\|_1 = \sum_{i \neq j} |S_{ij}|$ (Note that we do not penalize the diagonal entries of S), and nuclear norm of matrix L is defined as $\|L\|_* = \text{trace} \sqrt{L^T L}$. Recall that both S and L are symmetric matrices. The entries of matrix S can either be positive or negative. Note that we have imposed the diagonal entries of the matrix X to be zeros. Given that $L = F^T D F$ where matrix D is diagonal with nonnegative diagonal entries, it is easy to see that matrix L is positive semidefinite; which consequently leads to $\|L\|_* = \text{trace}(L)$, which is a linear functional to the matrix L . The nuclear norm of L mimicks the number of nonzero eigenvalues of L , which is the same as the rank of L . The regularization based on the nuclear norm was proposed in [7] and its statistical properties are studied in [2].

After we have obtained \hat{S} in (10), we can uncover the graphical model by investigating non-zero entries in \hat{S} . On the other hand when we have calculated \hat{L} , we may not be able to find binary matrix F and nonnegative diagonal matrix D such that $\hat{L} = F^T D F$. This is the price we have to pay for an amenable computational approach. The rank of estimated \hat{L} will be our

estimate of the number of factors (i.e., the number of underlying common topics). For the issue on assigning the community membership of each node i , we will discuss this later in Section 6.

4 | NON-ASYMPTOTIC ERROR BOUND OF THE ESTIMATOR

In this section, we focus on investigating the behaviour of non-asymptotic error bound of our estimator in the context where the number of papers in network is explicitly tracked. We are interested in solving the following optimization problem :

$$\min_{\substack{\alpha \in \mathbb{R}, S=S^T \\ L \geq 0}} -\frac{1}{n} \log \prod_{1 \leq i, j \leq n} \frac{\exp(X_{ij}(\alpha + L_{ij} + S_{ij}))}{1 + \exp(\alpha + L_{ij} + S_{ij})} + \delta \|L\|_* + \gamma \|S\|_1 \quad (11)$$

For the convenience of theoretical investigation, we slightly modify the first term in the objective function summing over all (i, j) pairs. After scaling, due to symmetry of X, L , and S , the only difference between (10) and (11) is in the inclusion of terms in diagonal pairs $(i, i), \forall i = 1, \dots, n$. This slight modification leads to neither theoretical consequence nor noticeable difference in practice. Let $(\hat{\alpha}, \hat{L}, \hat{S})$ be the solution to (11), and (α^*, L^*, S^*) be the ground truth, which governs the data generating process. Let $\hat{\Theta}, \Theta^*$ be defined respectively as $\hat{\Theta} = \hat{\alpha} \mathbb{1} \mathbb{1}^T + \hat{L} + \hat{S}$ and $\Theta^* = \alpha^* \mathbb{1} \mathbb{1}^T + L^* + S^*$. And denote the error term for each parameter as $\hat{\Delta}^\Theta = \hat{\Theta} - \Theta^*, \hat{\Delta}^\alpha = \hat{\alpha} - \alpha^*, \hat{\Delta}^L = \hat{L} - L^*, \hat{\Delta}^S = \hat{S} - S^*$. Throughout the discussion, let $P^* = \left\{ \frac{\exp(\Theta_{ij}^*)}{1 + \exp(\Theta_{ij}^*)} \right\}_{1 \leq i, j \leq n} \in \mathbb{R}^{n \times n}$. We then impose several assumptions for theoretical guarantees of our estimator.

Assumption 1. (Strong convexity) For any $\Theta \in \mathbb{R}^{n \times n}$, define the log-likelihood in (11), $h(\Theta) = -\frac{1}{n} \sum_{i,j} \{X_{ij} \Theta_{ij} - \log(1 + \exp(\Theta_{ij}))\}$. We assume that $h(\Theta)$ is τ -strongly convex in a sense that lowest eigenvalue of Hessian matrix of likelihood function is bounded away from zero ($\tau > 0$):

$$\nabla^2 h(\Theta) = \text{diag} \left(\text{vec} \left(\frac{\exp(\Theta)}{n(1 + \exp(\Theta))^2} \right) \right) \geq \tau I_{n^2 \times n^2}$$

For any vector a , $\text{diag}(a)$ is the diagonal matrix with elements of a on its diagonal. For any square matrix A and B , $A \geq B$ if and only if $A - B$ is positive semi-definite.

Assumption 2. (Identifiability of $\alpha \mathbb{1} \mathbb{1}^T$ and L , Spikiness of L) To ensure the identifiability of $\alpha \mathbb{1} \mathbb{1}^T$ and L , we assume the latent variables are centered, that is $JL = L$, where $J = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T$, $\mathbb{1}$ denotes all one vector in \mathbb{R}^n . We impose a spikiness condition $\|L\|_\infty \leq \frac{\kappa}{\sqrt{n \times n}}$ on L , to ensure the separation of L and S matrix [1]. We would also like to note that the constraint $|\alpha| \leq C\kappa$, for an absolute constant C , is included partially for obtaining theoretical guarantees.

With these assumptions, we present the behavior of non-asymptotic error bound of our estimator through following theorem. In our result, we measure error using squared Frobenius norm summed across three matrices:

$$e^2(\hat{\alpha} \mathbb{1} \mathbb{1}^T, \hat{L}, \hat{S}) := \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \|\hat{\Delta}^L\|_F^2 + \|\hat{\Delta}^S\|_F^2$$

Theorem 1. Under the assumptions 1 and 2, if we solve the convex problem (11) with a pair of regularization parameter (δ, γ) satisfying

$$\delta \geq 2 \left\| \frac{1}{n} (X - P^*) \right\|_{op} \quad \text{and} \quad \gamma \geq 2 \left\| \frac{1}{n} (X - P^*) \right\|_\infty + 4\kappa\tau \left(\frac{Cn+1}{n} \right) \quad (12)$$

There exist universal constants $c_j, j = 1, 2, 3$, for all integers $k = 1, 2, \dots, n$ and $s = 1, 2, \dots, n^2$, we have following upper bound on $e^2(\hat{\alpha} \mathbb{1} \mathbb{1}^T, \hat{L}, \hat{S})$

$$e^2(\hat{\alpha} \mathbb{1} \mathbb{1}^T, \hat{L}, \hat{S}) \leq \underbrace{c_1 \frac{\delta^2}{\tau^2}}_{\mathcal{K}_{\alpha^*}} + \underbrace{c_2 \frac{\delta^2}{\tau^2} \left\{ k + \frac{\tau}{\delta} \sum_{j=k+1}^n \sigma_j(L^*) \right\}}_{\mathcal{K}_{L^*}} + \underbrace{c_3 \frac{\gamma^2}{\tau^2} \left\{ s + \frac{\tau}{\gamma} \sum_{(i,j) \notin M} |S_{ij}^*| \right\}}_{\mathcal{K}_{S^*}} \quad (13)$$

where M is an arbitrary subset of matrix indices of cardinality at most s .

We would first like to note readers that the result presented in Theorem 1 can be thought of as an extension to Theorem 1 presented in paper [1] to generalized linear model. Our work considers a logistic loss function whose parameters are characterized by sparse matrix plus low rank matrix, whereas [1] works on quadratic loss of noisy realization of linear transformation of sum of low rank and sparse matrices. We show that similar proof techniques can be applied in our case, giving us similar results.

