

9 Community detection and the Stochastic Block Model

9.1 Community Detection

Community detection in a network is a central problem in data science. A few lectures ago we discussed clustering and gave a performance guarantee for spectral clustering (based on Cheeger's Inequality) that was guaranteed to hold for any graph. While these guarantees are remarkable, they are worst-case guarantees and hence pessimistic in nature. In what follows we analyze the performance of a convex relaxation based algorithm on typical instances of the community detection problem (where typical is defined through some natural distribution of the input).

We focus on the problem of minimum graph bisection. The objective is to partition a graph in two equal-sized disjoint sets (S, S^c) while minimizing $\text{cut}(S)$ (note that in the previous lecture, for the Max-Cut problem, we were maximizing it instead!).

9.2 Stochastic Block Model

We consider a random graph model that produces graphs that have a clustering structure. Let n be an even positive integer. Given two sets of $m = \frac{n}{2}$ nodes consider the following random graph G : For each pair (i, j) of nodes, (i, j) is an edge of G with probability p if i and j are in the same set, and with probability q if they are in different sets. Each edge is drawn independently and $p > q$. This is known as the Stochastic Block Model on two communities.

(Think of nodes as habitants of two different towns and edges representing friendships, in this model, people leaving in the same town are more likely to be friends)

The goal will be to recover the original partition. This problem is clearly easy if $p = 1$ and $q = 0$ and hopeless if $p = q$. The question we will try to answer is for which values of p and q is it possible to recover the partition (perhaps with high probability). As $p > q$, we will try to recover the original partition by attempting to find the minimum bisection of the graph.

9.3 What does the spike model suggest?

Motivated by what we saw in previous lectures, one approach could be to use a form of spectral clustering to attempt to partition the graph.

Let A be the adjacency matrix of G , meaning that

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E(G) \\ 0 & \text{otherwise.} \end{cases} \quad (78)$$

Note that in our model, A is a random matrix. We would like to solve

$$\begin{aligned} \max \quad & \sum_{i,j} A_{ij} x_i x_j \\ \text{s.t.} \quad & x_i = \pm 1, \forall_i \\ & \sum_j x_j = 0, \end{aligned} \quad (79)$$

The intended solution x takes the value $+1$ in one cluster and -1 in the other.

Relaxing the condition $x_i = \pm 1, \forall_i$ to $\|x\|_2^2 = n$ would yield a spectral method

$$\begin{aligned} \max \quad & \sum_{i,j} A_{ij} x_i x_j \\ \text{s.t.} \quad & \|x\|_2 = \sqrt{n} \\ & \mathbf{1}^T x = 0 \end{aligned} \tag{80}$$

The solution consists of taking the top eigenvector of the projection of A on the orthogonal of the all-ones vector $\mathbf{1}$.

The matrix A is a random matrix whose expectation is given by

$$\mathbb{E}[A] = \begin{cases} p & \text{if } (i, j) \in E(G) \\ q & \text{otherwise.} \end{cases}$$

Let g denote a vector that is +1 in one of the clusters and -1 in the other (note that this is the vector we are trying to find!). Then we can write

$$\mathbb{E}[A] = \frac{p+q}{2} \mathbf{1}\mathbf{1}^T + \frac{p-q}{2} gg^T,$$

and

$$A = (A - \mathbb{E}[A]) + \frac{p+q}{2} \mathbf{1}\mathbf{1}^T + \frac{p-q}{2} gg^T.$$

In order to remove the term $\frac{p+q}{2} \mathbf{1}\mathbf{1}^T$ we consider the random matrix

$$\mathcal{A} = A - \frac{p+q}{2} \mathbf{1}\mathbf{1}^T.$$

It is easy to see that

$$\mathcal{A} = (A - \mathbb{E}[\mathcal{A}]) + \frac{p-q}{2} gg^T.$$

This means that \mathcal{A} is a superposition of a random matrix whose expected value is zero and a rank-1 matrix, i.e.

$$\mathcal{A} = W + \lambda vv^T$$

where $W = (A - \mathbb{E}[\mathcal{A}])$ and $\lambda vv^T = \frac{p-q}{2} n \left(\frac{g}{\sqrt{n}}\right) \left(\frac{g}{\sqrt{n}}\right)^T$. In previous lectures we saw that for large enough λ , the eigenvalue associated with λ pops outside the distribution of eigenvalues of W and whenever this happens, the leading eigenvector has a non-trivial correlation with g (the eigenvector associated with λ).

