# GRAPH CLUSTERING AND THE NUCLEAR WASSERSTEIN METRIC.

DANIEL DE ROUX AND MAURICIO VELASCO

ABSTRACT. We study the problem of discovering the cluster structure of a random graph $\mathcal{G}$ from an independent sample of size $N$. We propose a Wasserstein robust formulation of this optimization problem and prove that it can be reformulated as a tractable convex optimization problem. We give theoretical performance guarantees for these problems when the Wasserstein metric is induced by the nuclear norm and $\mathcal{G}$ is distributed according to the stochastic block model explaining its good practical behavior. Finally we present our Julia implementation of the proposed algorithm and use it to analyze the voting patterns of the colombian senate in ♣♣♣ Mauricio: [period?].

## 1. INTRODUCTION

Let $\mathcal{G}$ be a random graph with $n$ vertices. By a deterministic summary of $\mathcal{G}$ we mean a (deterministic) graph $H^*$ which, on average, differs from $\mathcal{G}$ by as few edges as possible. In this article we study the problem of finding deterministic summaries *from an independent sample* of $\mathcal{G}$ of size $N$. More precisely we will address the following problem:

**Problem 1.1.** *Given adjacency matrices $B_1, \ldots, B_N$ of an independent sample with the same distribution as $\mathcal{G}$ find a symmetric matrix $A^*$ in $\arg\min_A \mathbb{E}_{B \sim \mathcal{G}}[\|A - B\|_1]$.*

Special cases of this problem arise in cluster detection and in data summarization, both heavily studied in the literature (♣♣♣ Mauricio: [refs?]).

A possible approach to problem 1.1 is to use the samples to use the empirical measure $\hat{\mu} := \sum_{i=1}^{N} \frac{1}{N} \delta_{B_i}$ as an approximation of the distribution of $\mathcal{G}$ and to find a minimizer $A$ of the resulting empirical risk

$$\mathbb{E}_{B \sim \mu}[\|A - B\|_1] = \frac{1}{N} \sum_{i=1}^{N} \|A - B_i\|_1.$$

When the sample size $N$ is not sufficiently large for $\hat{\mu}$ to be a good approximation for the distribution of $\mathcal{G}$ this approach leads to overfitting. To mitigate this we propose a robust version of the summarization problem by minimizing the worst-case risk when the distribution of $B$ is allowed to vary in a ball $\mathcal{N}_\delta(\hat{\mu})$ of radius $\delta > 0$ centered at $\hat{\mu}$ in a suitable metric, leading to the following robust optimization problem:

---

**Problem 1.2.** *Given adjacency matrices $B_1, \ldots, B_N$ of an independent sample with the same distribution as $\mathcal{G}$ find a symmetric matrix $\overline{A}$ which minimizes*

$$R(A) := \left( \sup_{\nu \in \mathcal{N}_\delta(\hat{\mu})} \mathbb{E}_{B \sim G}[\|A - B\|_1] \right).$$

The seminal work of Kuhn et. al. ♣♣♣ Mauricio: [ref?] shows that using a Wasserstein metric among probability distributions in the problem above leads to a tractable convex optimization problem. Our first result is that this is also true for the robust graph summarization problem when the Wasserstein metric is induced by a semidefinitely representable metric.

**Theorem 1.3.** *Let $K$ be the set of $n \times n$ symmetric matrices with entries in $[0, 1]$. If $\delta > 0$ then problem 1.2 is equivalent to:*

$$\min_{(A, s_1, \ldots, s_N, \lambda) \in \mathcal{T}} \left( \lambda\delta + \frac{1}{N} \sum_{i=1}^{N} s_i + \frac{1}{N} \sum \|A - B_i\|_1 \right)$$

*where $\mathcal{T}$ is the set of $(A, \vec{s}, \lambda)$ satisfying the inequalities $\lambda \geq 0$ and $\eta_{\epsilon,i}(\lambda) \leq s_i$ as $i = 1, \ldots, N$ and $\epsilon$ ranges over all $\{0, 1\}$-symmetric matrices and*

$$\eta_{\epsilon,i}(\lambda) := \sup_{Y \in K} \left( \langle \epsilon, A - Y \rangle - \lambda \|B_i - Y\| \right)$$

The formulation of 1.2 in Theorem 1.3 is a finite-dimensional convex optimization problem which unfortunately contains an exponential number of constraints. We therefore introduce:

(1) A more tractable simplification which agrees with the original problem whenever the optimal $A$ occurs at a matrix with entries in $\{0, 1\}$.
(2) A relaxation which takes the form of a regularization of empirical risk minimization.

More specifically, we prove the following

**Theorem 1.4.** *If $A$ is a symmetric $\{0, 1\}$-matrix then the following statements hold:*

(1) *The number $\sup_{\nu \in B_\delta(\hat{\mu})} \mathbb{E}[\|A - B\|_1]$ is equal to the optimal value of the problem*

$$\min_{\mathcal{H}} \left( \lambda\delta + \frac{1}{N} \sum_{i=1}^{N} s_i + \frac{1}{N} \sum_{i=1}^{N} \|A - B_i\|_1 \right)$$

*where $\mathcal{H}$ is the set of $(\lambda, s_1, \ldots, s_N, \Lambda, W) \in \mathbb{R} \times \mathbb{R}^N \times \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n}$ which satisfy the inequalities:*

$$\|W\|_* \leq \lambda$$
$$\Lambda \geq 0$$
$$2A - 11^t - W + \Lambda \geq 0$$
$$\|\Lambda\|_1 \leq s_i - \langle 2A - 11^t - W, B_i \rangle \text{ for } i = 1, \ldots, N$$