Result of the Theorem 1 provides a family of upper-bounds, one for each indexed by a specific choice of model subspace M , and rank parameter k , helping us to understand the behavior of our proposed estimator both in practice and in ideal cases. In real-world, large scale citation network analysis, it is difficult to expect network with L^* matrix with exactly $k(\ll n)$ non-zero eigenvalues, but with a matrix L^* with approximately low rank. In this case, our theorem provides the intuition on the rate on how fast the error bound \mathcal{K}_{L^*} will grow. In ideal cases where L^* is an exact low rank matrix with rank k (i.e., $\text{rank}(L^*) = k$) and S^* is a sparse matrix, which lies within model subspace M (i.e., $|\text{supp}(S^*)| = s$), we can easily verify “approximation error” terms in \mathcal{K}_{L^*} (i.e., $\tau \sum_{j=k+1}^n \sigma_j(L^*)$) and in \mathcal{K}_{S^*} (i.e., $\tau \sum_{(i,j) \notin M} |S_{ij}^*|$) disappear, giving us Frobenius error bound as follows:

$$e^2(\hat{\alpha} \mathbb{1}^T, \hat{L}, \hat{S}) \lesssim \delta^2(k+1) + \gamma^2 s$$

where \lesssim denotes that we ignore constant factors.

5 | COMPUTATION

We propose a method that takes advantage of the special structure of the L_1 and the nuclear norm by means of the alternating direction method of multiplier (ADMM), which is a method that has recently gained momentum. An examination of the objective function in (10) unveils that terms

$$\alpha \sum_{1 \leq i < j \leq n} X_{ij} + \frac{1}{2} X \bullet L + \frac{1}{2} X \bullet S$$

are linear in α , L , and S . The term

$$\sum_{1 \leq i < j \leq n} \log(1 + e^{\alpha + L_{ij} + S_{ij}})$$

is convex with respect to α , L , and S . Functions $\|S\|_1$ and $\|L\|_*$ are known to be convex functions. Therefore, the objective function in (10) is convex. The above convex optimization problem can be solved via ADMM as follows.

5.1 | ADMM approach

We give a review of the alternating direction method of multiplier (ADMM). Consider two closed convex functions

$$f : \mathcal{X}_f \rightarrow \mathbb{R} \text{ and } g : \mathcal{X}_g \rightarrow \mathbb{R},$$

where the domain \mathcal{X}_f and \mathcal{X}_g of functions f and g are closed convex subsets of \mathbb{R}^d , and $\mathcal{X}_f \cap \mathcal{X}_g$ is nonempty. Both f and g are possibly non-differentiable. The alternating direction method of multiplier is an iterative algorithm that solves the following generic optimization problem:

$$\min_{x \in \mathcal{X}_f \cap \mathcal{X}_g} \{f(x) + g(x)\},$$

or equivalently

$$\begin{aligned} \min_{x \in \mathcal{X}_f, z \in \mathcal{X}_g} \{f(x) + g(z)\}, \\ \text{subject to} \quad x = z. \end{aligned} \tag{14}$$

To describe the algorithm, we will need the following proximal operators

- $\mathbf{P}_{\lambda, f} : \mathbb{R}^d \rightarrow \mathcal{X}_f$ as

$$\mathbf{P}_{\lambda, f}(v) = \arg \min_{x \in \mathcal{X}_f} \left\{ f(x) + \frac{1}{2\lambda} \|x - v\|_2^2 \right\}$$

- and $\mathbf{P}_{\lambda, g} : \mathbb{R}^d \rightarrow \mathcal{X}_g$ as

$$\mathbf{P}_{\lambda, g}(v) = \arg \min_{x \in \mathcal{X}_g} \left\{ g(x) + \frac{1}{2\lambda} \|x - v\|_2^2 \right\},$$

where $\|\cdot\|_2$ is the usual Euclidean norm on \mathbb{R}^d and λ is a scale parameter that is a fixed positive constant.

The algorithm starts with some initial values $x^0 \in \mathcal{X}_f, z^0 \in \mathcal{X}_g, u^0 (= \lambda y^0) \in \mathbb{R}^d$. At the $(m+1)$ th iteration, (x^m, z^m, u^m) is updated according to the following steps until convergence

- Step 1: $x^{m+1} = \mathbf{P}_{\lambda, f}(z^m - u^m)$,

- Step 2: $z^{m+1} = \mathbf{P}_{\lambda,g}(x^{m+1} + u^m)$,
- Step 3: $u^{m+1} = u^m + x^{m+1} - z^{m+1}$.

The convergence properties of the algorithm are summarized in the following result in [3]. Let p^* be the minimal value in (14).

Theorem 2 (Boyd et al., 2011). Assume functions $f : \chi_f \rightarrow \mathbb{R}$ and $g : \chi_g \rightarrow \mathbb{R}$ are closed convex functions, whose domains χ_f and χ_g are closed convex subsets of \mathbb{R}^d and $\chi_f \cap \chi_g \neq \emptyset$. Assume the Lagrangian of (14)

$$L(x, z, y) = f(x) + g(z) + y^T(x - z)$$

has a saddle point, that is, there exists (x^*, z^*, y^*) (not necessarily unique) that $x^* \in \chi_f$ and $z^* \in \chi_g$, for which

$$L(x^*, z^*, y) \leq L(x^*, z^*, y^*) \leq L(x, z, y^*), \quad \forall x, z, y \in \mathbb{R}^d.$$

Then the ADMM has the following convergence properties.

1. Residual convergence. $x^m - z^m \rightarrow 0$ as $m \rightarrow \infty$; i.e., the iterates approach feasibility.
2. Objective convergence. $f(x^m) + g(z^m) \rightarrow p^*$ as $m \rightarrow \infty$; i.e., the objective function of the iterates approaches the optimal value.
3. Dual variable convergence. $y^m \rightarrow y^*$ as $m \rightarrow \infty$, where y^* is a dual optimal point.

Now we see how ADMM can be adopted to solve for our penalized likelihood estimation (10). We reparameterize $M = L + S$ and let $x = (\alpha, M, L, S)$ (viewed as a vector). We define the following:

$$\begin{aligned} \chi_f &= \{(\alpha, M, L, S) : \alpha \in \mathbb{R}, M, L, S \in \mathbb{R}^{n \times n}, L \text{ is positive semidefinite}, S \text{ is symmetric}\}, \\ f(x) &= -\frac{\alpha}{n} \sum_{1 \leq i < j \leq n} X_{ij} - \frac{1}{2n} X \bullet M + \frac{1}{n} \sum_{1 \leq i < j \leq n} \log(1 + e^{\alpha + M_{ij}}) + \gamma \|S\|_1 + \delta \|L\|_*, \\ \chi_g &= \{(\alpha, M, L, S) : \alpha \in \mathbb{R}, M, L, S \in \mathbb{R}^{n \times n}, M \text{ is symmetric and } M = L + S\}, \text{ and} \\ g(x) &= 0, \text{ for } x \in \chi_g. \end{aligned}$$

One can verify that (10) can be written as

$$\min_{x \in \chi_f \cap \chi_g} \{f(x) + g(x)\}.$$

We now present each of the three steps of the ADMM algorithm and show that the proximal operators $\mathbf{P}_{\lambda,f}$ and $\mathbf{P}_{\lambda,g}$ are easy to evaluate. Let

$$x^m = (x_\alpha^m, x_M^m, x_L^m, x_S^m), \quad z^m = (z_\alpha^m, z_M^m, z_L^m, z_S^m), \quad u^m = (u_\alpha^m, u_M^m, u_L^m, u_S^m).$$