Note that since to obtain \mathcal{A} we simply subtracted a multiple of $\mathbf{1}\mathbf{1}^T$ from A , problem (80) is equivalent to

$$\begin{aligned} \max \quad & \sum_{i,j} \mathcal{A}_{ij} x_i x_j \\ \text{s.t.} \quad & \|x\|_2 = \sqrt{n} \\ & \mathbf{1}^T x = 0 \end{aligned} \tag{81}$$

Now that we removed a suitable multiple of $\mathbf{1}\mathbf{1}^T$ we will even drop the constraint $\mathbf{1}^T x = 0$, yielding

$$\begin{aligned} \max \quad & \sum_{i,j} \mathcal{A}_{ij} x_i x_j \\ \text{s.t.} \quad & \|x\|_2 = \sqrt{n}, \end{aligned} \tag{82}$$

whose solution is the top eigenvector of \mathcal{A} .

Recall that if $\mathcal{A} - \mathbb{E}[\mathcal{A}]$ was a Wigner matrix with i.i.d entries with zero mean and variance σ^2 then its empirical spectral density would follow the semicircle law and it will essentially be supported in $[-2\sigma\sqrt{n}, 2\sigma\sqrt{n}]$. We would then expect the top eigenvector of \mathcal{A} to correlate with g as long as

$$\frac{p-q}{2}n > \frac{2\sigma\sqrt{n}}{2}. \tag{83}$$

Unfortunately $\mathcal{A} - \mathbb{E}[\mathcal{A}]$ is not a Wigner matrix in general. In fact, half of its entries have variance $p(1-p)$ while the variance of the other half is $q(1-q)$.

If we were to take $\sigma^2 = \frac{p(1-p)+q(1-q)}{2}$ and use (83) it would suggest that the leading eigenvector of \mathcal{A} correlates with the true partition vector g as long as

$$\frac{p-q}{2} > \frac{1}{\sqrt{n}} \sqrt{\frac{p(1-p)+q(1-q)}{2}}, \tag{84}$$

However, this argument is not necessarily valid because the matrix is not a Wigner matrix. For the special case $q = 1-p$, all entries of $\mathcal{A} - \mathbb{E}[\mathcal{A}]$ have the same variance and they can be made to be identically distributed by conjugating with gg^T . This is still an impressive result, it says that if $p = 1-q$ then $p-q$ needs only to be around $\frac{1}{\sqrt{n}}$ to be able to make an estimate that correlates with the original partitioning!

An interesting regime (motivated by friendship networks in social sciences) is when the average degree of each node is constant. This can be achieved by taking $p = \frac{a}{n}$ and $q = \frac{b}{n}$ for constants a and b . While the argument presented to justify condition (84) is not valid in this setting, it nevertheless suggests that the condition on a and b needed to be able to make an estimate that correlates with the original partition is

$$(a-b)^2 > 2(a+b). \tag{85}$$

Remarkably this was posed as conjecture by Decelle et al. [DKMZ11] and proved in a series of works by Mossel et al. [MNS14b, MNS14a] and Massoulié [Mas14].

9.3.1 Three of more communities

The stochastic block model can be similarly defined for any $k \geq 2$ communities: G is a graph on $n = km$ nodes divided on k groups of m nodes each. Similarly to the $k = 2$ case, for each pair (i, j) of nodes, (i, j) is an edge of G with probability p if i and j are in the same set, and with probability q if they are in different sets. Each edge is drawn independently and $p > q$. In the sparse regime, $p = \frac{a}{n}$ and $q = \frac{b}{n}$, the threshold at which it is possible to make an estimate that correlates with the original partition is open.

Open Problem 9.1 Consider the balanced Stochastic Block Model for $k > 3$ (constant) communities with inner probability $p = \frac{a}{n}$ and outer probability $q = \frac{b}{n}$, what is the threshold at which it becomes possible to make an estimate that correlates with the original partition is open (known as the partial recovery or detection threshold). We refer the reader to [DKMZ11, ZMZ14, GZC⁺15] for more information on this and many other interesting conjectures often motivated from statistical physics.

9.4 Exact recovery

We now turn our attention to the problem of recovering the cluster membership of every single node correctly, not simply having an estimate that correlates with the true labels. We'll restrict to two communities for now. If the probability of intra-cluster edges is $p = \frac{a}{n}$ then it is not hard to show that each cluster will have isolated nodes making it impossible to recover the membership of every possible node correctly. In fact this is the case whenever $p \ll \frac{2 \log n}{n}$. For that reason we focus on the regime

$$p = \frac{\alpha \log(n)}{n} \text{ and } q = \frac{\beta \log(n)}{n}, \quad (86)$$

for some constants $\alpha > \beta$.

