Covariate Regularized Community Detection in Sparse Graphs

Bowei Yan and Purnamrita Sarkar University of Texas at Austin

Abstract

In this paper, we investigate community detection in networks in the presence of node covariates. In many instances, covariates and networks individually only give a partial view of the cluster structure. One needs to jointly infer the full cluster structure by considering both. In statistics, an emerging body of work has been focused on combining information from both the edges in the network and the node covariates to infer community memberships. However, so far the theoretical guarantees have been established in the dense regime, where the network can lead to perfect clustering under a broad parameter regime, and hence the role of covariates is often not clear. In this paper, we examine sparse networks in conjunction with finite dimensional sub-gaussian mixtures as covariates under moderate separation conditions. In this setting each individual source can only cluster a non-vanishing fraction of nodes correctly. We propose a simple optimization framework which provably improves clustering accuracy when the two sources carry partial information about the cluster memberships, and hence perform poorly on their own. Our optimization problem can be solved using scalable convex optimization algorithms. Using a variety of simulated and real data examples, we show that the proposed method outperforms other existing methodology.

Keywords: stochastic block models, kernel method, semidefinite programming, sub-gaussian mixture, asymptotic analysis

1 Introduction

Community detection in networks is a fundamental problem in machine learning and statistics. A variety of important practical problems like analyzing socio-political ties among leading politicians (Gil-Mendieta & Schmidt 1996), understanding brain graphs arising from diffusion MRI data (Binkiewicz et al. 2017), investigating ecological relationships between different tiers of the food chain (Jacob et al. 2011) can be framed as community detection problems. Much attention has been focused on developing models and methodology to recover latent community memberships. Among generative models, the stochastic block model (Holland et al. 1983) and its variants (Airoldi et al. (2008) etc.) have attracted a lot of attention, since their simplicity facilitates efficient algorithms and asymptotic analysis (Rohe et al. 2011, Amini et al. 2013, Chen & Xu 2016).

Although most real world network datasets come with covariate information associated with nodes, existing approaches are primarily focused on using the network for inferring the hidden community memberships or labels. Take for example the Mexican political elites network (described in detail in Section 4). This dataset comprises of 35 politicians (military or civilian) and their connections. The associated covariate for each politician is the year when one came into power.

After the military coup in 1913, the political arena was dominated by the military. In 1946, the first civilian president since the coup was elected. Hence those who came into power later are more likely to be civilians. Politicians who have similar number of connections to the military and civilian groups are hard to classify from the network alone. Here the temporal covariate is crucial in resolving which group they belong to. On the other hand, politicians who came into power around 1940's, are ambiguous to classify using covariates. Hence the number of connections to the two groups in the network helps in classifying these nodes. Our method can successfully classify these politicians and has higher classification accuracy than existing methods (Binkiewicz et al. 2017, Zhang et al. 2016).

In Statistics literature, there has been some interesting work on combining covariates and dense networks (average degree growing faster than logarithm of the number of nodes). In Binkiewicz et al. (2017), the authors present assortative covariate-assisted spectral clustering (ACASC) where one does Spectral Clustering on the the gram matrix of the covariates plus the regularized graph Laplacian weighted by a tuning parameter. A joint criterion for community detection (JCDC) with covariates is proposed by Zhang et al. (2016), which could be seen as a covariate reweighted Newman-Girvan modularity. This approach enables learning different influence on each covariate. In concurrent work Weng & Feng (2016) provide a variational approach for community detection.

All of the above works are carried out in the dense regime with strong separability conditions on the linkage probabilities. ACASC also requires the number of dimensions of covariates to grow with the number of nodes for establishing consistency.

In contrast, we prove our result for sparse graphs where the average degree is constant and the the covariates are finite dimensional sub-gaussian mixtures with moderate separability conditions. In our setting, neither source can yield consistent clustering in the limit. We show that combining the two sources leads to improved upper bounds on clustering accuracy under weaker conditions on separability on each individual source.

Leveraging information from multiple sources have been long studied in Machine learning and Data mining under the general envelop of multi-view clustering methods. Kumar et al. (2011) use a regularization framework so that the clustering adheres to the dissimilarity of clustering from each view. Liu et al. (2013) optimize the nonnegative matrix factorization loss function on each view, plus a regularization forcing the factors from each view to be close to each other. The only provable method is by Chaudhuri et al. (2009), where the authors obtain guarantees where the two views are mixtures of Log-concave distributions. This algorithm does not apply to networks.

In this paper, we propose a penalized optimization framework for community detection when node covariates are present. We take the sparse degree regime of Stochastic Blockmodels, where one can only correctly cluster a non-vanishing fraction of nodes. Similarly, for covariates, we assume that the covariates are generated from a finite dimensional sub-gaussian mixture with moderate separability conditions. We prove that our method leads to an improved clustering accuracy under weaker conditions on the separation between clusters from each source. As byproducts of our theoretical analysis we obtain new asymptotic results for sparse networks under weak separability conditions and kernel clustering of finite dimensional mixture of sub-gaussians. Using a variety of real world and simulated data examples, we show that our method often outperforms existing methods. Using simulations, we also illustrate that when the two sources only have partial and in some sense orthogonal information about the clusterings, combining them leads to better clustering than using the individual sources.

In Section 2, we introduce relevant notation and present our optimization framework. In Sec-

tion 3, we present our main results, followed by experimental results on simulations and real world networks in Section 4. Majority of the proofs are presented in the appendix, with details deferred to the supplementary.

2 Problem Setup

In this section, we introduce our model and set up the convex relaxation framework. For clarity, we list all definitions and notations that will be used later in Table 3.

Assume (C_1, \dots, C_r) represent a r-partition for n nodes $\{1, \dots, n\}$. Let $m_i = |C_i|$ be the size of cluster i, and let m_{\min} and m_{\max} be the minimum and maximum cluster sizes respectively. We use $\pi_i := \frac{m_i}{n}$, $\pi_{\min} = \frac{m_{\min}}{n}$ and $\alpha = m_{\max}/m_{\min}$. We denote by A the $n \times n$ binary adjacency matrix and by Y the $n \times d$ matrix of d dimensional covariates. The generation of A and Y share the true and unknown membership matrix $Z = \{0,1\}^{n \times r}$. We define the graph model as:

(Graph Model)
$$P(A_{ij} = 1|Z) = Z_i^T B Z_j$$
 For $i \neq j$ (1)

B is a $r \times r$ matrix of within and across cluster connection probabilities. Furthermore $A_{ii} = 0, \forall i \in [n]$. We consider the sparse regime where $n \max_{k\ell} B_{k\ell}$ is a constant and hence average expected degree is also a constant w.r.t n. Amini et al. (2018) define two different classes of block models in terms of separability properties of B. We state this below.

Definition 1. A stochastic block model is called strongly assortative if $\min_k B_{kk} > \max_{k \neq \ell} B_{k\ell}$. It is called weakly assortative if $\forall k \neq \ell$, $B_{kk} > B_{k\ell}$.

This distinction is important because the weakly assortative class of blockmodels is a superset of strongly assortative models, and most of the analysis are done in the stronger setting. To our knowledge, there has not been any work on weakly assortative blockmodels in the sparse setting. For the covariates, we define,

(Covariate Model)
$$Y_i = \sum_{a=1}^r Z_{ia}\mu_a + W_i$$
 (2)

 W_i are mean zero d dimensional sub-gaussian vectors with spherical covariance matrices $\sigma_k^2 I_d$ and sub-gaussian norm ψ_k (for $i \in C_k$). Standard definitions of sub-gaussian random variables (for more detail see Vershynin (2010)) are provided in the Supplementary material. We define the distance between clusters C_k and C_ℓ as $d_{k\ell} = \|\mu_k - \mu_\ell\|$ and the separation as $d_{\min} = \min_{k \neq \ell} d_{k\ell}$.

Notation	Mathematical Definition	Explanation
$A \in \{0,1\}^{n \times n}$	$A_{ij} i \in C_k, j \in C_\ell \sim Ber(B_{k\ell})$	Adjacency matrix (Symmetric)
$Y_i \in \mathbb{R}^d$		Covariate observation for <i>i</i> th point
$K \in [0,1]^{n \times n}$	$K(i,j) = f(Y_i - Y_j _2^2)$	Kernel matrix, symmetric and positive definite

Table 2: Random variables used in the paper

Notation	Mathematical Definition	Explanation		
n, d		Number of nodes, dimensionality of covariates		
I_d		identity matrix of size $d \times d$		
$\operatorname{diag}(v_1,\ldots,v_k)\in\mathbb{R}^{k\times k}$		Diagonal matrix with diagonal (v_1, \ldots, v_k)		
r	$\Theta(1)$	Number of clusters		
$B \in [0,1]^{r \times r}$	$\Theta(1/n)$	Symmetric Probability matrix in SBM		
$Z \in \{0,1\}^{n \times r}$		Latent class memberships		
m_i	$\sum_{j} Z(j,i)$	Number of points in <i>i</i> th cluster		
π_i	$\frac{m_i}{n}$	Proportion of points in <i>i</i> th cluster		
$m_{\rm max}$	$\max_k m_k, \Theta(n)$	Largest cluster size		
m_{\min}	$\min_k m_k, \Theta(n)$	Smallest cluster size		
α	$m_{\mathrm{max}}/m_{\mathrm{min}},\Theta(1)$	Ratio between largest and smallest clusters		
C_k	$\{j: Z(j,i) = 1\}$	Point set for kth cluster		
$X_0 \in \mathbb{R}^{n \times n}$	$Z \mathrm{diag}(1/m_1,\ldots,1/m_r)Z^T$	Ground truth clustering matrix		
$a_k, b_k = \Theta(1)$	$a_k = nB_{kk}, b_k = n \max_{\ell \neq k} B_{k\ell}$	Rescaled probabilities		
$g \in \mathbb{R}$	$\frac{2}{n-1}\sum_{i\leq j} Var(A_{ij}), \Theta(1)$	Average variance of Graph edges		
$\mu_k, \sigma_k I_d$		Mean, covariance matrix for Y_i if $i \in C_k$		
ψ_k		subgaussian norm for Y_i if $i \in C_k$		
$d_{k\ell}$	$\ \mu_k-\mu_\ell\ $	Distance between cluster centers for the covariates		
K_I	Eq. (11)	Reference matrix for the kernel		
$ u_k$	Eq. (6)	Separation in K_I		
γ	$\min_k (a_k - b_k + \lambda \nu_k), \Theta(1)$	Separation of $ZBZ^T + \lambda K$		

Table 1: Population quantities used in the paper

Notation For a matrix $M \in \mathbb{R}^{n \times n}$, we use $||M||_F$ and ||M|| to denote the Frobenius and operator norms of M respectively. The ℓ_{∞} norm is defined as: $||M||_{\infty} = \max_{i,j} |M_{ij}|$. For two matrices $M, Q \in \mathbb{C}^{m \times n}$, their inner product is $\langle M, Q \rangle = \operatorname{trace}(M^TQ)$. The $\ell_{\infty} \to \ell_1$ norm of a matrix M is defined as $||M||_{\ell_{\infty} \to \ell_1} = \max_{\|s\|_{\infty} \le 1} ||Ms||_1$. From now on we use I_n to denote the identity matrix of size n, $\mathbf{1}_n$ to represent the all one n-vector and $E_n, E_{n,k}$ to represent the all one matrix with size $n \times n$ and $n \times k$ respectively. We use standard order notations O, o, Ω, ω , etc. For example, we use $t(n) = \Theta(1/n)$ to denote that $t(n) \times n$ is a constant w.r.t n. We also use \tilde{O} notation to exclude multiplicative terms that are logarithmic in n.