(2) *The following inequality holds:*

$$\sup_{\nu \in B_\delta(\hat{\mu})} \mathbb{E}[\|A - B\|_1] \leq \frac{1}{N} \sum_{i=1}^{N} \|A - B_i\|_1 + \delta\|2A - 11^t\|_*$$

*(3) If $\|\cdot\|$ is a semidefinitely representable function then the problems:*

$$(1) \qquad \min_A \min_{\mathcal{H}} \left( \lambda \delta + \frac{1}{N} \sum_{i=1}^{N} s_i + \frac{1}{N} \sum_{i=1}^{N} \|A - B_i\|_1 \right) \quad and$$

$$(2) \qquad \min_A \left( \frac{1}{N} \sum_{i=1}^{N} \|A - B_i\|_1 + \delta \|2A - 11^t\|_* \right)$$

*are semidefinite programming problems.*

Solving the optimization problems in part (3) of Theorem 1.4 leads to a new algorithm for estimating deterministic graph summaries which we call *Wasserstein robust graph summarization.* We applied this algorithm to a variety of graphs $G$ distributed according to the stochastic block model and noted that the regularized problem with the Wasserstein metric induced by the nuclear norm, henceforth *Wasserstein robust nuclear norm summarization,* was able to recover the correct cluster structure using only very few samples outperforming all other methods known to us. Our next result is a performance guarantee which shows that exact recovery occurs for suitable $\delta > 0$ with overwhelming probability explaining the good practical performance of Wasserstein nuclear norm summarization for cluster detection.

Recall that a random graph $\mathcal{G}$ has distribution given by the stochastic block model in $n$ vertices if there is a partition of $[n]$ into disjoint subsets $C_1, \ldots, C_k$, real numbers $0 \le q < \frac{1}{2} < p_1, \ldots, p_k \le 1$ and edges are added independently with probability $p_{ij}$ of joining vertices $i, j$ by an edge where

$$p_{ij} := \begin{cases} p_t, \text{ if } \{i, j\} \subseteq C_t \\ q, \text{ else.} \end{cases}$$

Note that a random graph $\mathcal{G}$ has a unique deterministic summary given by $A^*$ with entries given by

$$A_{ij}^* := \begin{cases} 1, \text{ if } \{i, j\} \in C_t \text{ for some } t \\ 0, \text{ otherwise.} \end{cases}$$

**Theorem 1.5.** *Suppose $B_1, \ldots, B_N$ are independent and have the same distribution as $\mathcal{G}$. If $\alpha = \min(|p_t - \frac{1}{2}|, |q - \frac{1}{2}|)$ and $\delta^*$ is the maximum of $a(\delta) := \frac{\delta \left( \alpha - \frac{\delta}{n} \right)^2}{\left( 1 + \frac{2\delta}{n} \right)}$ in $[0, \alpha n]$ then the probability that $A^*$ is not a minimizer of (2) is bounded above by*

$$\exp \left( -\frac{2N(n-1)^3(\alpha n - \delta^*)^2}{n} \right) + e^{-N(\delta^* a(\delta^*))} \prod_{i \ne j} \left( 1 + (e^{-4a}\widetilde{p_{ij}} + \widetilde{q_{ij}})^{\frac{N}{2}} \right).$$

*Moreover this quantity decreases exponentially with the sample size $N$.*

The key point of the Theorem, discussed at length in Section **??**, is that we were able to understand why the nuclear norm is a good regularizer for graph summarization problems, namely because the subdifferential of the regularizer at $A^*$ is sufficiently rich so as to contain enough transportation matrices.

Next we focus on the practical performance of the proposed algorithm. Problems (1) and (2) require solving large semidefinite programs which are beyond the

capacity of standard off-the-shelf software even for relatively small graphs (of say 40 vertices with $N = 4$). One possible reason is that off-the-shelf solvers often use interior point methods. A better alternative, especially well suited for solving **??** is to use first order numerical optimization methods such as ADMM. In Section **??** we adapt the ADMM algorithm for our regularized problem and present our open source Julia implementation allowing us to run the algorithm in graphs of up to 10000 vertices on a common laptop. In Section **??** we use this algorithm to find a deterministic summary of voting patterns in the colombian senate ♣♣♣ Mauricio: [Period].

1.1. **Acknowledgments.** We wish to thank Fabrice Gamboa, Mauricio Junca and Thierry Klein for useful conversations during the completion of this project. ♣♣♣ Mauricio: [CienciaPolitica]. D. De Roux was partially supported by a grant from Facultad de Ciencias, Universidad de los Andes. M. Velasco was partially supported by research funds from Universidad de los Andes.

## 2. PRELIMINARIES

2.1. **Preliminaries on graphs.** By a graph $G$ we mean a finite loopless undirected graph. We say that $G$ is weighted if it is endowed with a function $w : E(G) \to \mathbb{R}$ which assigns to every edge a real number in $[0, 1]$. If $G$ has $n$ vertices then it is completely specified by its adjacency matrix $A \in \{0, 1\}^{n \times n}$ defined by $A_{ij} = 1$ if and only if vertices $i, j$ are connected. If $G$ is weighted then we use the term adjacency matrix of $G$ to denote the matrix with entries $A_{i,j} = w(i, j)$.