Step 1. We solve $x^{m+1} = \mathbf{P}_{\lambda,f}(z^m - u^m)$. Due to the special structure of $f(\cdot)$, x_α^{m+1} , x_M^{m+1} , x_L^{m+1} , and x_S^{m+1} can be updated separately. More precisely, we have

$$\begin{aligned} x_\alpha^{m+1}, x_M^{m+1} &= \arg \min_{\alpha, M} -\frac{\alpha}{n} \sum_{1 \leq i < j \leq n} X_{ij} - \frac{1}{2n} X \bullet M + \frac{1}{n} \sum_{1 \leq i < j \leq n} \log(1 + e^{\alpha + M_{ij}}) \\ &\quad + \frac{1}{2\lambda} [\alpha - (z_\alpha^m - u_\alpha^m)]^2 + \frac{1}{2\lambda} \|M - (z_M^m - u_M^m)\|_F^2, \end{aligned} \quad (15)$$

$$\begin{aligned} x_L^{m+1} &= \arg \min_L \delta \|L\|_* + \frac{1}{2\lambda} \|L - (z_L^m - u_L^m)\|_F^2, \\ &\text{subject to } L \text{ is positive semidefinite;} \end{aligned} \quad (16)$$

$$\begin{aligned} x_S^{m+1} &= \arg \min_S \gamma \|S\|_1 + \frac{1}{2\lambda} \|S - (z_S^m - u_S^m)\|_F^2, \\ &\text{subject to } S \text{ is symmetric,} \end{aligned} \quad (17)$$

where $\|\cdot\|_F$ is the matrix Frobenius norm, defined as $\|M\|_F^2 = \sum_{i,j} m_{ij}^2$ for a matrix $M = (m_{ij})$. The problem in (15) may not have a closed-form solution. We use a simple gradient descent to solve in this step, setting the step size equal to 0.05 and stopping criteria as $\max(|x_{\alpha,m}^{(t+1)} - x_{\alpha,m}^{(t)}|, \|x_{M,m}^{(t+1)} - x_{M,m}^{(t)}\|_\infty) \leq 10^{-9}$. Note that there are close-form solutions to (16) and (17), while (15) is a unconstrained convex optimization problem. More specifically, in (16), suppose the eigenvalue decomposition of the symmetric matrix $(z_L^m - u_L^m)$ can be written as

$$z_L^m - u_L^m = T \Lambda T^T,$$

where T is orthogonal ($TT^T = I_n$). Then, for $J = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T$, we have

$$x_L^{m+1} = J(T \text{diag}(\Lambda - \lambda\delta)_+ T^T) J^T,$$

and $\text{diag}(\Lambda - \lambda\delta)_+$ is a diagonal matrix with the j th diagonal entry being

$$(\Lambda_{jj} - \lambda\delta)_+ = \begin{cases} 0, & \text{if } \Lambda_{jj} < \lambda\delta, \\ \Lambda_{jj} - \lambda\delta, & \text{if } \Lambda_{jj} \geq \lambda\delta. \end{cases}$$

In (17), we have, for $i \neq j$,

$$S_{ij} = \begin{cases} 0, & \text{if } |(z_S^m - u_S^m)_{ij}| < \lambda\gamma, \\ (z_S^m - u_S^m)_{ij} - \lambda\gamma, & \text{if } (z_S^m - u_S^m)_{ij} > \lambda\gamma, \\ (z_S^m - u_S^m)_{ij} + \lambda\gamma, & \text{if } (z_S^m - u_S^m)_{ij} < -\lambda\gamma. \end{cases}$$

Step 2. We solve $z^{m+1} = \mathbf{P}_{\lambda,g}(x^{m+1} + u^m)$. A close-form solution exists here. Denote $\bar{\alpha} = x_\alpha^{m+1} + u_\alpha^m$, $\bar{M} = x_M^{m+1} + u_M^m$, $\bar{L} = x_L^{m+1} + u_L^m$, and $\bar{S} = x_S^{m+1} + u_S^m$, then evaluating $\mathbf{P}_{\lambda,g}(x^{m+1} + u^m)$ becomes

$$\begin{aligned} \min_{\alpha, M, L, S} \quad & \frac{1}{2}[\alpha - \bar{\alpha}]^2 + \frac{1}{2}\|M - \bar{M}\|_F^2 + \frac{1}{2}\|L - \bar{L}\|_F^2 + \frac{1}{2}\|S - \bar{S}\|_F^2 \\ \text{subject to} \quad & M \text{ is symmetric and } M = L + S. \end{aligned}$$

The above optimization problem has a close-form solution, which is as follows:

$$\begin{aligned} z_\alpha^{m+1} &= \bar{\alpha}, \\ z_M^{m+1} &= \frac{1}{3}\bar{M} + \frac{1}{3}\bar{M}^T + \frac{1}{3}\bar{L} + \frac{1}{3}\bar{S}, \\ z_L^{m+1} &= \frac{1}{6}\bar{M} + \frac{1}{6}\bar{M}^T + \frac{2}{3}\bar{L} - \frac{1}{3}\bar{S}, \quad \text{and} \\ z_S^{m+1} &= \frac{1}{6}\bar{M} + \frac{1}{6}\bar{M}^T - \frac{1}{3}\bar{L} + \frac{2}{3}\bar{S}. \end{aligned}$$

Step 3. We solve $u^{m+1} = u^m + x^{m+1} - z^{m+1}$, which is a simple arithmetic.

The most important implementation details of this algorithm are the choice of λ and stopping criterion. In this work, we simply choose $\lambda = 0.5$. We terminate the algorithm when in the m th iteration, we have $\|x_M^m - x_L^m - x_S^m\|_F \leq \delta$, with $\delta = 10^{-7}$.

6 | NUMERICAL ANALYSIS AND APPLICATIONS

Section 6 is divided into two parts. In subsection 6.1, we conduct an empirical study of measuring the performance of our proposed method with artificially specified graphical structures. In subsection 6.2, we perform a real data analysis with citation network for statisticians.

Remark 1. Throughout the Section, we add $*$ to the superscript of parameters when they govern a data-generating process. α^* is an intercept term in logistic regression model. f_j^* is the j th column of factor loading matrix $F^* \in \mathbb{R}^{K \times n}$, where K refers to the number of topics embedded in network, and n denotes the number of papers in a network. Matrix $D^* \in \mathbb{R}^{K \times K}$ is a diagonal matrix whose entries denote weight coefficient of factors. S_{ij}^* denotes an ad-hoc citational relation between i th and j th paper. We use the notation $|S^*|$ to denote the cardinality of non-zero components of upper-triangular part of S^* matrix (i.e., $|S^*| = |\{(i, j) : S_{ij}^* \neq 0, 1 \leq i < j \leq n\}|$). We will use $X \sim \text{Unif}[a, b]$ to denote X is drawn from uniform distribution supported on the interval $[a, b]$.

6.1 | Numerical experiments with synthetic data

In this subsection, after a brief introduction on several synthetic scenarios that we want to explore, we will discuss how to set ground truth parameters for generating synthetic networks from pre-specified scenarios. Then, three model selection criteria and related evaluation metrics for the selected model will be introduced. Lastly, several interesting findings from the experiments will be presented.

6.1.1 | Synthetic Setting

We consider following three scenarios: $\{(n^{(i)}, K^{(i)}, nnz^{(i)})\}_{i=1}^3 = \{(30, 3, 10), (80, 4, 20), (120, 5, 40)\}$.