2.1 Optimization Framework

We now present our optimization framework. There are many available semidefinite programming (SDP) relaxations for clustering blockmodels (Amini et al. 2018, Cai & Li 2015, Chen & Xu 2016). The common element in all of these is maximizing the inner product between A and X, for a positive semidefinite matrix X. Here X is a stand-in for the clustering matrix ZZ^T . Unequal-sized clusters is usually tackled with an extra regularization term added to the objective function (see Hajek et al. (2016), Perry & Wein (2017), Cai & Li (2015) among others). While the above consistency results are for dense graphs, Guédon & Vershynin (2015), Montanari & Sen (2016) show that in the sparse regime one can use this method to obtain an error rate which is a constant w.r.t n and

Notation	Mathematical Definition	Explanation	
1_n		All one vector of length n	
E_n	$\mid 1_n 1_n^T$	All ones matrix of size $n \times n$	
I_d		Identity matrix of size $d \times d$	
K_G	≤ 1.783	Grothendieck's constant	
$f(x): \mathbb{R}_+ \to [0,1]$		Kernel function	
\mathcal{F}	$\begin{cases} \{X \succeq 0, 0 \le X \le \frac{1}{m_{\min}}, \\ X1_n = 1_n, \operatorname{trace}(X) = r \} \end{cases}$	Feasible set of the SDP	
X_M	$ \arg \max_X \langle M, X \rangle $ s.t. $X \in \mathcal{F}$	Solution matrix of the SDP	
$\theta_i(M)$		i-th eigenvalue of M	
λ_n, λ_0	$\lambda_n = \lambda_0/n, \lambda_0 = \Theta(1)$	Tuning parameter between graph and covariates	

Table 3: Useful notations and definitions

depends on the gap between the within and across cluster probabilities.

SDPs are not only limited to network clustering. Several convex relaxations for k-means type loss are proposed in the literature (see Peng & Wei (2007), Mixon et al. (2017), Yan & Sarkar (2016) for more references). In particular in these settings one maximizes $\langle W, X \rangle$, for some positive semidefinite matrix X, where W is a matrix of similarities between pairwise data points. For classical k-means W_{ij} can be $Y_i^T Y_j$ whereas for k-means in the kernel space one uses a suitably defined kernel similarity function between the ith and jth covariates. We analyze the widely-used Gaussian kernel to allow for non-linear boundaries between clusters. Let K be the $n \times n$ kernel matrix whose (i,j)-th entry is $K(i,j) = f(\|Y_i - Y_j\|_2^2)$, where $f(\cdot)$, where $f(x) = \exp(-\eta x)$ for $x \ge 0$. This kernel function is upper bounded by 1 and is Lipschitz continuous w.r.t. the distance between two observations. Furthermore, in contrast to network based SDPs, the above uses X as a stand in for the normalized variant of the clustering matrix ZZ^T , i.e. the desired solution is $(X_0)_{ij} = \frac{1(k=\ell)}{m_k}$, if $i \in C_k$, $j \in C_\ell$. It can be seen that $\|X_0\|_F^2 = r$.

In our optimization framework, we propose to add a k-means type regularization term to the network objective, which enforces that the estimated clusters are consistent with the latent memberships in the covariate space.

$$X = \arg\max_{X} \langle A + \lambda_n K, X \rangle \quad s.t. \quad X \in \mathcal{F}, \tag{3}$$

where λ_n is a tuning parameter (possibly depending on n) and the constraint set $\mathcal{F} = \{X \succeq 0, 0 \le X \le \frac{1}{m_{\min}}, X \mathbf{1}_n = \mathbf{1}_n, \operatorname{trace}(X) = r\}$ is similar to Peng & Wei (2007). The m_{\min} in the constraint can be replaced by any lower bound on the smallest cluster size, and is mainly of convenience for the analysis. In the implementation, it suffices to enforce the elementwise positivity constraints, and other linear constraints. For ease of exposition, we define

$$X_M = \arg\max_{X} \langle M, X \rangle \quad s.t. \quad X \in \mathcal{F},$$
 (4)

When $K(i, j) = Y_i^T Y_j$, then the non-convex variant of the objective function naturally assumes a form similar to the work of ACASC (modulo normalization of A).

3 Main Results

Typically in existing SDP literature for sparse networks or subgaussian mixtures (Guédon & Vershynin 2015, Mixon et al. 2017), one obtains a relative error bound of the deviation of X_M (the

solution of the SDP) from the ideal clustering matrix X_0 . This relative error is typically proportional to the ratio of the observed matrix with a suitably defined reference matrix, and some quantity which measures the separation between the different clusters. Our theoretical result shows that the relative error of the solution to the combined SDP is proportional to the ratio of the observed $A + \lambda_n K$ matrix to a suitably defined reference matrix to a quantity which measures separation between clusters. This quantity is a non-linear combination of the separations stemming from the two sources. We first present an informal version of the main result. Main theorem (informal): Let $X_{A+\lambda_n K}$ be the solution of SDP (4). Let s_G^k and s_C^k be constants denoting the separations of cluster k from the other clusters defined in terms of the model parameters of the network and the covariates respectively. If the tuning parameter $\lambda_n = \lambda_0/n$ for some constant λ_0 , then

$$||X_{A+\lambda_n K} - X_0||_F^2 \le \frac{c_G + \lambda_0 c_C}{\min_k \left(s_G^k + \ell s_C^k\right)},$$

where c_G and c_C are constants representing the error corresponding to the graph and the covariates.

Note that in SBM, the separation is well-defined, i.e. when M=A, a natural choice of the reference matrix is E[A|Z] which is blockwise constant. In this case, the separation is given by $\min_k(B_{kk} - \max_\ell B_{k\ell})$, and leads to a result on weakly assortative sparse block models which we present in more details in Section 3.1. However, for the kernel matrix K, the main difficulty is that one cannot achieve element-wise or operator norm concentration of K (also discussed in Von Luxburg et al. (2008)). This makes the choice of the reference matrix difficult. To better understand the role of the separation parameter, we first present a key technical lemma bounding $||X_M - X_0||_F$. The main goal of this lemma is to establish an upper bound on the frobenius norm difference between the solution to an SDP with input matrix M to the ideal clustering matrix.

Lemma 1. Let X_M be defined by Eq (4) for some input matrix M. Also let Q be a reference matrix where $Q_{ij} = \beta_k^{(in)}, \forall i, j \in C_k$, and $\beta_k^{(out)} \ge Q_{ij} \ge 0, \forall i \in C_k, j \in C_\ell, k \ne \ell$. If $\min_k(\beta_k^{(in)} - \beta_k^{(out)}) \ge 0$, then

$$||X_M - X_0||_F^2 \le 2 \frac{\langle M - Q, X_M - X_0 \rangle}{m_{\min} \min_k (\beta_k^{(in)} - \beta_k^{(out)})}$$
(5)

Remark 1. The key to the above lemma is to find a suitable reference matrix Q which satisfies some separation conditions between the blocks. The deviation between X_M and X_0 is small if M-Q is small, and large if the separation between blocks in Q is small. While the proof technique is inspired by Guédon & Vershynin (2015), the details are different because of our use of different constraints and because our reference matrix Q does not have to be blockwise constant and can be weakly assortative instead of strongly assortative.

The results on networks, covariates and the combination of the two essentially reduces to identifying good reference matrices (Q) for the input matrices A, K, and $A + \lambda K$, which

- 1. Satisfies the properties of Q in the above lemma.
- 2. Has a large separation $\min_k(\beta_k^{(in)} \beta_k^{(out)})$ increasing the denominator of Eq. (5).
- 3. Has a small deviation from M, thereby reducing the numerator of Eq (5).

Now the main work is to choose the reference matrix Q for $A + \lambda K$. As pointed out before, a common choice for reference matrix of A is $\mathbb{E}[A|Z]$. For the covariates, we divide the nodes into "good" nodes $S_k := \{i \in C_k : ||Y_i - \mu_k|| \le \Delta_k\}$ and the rest. Also define $S = \bigcup_{k=1}^r S_k$. Δ_k will be defined such that the kernel matrix induced by the rows and columns in S is weakly assortative, and $3\Delta_k + \Delta_\ell \le d_{k\ell}$. Define

$$r_k := f(2\Delta_k), \quad s_k := \max_{\ell \neq k} f(d_{k\ell} - \Delta_k - \Delta_\ell), \qquad \nu_k = r_k - s_k \tag{6}$$

A simple use of triangle inequality gives $\min_{i,j\in\mathcal{S}_k} K_{ij} \geq r_k$ and $\max_{i\in\mathcal{S}_k,j\in\mathcal{S}_\ell,\ell\neq k} K_{ij} \leq s_k$. Hence the separation for cluster k is $\nu_k := r_k - s_k$. We define the reference matrix K_I as:

$$(K_I)_{ij} = \begin{cases} f(2\Delta_k), & \text{if } i, j \in C_k \\ \min\{f(d_{k\ell} - \Delta_k - \Delta_\ell), K_{ij}\}, & \text{if } i \in C_k, j \in C_\ell, k \neq \ell \end{cases}$$
 (7)

The choice of Δ_k is crucial. A large Δ_k makes the size of non-separable nodes \mathcal{S}^c small, but drives down the separation ν_k .

We are now ready to present our main result. As we will show in the proof, the new separation is $\gamma = \min_k \frac{(a_k - b_k) + \lambda_0 \nu_k}{n}$. Typically, in the general case with unequal sub-gaussian norms, one should benefit from using different Δ_k 's for different clusters. For example for a cluster with a large $a_k - b_k$, we can afford to have a small ν_k . To think in terms of Δ_k , for this cluster one can have a large Δ_k , which will make $|\mathcal{S}_k|$ larger than before, but will not affect the separation $(a_k - b_k) + \lambda_0 \nu_k$ of cluster k very detrimentally. We now present our first main theorem.

Theorem 1. Let $a_k = nB_{kk}$, $b_k = n\max_{\ell \neq k} B_{k\ell}$, $g := \frac{2}{n-1}\sum_{i < j} Var(a_{ij}) \geq 9$. Take $\lambda_n = \lambda_0/n$, $m_k = n\pi_k$, $m_{\min} = n\pi_{\min}$, and $\pi_0 := \sum_k (m_k \exp(-\Delta_k^2/5\psi_k^2) + \sqrt{m_k \log m_k/2})/n$. Let $X_{A+\lambda_n K}$ be defined as in Eq (4). If $\pi_{\min} = \Theta(1)$ and $\min_k (a_k - b_k + \lambda_0 \nu_k) > 0$, then, with probability tending to one,

$$||X_{A+\lambda K} - X_0||_F^2 \le 2K_G \frac{6\sqrt{g} + \lambda_0 \left(2\pi_0 + \sum_k \pi_k^2 (1 - f(2\Delta_k))\right)}{\pi_{\min}^2 \min_k (a_k - b_k + \lambda_0 \nu_k)},$$

where $\nu_k = f(2\Delta_k) - \max_{\ell \neq k} f(d_{k\ell} - \Delta_k - \Delta_\ell)$ for some $\Delta_k, \Delta_\ell \geq 0$ and $\max(\Delta_k, \Delta_\ell) \leq d_{k\ell}/4$.