If $A$ is a matrix then we use $\| \bullet \|$, $\| \bullet \|_1$ to denote its operator norm and $\ell^1$-norm respectively.

## 3. A DESCRIPTION OF THE PROBLEM

By a random graph $B$ we mean a random variable $B$ taking values on the set of adjacency matrices of graphs (i.e. symmeric matrices in $\{0, 1\}^{n \times n}$ with zero diagonal). By a random weighted graph we mean a random variable taking values in the adjacency matrices of weighted graphs (i.e. symmetric matrices with 0 diagonal all of whose off-diagonal entries lie in $[0, 1]$).

**Definition 3.1.** *Let $B$ be a random weighted graph and let $A$ be a (weighted) adjacency matrix. Define the risk of choosing $A \in \{0, 1\}^n$ as a deterministic summary of $B$ as*

$$R(A) := \mathbb{E}[\|A - B\|_1].$$

*We say that $A^*$ is an optimal summary of a random graph $B$ in the set $S$ if it is a minimizer of the optimization problem $\min_{A \in S} R(A)$.*

## 4. RECOVERING CLUSTER STRUCTURES IN THE STOCHASTIC BLOCK MODEL

Suppose $B$ is a random (undirected loopless) graph on $n$ vertices generated by the stochastic block model. This means that we fix a set partition $C_1, \ldots, C_l$ of $[n]$ into sets we call clusters and real numbers $0 \le p_i, \bar{p} \le 1$ for $i = 1, \ldots, l$. The edges of $B$ are independent random variables and an edge joins vertices $i, j$ with

probability $p_t$ if $\{i, j\} \subseteq C_t$ for some cluster $C_t$ and with probability $\bar{p}$ if $\{i, j\}$ is not contained in any $C_t$. We let $O \subseteq [n] \times [n]$ be the set of pairs of vertices which are not simultaneously contained in any cluster.

For an integer $N$ let $B_1, \ldots, B_N$ be an independent sample of $N$ graphs with the distribution of $B$. Let $A^*$ be the $n \times n$ matrix with entries in $\{0, 1\}$ which captures the underlying cluster structure, namely $A^*_{ij} = 1$ iff there is a cluster $C_t$ which contains both $ij$. In this section we study the probability, as a function of $\delta$ that the optimization problem $\min_A \Delta(A)$

$$\Delta(A) = \delta \|2A - 11^t\|_* + \frac{1}{N} \sum_{k=1}^N \|A - B_k\|_1$$

has the correct cluster structure $A^*$ as a minimizer. For vertices $i, j \in [n]$ define $n_1(ij)$ (resp. $n_0(ij)$) the random variables which count the number of times that a given pair is (resp is not) an edge of some $B_j$, $j = 1, \ldots N$. Note that the $n_q(ij)$ for $q = 0, 1$ are binomial random variables.

**Lemma 4.1.** *The following statements hold:*

(1) *The subdifferential of $\frac{1}{N} \sum_{k=1}^N \|A - B_k\|_1$ at $A^*$ is the set of symmetric matrices $C$ satisfying the inequalities*

$$\frac{n_0(ij) - n_1(ij)}{N} \le C_{ij} \le 1, \text{ if } \{i, j\} \subseteq C_t \text{ for some } t \text{ and}$$

$$-1 \le C_{ij} \le \frac{n_0(ij) - n_1(ij)}{N} \text{ if } \{i, j\} \text{ does not belong to any cluster.}$$

(2) *The subdifferential of $\delta \|2A - 11^t\|_*$ at $A^*$ is given by the set of symmetric matrices of the form $2\delta C$ where $C$ has spectral norm $\|C\| \le 1$ and satisfies $\langle C, 2A - 11^t \rangle = n$.*

*Proof.* (1) Since the subdifferential is additive it suffices to understand the subdifferential of the absolute value. If $i, j \in C_t$ then $A^*_{ij} = 1$ and the entry $ij$ of the subdifferential of the sum at $A^*$ is $[-1, 1]$ for each $B_i$ containing the edge and it is 1 for each $B_i$ for which $(ij)$ is not an edge. If $i, j$ is not contained in any cluster then $A^*_{ij} = 0$ and the entry $ij$ of the subdifferential of the sum at $A^*$ is $-1$ for each $B_i$ which contains the edge $ij$ and $[-1, 1]$ for each $B_i$ which does not, proving the claim. (2) It is easy to prove that the subdifferential of any norm $\| \bullet \|$ at a point $X$ is given by those $C$ for which the dual norm $\|C\|_* \le 1$ and $\langle C, X \rangle = \|X\|$. Claim (2) follows because $\|2A - 11^t\| = Tr(2A - 11^t) = n$ where the first equality holds since $2A - 11^t$ is positive semidefinite.♣♣♣ Mauricio: [Esto es obvio con solo dos clusters (la matriz es $uu^t$ donde $u$ es el vector con 1's en un cluster y $-1$'s en el complemento pero hay que demostrarlo para tres o mas)]. □

**Lemma 4.2.** *If $\delta = 0$ then*

$$\mathbb{P}\{A^* \in \arg\min \Delta\} = \prod_{ij \in O} \mathbb{P}\{n_1(ij) \le n_0(ij)\} \prod_{t=1}^k \left( \prod_{(ij) \in C_t} \mathbb{P}\{n_0(ij) \le n_1(ij)\} \right)$$

*where $\mathbb{P}\{n_0(ij) \le n_1(ij)\}$ is given by the following formula* ♣♣♣ Mauricio: [Ejercicio para Daniel: Encontrar una fórmula, es la probabilidad de que haya mas caras que

sellos en $N$ lanzamientos de una moneda trucada donde las probabilidades de la moneda dependen sólo de la arista $ij$].