6.1.2 | Data Generation

Following the notation of our paper, each edge X_{ij} follows Bernoulli distribution, whose parameter is parameterized by the probability, $P_{ij}^* = \frac{\exp(\alpha^* + F_i^{*T} D^* F_j^* + S_{ij}^*)}{1 + \exp(\alpha^* + F_i^{*T} D^* F_j^* + S_{ij}^*)}$. We set the model parameters randomly following these steps:

1. Generate α^* from $\text{Unif}[-3, -2]$ (i.e., $\alpha^* \sim \text{Unif}[-3, -2]$).
2. For the binary factor loading matrix $F^* (\in \mathbb{R}^{n \times K})$, first fill in each row of the matrix with $\lfloor \frac{n}{K} \rfloor$ ones, and make sure that each column of F^* has only one 1. In last row of F^* , fill in the last remaining $n - K \lfloor \frac{n}{K} \rfloor$ entries with 1's. Then, we randomly choose one of columns, and fill that column with ones. Lastly, set $F^* = F^* J$, where $J = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T$.
3. Each element of diagonal matrix D^* is generated from uniform distribution supported on $[7, 8]$.
4. For the sparse component S^* , first draw $|S^*|$ sparsity pattern first draw nnz distinct random pairs of indices that are uniformly distributed over the indices of upper triangular part (off diagonal) of matrix S^* . For each drawn pair of indices (i, j) , draw random number from $\text{Unif}[7, 8]$ (i.e., $S_{ij}^* \sim \text{Unif}[7, 8]$).

6.1.3 | Choosing the tuning parameters and evaluation criteria

We observe that choosing a good pair of tuning parameters is an important yet challenging issue in our setting. Here we present a heuristic approach. Following the approach in Ji and Jin [11], we plot the largest 20 eigenvalues of adjacency matrix X , which can tell us the number of communities embedded in network. Then, we record the rank of \hat{L} and number of non-zero elements in \hat{S} for each tuning parameter pair on a given grid, and choose a proper pair that give us interesting interpretations of data.

One might wonder how traditional model selection methods, such as the Bayes Information Criterion (BIC; [18]) and the AIC, work where we recall that BIC and AIC are defined as follows:

$$\text{BIC}(M) = -2\mathbb{L}_n(\hat{\beta}(M)) + |M| \log \left(\frac{n(n-1)}{2} \right),$$

and

$$\text{AIC}(M) = -2\mathbb{L}_n(\hat{\beta}(M)) + 2|M|,$$

where M is the current model, $\mathbb{L}_n(\hat{\beta}(M))$ is the maximal log-likelihood for a given model M , and $|M|$ is the number of free parameters in M , which is determined by the number of non-zeros in S and the low-rank matrix L . If we have $\text{rank}(L) = K$, we can establish the following

$$|M| = \sum_{i < j} 1_{\{S_{ij} > 0\}} + nK - \frac{K(K-1)}{2} + 1;$$

since the number of free parameters in L is K plus $nk - K(K+1)/2$, which is the number of free parameters in determining K orth-normal vectors. Additionally 1 in the last is due to α . We want to find an M , which minimizes $\text{BIC}(M)$ or $\text{AIC}(M)$ as a function of M . However these traditional methods seem to entail a serious problem. Our experience shows that both BIC and AIC choose the most parsimonious model with the smallest $K (= \text{rank}(\hat{L}))$ and the smallest cardinality $(= |\text{supp}(\hat{S})|/2)$ of ad-hoc dependency among all the estimated (\hat{L}, \hat{S}) 's. We evaluate the models that are selected via our heuristic approach, BIC, and AIC by using following three evaluation metrics.

$$\begin{aligned} M_1 &= \mathbb{1}\{\text{rank}(\hat{L}) = \text{rank}(L^*)\}, \\ M_2 &= \frac{|\{(i, j) : i < j : S_{i,j}^* \neq 0 \ \& \ \hat{S}_{i,j} \neq 0\}|}{|\{(i, j) : i < j : S_{i,j}^* \neq 0\}|}, \\ M_3 &= \frac{|\{(i, j) : i < j : S_{i,j}^* = 0 \ \& \ \hat{S}_{i,j} \neq 0\}|}{|\{(i, j) : i < j : S_{i,j}^* = 0\}|}, \end{aligned}$$

where M_1 is a metric on whether the selected model recovers true low rank structure of network, M_2 evaluates the positive selection rate of the sparse ad-hoc structure in network, and M_3 evaluates the false discovery rate of ad-hoc edges. With properly selected tuning parameter, M_1 will be 1, M_2 will be close to 1, and M_3 will get close to 0. We present the evaluation results of three scenarios via these three criteria in Table 1 .

	Case 1			Case 2			Case 3		
	Heuristic	AIC	BIC	Heuristic	AIC	BIC	Heuristic	AIC	BIC
M_1	1	0	0	1	0	0	1	0	0
M_2	0.5	0	0	0.842	0	0	0.725	0	0
M_3	0.007	0	0	0.008	0	0	0.0113	0	0

TABLE 1 For three scenarios, our heuristic method chooses the model with \hat{L} with true rank, \hat{S} whose M_2 value is close to 1, and M_3 value is close to 0. Note that, for all three scenarios, both AIC and BIC choose the model with no ad-hoc structure (i.e., $\hat{S} = 0$), and \hat{L} with rank 1 (i.e., $\text{rank}(\hat{L}) = 1$).

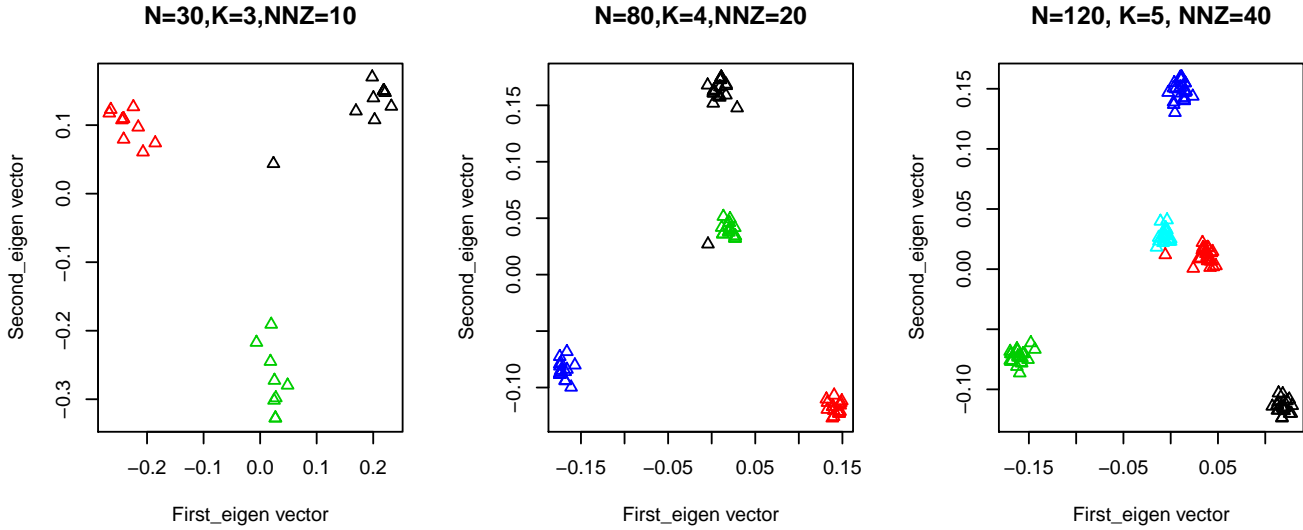


FIGURE 1 Illustrations of assignment results via k -means algorithm on K eigenvectors of the estimated \hat{L} matrix.

6.1.4 | Findings

Given a network data X , after fitting the model with a proper pair of tuning parameters, (γ, δ) , we need to determine whether node i belongs to topic k . We applied a simple k -means clustering on the fitted L matrix's K eigenvectors. For the illustration of performance of our fitted model after applying k -means algorithm on three synthetic scenarios, in Fig1 , we plot the coordinates of each rows of first two eigen-vectors on 2D plane.