Here K_G is the Grothendieck's constant. The best value of K_G is still unknown, and the best known bound is $K_G \leq 1.783$ (Braverman et al. 2013). First note that in the sparse case, we take $\lambda_n = \lambda_0/n$ for some constant λ_0 . In general the upper bound depends on several parameters such as λ_n and the scale parameter η in the gaussian kernel. We provide procedures for tuning λ_n and η in Section 4. The Δ_k 's show up in the numerator as well as the denominator. Finding the optimal Δ_k is cumbersome in the general case with unequal ψ_k 's. In Section 3.2 we derive an upper bound for equal Δ_k 's for concreteness.

Now we present two natural byproducts of our analysis, namely the result on graphs, i.e. bounds on $||X_0 - X_A||_F$ and the result on covariate clustering i.e. bounds on $||X_0 - X_K||_F$.

3.1 Result on Sparse Graph

While most dense network-based community detection schemes give perfect clustering in the limit (Amini et al. 2013, 2018, Cai & Li 2015, Chen & Xu 2016, Yan, Sarkar & Cheng 2017), in the sparse case no algorithm is consistent; however semidefinite relaxations (among others) can achieve an error rate governed by the within and across cluster probabilities (Guédon & Vershynin 2015, Montanari & Sen 2016). The sparse network analysis is done under strongly assortative settings.

Proposition 1 (Analysis for graph). Let a_k, b_k defined as in Theorem 1 are positive constants and $g \geq 9$. Then with probability tending to 1,

$$\frac{\|X_A - X_0\|_F}{\|X_0\|_F} \le \epsilon,$$

if
$$\min_k (a_k - b_k) \ge \frac{23\alpha^2 r \sqrt{g}}{\epsilon^2}$$
 where $\alpha := m_{\max}/m_{\min}$.

Note that in the above result, in order to have the error rate ϵ to go to zero, one would require $a_k - b_k$ to go to infinity, whereas by definition a_k, b_k are constants. Therefore one can only hope for a small albeit constant ϵ . In addition, both number of clusters r and the ratio between largest and smallest cluster sizes α needs to be constant order w.r.t n in order to guarantee the error rate does not increase when the network grows.

Remark 2 (Comparison with prior work). In contrast to having $\min_k a_k - \max_k b_k$ (strong assortativity) in the denominator like Guédon & Vershynin (2015), we have $\min_k (a_k - b_k)$ (weak assortativity), which allows for a much broader parameter regime.

3.2 Result on Covariates

We present a result for covariates analogous to the sparse graph setting, which establishes that, while SDP with covariates is not consistent with finite signal-to-noise ratio, it achieves a small error rate if the cluster centers are further apart. But before delving into our analysis, we provide a brief overview of existing work.

For covariate clustering, it is common to make distributional assumptions; usually a mixture model with well-separated centers suffices to show consistency. The most well-studied model is Gaussian mixture models, which can be inferred by Expectation-Maximization algorithm, for which recently there has been some local convergence results (Balakrishnan et al. 2017, Yan, Yin & Sarkar 2017) and its variants (Dasgupta & Schulman 2007). The condition required for provable recovery on the separation is usually the minimum distance between clusters is greater than some multiple of the square root of dimension (or effective dimension).

Another popular technique is based on SDP relaxations. For example, Peng & Wei (2007), Mixon et al. (2017) propose a SDP relaxation for k-means type clustering. To make the analysis concrete, for Proposition 2, we use $\Delta_k = \Delta$.

Proposition 2 (Analysis for Covariates). Let K be the kernel matrix generated from kernel function f. Denote ν_k as in Eq (6). If $\frac{d_{\min}}{\psi_{\max}} > \max\left\{\sqrt{d}, \frac{180}{\sqrt{d}}\right\}$, then with properly chosen η , with probability at least $1 - \sum_k \frac{1}{m_k}$,

$$\frac{\|X_K - X_0\|_F^2}{\|X_0\|_F^2} \le C\alpha^2 d \frac{\psi_{\max}^2}{d_{\min}^2} \max \left\{ \log \left(\frac{d_{\min}}{\psi_{\max} \sqrt{d}} \right), r \right\}$$

Remark 3 (Comparison with prior work). In recent work, Mixon et al. (2017) show the effectiveness of SDP relaxation with k-means clustering for sub-gaussian mixtures, provided the minimum distance between centers is greater than the standard deviation of the sub-gaussian times the number of clusters r. We provide a dimensionality reduction scheme, which also shows that the separation condition requires that $d_{\min} = \Omega(\sqrt{\min(r,d)})$. Our proof technique is new and involves carefully constructing a reference matrix for Lemma 1.

3.3 Analysis of Covariate Clustering when $d \gg r$

In high dimensional statistical problems, the signal is often assumed to lie in a low dimensional subspace or manifold. This is why much of Gaussian Mixture modeling literature first computes some projection of the data onto a low dimensional subspace (Vempala & Wang 2004). To reduce the dimensionality of the raw data, one could do a feature selection for the covariates (e.g. Jin et al. (2017), Verzelen et al. (2017)). In contrast, here we propose a much simpler dimensionality reduction step, which does not distort the pairwise distances between cluster means too much. The intuition is that, for clustering a subgaussian mixture, if $d \gg r$, the effective dimensionality of the data is r since the cluster means lie in an at most r-dimensional subspace.

Hence we propose the following simple dimensionality reduction algorithm when $d \gg r$ in a spirit similar to Chaudhuri et al. (2009). We split up the sample into two random subsets P_1 and P_2 of sizes n_1 and $n-n_1$ and compute the top r-1 eigenvectors U_{r-1} of the matrix $\hat{S} = \frac{\sum_{i \in P_1} (Y_i - \bar{Y})(Y_i - \bar{Y})^T}{n_1} \in \mathbb{R}^{d \times d}$, where $\bar{Y} = \frac{\sum_{i \in P_1} Y_i}{n_1}$. Now we project the covariates from subset P_2 onto this lower dimensional subspace as $Y_i' = U_{r-1}^T Y_i$ to get the low dimensional projections. We take $n_1 = n/\log n$.

Lemma 2. Let $M:=\sum_k \pi_k \mu_k \mu_k^T$. If $\sum_k \pi_k \mu_k = 0$, and the smallest eigenvalue of M satisfies $\theta_{r-1}(M) \geq 5\psi_{\max}^2 + C\sqrt{\frac{d\log^2 n}{n}}$ for some constant C, the projected Y_i' are also independent data points generated from an isotropic sub-gaussian mixture in r-1 dimensions. Furthermore the minimum distance between the means in the r-1 dimensional space is at least $d_{\min}/2$ with probability at least $1-\tilde{O}(r^2n^{-d})$, where d_{\min} is the separation in the original space.

The proof of this lemma is deferred to the supplementary material. We believe the proof can be generalized to non-spherical cases as long as the largest eigenvalue of covariance matrix for each cluster is bounded. Typically $\theta_{r-1}(M)$ signifies the amount of signal. For example, for the simple case of mixture of two gaussians with $\pi_1 = 1/2$, and $\mu_2 = -\mu_1$, $\theta_{r-1}(M) = \|\mu_1\|^2$, which is essentially $d_{\min}^2/4$. Hence the condition on $\theta_{r-1}(M)$ essentially translates to a lower bound on the signal to noise ratio, i.e. $d_{\min}^2 \geq 48\psi_{\max}^2 + C'\sqrt{\frac{d\log^2 n}{n}}$ for some constant C'. When d > r, if one applies Lemma 2 on the r-1 dimensional space, then as long as $d_{\min}^2 = \Omega(\psi_{\max}^2 r)$, the separation in the low dimensional space also satisfies the separation condition in Proposition 2. Thus the dimensionality reduction brings down the separation condition in Proposition 2 from $\Omega(\psi_{\max}\sqrt{d})$ to $\Omega(\psi_{\max}\sqrt{\min(r,d)})$.

The sample splitting is merely for theoretical convenience which ensures that the projection matrix and the projected data are independent, resulting in the fact that the final projection is also an independent sample from a sub-gaussian mixture. To be concrete, the labels of P_1 do not matter asymptotically, since they incur a relative error in $||X_0 - X_K||_F / ||X_0||_F$ less than $\sqrt{n^2/(m_{\min}^2 \log n)/\sqrt{r}} \le \sqrt{\alpha^2 r/\log n}$, where α and r are both constants. In our setting, the relative error in Proposition 2 is a small but non-vanishing constant, and so this additional vanishing error term does not affect it. However this sample splitting step is not necessary in practice (Chaudhuri et al. 2009), and so we do not pursue this further.

We now present the tuning procedure, and experimental results.

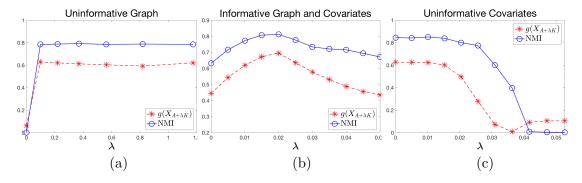


Figure 1: Tuning: (a) $B = 0.005E_3, n = 1000, d = 6, d_{\min} = 15\sigma$; (b) $d = 6, d_{\min} = 1.3, \sigma = (1, 1, 5), B = \text{diag}(0.004, 0.024, 0.024) + 0.004E_3$; (c) $d = 6, d_{\min} = 0, B = 0.0144I_3 + 0.0016E_3$.

4 Experiments

In this section, we present results on real and simulated data. The cluster labels in our method are obtained by spectral clustering of the solution matrix returned by the SDP. We will use SDP-comb, SDP-net, SDP-cov to represent the labels estimated from $X_{A+\lambda_n K}$, X_A and X_K respectively. Performance of the clustering is measured by normalized mutual information (NMI), which is defined as the mutual information of the two distributions divided by square root of the product of their entropies. We have also calculated classification accuracy and they show similar trends, so only NMI is reported in this section. For real and simulated data, we compare: (1) Covariate-assisted spectral clustering (ACASC) (Binkiewicz et al. 2017); (2) JCDC (Zhang et al. 2016), (3) SDP-comb, (4) SDP-net and (5) SDP-cov. The last two are used as references of graph-only and covariate-only clustering respectively.

4.1 Implementation and computational cost

Solving semidefinite programming with linear and non-linear constraints has been a challenging problems in numerical optimization community. Many SDPs proposed in statistical literature (Cai & Li 2015, Chen & Xu 2016, Amini et al. 2018) are solved by the alternating descent method of multipliers (ADMM) algorithm (Boyd et al. 2011). Although ADMM is tractable for middle-sized problems and reasonable numerical behavior, whether it convergences in presence of non-negative constraints, which is prevalent in network literatures, remains an open problem. Recently, Yang et al. (2015) propose a majorized semismooth Newton-CG augmented Lagrangian method, called SDPNAL+, which is provably convergent. We solve the SDP using the matlab package of SDPNAL+ in all our experiments¹. The package provides an efficient implementation of the algorithm. Solving the SDP for matrix of size 1000 × 1000 takes less than a minute on a Macbook with a 1.1 GHz Intel Core M processor.