*Proof.* The matrix $A^*$ is a minimizer of the above convex function if and only if its subdifferential at $A^*$ contains the matrix 0. By part (1) of the previous Lemma this occurs if and only if $\frac{n_0(ij)-n_1(ij)}{N} \leq 0$ for $(ij)$ in a cluster and $\frac{n_0(ij)-n_1(ij)}{N} \geq 0$ for $(ij)$ in $O$. Independence of the edges then implies the above formula.  $\square$

*Remark* 4.3. Could the set of minimizers be larger? If $A'$ has at least one entry $A_{ij} \in (0,1)$ the corresponding component in the subgradient is the constant $n_0(ij)-n_1(ij)$ and this equals zero with much smaller probability, precisely when both terms equal to $\frac{N}{2}$. In particular it is impossible if $N$ is odd and in this case $A^*$ is the only minimizer.

Next we ask whether it is possible to increase the probability of correct recovery by allowing $\delta > 0$. By Lemma 4.1 $A^* \in \arg\min \Delta$ if and only if there exists $S$ in the subdifferential of $\delta\|2A-11^t\|_*$ at $A = A^*$ such that $-S$ belongs to the subdifferential of $\frac{1}{N}\sum_{k=1}^N \|A - B_k\|_1$ at $A^*$.

The most immediate way to do this would be to find an $S$ with $S_{ij} \leq 0$ negative on edges $ij \in C_t$ and $S_{ij} \geq 0$ on edges of $O$. More precisely we would like to find a best such $C$ by solving the optimization problem:

$$\min\left(\sum_t \sum_{(ij)\in C_t} C_{ij}\right) - \sum_{(ij)\in O} C_{ij} \text{ s.t. } \|B\| \leq 1, \langle C, 2A^* - 11^t \rangle = n$$

However it is easy to see that we cannot do this improvement simultaneously in all components, because $n = \langle S, 2A - 11^t \rangle = Tr(S) + \sum_{ij \in O^c} S_{ij} - \sum_{ij \in O} S_{ij} \leq n + u$ where $u$ is the objective function in the problem above. We conclude that $u$ must be nonnegative so we cannot improve simultaneously in all directions at once. In the following section we discuss how tradeoffs between components explain the improved recovery probability induced by the spectral norm.

## 5. ESTIMATING RECOVERY PROBABILITIES

Let $\Gamma$ be the symmetric matrix with zero diagonal and off-diagonal entries given by $\Gamma_{ij} = \frac{n_0(ij)-n_1(ij)}{N}$. By Lemma 4.1 a symmetric matrix $C$ lies in the subdifferential if and only if it satisfies the inequalities

$$\Gamma_{ij} \leq C_{ij} \leq 1, \text{ if } \{i,j\} \subseteq C_t \text{ for some } t \text{ and}$$

$$-1 \leq C_{ij} \leq \Gamma_{ij} \text{ if } \{i,j\} \text{ does not belong to any cluster.}$$

To simplify these inequalities we define a linear operator $\widetilde{\bullet}$ on symmetric matrices by the formula

$$\widetilde{A} = \begin{cases} A_{ij} \text{ if } i = j \text{ or } ij \in I \text{ and} \\ -A_{ij} \text{ if } ij \in O. \end{cases}$$

in this language $C_{ij}$ belongs to the subdifferential if and only if $\widetilde{\Gamma}_{ij} \leq \widetilde{C}_{ij}$ for $i \neq j$. The following key result gives sufficient conditions for the true cluster structure $A^*$

to be a minimizer of the proposed optimization problem. In order to describe it we introduce the following notation.

**Definition 5.1.** *Let $\delta$ be a positive real number. For a symmetric matrix $\Gamma$ define the quantities*

$$b(\Gamma, \delta) := \sum_{i \neq j} \max \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n}, 0 \right) \ \text{and} \ a(\Gamma, \delta) := \sum_{i \neq j} \max \left( -\widetilde{\Gamma_{ij}} - \frac{2\delta}{n}, 0 \right)$$

The quantity $b(\Gamma, \delta)$ (resp. $a(\Gamma, \delta)$) measures the total amount by which the matrix $-\frac{2\delta}{n} 11^t$ fails (resp. succeeds) to be in the subdifferential of Lemma 4.1 in the sense that it sums over all $ij$ the amount by which the inequalities $\widetilde{\Gamma_{ij}} \leq \frac{2\delta}{n} 11^t$ fail (resp. succeed). The key point of the following Theorem is that if the inequality fails by less than it succeeds then the subdifferential of the spectral norm is sufficiently rich so as to allow us to redistribute these quantities. In this sense the following Theorem explains the success of the spectral norm in cluster recovery algorithms.