Let $x \in \mathbb{R}^n$ with $x_i = \pm 1$ representing the partition (note there is an ambiguity in the sense that x and $-x$ represent the same partition). Then, if we did not worry about efficiency then our guess (which corresponds to the Maximum Likelihood Estimator) would be the solution of the minimum bisection problem (79).

In fact, one can show (but this will not be the main focus of this lecture, see [ABH14] for a proof) that if

$$\sqrt{\alpha} - \sqrt{\beta} > \sqrt{2}, \quad (87)$$

then, with high probability, (79) recovers the true partition. Moreover, if

$$\sqrt{\alpha} - \sqrt{\beta} < \sqrt{2},$$

no algorithm (efficient or not) can, with high probability, recover the true partition.

We'll consider a semidefinite programming relaxation algorithm for SBM and derive conditions for exact recovery. The main ingredient for the proof will be duality theory.

9.5 The algorithm

Note that if we remove the constraint that $\sum_j x_j = 0$ in (79) then the optimal solution becomes $x = \mathbf{1}$. Let us define $B = 2A - (\mathbf{1}\mathbf{1}^T - I)$, meaning that

$$B_{ij} = \begin{cases} 0 & \text{if } i = j \\ 1 & \text{if } (i, j) \in E(G) \\ -1 & \text{otherwise} \end{cases} \quad (88)$$

It is clear that the problem

$$\begin{aligned}
& \max \sum_{i,j} B_{ij}x_ix_j \\
& \text{s.t. } x_i = \pm 1, \forall_i \\
& \sum_j x_j = 0
\end{aligned} \tag{89}$$

has the same solution as (79). However, when the constraint is dropped,

$$\begin{aligned}
& \max \sum_{i,j} B_{ij}x_ix_j \\
& \text{s.t. } x_i = \pm 1, \forall_i,
\end{aligned} \tag{90}$$

$x = \mathbf{1}$ is no longer an optimal solution. Intuitively, there is enough “−1” contribution to discourage unbalanced partitions. In fact, (90) is the problem we’ll set ourselves to solve.

Unfortunately (90) is in general NP-hard (one can encode, for example, **Max-Cut** by picking an appropriate B). We will relax it to an easier problem by the same technique used to approximate the **Max-Cut** problem in the previous section (this technique is often known as *matrix lifting*). If we write $X = xx^T$ then we can formulate the objective of (90) as

$$\sum_{i,j} B_{ij}x_ix_j = x^T Bx = \text{Tr}(x^T Bx) = \text{Tr}(Bxx^T) = \text{Tr}(BX)$$

Also, the condition $x_i = \pm 1$ implies $X_{ii} = x_i^2 = 1$. This means that (90) is equivalent to

$$\begin{aligned}
& \max \quad \text{Tr}(BX) \\
& \text{s.t.} \quad X_{ii} = 1, \forall_i \\
& \quad \quad X = xx^T \text{ for some } x \in \mathbb{R}^n.
\end{aligned} \tag{91}$$

The fact that $X = xx^T$ for some $x \in \mathbb{R}^n$ is equivalent to $\text{rank}(X) = 1$ and $X \succeq 0$. This means that (90) is equivalent to

$$\begin{aligned}
& \max \quad \text{Tr}(BX) \\
& \text{s.t.} \quad X_{ii} = 1, \forall_i \\
& \quad \quad X \succeq 0 \\
& \quad \quad \text{rank}(X) = 1.
\end{aligned} \tag{92}$$

We now relax the problem by removing the non-convex rank constraint

$$\begin{aligned}
& \max \quad \text{Tr}(BX) \\
& \text{s.t.} \quad X_{ii} = 1, \forall_i \\
& \quad \quad X \succeq 0.
\end{aligned} \tag{93}$$

This is an SDP that can be solved (up to arbitrary precision) in polynomial time [VB96].

Since we removed the rank constraint, the solution to (93) is no longer guaranteed to be rank-1. We will take a different approach from the one we used before to obtain an approximation ratio for Max-Cut, which was a worst-case approximation ratio guarantee. What we will show is that, for some values of α and β , with high probability, the solution to (93) not only satisfies the rank constraint but it coincides with $X = gg^T$ where g corresponds to the true partition. After X is computed, g is simply obtained as its leading eigenvector.

9.6 The analysis

Without loss of generality, we can assume that $g = (1, \dots, 1, -1, \dots, -1)^T$, meaning that the true partition corresponds to the first $\frac{n}{2}$ nodes on one side and the other $\frac{n}{2}$ on the other.

9.6.1 Some preliminary definitions

Recall that the degree matrix D of a graph G is a diagonal matrix where each diagonal coefficient D_{ii} corresponds to the number of neighbours of vertex i and that $\lambda_2(M)$ is the second smallest eigenvalue of a symmetric matrix M .