Ad-hoc links. Ad-hoc links of synthetic network data can be thought of as “across edges” between clusters of nodes. Since indices of non-zero entries of S^* are randomly chosen, there might be some indices of positive S_{ij}^* entries where L_{ij}^* is also positive. In this case, this makes the edges indistinguishable if they come from S_{ij}^* or L_{ij}^* . So the total number of “across edges” might not be exactly same as “nnz” as we set in our data generation setting.

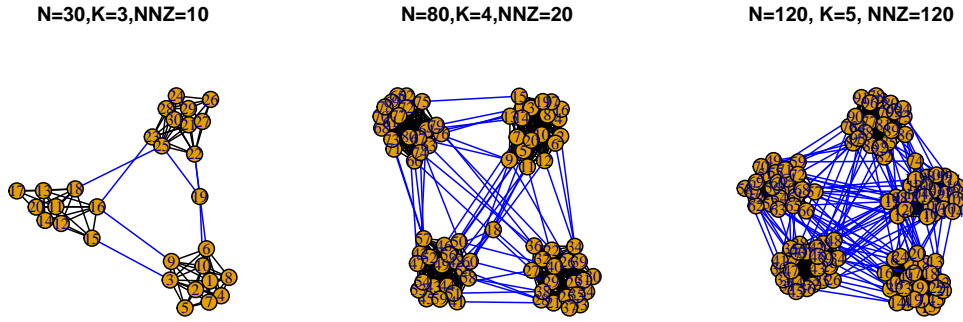


FIGURE 2 Edges colored with blue correspond to non-zero entries of estimated \hat{S} matrix. For all three cases, our algorithm correctly separates the “across edges” from edges coming from low-rank matrix structure.

6.2 | Citation networks for statisticians

Recently, Ji and Jin [11] published an interesting data set on citation and co-authorship networks of statisticians. Citational relationships over 3000 papers from 4 statistical journals were collected, and collaborative network among authors of the papers is also given. In our work, we only analyze the citation network of papers. This dataset is based on all papers published from

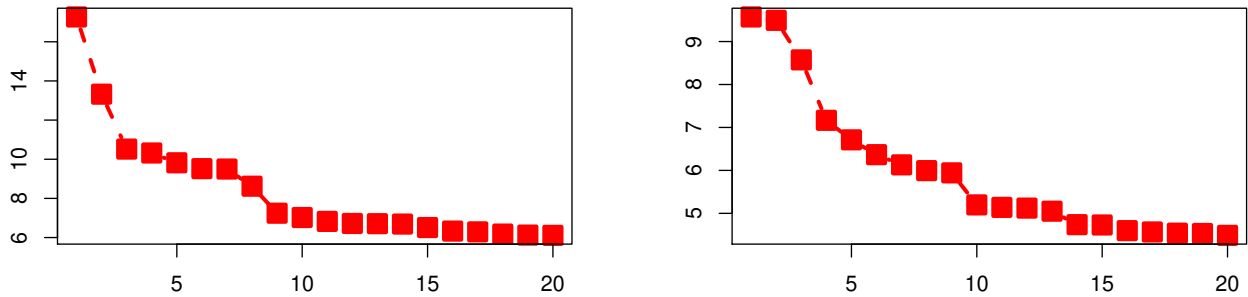


FIGURE 3 Scree plots. From left to right: 543-Citation network, Mixed Topics-Citation network.

2003 to first half of 2012, in 4 of the top statistical journals: Annals of Statistics, Biometrika, Journal of American Statistical Association and Journal of Royal Statistical Society (Series B). Among 3248 papers in total, we restrict our attention on papers, which have greater than or equal to 7 citations from the collected papers. Consequently, we have 543 papers in total. The elbow points of the eigen-value scree plot from may be at the 3rd, 7th, or 9th largest eigenvalue, suggesting that there might be from 2 to 8 topics embedded in network (See Fig3). In light of this, we performed our analysis under the assumption that there might exist 2 distinct topics and one giant component, in which 3 or 4 topics are mixed up.

A pair of parameters, $(\gamma, \delta) = (0.000912, 0.0097)$, gives us \hat{L} with rank 3, and \hat{S} with $|\hat{S}| = 197$. We apply k -means algorithm on 3 eigenvectors of estimated \hat{L} matrix, and visualize the result of clustering by plotting the rows of first two eigenvectors of \hat{L} on plane in Fig 4 . 404 papers colored in green are densely clustered nearby the origin with a short tail, (Fig4), and k -means classifies these papers as third topic. We list the first two topics discovered through our analysis.

- Variable selection (85 papers)
- Multiple Hypothesis Testing (54 papers)

These two topics are classical research topics in field of statistics. First group talks about Variable Selection in high-dimensional data. Majority of papers in second community study about Controlling False Discovery Rate in various statistical settings. As expected, third group is quite hard to interpret and seems to have substructures. For further investigation on this group, after obtaining the sub-network which is comprised only with the papers in the third group, we performed same analysis once again under the assumption that $K = 4$ (See Fig3 ,right). We obtain four sub-communities as follows:

- Non-parametric Bayesian Statistics (17 papers)
- Functional Data analysis (35 papers)
- Dimension Reduction (23 papers)
- Mixed topics (329 papers)

Out of this one big chunk, we got three small, but meaningful topics: Bayesian Statistics, Functional / Longitudinal Data Analysis, Dimension Reduction. Due to small volume of each communities, we could manually check that false discoveries for each community are all zero.

Sub-network structure has also a big chunk of papers which we refer to as Mixed Topics. Not only could we see the papers with topics on Learning Theory, Non-parametric/Semi-parametric Statistics, Spatial Statistics, Theoretical Machine Learning, which does not seem to belong to any of the five communities listed above, but also we could identify the papers with combinations of two or three topics. Papers, such as The Bayesian Lasso (T. Park, et al. 2008), Coordinate-independent sparse sufficient dimension reduction and variable selection (X. Chen, et al. 2010), are the examples of these papers. It is also interesting to think about a reason on papers which seem to have obvious membership in one of 5 communities other than Mixed topic classified as Mixed topic. For instance, "On the "degrees of freedom" on the LASSO (H. Zou, et al. 2007)" is classified as Mixed Topic paper. We can simply guess model selection has lots of applications in other topics, so it might cite or have been cited by many papers in other communities. Actually, out of 19 citation relationships it has with other papers, 12 of them came from the relationships with papers from Mixed topics.

Furthermore, it is worth to note that our result is consistent with that from Ji and Jin [11] in a sense that they also recover topics such as Multiple Testing, Variable Selection and Non-parametric Bayesian statistics. This is interesting since even though the dataset we consider is different from theirs, the results are consistent. They focus their attentions on analyzing "weakly connected giant component" for a citation network of each paper's authors (i.e., each node is an author). We consider the citation network of papers whose node degree is greater than or equal to 7.

Non-zero components of \hat{S} capture the citation relationships among papers that are not attributable to the common topics. The model we select has 197 sparse edges (9% of total edges), and all of them are positive edges. In Table 6.2, we provide 15 pairs of papers that have the largest estimated \hat{S}_{ij} . All the 15 edges come from pairs of papers from different communities. For instance, the first pair of papers comes from the Functional Analysis community (denoted by *FuncAn*) and Variable Selection (*VarSel*) community. *FuncAn* paper cites *VarSel* paper for borrowing a mathematical representation to build a theorem. Though it might seem to be a crucial step for building a theorem in their paper, we cannot say that two papers are closely related in terms of topic. Second pair comes from Multiple Testing (denoted as *MulT*) and Variable Selection (*VarSel*) community. *MulT* paper briefly mentions about *VarSel* paper in future work section, suggesting a possible way of combining their work and work in *VarSel* paper. And we also observe that as the estimated weight of sparse edges (i.e., \hat{S}_{ij}) decreases, the number of pairs of papers that come from Mixed topics community increases.