**Theorem 5.2.** *Assume there are only two clusters. If $b(\Gamma, \delta) \leq \min(\delta, a(\Gamma, \delta))$ then $A^*$ is a minimizer of the optimization problem $\min_A \Delta(A)$.*

*Proof.* We will show that there exists a matrix $C_{ij}$ such that $-C_{ij} \in \partial \left( \delta \| 2A - 11^t \| \right) (A^*)$ for which $\widetilde{\Gamma_{ij}} \leq \widetilde{C_{ij}}$ for $i \neq j$. It will then follow that $0 = C - C$ belongs to the subdifferential of $\Delta(A)$ at $A^*$ and thus $A^*$ is a minimizer as claimed.

Recall that $-C \in \partial \left( \delta \| 2A - 11^t \| \right) (A^*)$ if and only if it satifies the conditions

$$\langle H^*, C \rangle = -2\delta n \text{ and } \|C\| \leq 2\delta$$

where $H^* := 2A^* - 11^t$ and $\| \bullet \|$ is the spectral norm. Both of these conditions are satisfied by setting $C = -\frac{2\delta}{n} H^*$. However this choice of $C$ will not, in general, satisfy the inequalities $\widetilde{\Gamma_{ij}} \leq -\frac{2\delta}{n} H^*_{ij} = -\frac{2\delta}{n} 11^t_{ij}$ for $i \neq j$.

We will adjust our candidate for $\widetilde{C}$ by adding to it a transportation matrix $\widetilde{K}$ that will guarantee that all these inequalities are satisfied when $b(\Gamma, \delta) \leq a(\Gamma, \delta)$. Crucially we will choose $\widetilde{K}$ so that $K$ satisfies $KH^* = H^*K = 0$ allowing us to control the spectral norm of $C$.

Since $b(\Gamma, \delta) \leq a(\Gamma, \delta)$ there exists a way to redistribute the quantity $b(\Gamma, \delta)$ by substracting it from the $ij$ for which $-\frac{\delta}{n} \leq \widetilde{\Gamma}_{ij}$ and adding it into those $st$ for which $\widetilde{\Gamma}_{st} \leq -\frac{\delta}{n}$. More specifically, if $b = \widetilde{\Gamma}_{ij} + \frac{2\delta}{n} > 0$ then there exists a set $(i_1, j_1), \ldots, (i_t, j_t)$ of non-diagonal entries and nonnegative constants $\gamma_{i_s, j_s}$ summing to one such that $\widetilde{\Gamma}_{i_s j_s} + \frac{2\delta}{n} + b\gamma_{i_s, j_s} \leq 0$. Define the off-diagonal elements of $\widetilde{C}$ by

$$\widetilde{C} := -\frac{2\delta}{n} 11^t + \sum_{i_s, j_s} b\gamma_{i_s, j_s}(-e_{ij} + e_{i_s, j_s})$$

where $e_{ij}$ is the symmetric matrix with one in positions $i, j$ and $j, i$ and zeroes otherwise. More generally, let $B$ be the set of paris $(i, j)$ with $i < j$ such that $\widetilde{\Gamma}_{ij} + \frac{2\delta}{n} > 0$. If $b(\Gamma, \delta) \leq a(\Gamma, \delta)$ then for each $ij \in B$ there exist nonnegative

constants $\gamma_{st}^{(ij)}$ such that $\sum_{s<t} \gamma_{st}^{(ij)} = 1$ and for which the matrix

$$\widetilde{C} := -\frac{2\delta}{n} 11^t + \sum_{ij \in B} \sum_{s<t} \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n} \right) \gamma_{st}^{(ij)} (-e_{ij} + e_{st})$$

satisfies $\widetilde{\Gamma}_{ab} \leq \widetilde{C}_{ab}$ for all $a \neq b$. We will show that if $b(\Gamma, \delta) \leq 2\delta$ then the matrix $C$ with off-diagonal entries given by

$$C_{ab} = -\frac{2\delta}{n} H_{ab}^* + \sum_{ij \in B} \sum_{s<t} \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n} \right) \gamma_{st}^{(ij)} (\widetilde{-e_{ij} + e_{st}})_{ab}$$

lies in $-\partial \left( \delta \| 2A - 11^t \|_* \right)(A^*)$ for some choice of diagonal, proving the Theorem. For a pair of indices $i < j$ let $\epsilon_{ij} = 1$ if $ij \in I$ and $\epsilon_{ij} = -1$ if $ij \in O$. Define the matrix $t_{ij}$ by $t_{ij} = \epsilon_{ij} e_{ij} - e_{ii} - e_{jj}$ and note that $t_{ij}$ satisfies the following three properties:

(1) The equalities $H^* t_{ij} = 0 = t_{ij} H^*$ hold. If $ij \in O$ this happens only when there are exactly two clusters and this is the only point in the proof where this assumption is used.
(2) The off-diagonal entries of $\widetilde{t_{ij}}$ are equal to those of $e_{ij}$. In particular the off-diagonal entries of $(\widetilde{-e_{ij} + e_{st}})_{ab}$ are always equal to those of $\widetilde{-t_{ij} + t_{st}}$.
(3) The inequality $\| - t_{ij} + t_{st} \| \leq 2$ holds for all $ij$ and $st$. This is immediate via direct calculation (the inequality is strict only if $|\{i, j\} \cap \{s, t\}| \geq 1$ and in this case it can take values of $\sqrt{3}$ and $0$).