Definition 9.1 Let \mathcal{G}_+ (resp. \mathcal{G}_-) be the subgraph of G that includes the edges that link two nodes in the same community (resp. in different communities) and A the adjacency matrix of G . We denote by $D_{\mathcal{G}}^+$ (resp. $D_{\mathcal{G}}^-$) the degree matrix of \mathcal{G}_+ (resp. \mathcal{G}_-) and define the Stochastic Block Model Laplacian to be

$$L_{SBM} = D_{\mathcal{G}}^+ - D_{\mathcal{G}}^- - A$$

9.7 Convex Duality

A standard technique to show that a candidate solution is the optimal one for a convex problem is to use convex duality.

We will describe duality with a game theoretical intuition in mind. The idea will be to rewrite (93) without imposing constraints on X but rather have the constraints be implicitly enforced. Consider the following optimization problem.

$$\max_X \min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \text{Tr}(BX) + \text{Tr}(QX) + \text{Tr}(Z(I_{n \times n} - X)) \quad (94)$$

Let us give it a game theoretical interpretation. Suppose that is a primal player (picking X) whose objective is to maximize the objective and a dual player, picking Z and Q after seeing X , trying to make the objective as small as possible. If the primal player does not pick X satisfying the constraints of (93) then we claim that the dual player is capable of driving the objective to $-\infty$. Indeed, if there is an i for which $X_{ii} \neq 1$ then the dual player can simply pick $Z_{ii} = -c \frac{1}{1-X_{ii}}$ and make the objective as small as desired by taking large enough c . Similarly, if X is not positive semidefinite, then the

dual player can take $Q = cvv^T$ where v is such that $v^T X v < 0$. If, on the other hand, X satisfy the constraints of (93) then

$$\text{Tr}(BX) \leq \min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \text{Tr}(BX) + \text{Tr}(QX) + \text{Tr}(Z(I_{n \times n} - X)),$$

since equality can be achieved if, for example, the dual player picks $Q = 0_{n \times n}$, then it is clear that the values of (93) and (94) are the same:

$$\max_{\substack{X, \\ X_{ii} \leq \forall_i \\ X \succeq 0}} \text{Tr}(BX) = \max_X \min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \text{Tr}(BX) + \text{Tr}(QX) + \text{Tr}(Z(I_{n \times n} - X))$$

With this game theoretical intuition in mind, it is clear that if we change the “rules of the game” and have the dual player decide their variables before the primal player (meaning that the primal player can pick X knowing the values of Z and Q) then it is clear that the objective can only increase, which means that:

$$\max_{\substack{X, \\ X_{ii} \leq \forall_i \\ X \succeq 0}} \text{Tr}(BX) \leq \min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \max_X \text{Tr}(BX) + \text{Tr}(QX) + \text{Tr}(Z(I_{n \times n} - X)).$$

Note that we can rewrite

$$\text{Tr}(BX) + \text{Tr}(QX) + \text{Tr}(Z(I_{n \times n} - X)) = \text{Tr}((B + Q - Z)X) + \text{Tr}(Z).$$

When playing:

$$\min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \max_X \text{Tr}((B + Q - Z)X) + \text{Tr}(Z),$$

if the dual player does not set $B + Q - Z = 0_{n \times n}$ then the primal player can drive the objective value to $+\infty$, this means that the dual player is forced to choose $Q = Z - B$ and so we can write

$$\min_{\substack{Z, Q \\ Z \text{ is diagonal} \\ Q \succeq 0}} \max_X \text{Tr}((B + Q - Z)X) + \text{Tr}(Z) = \min_{\substack{Z \\ Z \text{ is diagonal} \\ Z - B \succeq 0}} \max_X \text{Tr}(Z),$$

which clearly does not depend on the choices of the primal player. This means that

$$\max_{\substack{X, \\ X_{ii} \leq \forall_i \\ X \succeq 0}} \text{Tr}(BX) \leq \min_{\substack{Z \\ Z \text{ is diagonal} \\ Z - B \succeq 0}} \text{Tr}(Z).$$

This is known as weak duality (strong duality says that, under some conditions the two optimal values actually match, see, for example, [VB96], recall that we used strong duality when giving a sum-of-squares interpretation to the Max-Cut approximation ratio, a similar interpretation can be given in this problem, see [Ban15b]).

Also, the problem

$$\begin{aligned}
\min \quad & \text{Tr}(Z) \\
\text{s.t.} \quad & Z \text{ is diagonal} \\
& Z - B \succeq 0
\end{aligned} \tag{95}$$

is called the dual problem of (93).