4.2 Choice of Tuning Parameters

As we pointed out earlier, the elementwise upper bound $\frac{1}{m_{\min}}$ is only for convenience of theoretical analysis. In the implementation, we do not enforce this constraint. So the main tuning parameters

¹The code used for the experiment can be found at https://github.com/boweiYan/SDP_SBM_unbalanced_size.

would be the scale parameter in the kernel matrix η and the tradeoff parameter between graph and covariates λ_n . In most of our experiments the number of clusters is assumed known. In this section, we also provide a practical way to choose among candidates of r when it is not given.

Choice of η We use the method proposed in Shi et al. (2009) to select the scale parameter. The intuition is to keep enough (say 10%) of the data points in the "range" of the kernel for most (say 95%) data points. Given the covariates, we first compute the pairwise distance matrix. Then for each data point Y_i , compute q_i as 10% quantile of $d(Y_i, Y_j), \forall j \in [n]$. The bandwidth is defined as

$$w = \frac{95\% \text{ quantile of } q_i}{\sqrt{95\% \text{ quantile of } \chi_d^2}}$$

and scale parameter $\eta = \frac{1}{2w^2}$.

Note when the data is high-dimensional, we will first conduct dimensionality reduction as in Section 3.3, then use the intrinsic dimension to tune the scale parameter.

Choice of λ_n As λ_n increases, the resulting $X_{A+\lambda_n K}$ clustering gradually changes from X_A clustering to X_K clustering. Our theoretical results show that, with the right λ_n , $X_{A+\lambda_n K}$ and X_0 should be close, and hence also have similar eigenvalues. Let $\theta_i(M)$ be the *i*-th eigenvalue of matrix M. Define the eigen gap function for clustering matrices $g(X) := (\theta_r(X) - \theta_{r+1}(X))/\theta_r(X)$. Using Weyl's inequality and the fact that $\|X_{A+\lambda_n K} - X_0\|_{\text{op}} \leq \|X_{A+\lambda_n K} - X_0\|_F$, we have: $\theta_r(X_0) - \|X_{A+\lambda_n K} - X_0\|_F \leq \theta_r(X_{A+\lambda_n K}) \leq \theta_r(X_0) + \|X_{A+\lambda_n K} - X_0\|_F$. Since $g(X_0) = 1$, we pick the λ_n maximizing $g(X_{A+\lambda_n K})$. In Figure 1 (a)-(c), figures from left to right represent the situation where graph is uninformative (Erdős-Rényi), both are informative and covariates are uninformative. We plot $g(X_{A+\lambda_n K})$ and NMI of the clustering from $X_{A+\lambda_n K}$ with the true labels against λ_n . Figure 1 shows that $g(X_{A+\lambda_n K})$ and NMI of the predicted clustering have a similar trend, justifying the effectiveness of the tuning procedure.

Unknown number of clusters In many real world settings, it is generally hard to possess the knowledge of number of clusters. Methods are proposed for selecting number of blocks under sparse stochastic block models (Le & Levina 2015), but most of these methods are designed specific for graph adjacency matrix and cannot be generalized to continuous matrix scenarios. We observe that the eigen gap acts as an informative indicator for picking the number of clusters. So when the number of clusters is unknown, we run the SDP over a grid of λ_n, k , and choose the pair that maximizes the eigen gap. As we show in Figure 2, we construct two settings and test the performance of using eigen gap to select r. In the first setting, the true model has 3 clusterings

with proportion 3:4:5, the probability matrix is $B = 0.01 * \begin{bmatrix} 1.6 & 1.2 & 0.16 \\ 1.2 & 1.6 & 0.02 \\ 0.16 & 0.02 & 1.2 \end{bmatrix}$. And the covariates

are high dimensional Gaussian centered at $\mu_1 = (0, 2, 0 \cdots, 0)$, $\mu_2 = (-1, -0.8, 0 \cdots, 0)$, $\mu_3 = (1, -0.8, 0 \cdots, 0)$. We sample n = 800 data points, and run SDP on top of it with different choice of λ_n and specified number of clusters k. For each pair of parameter, we compute the NMI and eigengap and plot them on the upper and lower panel of Figure 2(a). As we can see, the eigengap presents a similar trend as the NMI, hence picking the pair that optimizes eigengap will have a relatively high NMI as well. Note here the mis-specified k = 2 has a higher NMI than that of the true value of r. This tells us even the number of clusters is mis-specified, the SDP is still able

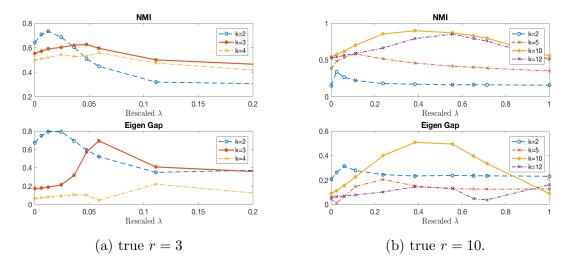


Figure 2: NMI and eigen gap for various choice of r.

to find structure that correlates with the underlying model. This phenomenon is also observed in several other works (Yan, Sarkar & Cheng 2017, Perry & Wein 2017).

In the second scenario, we generate a planted partition model with 10 equal-sized clusters, where $B = 0.046I_{10} + 0.004E_{10}$, along with Gaussian covariates centered at $[3 * I_{10} | \mathbf{0}_{3,90}]$. We conduct the same type of experiment as above and plot the NMI and eigengap. In this case, the eigen gap succussfully recovered the true number of clusters.

4.3 Simulation Studies

In this part we consider two simulation settings. In the first setting, we generate three clusters with sizes 3:4:5, with n=800. The probability matrix is $B=0.01*\begin{bmatrix} 1.6 & 1.2 & 0.16 \\ 1.2 & 1.6 & 0.02 \\ 0.16 & 0.02 & 1.2 \end{bmatrix}$, and the covariates

for each cluster are generated with 100 dimensional unit variance isotropic Gaussians, whose centers are only non-zero on the first two dimensions with $\mu_1 = (0, 2, 0 \cdots, 0)$, $\mu_2 = (-1, -0.8, 0 \cdots, 0)$, $\mu_3 = (1, -0.8, 0 \cdots, 0)$. This is the same setting as in the first simulation for unknown r. In this example, the network cannot separate out clusters one and two well, whereas the covariates can. On the other hand, clusters two and three are not well separated in the covariate space, while they are well separated using the network parameters. The experiments are repeated on 10 independently generated samples and the box plot for NMI is shown as in Figure 3(c). In the second row of Figure 3, we examine covariates with nonlinear cluster boundaries. The graph used here is the same as above, and the covariates are 2-dimensional, whose scatter plot is shown in Figure 3(e). In this case, the kernel matrix is able to pick up local similarities hence performs better than combination via inner product similarity as used in ACASC. In both simulations, SDP-comb outperforms others.

4.4 Real World Networks

Now we present results on a real world social network and an ecological network. The performance of clustering is evaluated by NMI with the ground truth labels.

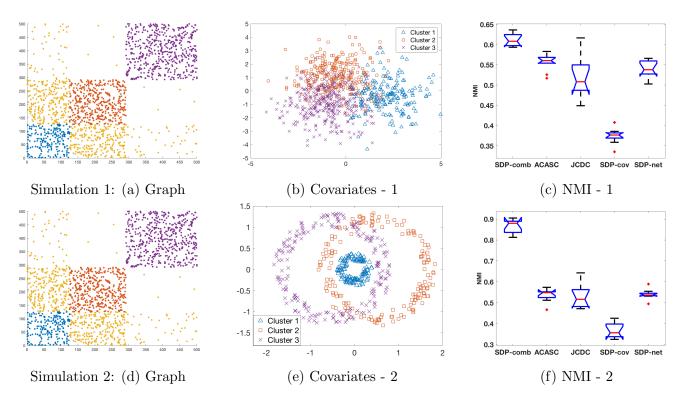


Figure 3: The first and second rows have results for isotropic Gaussian covariates and covariates lies on a nonlinear manifold respectively. We plot the adjacency matrix A in (a) and (b), where blue, red and purple points represent within cluster edges for 3 ground truth clusters respectively and yellow points represent inter-cluster edges. In (b) and (e) we plot covariates; different shapes and colors imply different clusters. (c) and (f) show the box plots for NMI.

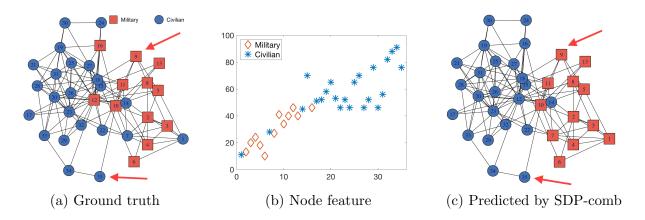


Figure 4: Mexican political network.

Mexican political elites As discussed before, this network (Gil-Mendieta & Schmidt 1996) depicts the political, kinship, or business interactions between 35 Mexican presidents and close collaborators, etc. The two ground truth clusters consist of the military and the civilians, indicating

the background of the politician. The year in which a politician first held a significant governmental position, is used as a covariate. Figure 4(b) shows that the covariate gives a good indication of the labels. This is because the military dominated the political arena after the revolution in the beginning of the twentieth century, and were succeeded by the civilians.

Table 4 shows the NMI of all methods, where our method outperforms other covariate-assisted approaches. From Figure 4(a, c), for example, node 35 has exactly one connection to each of the military and civilian groups, but seized power in the 90s, which strongly indicates a civilian background. On the other hand, node 9 took power in 1940, a year when civilian and military had almost equal presence in politics, making it hard to detect node 9's political affiliation. However, this node has more edges to the military group than the civilian group. By taking the graph structure into consideration, we can correctly assign the military label to it.

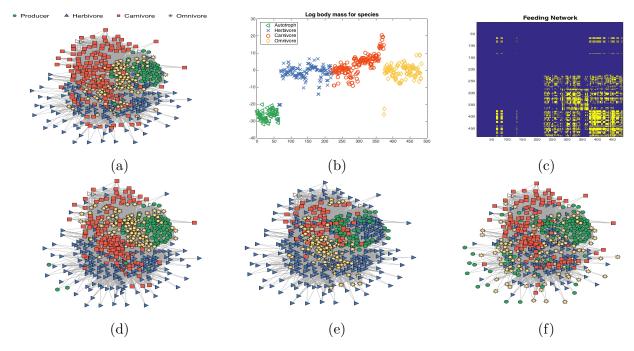


Figure 5: Weddell sea network: (a) True labels; (b) Log body mass; (c) Constructed adjacency matrix A_{τ} ; we show labels from (d) SDP-comb; (e) SDP-net; (f) SDP-cov.

Weddell sea trophic dataset The next example we consider is an ecological network collected by Jacob et al. (2011) describing the marine ecosystem of Weddell Sea, a large bay off the coast of Antarctica. The dataset lists 489 marine species and their directed predator-prey interactions, as well as the average adult body mass for each of the species. We use a thresholded symmetrization of the directed graph as the adjacency matrix. Let G be the directed graph, the $(i, j)^{th}$ entry of GG^T captures the number of other species which i and j both feed on. We create binary matrices $A_{\tau} = 1(GG^T \geq \tau)$. Choosing different τ 's between 1 to 10 gives similar clustering. We use $\tau = 5$.