If $C$ denotes the matrix given by

$$C := -\frac{2\delta}{n} H^* + \sum_{ij \in B} \sum_{s<t} \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n} \right) \gamma_{st}^{(ij)} (-t_{ij} + t_{st})$$

then the following properties hold:

(1) The off-diagonal entries agree with those in our previous expression so $\widetilde{\Gamma}_{ab} \leq \widetilde{C}_{ab}$ for all $a \neq b$ and thus $C$ is in the subdifferential of Lemma 4.1.
(2) Since $H^* t_{ij} = 0 = t_{ij} H^*$ the equality $\langle H^*, C \rangle = \langle H^*, -\frac{2\delta}{n} H^* \rangle = -2\delta n$ holds and moreover

$$\|C\| = \max \left( \left\| -\frac{2\delta}{n} H^* \right\|, \left\| \sum_{ij \in B} \sum_{s<t} \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n} \right) \gamma_{st}^{(ij)} (-t_{ij} + t_{st}) \right\| \right).$$

The operator norm of the first term in the maximum equals $2\delta$ and that of the second term is bounded by $2b(\Gamma, \delta)$ by the triangle inequality and the definition of $b(\Gamma, \delta)$. We conclude that $\|C\|$ is bounded by $2\delta$ because $b(\Gamma, \delta) \leq \delta$. As a result $-C \in \partial \left( \| 2A - 11^t \| \right)(A^*)$ proving the Theorem.  □

♣♣♣ Mauricio: [Como extendemos este razonamiento al caso de tres o más clusters? Debe ser posible utilizar otras matrices de transporte para este caso que vivan en el kernel. Algo interesante es que con tres clusters es necesario que los promedios de una matriz aniquilada por $H^*$ dentro de cada uno de los bloques $C_i \times C_j$ deban ser cero.]

We will now prove the general case when there are more than two clusters. The proof will be similar the proof of the previous theorem. We will start with the

candidate matrix $\frac{-2\delta}{n}H^*$ and correct it by adding transport matrices that assure that the corrected matrix belong to the subdifferential of $\Delta(A)$ at $A^*$. The difficulty in applying the tools of the previous theorem to solve the general case is that the matrices $K$ that transport weight from one cluster to another do not, in general, satisfy the relation $KH^* = H^*K = 0$ when there are more than two clusters. This problem can be solved by splitting the matrix $H^*$ into matrices than only take into account 2 clusters, and using the previous theorem. We begin recalling the following simple result:

**Claim 5.3.** *Let $a_i, b_i \geq 0$, $i = 1, \ldots, p$ be two non-negative, finite sequences of real numbers such that*

$$\sum_{i=1}^{p} a_i \geq \sum_{i=1}^{p} b_i.$$

*Then there exist a finite sequence of reals $c_i$ such that*

- *$\sum_{i=1}^{p} c_i = 0$.*
- *$a_i \geq b_i + c_i \ \forall i$.*

Now we proceed to do the proof. For clusters $C_i \neq C_j$ define the matrix $H^{C_iC_j}$ whose entries are given by:

$$H_{uv}^{C_iC_j} = \begin{cases} 1 \text{ if } u, v \in C_i \text{ or } u, v \in C_j. \\ -1 \text{ if } u \in C_i, v \in C_j \text{ or } u \in C_i, v \in C_j. \\ 0 \text{ in any other case.} \end{cases}$$

Notice that this matrix has 4 blocs. Two with only 1 and two with only $-1$ moreover, its spectral norm is equal to $|C_i| + |C_j|$.

**Lemma 5.4.** *Assume there are $l$ clusters. Suppose that $b(\Gamma, \delta) \leq min(\delta, a(\Gamma, \delta))$. Then, $A^*$ is a minimizer of the optimization problem $\min_A \Delta(A)$.*

*Proof.* First of all, observe that

$$\frac{-2\delta}{n}H^* = \frac{-2\delta}{n}\frac{1}{l-1}\sum_{1\leq i<j\leq l} H^{C_iC_j}.$$

For each, $C_i \neq C_j$ construct the transport matrix $\Delta_{ij}$ as in the previous theorem, as to assure that $\Delta_{ij}H^{C_iC_j} = H^{C_iC_j}\Delta_{ij} = 0$. This can be done since $H^{C_iC_j}$ takes into account only two clusters. Recall that the total amount of weight to be corrected is $b(\Gamma, \delta)$. Let $w_{i,j}$ the weight to be distributed from cluster $i$ to cluster $j$. In the notation of the previous theorem, $w_{ij}$ is just the sum of the $\gamma_{i_sj_s}$ where $i_s \in i$ and $j_s \in j$.

Let

$$C := -\frac{2\delta}{n}H^* + \sum_{i<j}\Delta_{ij} = \frac{-2\delta}{n}\frac{1}{l-1}\sum_{1\leq i<j\leq l}H^{C_iC_j} + \sum_{i<j}\Delta_{ij}$$

Finally, assume that for each $i < j$, $\frac{1}{l-1}\|H^{C_iC_j}\| \geq \|\Delta_{ij}\|$. Then,

$$\|C\| = \left\|\frac{-2\delta}{n}\frac{1}{l-1}\sum_{1\leq i<j\leq l} H^{C_iC_j} + \sum_{i<j}\Delta_{ij}\right\|$$

$$\leq \frac{2\delta}{n(l-1)}\sum_{1\leq i<j\leq l}\left\|H^{C_iC_j} + \frac{n(l-1)}{2\delta}\Delta_{ij}\right\|$$

$$= \frac{2\delta}{n(l-1)}\sum_{1\leq i<j\leq l}\max(\|H^{C_iC_j}\|, \left\|\frac{n(l-1)}{2\delta}\Delta_{ij}\right\|)$$