The derivation above explains why the objective value of the dual is always larger or equal to the primal. Nevertheless, there is a much simpler proof (although not as enlightening): let X, Z be respectively a feasible point of (93) and (95). Since Z is diagonal and $X_{ii} = 1$ then $\text{Tr}(ZX) = \text{Tr}(Z)$. Also, $Z - B \succeq 0$ and $X \succeq 0$, therefore $\text{Tr}[(Z - B)X] \geq 0$. Altogether,

$$\text{Tr}(Z) - \text{Tr}(BX) = \text{Tr}[(Z - B)X] \geq 0,$$

as stated.

Recall that we want to show that gg^T is the optimal solution of (93). Then, if we find Z diagonal, such that $Z - B \succeq 0$ and

$$\text{Tr}[(Z - B)gg^T] = 0, \quad (\text{this condition is known as complementary slackness})$$

then $X = gg^T$ must be an optimal solution of (93). To ensure that gg^T is the unique solution we just have to ensure that the nullspace of $Z - B$ only has dimension 1 (which corresponds to multiples of g). Essentially, if this is the case, then for any other possible solution X one could not satisfy complementary slackness.

This means that if we can find Z with the following properties:

1. Z is diagonal
2. $\text{Tr}[(Z - B)gg^T] = 0$
3. $Z - B \succeq 0$
4. $\lambda_2(Z - B) > 0$,

then gg^T is the unique optima of (93) and so recovery of the true partition is possible (with an efficient algorithm).

Z is known as the dual certificate, or dual witness.

9.8 Building the dual certificate

The idea to build Z is to construct it to satisfy properties (1) and (2) and try to show that it satisfies (3) and (4) using concentration.

If indeed $Z - B \succeq 0$ then (2) becomes equivalent to $(Z - B)g = 0$. This means that we need to construct Z such that $Z_{ii} = \frac{1}{g_i} B[i, :]g$. Since $B = 2A - (\mathbf{1}\mathbf{1}^T - I)$ we have

$$Z_{ii} = \frac{1}{g_i} (2A - (\mathbf{1}\mathbf{1}^T - I))[i, :]g = 2\frac{1}{g_i} (Ag)_i + 1,$$

meaning that

$$Z = 2(D_{\mathcal{G}}^+ - D_{\mathcal{G}}^-) + I$$

is our guess for the dual witness. As a result

$$Z - B = 2(D_{\mathcal{G}}^+ - D_{\mathcal{G}}^-) - I - [2A - (\mathbf{1}\mathbf{1}^T - I)] = 2L_{SBM} + \mathbf{1}\mathbf{1}^T$$

It trivially follows (by construction) that

$$(Z - B)g = 0.$$

Therefore

Lemma 9.2 *If*

$$\lambda_2(2L_{SBM} + \mathbf{1}\mathbf{1}^T) > 0, \tag{96}$$

then the relaxation recovers the true partition.

Note that $2L_{SBM} + \mathbf{1}\mathbf{1}^T$ is a random matrix and so this boils down to “an exercise” in random matrix theory.

9.9 Matrix Concentration

Clearly,

$$\mathbb{E}[2L_{SBM} + \mathbf{1}\mathbf{1}^T] = 2\mathbb{E}L_{SBM} + \mathbf{1}\mathbf{1}^T = 2\mathbb{E}D_{\mathcal{G}}^+ - 2\mathbb{E}D_{\mathcal{G}}^- - 2\mathbb{E}A + \mathbf{1}\mathbf{1}^T,$$

and $\mathbb{E}D_{\mathcal{G}}^+ = \frac{n}{2} \frac{\alpha \log(n)}{n} I$, $\mathbb{E}D_{\mathcal{G}}^- = \frac{n}{2} \frac{\beta \log(n)}{n} I$, and $\mathbb{E}A$ is a matrix such with $4 \frac{n}{2} \times \frac{n}{2}$ blocks where the diagonal blocks have $\frac{\alpha \log(n)}{n}$ and the off-diagonal blocks have $\frac{\beta \log(n)}{n}$. We can write this as $\mathbb{E}A = \frac{1}{2} \left(\frac{\alpha \log(n)}{n} + \frac{\beta \log(n)}{n} \right) \mathbf{1}\mathbf{1}^T + \frac{1}{2} \left(\frac{\alpha \log(n)}{n} - \frac{\beta \log(n)}{n} \right) gg^T$