All species are labeled into four categories based on their prey types. Autotrophs (e.g. plants) do not feed on anything. Herbivores feed on autotrophs. Carnivores feed on animals that are not autotrophs, and the remaining are omnivores, which feed both on autotrophs and other animals (herbivore, carnivore, or omnivores). Since body masses of species vary largely from nanograms

Dataset	SDP-net	SDP-cov	SDP-comb	ACASC	JCDC
Mexican politicians	0.37	0.43	0.46	0.37	0.25
Weddell Sea	0.36	0.22	0.51	0.32	0.42

Table 4: NMI with ground truth for various methods

to tons, we work with the normalized logarithm of mass following the convention in Newman & Clauset (2016). Figure 5(b) illustrates the log body mass for species. Without loss of generality, we order the nodes as autotrophs, herbivores, carnivores and omnivores.

In Figures 5(c), we plot A_{τ} . Since the autotrophs do not feed on other species in this dataset, and since herbivores do not have too much overlap in the autotrophs they feed on, the upper left corner of the input network is extremely sparse. On the other side, the body sizes for autotrophs are much smaller than those of other prey types. Therefore the kernel matrix clearly separates them out.

We see that SDP-net (Figure 5(e)) heavily misclusters the autotrophs since it only replies on the network. SDP-net (Figure 5(f)) only takes the covariates into account and cannot distinguish herbivores from omnivores, since they possess similar body masses. However, SDP-comb (Figure 5(d)) achieves a significantly better NMI by combining both sources. Table 4 shows the NMI between predicted labels and the ground truth from SDP-comb, JCDC and ACASC. While JCDC and ACASC can only get as good as the the best of graph or covariates, our method achieves a higher NMI.

5 Discussion

In this paper, we propose a regularized convex optimization framework to infer community memberships jointly from sparse networks and finite dimensional covariates. We theoretically show that our framework can improve clustering accuracy of either source under weaker separation conditions. In particular, when each source only has partial information about the clustering, our methodology can lead to high clustering accuracy, when either source fails. We demonstrate the performance of our methodology on simulated and real networks, and show that it in general performs better than other state-of-the-art methods. While for ease of exposition we limit ourselves to two sources, our method can be easily generalized to multiple views or sources. Empirically, we demonstrate that our method works for covariates with nonlinear cluster boundaries; we intend to extend our theoretical analysis to this setting and non-isotropic covariates as well.

Acknowledgements

We thank Arash Amini and Yuan Zhang for generously sharing their code. We are grateful to Soumendu Mukherjee, Peter J. Bickel, Dave Choi and Harry Zhou for interesting discussions on our paper.

References

- Airoldi, E. M., Blei, D. M., Fienberg, S. E. & Xing, E. P. (2008), 'Mixed membership stochastic blockmodels', J. Mach. Learn. Res. 9, 1981–2014.
- Amini, A. A., Chen, A., Bickel, P. J., Levina, E. et al. (2013), 'Pseudo-likelihood methods for community detection in large sparse networks', *Ann. Statist.* **41**(4), 2097–2122.
- Amini, A. A., Levina, E. et al. (2018), 'On semidefinite relaxations for the block model', *The Annals of Statistics* **46**(1), 149–179.
- Balakrishnan, S., Wainwright, M. J. & Yu, B. (2017), 'Statistical guarantees for the em algorithm: From population to sample-based analysis', *Ann. Statist.* **45**(1), 77–120. URL: http://dx.doi.org/10.1214/16-AOS1435
- Binkiewicz, N., Vogelstein, J. T. & Rohe, K. (2017), 'Covariate-assisted spectral clustering', *Biometrika* **104**(2), 361–377.
- Boyd, S., Parikh, N., Chu, E., Peleato, B. & Eckstein, J. (2011), 'Distributed optimization and statistical learning via the alternating direction method of multipliers', Foundations and Trends® in Machine Learning 3(1), 1–122.
- Braverman, M., Makarychev, K., Makarychev, Y. & Naor, A. (2013), The grothendieck constant is strictly smaller than krivine's bound, in 'Forum of Mathematics, Pi', Vol. 1, Cambridge Univ Press, p. e4.
- Cai, T. T. & Li, X. (2015), 'Robust and computationally feasible community detection in the presence of arbitrary outlier nodes', *Ann. Statist.* **43**(3), 1027–1059.
- Chaudhuri, K., Kakade, S. M., Livescu, K. & Sridharan, K. (2009), Multi-view clustering via canonical correlation analysis, *in* 'Proceedings of the 26th annual international conference on machine learning', ACM, pp. 129–136.
- Chen, Y. & Xu, J. (2016), 'Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices', *Journal of Machine Learning Research* 17(27), 1–57.
- Dasgupta, S. & Schulman, L. (2007), 'A probabilistic analysis of em for mixtures of separated, spherical gaussians', *The Journal of Machine Learning Research* 8, 203–226.
- Gil-Mendieta, J. & Schmidt, S. (1996), 'The political network in mexico', *Social Networks* 18(4), 355–381.
- Guédon, O. & Vershynin, R. (2015), 'Community detection in sparse networks via grothendieck's inequality', *Probability Theory and Related Fields* pp. 1–25.
- Hajek, B., Wu, Y. & Xu, J. (2016), 'Achieving exact cluster recovery threshold via semidefinite programming', *IEEE Transactions on Information Theory* **62**(5), 2788–2797.
- Holland, P. W., Laskey, K. B. & Leinhardt, S. (1983), 'Stochastic blockmodels: First steps', *Social networks* 5(2), 109–137.

- Hsu, D., Kakade, S. M. & Zhang, T. (2012), 'A tail inequality for quadratic forms of subgaussian random vectors', *Electron. Commun. Probab* 17(52), 1–6.
- Jacob, U., Thierry, A., Brose, U., Arntz, W. E., Berg, S., Brey, T., Fetzer, I., Jonsson, T., Mintenbeck, K., Mollmann, C. et al. (2011), 'The role of body size in complex food webs: A cold case', Advances In Ecological Research 45, 181–223.
- Jin, J., Ke, Z. T., Wang, W. et al. (2017), 'Phase transitions for high dimensional clustering and related problems', *The Annals of Statistics* **45**(5), 2151–2189.
- Kumar, A., Rai, P. & Daume, H. (2011), Co-regularized multi-view spectral clustering, in 'Advances in Neural Information Processing Systems 24', pp. 1413–1421.
- Le, C. M. & Levina, E. (2015), 'Estimating the number of communities in networks by spectral methods', arXiv preprint arXiv:1507.00827.
- Liu, J., Wang, C., Gao, J. & Han, J. (2013), Multi-view clustering via joint nonnegative matrix factorization, in 'Proceedings of the 2013 SIAM International Conference on Data Mining', SIAM, pp. 252–260.
- Mixon, D. G., Villar, S. & Ward, R. (2017), 'Clustering subgaussian mixtures by semidefinite programming', *Information and Inference: A Journal of the IMA* p. iax001.
- Montanari, A. & Sen, S. (2016), Semidefinite programs on sparse random graphs and their application to community detection, *in* 'Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing', ACM, New York, NY, USA, pp. 814–827.
- Newman, M. E. & Clauset, A. (2016), 'Structure and inference in annotated networks', *Nature Communications* 7.
- Peng, J. & Wei, Y. (2007), 'Approximating k-means-type clustering via semidefinite programming', SIAM Journal on Optimization 18(1), 186–205.
- Perry, A. & Wein, A. S. (2017), A semidefinite program for unbalanced multisection in the stochastic block model, *in* 'Sampling Theory and Applications (SampTA), 2017 International Conference on', IEEE, pp. 64–67.
- Rohe, K., Chatterjee, S. & Yu, B. (2011), 'Spectral clustering and the high-dimensional stochastic blockmodel', *The Annals of Statistics* pp. 1878–1915.
- Shi, T., Belkin, M. & Yu, B. (2009), 'Data spectroscopy: Eigenspaces of convolution operators and clustering', *The Annals of Statistics* pp. 3960–3984.
- Vempala, S. & Wang, G. (2004), 'A spectral algorithm for learning mixture models', *Journal of Computer and System Sciences* **68**(4), 841–860.
- Vershynin, R. (2010), 'Introduction to the non-asymptotic analysis of random matrices', arXiv preprint arXiv:1011.3027.
- Verzelen, N., Arias-Castro, E. et al. (2017), 'Detection and feature selection in sparse mixture models', *The Annals of Statistics* **45**(5), 1920–1950.

- Von Luxburg, U., Belkin, M. & Bousquet, O. (2008), 'Consistency of spectral clustering', *The Annals of Statistics* pp. 555–586.
- Weng, H. & Feng, Y. (2016), 'Community detection with nodal information', arXiv preprint arXiv:1610.09735.
- Yan, B. & Sarkar, P. (2016), On robustness of kernel clustering, in 'Advances in Neural Information Processing Systems', pp. 3098–3106.
- Yan, B., Sarkar, P. & Cheng, X. (2017), 'Exact recovery of number of blocks in blockmodels', arXiv preprint arXiv:1705.08580.
- Yan, B., Yin, M. & Sarkar, P. (2017), Convergence of gradient em on multi-component mixture of gaussians, in 'Advances in Neural Information Processing Systems', pp. 6959–6969.
- Yang, L., Sun, D. & Toh, K.-C. (2015), 'Sdpnal + +: a majorized semismooth newton-cg augmented lagrangian method for semidefinite programming with nonnegative constraints', *Mathematical Programming Computation* 7(3), 331–366.
- Yu, Y., Wang, T. & Samworth, R. J. (2014), 'A useful variant of the davis–kahan theorem for statisticians', *Biometrika* **102**(2), 315–323.
- Zhang, Y., Levina, E. & Zhu, J. (2016), 'Community detection in networks with node features', Electron. J. Statist. 10(2), 3153–3178.

A Background materials on sub-gaussian random vectors

In this section, we present some properties of sub-gaussian random variables. A sub-gaussian random variable is defined by the following equivalent properties. More discussions on this topic can be found in Vershynin (2010).

Lemma 3 (Vershynin (2010)). The sub-gaussian norm of X is denoted by $\|X\|_{\psi_2} = \sup_{p\geq 1} p^{-1/2} (\mathbb{E}|X|^p)^{1/p}$. A random vector $X \in \mathbb{R}^n$ is defined to be sub-gaussian if the one-dimensional marginals $\langle X, x \rangle$ are sub-gaussian random variables for all $x \in \mathbb{R}^n$ with sub-gaussian norm $\|X\|_{\psi_2} = \sup_{x \in S^{n-1}} \|\langle X, x \rangle\|_{\psi_2}$. Every sub-gaussian random variable X satisfies:

- (1) $P(|X| > t) \le \exp(1 ct^2 / ||X||_{\psi_2}^2)$ for all $t \ge 0$;
- (2) (Rotation invariance) Consider a finite number of independent centered sub-gaussian random variables X_i . Then $\sum_i X_i$ is also a centered sub-gaussian random variable. Moreover, $\|\sum_i X_i\|_{\psi_2}^2 \leq C \sum_i \|X_i\|_{\psi_2}^2$.
- (3) Let X_1, \dots, X_n be independent centered sub-gaussian random variables. Then $X = (X_1, X_2, \dots, X_n)$ is a centered sub-gaussian random vector in \mathbb{R}^n and $\|X\|_{\psi_2} \leq C \max_i \|X_i\|_{\psi_2}$.