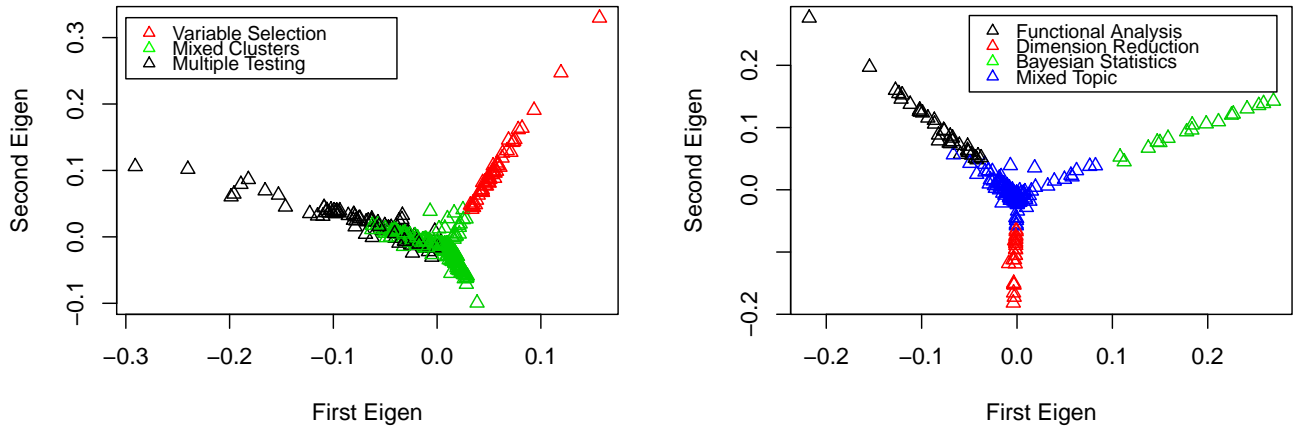


FIGURE 4 Illustration of first two eigenvectors of estimated \hat{L} matrix for original citation network with 543 papers (left) and sub-network graph (right). We present the results of assignment of each paper to each topic via k -means algorithm with different colors.

Pair	Weight	Community	Title
1	6.03	FuncAn	Properties of principal component methods for functional and longitudinal data analysis
2	3.08	VarSel	Nonconcave penalized likelihood with a diverging number of parameters
3	2.47	MulT	Innovated higher criticism for detecting sparse signals in correlated noise
4	2.15	VarSel	Regularized estimation of large covariance matrices
5	1.86	DimRed	Contour projected dimension reduction
6	1.77	VarSel	Factor profiled sure independence screening
7	1.71	VarSel	Factor profiled sure independence screening
8	1.66	DimRed	Sliced regression for dimension reduction
9	1.61	VarSel	Covariance regularization by thresholding
10	1.54	MulT	Adapting to unknown sparsity by controlling the false discovery rate
11	1.53	VarSel	Asymptotic properties of bridge estimators in sparse high-dimensional regression models
12	1.39	Mixed	Marginal asymptotics for the "large p small n" paradigm: with applications to microarray data
13	1.29	Mixed	A majorization-minimization approach to Variable Selection using spike and slab priors
14	1.20	VarSel	Empirical Bayes selection of wavelet thresholds
15	1.19	Mixed	Spades and mixture models
			Adapting to unknown sparsity by controlling the false discovery rate
			A constructive approach to the estimation of dimension reduction directions
			Factor profiled sure independence screening
			A majorization-minimization approach to variable selection using spike and slab priors
			Spike and slab variable selection: frequentist and Bayesian strategies
			variable selection in nonparametric additive models
			Nonparametric estimation of an additive model with a link function
			Nonparametric inferences for additive models
			Nonparametric independence screening in sparse ultra-high-dimensional additive models
			Bayesian variable selection in structured high-dimensional covariate spaces with applications in genomics
			Model selection and estimation in regression with grouped variables
			The sparsity and bias of the LASSO selection in high-dimensional linear regression
			Tests for high-dimensional covariance matrices
			Empirical dynamics for longitudinal data
			Variable selection in nonparametric varying-coefficient models for analysis of repeated measurements

TABLE 2 15 pairs of papers that have the largest \hat{S}_{ij} among the 197 non-zero component of \hat{S} . Note that all the 15 edges come from pairs of papers from different communities.

7 | CONCLUSION

We propose a new combined latent factor and sparse graphical model. We consider the regularized likelihood by means of L_1 and nuclear norm penalties. The computation of the regularized estimator is facilitated by developing an algorithm based on the alternating direction method of multiplier to optimize a non-smooth and convex objective function. The proposed method is applied to citation network of statisticians, and the estimated model renders good interpretative power. Specifically, our analysis on statistician's citation network sheds the new light on the interpretation of dataset.

Nonetheless, there are still several questions remaining to be answered. First of all, it remains unclear on how to choose the proper tuning parameters. Classical methods for choosing tuning parameters such as BIC or AIC did not work in our case, since they tend to choose the most parsimonious models. We also do not have systematic ways to do cross validation in network data. Not only because it is computationally expensive procedure, but also because if we partition the network data, we can lose fair amount of information on dependent structures among edges. This problem is also closely related with determining the number of communities in network. In lieu of using BIC or AIC, our analysis is heavily relying on heuristic approach when choosing the tuning parameter, and during this procedure, we use the screeplot for determining number of communities in network. Screeplot approach works well in general situation, but it does not necessarily always guarantee the correct estimate of number of communities. We need a more reliable and theoretically well understood way to determine K .

Secondly, we only consider the undirected case which is somewhat unrealistic assumption for citation network in real world, in a sense that it is not usual for two papers to cite each other at the same time. Since, in our research, we were interested in separating the low rank structure of edges and ad-hoc links in network, we did not take into account the directions of edges in our model. However, it would be interesting to consider incorporating the directed network into our matrix decomposition framework.

Last but not least, when we assign the memberships of each nodes, we use k -means clustering algorithm. However, k -means algorithm turns out that it tends to assign nodes conservatively to each communities. For example, in Fig4 (*left*), we can see that bunch of Multiple Testing papers are assigned as Mixed cluster, and in Fig4 (*right*) also, many papers which should have been classified to among three communities other than mixed topic, have been assigned as Mixed topic community. This is probably because k -means does not allow overlapping membership of each nodes. It would be interesting to see what happens if we apply some clustering methods which allow the overlapped memberships of nodes in network.

8 | APPENDIX

8.1 | Proof of Theorem 1

In proof, we use the notion of *decomposable regularizer* which is studied at length in the work [16]. Let (M, M^\perp) denote arbitrary subspace pair for which component-wise L_1 norm is decomposable. Throughout the proof, we adopt the convenient short-hand notation on projection of matrix P on subspace M as P_M . And also note that proof relies on following lemma:

Lemma 1. If a pair of regularization parameters (δ, γ) satisfies condition (12), then for $\mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S)$, we have

$$\mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S) \leq \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F + 3\mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + 4 \sum_{j=k+1}^n \sigma_j(L^*) + 4\frac{\gamma}{\delta} \|S_{M^\perp}^*\|_1$$

proof of Theorem 1.