This means that

$$\mathbb{E}[2L_{SBM} + \mathbf{1}\mathbf{1}^T] = ((\alpha - \beta) \log n) I + \left(1 - (\alpha + \beta) \frac{\log n}{n} \right) \mathbf{1}\mathbf{1}^T - (\alpha - \beta) \frac{\log n}{n} gg^T.$$

Since $2L_{SBM}g = 0$ we can ignore what happens in the span of g and it is not hard to see that

$$\lambda_2 \left[((\alpha - \beta) \log n) I + \left(1 - (\alpha + \beta) \frac{\log n}{n} \right) \mathbf{1}\mathbf{1}^T - (\alpha - \beta) \frac{\log n}{n} gg^T \right] = (\alpha - \beta) \log n.$$

This means that it is enough to show that

$$\|L_{SBM} - \mathbb{E}[L_{SBM}]\| < \frac{\alpha - \beta}{2} \log n, \tag{97}$$

which is a large deviations inequality. ($\|\cdot\|$ denotes operator norm)

We will skip the details here (and refer the reader to [Ban15c] for the details), but the main idea is to use an inequality similar to the ones presented in the lecture about concentration of measure (and, in particular, matrix concentration). The main idea is to separate the diagonal from the non-diagonal part of $L_{SBM} - \mathbb{E}[L_{SBM}]$. The diagonal part depends on in and out-degrees of each node and can be handled with scalar concentration inequalities for trinomial distributions (as it was in [ABH14] to obtain the information theoretical bounds). The non-diagonal part has independent entries and so its spectral norm can be controlled by the following inequality:

Lemma 9.3 (Remark 3.13 in [BvH15]) *Let X be the $n \times n$ symmetric matrix with independent centered entries. Then there exists a universal constant c' , such that for every $t \geq 0$*

$$\text{Prob}[\|X\| > 3\sigma + t] \leq ne^{-t^2/c'\sigma_\infty^2}, \quad (98)$$

where we have defined

$$\sigma := \max_i \sqrt{\sum_j \mathbb{E}[X_{ij}^2]}, \quad \sigma_\infty := \max_{ij} \|X_{ij}\|_\infty.$$

Using these techniques one can show (this result was independently shown in [Ban15c] and [HWX14], with a slightly different approach)

Theorem 9.4 *Let G be a random graph with n nodes drawn accordingly to the stochastic block model on two communities with edge probabilities p and q . Let $p = \frac{\alpha \log n}{n}$ and $q = \frac{\beta \log n}{n}$, where $\alpha > \beta$ are constants. Then, as long as*

$$\sqrt{\alpha} - \sqrt{\beta} > \sqrt{2}, \quad (99)$$

the semidefinite program considered above coincides with the true partition with high probability.

Note that, if

$$\sqrt{\alpha} - \sqrt{\beta} < \sqrt{2},$$

then exact recovery of the communities is impossible, meaning that the SDP algorithm is optimal. Furthermore, in this regime one can show that there will be a node on each community that is more connected to the other community than to its own, meaning that a partition that swaps them would have more likelihood. In fact, the fact that the SDP will start working essentially when this starts happening appears naturally in the analysis; the diagonal part corresponds exactly to differences between in and out-degrees and Lemma 9.3 allows one to show that the contributions of the off-diagonal part are of lower order.

Remark 9.5 *A simpler analysis (and seemingly more adaptable to other problems) can be carried out by using Matrix Bernstein's inequality [Tro12] (described in the lecture about Matrix Concentration). The idea is simply to write $L_{SBM} - \mathbb{E}[L_{SBM}]$ as a sum of independent matrices (where each matrix corresponds to a pair of nodes) and to apply Matrix Bernstein (see [ABH14]). Unfortunately, this only shows exact recovery of a suboptimal threshold (suboptimal essentially by a factor of 2).*

9.10 More communities

A natural question is to understand what is the exact recovery threshold for the Stochastic Block Model on $k \geq 2$ communities. Recall the definition: The stochastic block model can be similarly defined for any $k \geq 2$ communities: G is a graph on $n = km$ nodes divided on k groups of m nodes each. Similarly to the $k = 2$ case, for each pair (i, j) of nodes, (i, j) is an edge of G with probability p if i and j are in the same set, and with probability q if they are in different sets. Each edge is drawn independently and $p > q$. In the logarithmic degree regime, we'll define the parameters in a slightly different way: $p = \frac{\alpha' \log m}{m}$ and $q = \frac{\beta' \log m}{m}$. Note that, for $k = 2$, we roughly have $\alpha = 2\alpha'$ and $\beta = 2\beta'$, which means that the exact recovery threshold, for $k = 2$, reads as: for

$$\sqrt{\alpha'} - \sqrt{\beta'} > 1$$

recovery is possible (and with the SDP algorithm), and for $\sqrt{\alpha'} - \sqrt{\beta'} < 1$ exact recovery is impossible.