Now notice that

$$\delta \geq b(\Gamma,\delta) \Rightarrow (l-1)n \geq \frac{(l-1)n}{\delta}b(\Gamma,\delta) \Rightarrow \sum_{i<j}\|H^{C_iC_j}\| \geq \frac{(l-1)n}{\delta}\sum_{i<j}w_{i,j}$$

$$\geq \frac{(l-1)n}{2\delta}\sum_{i<j}\|\Delta_{i,j}\|.$$

By the claim, we can assume without loss of generality that for each $i < j$,

$$\|H^{C_i,C_j}\| \geq \frac{(l-1)n}{\delta}\|\Delta_{i,j}\|.$$

This implies that the last sum reduces to

$$\frac{2\delta}{n(l-1)}\sum_{1\leq i<j\leq l}\|H^{C_iC_j}\| = \sum_{1\leq i<j\leq l}\frac{2\delta(|C_i|+|C_j|)}{n(l-1)} = \frac{2\delta(l-1)}{n(l-1)}\sum_{1\leq i<j\leq l}|C_i|+|C_j| = 2\delta.$$

And so $\|C\| \leq 2\delta$.

$\square$

Using the previous Theorem we now estimate the probabilities of perfect recovery of the correct cluster structure.

5.1. **A bound for recovery probabilities.** In this section we will bound the probability that the correct $A^*$ is not an optimal solution of our proposed optimization problem. A key tool will be the following version of Hoeffding's inequality: If $X_1, \ldots, X_T$ are independent random variables with values in $[c_i, d_i]$ and $\Lambda_T := \sum_{i=1}^T X_i$ then the following inequality holds for all $t \geq 0$

$$\mathbb{P}\{\Lambda_T - \mathbb{E}[\Lambda_T] \geq t\} \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^T(d_i-c_i)^2}\right).$$

**Theorem 5.5.** *Suppose $B_1, \ldots, B_N$ are independent and have the same distribution as $\mathcal{G}$. If $\alpha = \min(|p_t - \frac{1}{2}|, |q - \frac{1}{2}|)$ and $\delta^*$ is the maximum of $a(\delta) := \frac{\delta(\alpha-\frac{\delta}{n})^2}{(1+\frac{2\delta}{n})}$ in $[0,\alpha n]$ then the probability that $A^*$ is not a minimizer of (2) is bounded above by*

$$\exp\left(-\frac{2N(n-1)^3(\alpha n - \delta^*)^2}{n}\right) + e^{-N(\delta^* a(\delta^*))}\prod_{i\neq j}\left(1 + (e^{-4a}p_{ij} + q_{ij})^{\frac{N}{2}}\right).$$

*Moreover this quantity decreases exponentially with the sample size $N$.*

*Proof.* By Lemma 5.4 and the union bound the probability that $A^*$ is not an optimal solution of problem (1.2) is bounded above by

$$(3) \qquad \mathbb{P}\{b(\Gamma, \delta) \geq a(\Gamma, \delta)\} + \mathbb{P}\{b(\Gamma, \delta) \geq \delta\}.$$

and we will find upper bounds for the individual terms in (3). For $t = 1, \dots, N$ and $i, j \in [n]$ define

$$Z_{ij}^{(t)} := \begin{cases} -1, & \text{if } (B_t)_{ij} = 1 \\ +1, & \text{if } (B_t)_{ij} = 0 \end{cases}$$

and note that for every $i \neq j$ the equality $\sum_{t=1}^{N} \frac{Z_{ij}^{(t)}}{N} = \Gamma_{ij}$ holds. As a result

$$b(\Gamma, \delta) - a(\Gamma, \delta) = \sum_{i \neq j} \left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n} \right) = \frac{2\delta n(n-1)}{n} + \sum_{t=1}^{N} \sum_{i \neq j} \frac{\widetilde{Z_{ij}^{(t)}}}{N}$$

and therefore if $M := \mathbb{E}\left[ \sum_{t=1}^{N} \sum_{i \neq j} \frac{\widetilde{Z_{ij}^{(t)}}}{N} \right]$ then the number $M$ is negative and is given by the formula

$$M = (2q - 1)2|O| + \sum_{i=1}^{l} 2 \binom{c_i}{2}(1 - 2p_i) \leq -2\alpha n(n-1)$$

We can therefore bound the probability in the first term with

$$\mathbb{P}\left\{ \frac{2\delta n(n-1)}{n} + \sum_{t=1}^{N} \sum_{i \neq j} \frac{\widetilde{Z_{ij}^{(t)}}}{N} \geq 0 \right\} = \mathbb{P}\left\{ \sum_{t=1}^{N} \sum_{i \neq j} \frac{\widetilde{Z_{ij}^{(t)}}}{N} - M \geq -2\delta(n-1) - M \right\} \leq$$

$$\leq \exp\left( -2 \frac{(-M - 2\delta(n-1))^2}{Nn(n-1)(\frac{2}{N})^2} \right) = \exp\left( -\frac{N}{2}\left(1 - \frac{1}{n}\right)\left(\frac{-M}{n-1} - 2\delta\right)^2 \right)$$

where the inequality follows from Hoeffding's inequality applied to the $Nn(n-1)$ independent random variables $\frac{\widetilde{Z_{ij}^{(t)}}}{N}$ which have values in $\left[-\frac{1}{N}, \frac{1}{N}\right]$. The inequality applies whenever $-\frac{M}{2(n-1)} > \delta > 0$. In particular whenever $0 < \delta < \alpha n$ we have

$$\mathbb{P}\{b(\Gamma, \delta) - a(\Gamma, \delta)\} \leq \exp\left( -\frac{2N(n-1)^3(\alpha n - \delta)^2}{n} \right).$$