A random variable is sub-exponential if the following equivalent properties hold with parameters $K_i > 0$ differing from each other by at most an absolute constant factor: (1) $P(|X| > t) \le \exp(1 - t/K_1)$ for all $t \ge 0$; (2) $(\mathbb{E}|X|)^{1/p} \le K_2 p$ for all $p \ge 1$; (3) $\mathbb{E}\exp(X/K_3) \le e$. The square of sub-gaussian random variable is sub-exponential.

Lemma 4 (Vershynin (2010)). A random variable X is sub-gaussian if and only if X^2 is sub-exponential. Moreover, $\|X\|_{\psi_2}^2 \leq \|X^2\|_{\psi_1} \leq 2\|X\|_{\psi_2}^2$.

B Proof of Lemma 1

We start with the following lemma, whose proof is in the Supplementary.

Lemma 5. For any X that satisfies $X \succeq 0, X \geq 0, X = 1$, we have $||X||_F^2 \leq trace(X)$.

Proof. We first show that for all such X, the eigenvalues of X are in [0,1]. Let v_i be the eigenvector of X corresponding to the i^{th} largest eigenvalue θ_i . Since X is positive semi-definite, $\theta_i \geq 0, \forall i$. Without loss of generality, let $i^* = \arg \max_i |v_1(i)|$, i.e. be the index of the entry with the largest absolute value of v_1 . Since $Xv_1 = \theta_1v_1$, and $\sum_i X_{ij} = 1, X_{ij} \geq 0$, we have:

$$|\theta_1 v_1(i^*)| = |\sum_j X_{i^*j} v_1(j)| \le \sum_j X_{i^*j} |v_1(j)| \le |v_1(i^*)|.$$

Therefore $|\theta_1| \leq 1$.

$$||X||_F^2 = \sum_i \theta_i^2 \le \sum_i \theta_i = \text{trace}(X)$$

Now we are in position to prove Lemma 1.

Proof of Lemma 1. Note that both X_0 and X_M are in the feasible set \mathcal{F} , by optimality, we have $\langle M, X_M \rangle \geq \langle M, X_0 \rangle$. We construct Q as stated in the lemma to obtain: $\langle Q, X_M - X_0 \rangle$, $\langle M - Q, X_M - X_0 \rangle \geq \langle Q, X_0 - X_M \rangle$. Note that Q is constant on diagonal blocks and upper bounded by q_k on off-diagonal blocks, with respect to the clustering of nodes. Using the fact that $|C_k| = m_k$, we have:

$$\langle M, X_0 - X_M \rangle = \sum_{k} \sum_{i \in C_k} \left(\beta_k^{(in)} \sum_{j \in C_k} \left(\frac{1}{m_k} - (X_M)_{ij} \right) + \sum_{\ell \neq k} \sum_{j \in C_\ell} Q_{ij} (0 - (X_M)_{ij}) \right)$$

$$\geq \sum_{k} \sum_{i \in C_k} \left(\beta_k^{(in)} \sum_{j \in C_k} \left(\frac{1}{m_k} - (X_M)_{ij} \right) - \beta_k^{(out)} \sum_{\ell \neq k} \sum_{j \in C_\ell} (X_M)_{ij} \right)$$

$$= \sum_{k} \sum_{i \in C_k} \left(\beta_k^{(in)} \left(1 - \sum_{j \in C_k} (X_M)_{ij} \right) - \beta_k^{(out)} \left(1 - \sum_{j \in C_k} (X_M)_{ij} \right) \right)$$

$$= \sum_{k} \sum_{i \in C_k} (\beta_k^{(in)} - \beta_k^{(out)}) \left(1 - \sum_{j \in C_k} (X_M)_{ij} \right) \geq \min_{k} (\beta_k^{(in)} - \beta_k^{(out)}) \sum_{k} \sum_{i \in C_k} \left(1 - \sum_{j \in C_k} (X_M)_{ij} \right)$$

The third line and last inequality uses the constraint that $\sum_{j} \hat{X}_{ij} = 1$, and $1 - \sum_{j \in C_k} \hat{X}_{ij} \ge 1 - \sum_{j} \hat{X}_{ij} = 0$. On the other hand,

$$||X_M - X_0||_F^2 = ||X_M||_F^2 - ||X_0||_F^2 + 2\langle X_0 - X_M, X_0 \rangle$$

By Lemma 5, and the fact that $||X_0||_F^2 = r$, we have $||X_M||_F^2 - ||X_0||_F^2 \le \operatorname{trace}(X_M) - r = 0$. Since $\min_k(\beta_k^{(in)} - \beta_k^{(out)}) \ge 0$,

$$||X_{M} - X_{0}||_{F}^{2} \leq 2\langle X_{0} - X_{M}, X_{0} \rangle = 2 \sum_{k} \sum_{i \in C_{k}} \sum_{j \in C_{k}} \frac{1}{m_{k}} \left(\frac{1}{m_{k}} - (X_{M})_{ij} \right)$$

$$= 2 \sum_{k} \sum_{i \in C_{k}} \frac{1}{m_{k}} \left(1 - \sum_{j \in C_{k}} (X_{M})_{ij} \right) \leq \frac{2}{m_{\min}} \sum_{k} \sum_{i \in C_{k}} \left(1 - \sum_{j \in C_{k}} (X_{M})_{ij} \right)$$

$$\leq \frac{2}{m_{\min} \min_{k} (\beta_{k}^{(in)} - \beta_{k}^{(out)})} \langle Q, X_{0} - X_{M} \rangle \leq \frac{2}{m_{\min} \min_{k} (\beta_{k}^{(in)} - \beta_{k}^{(out)})} \langle M - Q, X_{M} - X_{0} \rangle$$

C Proof of Proposition 1

We first introduce the following result on sparse graph with Grothendieck's inequality by Guédon & Vershynin (2015).

Lemma 6 (Guédon & Vershynin (2015)). Let $\mathcal{M}_{G}^{+} = \{X : X \succeq 0, diag(X) \preceq I_{n}\}$, $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ be a symmetric matrix whose diagonal entries equal 0, and entries above the diagonal are independent random variables satisfying $0 \leq a_{ij} \leq 1$. Let P = E[A|Z]. Assume that $\bar{p} := \frac{2}{n(n-1)} \sum_{i < j} Var(a_{ij}) \geq \frac{9}{n}$. Then, with probability at least $1 - e^3 5^{-n}$, we have $\max_{X \in \mathcal{M}_{G}^{+}} |\langle A - P, X \rangle| \leq K_G ||A - P||_{\ell_{\infty} \to \ell_{1}} \leq 3K_G \bar{p}^{1/2} n^{3/2}$, where K_G is the Grothendieck's constant, and its best know upper bound is 1.783.

Proof of Proposition 1. Notice that A and P := E[A|Z] has zero diagonals. Therefore,

$$\langle P - Q, X_A - X_0 \rangle = \sum_k \sum_{i \in C_k} a_k / n \left(\frac{1}{m_k} - (X_A)_{ii} \right)$$

$$\leq \sum_k p_k - p_{\min} \operatorname{trace}(X_A) \leq r(p_{\max} - p_{\min})$$
(8)

where $p_{\text{max}} = \max_k a_k/n$ and $p_{\text{min}} = \min_k a_k/n$. Thus by Lemma 1 and Eq (8),

$$||X_A - X_0||_F^2 \le \frac{2}{m_{\min}\min_k(a_k/n - b_k/n)} (\langle A - P, X_A - X_0 \rangle + r(p_{\max} - p_{\min}))$$

In sparse regime, both $m_{\min}X_0$ and $m_{\min}X_A$ belong to the set \mathcal{M}_G^+ . Let $g = n\bar{p} \geq 9$, applying Lemma 6 we get with probability at least $1 - e^3 5^{-n}$,

$$||X_A - X_0||_F^2 \le \frac{22\sqrt{n^2g}}{m_{\min}^2 \min_k(a_k/n - b_k/n)} + \frac{2r(p_{\max} - p_{\min})}{m_{\min} \min_k(a_k/n - b_k/n)}$$

Substituting $p_k = a_k/n$, $q_k = b_k/n$, and using the fact that

$$\frac{2r(p_{\max} - p_{\min})}{m_{\min}\min_k(p_k - q_k)} = \frac{2rm_{\min}(p_{\max} - p_{\min})}{m_{\min}^2\min_k(p_k - q_k)} \le \frac{2\max_k a_k}{m_{\min}^2\min_k(p_k - q_k)} = o(\sqrt{n^2g}),$$

Recall that $\alpha := m_{\text{max}}/m_{\text{min}}$, we get with probability tending to 1,

$$\frac{\|\hat{X} - X_0\|_F^2}{\|X_0\|_F^2} \le \frac{23n^2\sqrt{g}}{rm_{\min}^2 \min_k(a_k - b_k)} \le \frac{23\alpha^2 r\sqrt{g}}{\min_k(a_k - b_k)}.$$

D Proof of Proposition 2

Proof of Proposition 2. Recall that by definition, for $i \in C_k$, $Y_i - \mu_k$ is sub-gaussian random vector with sub-gaussian norm ψ_k . Using the following concentration inequality from Hsu et al. (2012) for sub-gaussian random vectors, we have:

For
$$i \in C_k$$
, $P(||Y_i - \mu_k||_2^2 > \psi_k^2(d + 2\sqrt{td} + 2t)) \le e^{-t}$

We take $t = c_k^2 d$ for $c_k \ge 1$. Since $1 + 2c_k + 2c_k^2 \le 5c_k^2$ for $c_k \ge 1$, we get $P(\|X - \mathbb{E}X\|^2 \le 5c_k^2 \psi_k^2 d) \ge 1 - \exp(-c_k^2 d)$. Let $\Delta_k = \sqrt{5}c_k \psi_k \sqrt{d}$, we can divide the nodes into "good nodes" (those close to their population mean) S_k and the rest as follows:

$$S_k = \{i \in C_k : ||Y_i - \mu_k|| \le \Delta_k\}, \qquad S = \bigcup_{k=1}^r S_k$$
(9)

Let $m_c^{(k)} = m_k - |\mathcal{S}_k|$. We want to bound $m_c^{(k)}$ with high probability. Note that $m_c^{(k)} = \sum_{i \in C_k} \mathbf{1}(||Y_i - \mu_k|| \ge \Delta_k)$ is a sum of i.i.d random variables. Therefore, using the Hoeffding bound we have:

$$P\left(m_c^{(k)} - m_k P(i \notin \mathcal{S}_k) \ge m_k \delta\right) \le \exp(-2m_k \delta^2)$$

Using $\delta = \sqrt{\log m_k/2m_k}$, we have:

$$P\left(m_c^{(k)} - m_k P(i \notin \mathcal{S}_k) \ge \sqrt{m_k \log m_k/2}\right) \le \frac{1}{m_k}$$

Since $P(i \notin S_k) \leq \exp(-c_k^2 d)$, we have:

$$P\left(m_c^{(k)} \ge m_k \exp(-c_k^2 d) + \sqrt{m_k \log m_k/2}\right) \le \frac{1}{m_k}$$