Proof. Since $\hat{\Theta}$ and Θ^* are optimal and feasible (respectively) for the convex program (??), we have

$$h(\hat{\Theta}) + \delta \|\hat{L}\|_* + \gamma \|\hat{S}\|_1 \leq h(\Theta^*) + \delta \|L^*\|_* + \gamma \|S^*\|_1 \quad (18)$$

Through the assumption of strong convexity on $h(\Theta)$, and by the Taylor expansion, we can get a following lower bound on the term $h(\hat{\Theta}) - h(\Theta^*)$:

$$h(\hat{\Theta}) - h(\Theta^*) \geq \langle \nabla_\Theta h(\Theta^*), \hat{\Theta} - \Theta^* \rangle + \frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2$$

By rearranging the term in (18) and plugging in above inequality relation, we get:

$$\frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2 \leq -\langle \nabla_\Theta h(\Theta^*), \hat{\Theta} - \Theta^* \rangle + \delta \|L^*\|_* + \gamma \|S^*\|_1 - \delta \|\hat{L}\|_* - \gamma \|\hat{S}\|_1 \quad (19)$$

Here, we introduce another notation, for any pair of positive tuning parameters (δ, γ) which is defined as the weighted combination of the two regularizers :

$$\mathbb{Q}(L, S) := \|L\|_* + \frac{\gamma}{\delta} \|S\|_1$$

Through the definition of \mathbb{Q} , we can rewrite (19) as follows:

$$\frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2 \leq -\langle \nabla_\Theta h(\Theta^*), \hat{\Theta} - \Theta^* \rangle + \delta \mathbb{Q}(L^*, S^*) - \delta \mathbb{Q}(\hat{\Delta}^L + L^*, \hat{\Delta}^S + S^*) \quad (20)$$

According to Agarwal et al [1]'s second element of lemma 1, the difference $\mathbb{Q}(L^*, S^*) - \mathbb{Q}(\hat{\Delta}^L + L^*, \hat{\Delta}^S + S^*)$ is upper-bounded by

$$\mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) - \mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S) + 2 \sum_{j=k+1}^n \sigma_j(L^*) + 2 \frac{\gamma}{\delta} \|S_{M^\perp}^*\|_1 \quad (21)$$

First, we want to control upper bound of the term $-\langle \nabla_\Theta h(\Theta^*), \hat{\Theta} - \Theta^* \rangle$ in (20).

$$\begin{aligned} -\langle \nabla_\Theta h(\Theta^*), \hat{\Theta} - \Theta^* \rangle &= \left\langle \frac{1}{n} (X - P^*), \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} + \hat{\Delta}^L + \hat{\Delta}^S \right\rangle \\ &\leq \frac{1}{n} \|X - P^*\|_{op} (\|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_* + \|\hat{\Delta}^L\|_*) + \frac{1}{n} \|X - P^*\|_\infty \|\hat{\Delta}^S\|_1 \\ &\leq \frac{1}{n} \|X - P^*\|_{op} (\|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_F + \|\hat{\Delta}_A^L\|_* + \|\hat{\Delta}_B^L\|_*) + \frac{1}{n} \|X - P^*\|_\infty (\|\hat{\Delta}_M^S\|_1 + \|\hat{\Delta}_{M^\perp}^S\|_1) \\ &\leq \frac{\delta}{2} (\|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_F + \|\hat{\Delta}_A^L\|_* + \|\hat{\Delta}_B^L\|_*) + \frac{\gamma}{2} (\|\hat{\Delta}_M^S\|_1 + \|\hat{\Delta}_{M^\perp}^S\|_1) \end{aligned} \quad (22)$$

Combining the inequalities (21) and (22), we can obtain the upper bound of RHS in (20) as follows:

$$\frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2 \leq \frac{\delta}{2} \|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_F + \frac{3\delta}{2} \mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + 2\delta \sum_{j=k+1}^n \sigma_j(L^*) + 2\gamma \|S_{M^\perp}^*\|_1 \quad (23)$$

Second, we wish to control the lower bound of the term $\frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2$ with respect to $\hat{\Delta}^\alpha, \hat{\Delta}^L, \hat{\Delta}^S$.

$$\begin{aligned} \|\hat{\Delta}^\Theta\|_F^2 &= \|\hat{\Theta} - \Theta^*\|_F^2 \\ &= \|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} + \hat{\Delta}^L + \hat{\Delta}^S\|_F^2 \\ &= \|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_F^2 + \|\hat{\Delta}^L + \hat{\Delta}^S\|_F^2 + 2\langle \hat{\Delta}^L + \hat{\Delta}^S, \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} \rangle \\ &= \|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_F^2 + \|\hat{\Delta}^L\|_F^2 + \|\hat{\Delta}^S\|_F^2 + 2\langle \hat{\Delta}^L + \hat{\Delta}^S, \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} \rangle + 2\langle \hat{\Delta}^L, \hat{\Delta}^S \rangle \end{aligned} \quad (24)$$

We want to get the further lower bound on trace inner product terms, $\langle \hat{\Delta}^L + \hat{\Delta}^S, \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} \rangle$, $\langle \hat{\Delta}^L, \hat{\Delta}^S \rangle$. To control the first trace inner product term, we use the relation $\hat{\Delta}^L \mathbb{1} = 0$, apply the definition of dual norm on inner product term, apply triangular inequality on $\hat{\Delta}^\alpha$, and lastly we apply the constraint imposed on $|\alpha|$ stated in Assumption 2.

$$\begin{aligned} |\langle \hat{\Delta}^L + \hat{\Delta}^S, \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} \rangle| &= |\langle \hat{\Delta}^S, \hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T} \rangle| \\ &\leq \|\hat{\Delta}^{\alpha \mathbb{1} \mathbb{1}^T}\|_\infty \|\hat{\Delta}^S\|_1 \\ &\leq (|\hat{\alpha}| + |\alpha^*|) \|\hat{\Delta}^S\|_1 \\ &\leq 2C\kappa \|\hat{\Delta}^S\|_1 \end{aligned} \quad (25)$$

To control the term $\langle \hat{\Delta}^L, \hat{\Delta}^S \rangle$, we first apply the definition of dual norm on trace inner product term, then apply triangular inequality on $\hat{\Delta}^L$ and spikiness condition.

$$\begin{aligned} |\langle \hat{\Delta}^L, \hat{\Delta}^S \rangle| &\leq \|\hat{\Delta}^L\|_\infty \|\hat{\Delta}^S\|_1 \\ &\leq (\|\hat{L}\|_\infty + \|L^*\|_\infty) \|\hat{\Delta}^S\|_1 \\ &\leq \left(\frac{2\kappa}{n}\right) \|\hat{\Delta}^S\|_1 \end{aligned} \quad (26)$$

We can combine the inequality (24), (25) and (26). Then applying the assumption on regularization parameter γ , and the fact $\|\hat{\Delta}^L\|_* \geq 0$ sequentially, we can get,

$$\begin{aligned}
\frac{\tau}{2} \|\hat{\Delta}^\Theta\|_F^2 &\geq \frac{\tau}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^L\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^S\|_F^2 - 2\kappa\tau\left(\frac{Cn+1}{n}\right) \|\hat{\Delta}^S\|_1 \\
&\geq \frac{\tau}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^L\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^S\|_F^2 - \frac{\gamma}{2} \|\hat{\Delta}^S\|_1 \\
&\geq \frac{\tau}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^L\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^S\|_F^2 - \frac{\delta}{2} \mathbb{Q}(\hat{\Delta}^L, \hat{\Delta}^S)
\end{aligned} \tag{27}$$

By combining the relations (23) and (27), applying triangular inequality, $\mathbb{Q}(\hat{\Delta}^L, \hat{\Delta}^S) \leq \mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + \mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S)$, and rearranging the term, we can get following inequality,

$$\frac{\tau}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^L\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^S\|_F^2 \leq \frac{\delta}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F + 2\mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + \frac{\delta}{2} \mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S) + 2\delta \sum_{j=k+1}^n \sigma_j(L^*) + 2\gamma \|S_{M^\perp}^*\|_1$$