Clearly, for any $k > 2$, if $\sqrt{\alpha'} - \sqrt{\beta'} < 1$ then exact recovery will also be impossible (simply imagine that an oracle tells us all of the community memberships except for those of two of the clusters, then the problem reduces to the $k = 2$ case). The remarkable fact is that, for $k = o(\log m)$ this is enough, not only for exact recovery to be possible, but also for an SDP based algorithm (very similar to the one above) to achieve exact recovery (see [AS15, ABKK15, HWX15, PW15]). However, for $k \approx \log n$, the situation is not understood.

Open Problem 9.2 *What is the threshold for exact recovery on the balanced symmetric Stochastic Block Model in $k \approx \log n$ communities and at what threshold does the SDP succeed at exactly determining the communities? (see [ABKK15]).*

9.11 Euclidean Clustering

The stochastic block model, although having fascinating phenomena, is not always an accurate model for clustering. The independence assumption assumed on the connections between pairs of vertices may sometimes be too unrealistic. Also, the minimum bisection of multisection objective may not be the most relevant in some applications.

One particularly popular form of clustering is k-means clustering. Given n points x_1, \dots, x_n and pairwise distances $d(x_i, x_j)$, the k-means objective attempts to partition the points in k clusters A_1, \dots, A_k (not necessarily of the same size) as to minimize the following objective³⁵

$$\min \sum_{t=1}^k \frac{1}{|A_t|} \sum_{x_i, x_j \in A_t} d^2(x_i, x_j).$$

A similar objective is the one in k-medians clustering, where for each cluster a center is picked (the center has to be a point in the cluster) and the sum of the distances from all points in the cluster to the center point are to be minimized, in other words, the objective to be minimized is:

$$\min \sum_{t=1}^k \min_{c_t \in A_t} \sum_{x_i \in A_t} d(x_i, c_t).$$

In [ABC⁺15] both a Linear Programming (LP) relaxation for k -medians and a Semidefinite Programming (SDP) relaxation for k -means are analyzed for a points in a generative model on which there are k disjoint balls in \mathbb{R}^d and, for every ball, points are drawn according to a isotropic distribution on each of the balls. The goal is to establish exact recovery of these convex relaxations requiring the least distance between the balls. This model (in this context) was first proposed and analyzed for k -medians in [NW13], the conditions for k -medians were made optimal in [ABC⁺15] and conditions for k -means were also given. More recently, the conditions on k -means were improved (made optimal for large dimensions) in [IMPV15a, IMPV15b] which also coined the term ‘‘Stochastic Ball Model’’.

For P the set of points, in order to formulate the k -medians LP we use variables y_p indicating whether p is a center of its cluster or not and z_{pq} indicating whether q is assigned to p or not (see [ABC⁺15] for details), the LP then reads:

³⁵When the points are in Euclidean space there is an equivalent more common formulation in which each cluster is assign a mean and the objective function is the sum of the distances squared to the center.

$$\begin{aligned}
\min \quad & \sum_{p,q} d(p,q) z_{pq}, \\
s.t. \quad & \sum_{p \in P} z_{pq} = 1, \quad \forall q \in P \\
& z_{pq} \leq y_p \\
& \sum_{p \in P} y_p = k \\
& z_{pq}, y_p \in [0, 1], \quad \forall p, q \in P.
\end{aligned}$$

the solution corresponds to an actual k-means solution if it is integral.

The semidefinite program for k-means is written in terms of a PSD matrix $X \in \mathbb{R}^{n \times n}$ (where n is the total number of points), see [ABC⁺15] for details. The intended solution is

$$X = \frac{1}{n} \sum_{t=1}^k \mathbf{1}_{A_t} \mathbf{1}_{A_t}^T,$$

where $\mathbf{1}_{A_t}$ is the indicator vector of the cluster A_t . The SDP reads as follows:

$$\begin{aligned}
\min_X \quad & \sum_{i,j} d(i,j) X_{ij}, \\
s.t. \quad & \text{Tr}(X) = k, \\
& X \mathbf{1} = \mathbf{1} \\
& X \succeq 0 \\
& X \succeq 0.
\end{aligned}$$

Inspired by simulations in the context of [NW13] and [ABC⁺15], Rachel Ward observed that the k-medians LP tends to be integral even for point configurations where no planted partition existed, and proposed the conjecture that k-medians is tight for typical point configurations. This was recorded as Problem 6 in [Mix15]. We formulate it as an open problem here:

Open Problem 9.3 *Is the LP relaxation for k-medians tight for a natural (random) generative model of points even without a clustering planted structure (such as, say, gaussian independent points)?*

Ideally, one would like to show that these relaxations (both the k-means SDP and the k-medians LP) are integral in instances that have clustering structure and not necessarily arising from generative random models. It is unclear however how to define what is meant by “clustering structure”. A particularly interesting approach is through stability conditions (see, for example [AJP13]), the idea is that if a certain set of data points has a much larger $k - 1$ -means (or medians) objective than a k -means (or medians) one, and there is not much difference between the k and the $k + 1$ objectives, then this is a good suggestion that the data is well explained by k clusters.