Finally, using union bound over all clusters we get:

$$P\left(m_c \ge \sum_k m_k e^{-c_k^2 d} + \sum_k \sqrt{m_k \log m_k/2}\right) \le \sum_k \frac{1}{m_k}$$

$$\tag{10}$$

Now define

$$(K_I)_{ij} = \begin{cases} f(2\Delta_k), & \text{if } i, j \in C_k \\ \min\{f(d_{k\ell} - \Delta_k - \Delta_\ell), K_{ij}\}, & \text{if } i \in C_k, j \in C_\ell, k \neq \ell \end{cases}$$
(11)

By Lemma 1, all diagonal blocks are blockwise constant and the off-diagonal blocks are upper bounded by $f(d_{k\ell} - \Delta_k - \Delta_\ell)$. Let $\nu_k = f(2\Delta_k) - \max_{\ell \neq k} f(d_{k\ell} - \Delta_k - \Delta_\ell)$, and $\gamma = \min_k \nu_k$. If $\nu_k \geq 0$, we have

$$||X_K - X_0||_F^2 \le \frac{2}{m_{\min}\gamma} \langle K - K_I, X_K - X_0 \rangle$$

Apply Grothendieck's inequality,

$$||X_K - X_0||_F^2 \le \frac{2K_G}{m_{\min}^2 \gamma} ||K - K_I||_{\ell_\infty \to \ell_1}$$
(12)

Now it remains to bound the $\ell_{\infty} \to \ell_1$ norm of $K - K_I$. Note that if $i \in S_k, j \in S_\ell, k \neq \ell$, then by a simple use of triangle inequality we have $K_{ij} \leq f(d_{k\ell} - \Delta_k - \Delta_\ell)$, so $K_{ij} = (K_I)_{ij}$; and if $i, j \in S_k$, then $K_{ij} \geq f(2\Delta_k)$.

$$||K - K_{I}||_{\ell_{\infty} \to \ell_{1}} = \max_{x,y \in \{\pm\}^{n}} \sum_{i,j} x_{i}y_{j} \left(K_{ij} - (K_{I})_{ij}\right)$$

$$\leq \max_{x,y \in \{\pm\}^{n}} \sum_{i,j \in \mathcal{S}} x_{i}y_{j} \left(K_{ij} - (K_{I})_{ij}\right) + \max_{x,y \in \{\pm\}^{n}} \sum_{i \notin \mathcal{S} \cup j \notin \mathcal{S}} x_{i}y_{j} \left(K_{ij} - (K_{I})_{ij}\right)$$

$$\stackrel{(i)}{\leq} \max_{x,y \in \{\pm\}^{n}} \sum_{i,j \in \mathcal{S}} x_{i}y_{j} \left(K_{ij} - (K_{I})_{ij}\right) + 2m_{c}n$$

$$\stackrel{(ii)}{=} \max_{x,y \in \{\pm\}^{n}} \sum_{k} \sum_{i,j \in \mathcal{S}_{k}} x_{i}y_{j} \left(K_{ij} - f(2\Delta_{k})\right) + 2m_{c}n$$

$$\leq \sum_{k} m_{k}^{2} (1 - f(2\Delta_{k})) + 2m_{c}n$$

$$(13)$$

where (i) is due to $|K_{ij} - (K_I)_{ij}| \le 1$, and (ii) comes from the definition of K_I . Now Eq 12 follows as

$$||X_{K} - X_{0}||_{F}^{2} \leq \frac{4K_{G}\left(\sum_{k} m_{k}^{2}(1 - f(2\Delta_{k})) + 2m_{c}n\right)}{m_{\min}^{2}\gamma}$$

$$= \frac{4K_{G}}{m_{\min}^{2}} \sum_{k} \left(m_{k}^{2} \frac{1 - f(2\Delta_{k})}{\gamma} + 2m_{k}ne^{-c_{k}^{2}d}/\gamma\right) + \frac{\sqrt{2}K_{G}n}{m_{\min}^{2}\gamma} \sum_{k} \sqrt{m_{k}\log m_{k}}$$
(14)

Recall that $f(x) = \exp(-\eta x^2)$, and $\gamma = \min_k \{f(2\Delta_k) - \max_{\ell \neq k} f(d_{k\ell} - \Delta_k - \Delta_\ell)\}$. For simplicity, we assume $c_k = c_0$. We take $c_0 = \sqrt{\log\left(\frac{d_{\min}^2}{\psi_{\max}^2 d}\right)/d}$ and the scale parameter $\eta = \frac{\phi}{20c_0^2\psi_{\max}^2 d}$, for some $\phi > 0$, which will be chosen later. Furthermore, we also define

$$\xi = \frac{d_{\min}}{2\sqrt{5}c_0\psi_{\max}\sqrt{d}} - 1. \tag{15}$$

If $\xi > 1$, then $d_{min} > 4\sqrt{5}c_0\psi_{max}\sqrt{d}$, and hence $\gamma > 0$. Also, since $\eta(d_{\min} - 2\sqrt{5}c_0\psi_{\max}\sqrt{d})^2 = \phi\xi^2$, $\forall k, \ell \in [r]$, if $d_{\min} := \min_{k\ell} d_{k\ell} > 4\sqrt{5}c_0\psi_{\max}\sqrt{d}$, then

$$\gamma \ge f(2\sqrt{5}c_0\psi_{\max}\sqrt{d}) - f(d_{\min} - 2\sqrt{5}c_0\psi_{\max}\sqrt{d}) = \exp(-\phi) - \exp(-\phi\xi^2).$$

and

$$1 - f(2\Delta_k) \le 1 - f(2\sqrt{5}c_0\psi_{\max}\sqrt{d}) = 1 - \exp(\phi)$$

Recall $\alpha = \frac{m_{\text{max}}}{m_{\text{min}}}$

$$||X_{K} - X_{0}||_{F}^{2}$$

$$\leq 4K_{G}r\alpha^{2} \cdot \frac{1 - f(2\sqrt{5}c_{0}\psi_{\max}\sqrt{d}) + 2r\exp(-c_{0}^{2}d)}{\gamma} + \frac{2\sqrt{2}K_{G}m_{\max}r^{2}\sqrt{m_{\max}\log m_{\max}}}{\gamma m_{\min}^{2}}$$

$$\leq \frac{4K_{G}r\alpha^{2}}{\gamma} \left(1 - \exp(-\phi) + \frac{2r\psi_{\max}^{2}\sqrt{d}}{d_{\min}^{2}} + r\sqrt{\log m_{\max}/2m_{\max}}\right)$$

$$\leq 4K_{G}r\alpha^{2} \left(\underbrace{\frac{(1 - \exp(-\phi) + 2r\psi_{\max}^{2}d/d_{\min}^{2})}{\exp(-\phi) - \exp(-\phi\xi^{2})}}_{A} + \underbrace{\frac{r\sqrt{\log m_{\max}/2m_{\max}}}{\exp(-\phi) - \exp(-\phi\xi^{2})}}_{B}\right)$$

$$(17)$$

We will first bound part (A).

$$(A) = \frac{\exp(\phi) - 1 + \exp(\phi) \frac{2r\psi_{\text{max}}^2 d}{d_{\text{min}}^2}}{1 - \exp(\phi - \phi\xi^2)} \stackrel{(i)}{\leq} \frac{\phi + \frac{\phi^2}{2} \exp(\phi) + \exp(\phi) \frac{2r\psi_{\text{max}}^2 d}{d_{\text{min}}^2}}{1 - \exp(\phi - \phi\xi^2)}$$
(18)

where (i) uses the Mean value theorem: for $e^x - 1 \le x + e^y x^2/2$ for $y \in [0, x]$. If $\frac{d_{\min}}{\psi_{\max} \sqrt{d}} > \max\{1, \frac{180}{d}\}$, using the fact that $\log x \le \sqrt{x}$, we have:

$$\frac{d_{\min}^2}{\psi_{\max}^2 d} > \frac{180}{d^2} \frac{d_{\min}}{\psi_{\max}} > \frac{180}{d} \log \left(\frac{d_{\min}^2}{\psi_{\max}^2 d} \right) = 180c_0^2.$$

Using Eq 15, we see that $\xi > \frac{\sqrt{180}}{2\sqrt{5}} - 1 = 2$, and hence $\gamma > 0$. Now we pick $\phi = \frac{\log \xi}{\xi^2}$.

Now we will use this to obtain a lower bound on $1 - \exp(\phi - \phi \xi^2)$. Since $\xi \ge 2$, we have $\xi^2/4 \ge 1$. Hence

$$1 - \exp(\phi - \phi \xi^2) \ge 1 - \exp(\phi \xi^2 / 4 - \phi \xi^2)$$
$$= 1 - \exp(-\phi 3\xi^2 / 4) = 1 - \exp(-3\log \xi / 4) = 1 - \xi^{-3/4}$$
$$> 1 - 2^{-3/4} = .4$$

Using the fact that the function $\frac{\log x}{x^2}$ is monotonically decreasing when x > 2, we see that $\phi < \log 2/2^2$ and $\exp(\phi) \le 1.2$. Furthermore,

$$\gamma \ge \exp(-\phi)(1 - \exp(\phi(1 - \xi^2))) \ge .3$$
 (19)

Now Eq. (18) yields:

$$\begin{split} (A) & \leq \frac{\phi + 1.2 \left(\frac{\phi^2}{2} + \frac{2r\psi_{\max}^2 d}{d_{\min}^2} \right)}{.4} \leq \frac{c \log \xi}{\xi^2} + \frac{3r\psi_{\max}^2 d}{d_{\min}^2} \\ & \leq \frac{c' \log(\xi + 1)}{(\xi + 1)^2} + \frac{3r\psi_{\max}^2 d}{d_{\min}^2} \leq c'' \frac{\psi_{\max}^2 d}{d_{\min}^2} \log \left(\frac{d_{\min}}{\psi_{\max} \sqrt{d}} \right) + \frac{3r\psi_{\max}^2 d}{d_{\min}^2}, \end{split}$$

for some constant c. To get (ii), note that

$$\frac{\log \xi}{\xi^2} \le \frac{\log(\xi+1)}{\xi^2} \le \frac{2.25 \log(\xi+1)}{(\xi+1)^2}, \forall \xi > 2$$

Finally, we bound (B) in Eq 17 using Eq 19.

$$(B) = \frac{r\sqrt{\log m_{\text{max}}/2m_{\text{max}}}}{\exp(-\phi) - \exp(-\phi\xi^2)} \le c_1 r \sqrt{\frac{\log m_{\text{max}}}{m_{\text{max}}}}$$

for some constant $c_1 > 0$. Putting pieces together, we have

$$\frac{\|X_K - X_0\|_F^2}{\|X_0\|_F^2} \le C\alpha^2 \max\left(\frac{\psi_{\max}^2 d}{d_{\min}^2} \max\left\{\log\left(\frac{d_{\min}}{\psi_{\max}\sqrt{d}}\right), r\right\}, r\sqrt{\frac{\log m_{\max}}{m_{\max}}}\right)$$

E Analysis for $X_{A+\lambda_n K}$

Proof of Theorem 1. Let K_I be defined as in Eq (11). Let $\gamma = \min_k (a_k/n - b_k/n + \lambda_n(f(2\Delta_k) - \max_{\ell \neq k} f(d_{k\ell} - \Delta_k - \Delta_\ell)))$. When $\gamma \geq 0$, Lemma 1 with $Q = ZBZ^T + \lambda_n K_I$, we have

$$||X_{A+\lambda_n K} - X_0||_F^2 \le \frac{2}{m_{\min} \gamma} \left(\langle A - P, X_{A+\lambda_n K} - X_0 \rangle + r(\max_k a_k/n - \min_k a_k/n) + \lambda_n \langle K - K_I, X_{A+\lambda_n K} - X_0 \rangle \right)$$