Further, by plugging in Lemma 1 to get an upper bound on $\mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S)$, we can rewrite the above inequality as follows:

$$\frac{\tau}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^L\|_F^2 + \frac{\tau}{2} \|\hat{\Delta}^S\|_F^2 - \frac{\delta}{2} \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F \leq \frac{7\delta}{2} \mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + 4\delta \sum_{j=k+1}^n \sigma_j(L^*) + 4\gamma \|S_{M^\perp}^*\|_1 \tag{28}$$

Noting that $\hat{\Delta}_A^L$ has rank at most $2k$ and that $\hat{\Delta}_M^S$ lies in the model space M , we find that

$$\begin{aligned}
\delta \mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) &\leq \sqrt{2k\delta} \|\hat{\Delta}_A^L\|_F + \Psi(M)\gamma \|\hat{\Delta}_M^S\|_F \\
&\leq \sqrt{2k\delta} \|\hat{\Delta}^L\|_F + \Psi(M)\gamma \|\hat{\Delta}^S\|_F
\end{aligned} \tag{29}$$

Here $\Psi(M)$ measures the compatibility between Frobenius norm and component-wise L_1 regularizer, where M is an arbitrary subset of matrix indices of cardinality at most s .

$$\Psi(M) := \sup_{U \in M, U \neq 0} \frac{\|U\|_1}{\|U\|_F}$$

Using Cauchy-Schwarz inequality, we can easily check the quantity $\Psi(M)$ is bounded by at most \sqrt{s} . Plugging in the relation (29) into (28) and rearranging the term relevant with $e^2(\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T, \hat{\Delta}^L, \hat{\Delta}^S)$ yield the claim. \square

8.2 | Proof of Lemma 1

Proof. Through the application of basic inequality by using optimality of $\hat{\Theta}$ and feasibility of Θ^* to convex program (??), we have

$$h(\hat{\Theta}) - h(\Theta^*) \leq \delta \mathbb{Q}(L^*, S^*) - \delta \mathbb{Q}(\hat{\Delta}^L + L^*, \hat{\Delta}^S + S^*) \tag{30}$$

By using convexity of $h(\Theta)$, we can write

$$\begin{aligned}
h(\hat{\Theta}) - h(\Theta^*) &\geq \langle \nabla_{\Theta} h(\Theta^*), \hat{\Theta} - \Theta^* \rangle \\
&= -\langle \frac{1}{n}(X - P^*), \hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T + \hat{\Delta}^L + \hat{\Delta}^S \rangle \\
&\geq -\frac{1}{n} \|X - P^*\|_{op} (\|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_* + \|\hat{\Delta}^L\|_*) + \frac{1}{n} \|X - P^*\|_\infty \|\hat{\Delta}^S\|_1 \\
&\geq -\frac{\delta}{2} (\|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F + \|\hat{\Delta}_A^L\|_* + \|\hat{\Delta}_B^L\|_*) - \frac{\delta}{2} (\|\hat{\Delta}_M^S\|_1 + \|\hat{\Delta}_{M^\perp}^S\|_1)
\end{aligned} \tag{31}$$

One more round of application on Agarwal et al [1]'s second element of lemma 1, we can get an upper bound of difference $\mathbb{Q}(L^*, S^*) - \mathbb{Q}(\hat{\Delta}^L + L^*, \hat{\Delta}^S + S^*)$ as follows:

$$\mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) - \mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S) + 2 \sum_{j=k+1}^n \sigma_j(L^*) + 2\frac{\gamma}{\delta} \|S_{M^\perp}^*\|_1 \tag{32}$$

By combining relations (31) and (32), we can get the upper bound of $\mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S)$:

$$\mathbb{Q}(\hat{\Delta}_B^L, \hat{\Delta}_{M^\perp}^S) \leq \|\hat{\Delta}^\alpha \mathbb{1} \mathbb{1}^T\|_F + 3\mathbb{Q}(\hat{\Delta}_A^L, \hat{\Delta}_M^S) + 4 \sum_{j=k+1}^n \sigma_j(L^*) + 4\frac{\gamma}{\delta} \|S_{M^\perp}^*\|_1$$

\square

Bibliography

- [1] Agarwal, A., S. Negahban, M. J. Wainwright, et al., 2012: Noisy matrix decomposition via convex relaxation: Optimal rates in high dimensions. *The Annals of Statistics*, **40**, no. 2, 1171–1197.
- [2] Bach, F. R., 2008: Consistency of trace norm minimization. *Journal of Machine Learning Research*, **9**, no. Jun, 1019–1048.
- [3] Boyd, S., N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al., 2011: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, **3**, no. 1, 1–122.
- [4] Candès, E. J., X. Li, Y. Ma, and J. Wright, 2011: Robust principal component analysis? *Journal of the ACM (JACM)*, **58**, no. 3, 11.
- [5] Chandrasekaran, V., P. A. Parrilo, and A. S. Willsky, 2010: Latent variable graphical model selection via convex optimization. *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, IEEE, 1610–1613.
- [6] Chen, Y., X. Li, J. Liu, and Z. Ying, 2016: A fused latent and graphical model for multivariate binary data. *arXiv preprint arXiv:1606.08925*.
- [7] Fazel, M., H. Hindi, S. P. Boyd, et al., 2001: A rank minimization heuristic with application to minimum order system approximation. *Proceedings of the American control conference*, Citeseer, volume 6, 4734–4739.
- [8] Gabay, D. and B. Mercier, 1975: *A dual algorithm for the solution of non linear variational problems via finite element approximation*. Institut de recherche d’informatique et d’automatique.
- [9] Glowinski, R. and A. Marrocco, 1975: On the solution of a class of nonlinear Dirichlet problems by a penalty-duality method and finite elements of order one. *Optimization Techniques IFIP Technical Conference*, Springer, 327–333.
- [10] Harman, H. H., 1960: *Modern factor analysis*. Univ. of Chicago Press.
- [11] Ji, P., J. Jin, et al., 2016: Coauthorship and citation networks for statisticians. *The Annals of Applied Statistics*, **10**, no. 4, 1779–1812.
- [12] Jöreskog, K. G., 1969: A general approach to confirmatory maximum likelihood factor analysis. *Psychometrika*, **34**, no. 2, 183–202.
- [13] — 1970: A general method for estimating a linear structural equation system. *ETS Research Bulletin Series*, **1970**, no. 2, i–41.
- [14] Lord, F. M. and M. R. Novick, 2008: *Statistical theories of mental test scores*. IAP.
- [15] McDonald, R. P., 2014: *Factor analysis and related methods*. Psychology Press.
- [16] Negahban, S. N., P. Ravikumar, M. J. Wainwright, B. Yu, et al., 2012: A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Statistical Science*, **27**, no. 4, 538–557.
- [17] Rasch, G., 1980: Probabilistic models for some intelligence and attainment tests. 1960. *Copenhagen, Denmark: Danish Institute for Educational Research*.
- [18] Schwarz, G. et al., 1978: Estimating the dimension of a model. *The annals of statistics*, **6**, no. 2, 461–464.
- [19] Zhou, Z., X. Li, J. Wright, E. Candès, and Y. Ma, 2010: Stable principal component pursuit. *2010 IEEE international symposium on information theory*, IEEE, 1518–1522.

How to cite this article: Suh Namjoon., Xiaoming Huo, Eric Heim, Timothy Van Slyke, and Lee Seversky (2019), Citation network, *journal.*, 2019;00:1–6.