Bounding the second term is more involved. Recall that

$$b(\Gamma, \delta) = \sum_{i \neq j} \max\left( \widetilde{\Gamma_{ij}} + \frac{2\delta}{n}, 0 \right).$$

Let $Y_{ij} := \widetilde{\Gamma_{ij}} + \frac{2\delta}{n}$ and let $X_{ij} := \max(Y_{ij}, 0)$. In order to prove a concentration inequality for the variables $X_{ij}$ we begin by studying their moment generating functions $m_{X_{ij}}(t)$. Note that for every real number $t$ the equality

$$\exp(tX_{ij}) = 1_{\{Y_{ij} \leq 0\}} + 1_{\{Y_{ij} \geq 0\}} \exp tY_{ij}$$

holds. Now $Y_{ij} = 1 + \frac{2\delta}{n} - 2\frac{n_{ij}}{N}$ where $n_{ij}$ is a binomial random variable with parameters $N$ and $p_{ij}$ given by

$$p_{ij} := \begin{cases} p_t, \text{ if } \{i, j\} \subseteq C_t \\ 1 - q, \text{ else} \end{cases}$$

As a result taking expected values on both sides of the expression above we conclude that

$$m_{X_{ij}}(t) \le \mathbb{P}\{Y_{ij} \le 0\} + e^{t\left(1 + \frac{2\delta}{n}\right)} \mathbb{E}\left(e^{-\frac{2t}{N}n_{ij}} 1_{\{Y_{ij} \ge 0\}}\right).$$

Using the Cauchy-Schwartz inequality and the known formula for the moment generating function of a binomial random variable it follows that

$$m_{X_{ij}}(t) \le \mathbb{P}\{Y_{ij} \le 0\} + e^{t\left(1 + \frac{2\delta}{n}\right)} \mathbb{E}\left(e^{-\frac{4t}{N}n_{ij}}\right)^{\frac{1}{2}} \mathbb{P}\{Y_{ij} \ge 0\}^{\frac{1}{2}} =$$

$$= \mathbb{P}\{Y_{ij} \le 0\} + e^{t\left(1 + \frac{2\delta}{n}\right)} \left(e^{-\frac{4t}{N}}p_{ij} + q_{ij}\right)^{\frac{N}{2}} \mathbb{P}\{Y_{ij} \ge 0\}^{\frac{1}{2}}$$

where $q_{ij} := 1 - p_{ij}$. By Hoeffding's inequality on Bernoulli random variables we know that

$$\mathbb{P}\{Y_{ij} \ge 0\} \le \exp\left(-\frac{N}{2}\left(-\frac{2\delta}{n} - (1 - 2p_{ij})\right)^2\right) \le \exp\left(-2N(\alpha - \delta/n)^2\right)$$

so if $t = aN$ the inequality

$$e^{t\left(1 + \frac{2\delta}{n}\right)} \exp\left(-N(\alpha - \delta/n)^2\right) \le 1$$

holds whenever $a \le \frac{(\alpha - \delta/n)^2}{\left(1 + \frac{2\delta}{n}\right)}$ and for all such $a$ we have

$$m_{X_{ij}}(aN) \le \left(1 + (e^{-4a}p_{ij} + q_{ij})^{\frac{N}{2}}\right)$$

We define $a(\delta) := \frac{(\alpha - \delta/n)^2}{\left(1 + \frac{2\delta}{n}\right)}$ and will use it to prove a moment concentration inequality for $b(\Gamma, \delta)$ which will give us a bound on the second term in (3). For every $t > 0$ we have

$$\mathbb{P}\{b(\Gamma, \delta) \ge \delta\} = \mathbb{P}\left\{\exp\left(t \sum_{i \ne j} X_{ij}\right) \ge e^{t\delta}\right\} \le e^{-t\delta} \prod_{i \ne j} \mathbb{E}[e^{tX_{ij}}] = e^{-t\delta} m_{X_{ij}}(t)$$

Choosing $t = a(\delta)N$ and using the previous inequality we see that

$$\mathbb{P}\{b(\Gamma, \delta) \ge \delta\} \le e^{-N\delta a(\delta)} \prod_{i \ne j} \left(1 + (e^{-4a}p_{ij} + q_{ij})^{\frac{N}{2}}\right)$$

Which decreases exponentially in $N$ for any $0 \le \delta \le n\alpha$. The rate of decrease of the first term is controlled by the positive factor

$$\delta a(\delta) = \frac{\delta \left(\alpha - \frac{\delta}{n}\right)^2}{\left(1 + \frac{2\delta}{n}\right)}.$$

and we let $\delta^*$ be a maximizer of this function.                                   $\square$

Departamento de matemáticas, Universidad de los Andes, Carrera $1^{\text{ra}}\#18A-12$, Bogotá, Colombia

*E-mail address*: d.de1033@uniandes.edu.co

Departamento de matemáticas, Universidad de los Andes, Carrera $1^{\text{ra}}\#18A-12$, Bogotá, Colombia

*E-mail address*: mvelasco@uniandes.edu.co