Now by Grothendieck's inequality on both $\langle A-P, X_{A+\lambda_n K}-X_0 \rangle$ and $\langle K-K_I, X_{A+\lambda_n K}-X_0 \rangle$, one gets,

$$||X_{A+\lambda_n K} - X_0||_F^2 \le \frac{2K_G}{m_{\min}^2 \gamma} \left(2||A - P||_{\ell_{\infty} \to \ell_1} + r(\max_k a_k/n - \min_k a_k/n) + 2\lambda_n ||K - K_I||_{\ell_{\infty} \to \ell_1} \right)$$

By Lemma 6 and Eq (13),

$$||X_{A+\lambda_n K} - X_0||_F^2 \le \frac{4K_G}{m_{\min}^2 \gamma} \left(6\sqrt{n^3 \bar{p}} + \lambda_n \left(2m_c n + \sum_k m_k^2 (1 - f(2\Delta_k)) \right) \right)$$

Using $\lambda_n = \lambda_0/n$, $m_k = n\pi_k$, $m_{\min} = n\pi_{\min}$, and $\pi_0 := \sum_k (m_k \exp(-\Delta_k^2/(5\psi_k^2)) + \sqrt{m_k \log m_k/2})/n$ in conjunction with Eq (10), we get with probability tending to 1,

$$||X_{A+\lambda_n K} - X_0||_F^2 \le 4K_G \frac{6\sqrt{g} + \lambda_0 \left(2\pi_0 + \sum_k \pi_k^2 (1 - f(2\Delta_k))\right)}{\pi_{\min}^2 \min_k (a_k - b_k + \lambda_0 \nu_k)}$$

F Analysis of covariate clustering when $d \gg r$

Before proving Lemma 2, we clearly state our assumptions and other useful lemmas.

Assumption 1. We assume that M is of rank r-1, i.e. the means are not collinear, or linearly dependent, other than the fact that they are centered.

Lemma 7. Let $M = \sum_k \pi_k \mu_k \mu_k^T$ and S be the covariance matrix of n data points from a subgaussian mixture, then $S = M + \sum_i \pi_i \sigma_i^2 I_d$. Let \hat{S} be the sample covariance matrix $\hat{S} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(Y_i - \bar{Y})^T}{n}$. We have $\|\hat{S} - S\| \le C\sqrt{\frac{d \log n}{n}}$ for some constant C with probability bigger than $1 - O(n^{-d})$.

This is a direct consequence of Corollary 5.50 from Vershynin (2010). The main ingredient of the proof is provided below.

Lemma 8. Let U_{r-1} be the top r-1 eigenvectors of \hat{S} estimated using P_1 , and λ be the smallest positive eigenvalue of M. For any vector v in the span of $\{\mu_i\}_{i=1}^r$, as long as $\lambda > 5\left(\psi_{\max}^2 + C\sqrt{\frac{d\log^2 n}{n}}\right)$ we have $\|U_{r-1}^Tv\| \ge \|v\|/2$ with probability at least $1 - \tilde{O}(n^{-d})$.

Proof. Take $n_1 = \frac{n}{\log n}$ and v to be a vector in the span of $\{\mu_i\}_{i=1}^r$. By definition, we have $\|Mv\| \ge \lambda \|v\|$. Let $R = \hat{S} - S$. Denote $\bar{\sigma}^2 = \sum_i \pi_i \sigma_i^2$, by Lemma 7, $S = M + \bar{\sigma}^2 I_d$. We also know that $\bar{\sigma}^2 \le \sigma_{\max}^2 \le \psi_{\max}^2$ by the property of sub-gaussian distributions. Since S is estimated from P_1 with n_1 points, applying Lemma 7 with n_1 we get $\|R\| \le \epsilon = C\sqrt{\frac{d \log n_1}{n_1}}$. By Weyl's inequality, $\|\hat{S}v\| = \|(M + R + \sum_i \sigma_i^2 I_d)v\| \ge (\lambda - \sigma_{\max}^2 - \epsilon)\|v\|$. Let $U_{r:d}$ be the eigenspace orthogonal to U_{r-1} . Assume the contradiction that $\|U_{r-1}^Tv\| < \|v\|/2$. Then there has to be a unit d dimensional vector $u \in \text{span}(U_{r:d})$, such that $\|u^Tv\| > \|v\|/2$. On one hand, if we write $u = c\frac{v}{\|v\|} + \sqrt{1 - c^2}v^{\perp}$, for |c| > 1/2 and some unit vector v^{\perp} orthogonal to v, we have $\|\hat{S}u\| \ge \frac{\lambda - \sigma_{\max}^2 - \epsilon}{2} - \sqrt{1 - c^2}\|\hat{S}v^{\perp}\|$. Note $\|\hat{S}v^{\perp}\| = \|(M + R + \bar{\sigma}^2 I_d)v^{\perp}\|$. Since v^{\perp} is orthogonal to the span of M, $\|\hat{S}v^{\perp}\| \le (\sigma_{\max}^2 + \epsilon)$. Hence

$$\|\hat{S}u\| \ge \frac{\lambda - 3(\sigma_{\max}^2 + \epsilon)}{2}.\tag{20}$$

On the other hand, since $u \in \text{span}(U_{r:d})$, by Weyl's inequality, $\|\hat{S}u\| \leq |\lambda_k(\hat{S})| \leq \sigma_{\max}^2 + \epsilon$. This contradicts with Eq. (20) since we assume $\lambda > 5(\psi_{\max}^2 + \epsilon) \geq 5(\sigma_{\max}^2 + \epsilon)$. The result is proven by contradiction.

Remark 4. Note that the result can be generalized to non-spherical case as long as the largest eigenvalue of covariance matrix for each cluster is bounded.

We are now ready to prove Lemma 2.

Proof of Lemma 2. Recall that $Y_i' = U_{r-1}^T Y_i$ where U_{r-1} and Y_i are from two different partitions and hence independent. Let $Z_i \in [r]$ denote that latent variable associated with i. Thus, $E[Y_i'|Z_i = a, P_2] = U_{r-1}^T E[Y_i|Z_i = a] = U_{r-1}^T \mu_a$. Thus the means of the new mixture are $\mu_a' := U_{r-1}^T \mu_a$ and the covariance matrix is isotropic, i.e. $E[(Y_i' - \mu_a')(Y_i' - \mu_a')^T | P_2, Z_i = a] = \sigma_a^2 I_{r-1}$. Furthermore, using Lemma 8 we have $\min_{k \neq \ell} \|\mu_k' - \mu_\ell'\| = \min_{k \neq \ell} \|U_{r-1}^T (\mu_k - \mu_\ell)\| \ge \|d_{\min}\|/2$. Since this requires an application of Lemma 8 to each of the vectors $\mu_k - \mu_\ell$, $k, \ell \in [r]$, the success probability is at least $1 - \tilde{O}(r^2 n^{-d})$ by union bound.

G From X to cluster labels

From some solution matrix \ddot{X} , we can apply spectral clustering on it to get the cluster labels. Below we present a theorem that bounds the misclassification error by the Frobenius norm of matrix difference. The proof technique is inspired by those in Rohe et al. (2011), Yan & Sarkar (2016).

Theorem 2. The number of misclassification nodes is bounded by $64m_{\text{max}}\|\hat{X} - X_0\|_F^2$.

Proof. Let \hat{U} be the top r eigenvectors of \hat{X} , $U \in \mathbb{R}^{n \times r}$ be the top r eigenvector of X_0 . Let $\nu \in \mathbb{R}^{r \times r}$ be the population value of the eigenvector corresponding to each cluster, $U = Z\nu$. By Davis-Kahan theorem Yu et al. (2014), we have

$$\|\hat{U} - UO\|_F^2 \le \frac{8\|\hat{X} - X_0\|_F^2}{(\theta_r(X_0) - \theta_{r+1}(X_0))^2} = 8\|\hat{X} - X_0\|_F^2$$
(21)

Define $\mathcal{M} = \{i : ||c_i - Z_i \nu O|| \ge \frac{1}{\sqrt{2m_{\text{max}}}}\}$. We now prove that \mathcal{M} is a superset of all misclassified nodes by the above procedure, and its cardinality is bounded as in the theorem statement. U is a unit basis so we know $I = U^T U = \nu^T Z^T Z \nu = \nu^T \mathrm{diag}(m_1, \dots, m_r) \nu$. So $\theta_{\min}(\nu^T \nu) \geq \frac{1}{m_{\max}}$. Define $\mathcal{C} = \{M \in \mathbb{R}^{n \times r} : M \text{ has no more than } r \text{ unique rows}\}$. Then minimizing the k-means

objective for \hat{U} is equivalent to

$$\min_{\{s_1, \dots, s_r\} \subset \mathbb{R}^r} \sum_{i} \min_{g} \|\hat{u}_i - s_g\|_2^2 = \min_{M \in \mathcal{C}} \|\hat{U} - M\|_F^2$$

So $C = [c_1, \dots, c_n] = \arg\min_{M \in \mathcal{C}} \|\hat{U} - M\|_F^2$ and $\|C - \hat{U}\| \le \|Z\nu O - \hat{U}\|$. c_i is the center assigned to point i by running k-means on U.

Now we prove all points lying outside of M is correctly labeled, or equivalently, $||c_i - Z_i \nu O|| <$ $||c_i - Z_j \nu O||_2$ for all $Z_j \neq Z_i$. To see this, note for $\forall i, j \in [n]$, when $Z_i \neq Z_j$,

$$||Z_i \nu - Z_j \nu|| = ||(Z_i - Z_j)\nu|| \ge \sqrt{2} \min_{x:||x||^2 = 1} \sqrt{x^T \nu^T \nu x} \ge \sqrt{\frac{2}{m_{\max}}}$$

So

$$||c_i - Z_j \nu O||_2 \ge ||Z_i \nu - Z_j \nu|| - ||c_i - Z_i \nu O|| \ge \sqrt{\frac{2}{m_{\text{max}}}} - \sqrt{\frac{1}{2m_{\text{max}}}} = \sqrt{\frac{1}{2m_{\text{max}}}}$$
(22)

Therefore when $Z_i \neq Z_j$, $||c_i - Z_i \nu O|| < \sqrt{\frac{r}{2n}} \Rightarrow ||c_i - Z_i \nu O||_2 < ||c_i - Z_j \nu O||_2$, which means node *i* is correctly clustered.

Below we bound the cardinality of \mathcal{M} . By Markov's inequality,

$$|\mathcal{M}| \le 2m_{\max} \sum_{i \in [n]} ||c_i - Z_i \nu O||_F^2$$

$$= 2m_{\max} ||C - UO||_F^2$$

$$\le 2m_{\max} (||C - \hat{U}||_F + ||\hat{U} - UO||_F)^2$$

Note

$$||C - \hat{U}||_F^2 \le ||\hat{U} - UO||_F^2$$

Therefore, we have

$$|\mathcal{M}| \le 8m_{\text{max}} \|\hat{U} - UO\|_F^2 \tag{23}$$

Combining with Eq. (21), we have

$$|\mathcal{M}| \le 64m_{\max} ||\hat{X} - X_0||_F^2$$