Open Problem 9.4 *Given integrality conditions to either the k-medians LP or the k-means SDP based on stability like conditions, as described above.*

9.12 Probably Certifiably Correct algorithms

While the SDP described in this lecture for recovery in the Stochastic Block Model achieves exact recovery in the optimal regime, SDPs (while polynomial time) tend to be slow in practice. There are faster (quasi-linear) methods that are also able to achieve exact recovery at the same threshold.

However, the SDP has an added benefit of producing a posteriori certificates. Indeed, if the solution from the SDP is integral (rank 1) then one is (a posteriori) sure to have found the minimum bisection. This means that the SDP (above the threshold) will, with high probability, not only find the minimum bisection but will also produce a posteriori certificate of such. Such an algorithm is referred to as Probably Certifiably Correct (PCC) [Ban15b]. Fortunately, one can get (in this case) get the best of both worlds and get a fast PCC method for recovery in the Stochastic Block Model essentially by using a fast method to find the solution and then using the SDP to only certify, which can be done considerably faster (see [Ban15b]). More recently, a PCC algorithm was also analyzed for k-means clustering (based on the SDP described above) [IMPV15b].

9.13 Another conjectured instance of tightness

The following problem is posed, by Andrea Montanari, in [Mon14], a description also appears in [Ban15a]. We briefly describe it here as well:

Given a symmetric matrix $W \in \mathbb{R}^{n \times n}$ the positive principal component analysis problem can be written as

$$\begin{aligned} \max \quad & x^T W x \\ \text{s. t.} \quad & \|x\| = 1 \\ & x \geq 0 \\ & x \in \mathbb{R}^n. \end{aligned} \tag{100}$$

In the flavor of the semidefinite relaxations considered in this section, (100) can be rewritten (for $X \in \mathbb{R}^{n \times n}$) as

$$\begin{aligned} \max \quad & \text{Tr}(WX) \\ \text{s. t.} \quad & \text{Tr}(X) = 1 \\ & X \geq 0 \\ & X \preceq 0 \\ & \text{rank}(X) = 1, \end{aligned}$$

and further relaxed to the semidefinite program

$$\begin{aligned} \max \quad & \text{Tr}(WX) \\ \text{s. t.} \quad & \text{Tr}(X) = 1 \\ & X \geq 0 \\ & X \preceq 0. \end{aligned} \tag{101}$$

This relaxation appears to have a remarkable tendency to be tight. In fact, numerical simulations suggest that if W is taken to be a Wigner matrix (symmetric with i.i.d. standard Gaussian entries), then the solution to (101) is rank 1 with high probability, but there is no explanation of this phenomenon. If the Wigner matrix is normalized to have entries $\mathcal{N}(0, 1/n)$, it is known that the typical value of the rank constraint problem is $\sqrt{2}$ (see [MR14]).

This motivates the last open problem of this section.

Open Problem 9.5 *Let W be a gaussian Wigner matrix with entries $\mathcal{N}(0, 1/n)$. Consider the fol-*

lowing Semidefinite Program:

$$\begin{aligned} \max \quad & \text{Tr}(WX) \\ \text{s. t.} \quad & \text{Tr}(X) = 1 \\ & X \succeq 0 \\ & X \preceq 0. \end{aligned} \tag{102}$$

Prove or disprove the following conjectures.

1. The expected value of this program is $\sqrt{2} + o(1)$.
2. With high probability, the solution of this SDP is rank 1.

Remark 9.6 The dual of this SDP motivates a particularly interesting statement which is implied by the conjecture. By duality, the value of the SDP is the same as the value of

$$\min_{\Lambda \succeq 0} \lambda_{\max}(W + \Lambda),$$

which is thus conjectured to be $\sqrt{2} + o(1)$, although no bound better than 2 (obtained by simply taking $\Lambda = 0$) is known.

MIT OpenCourseWare
<http://ocw.mit.edu>

18.S096 Topics in Mathematics of Data Science
Fall 2015

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.