The condensation threshold in stochastic block models

Joe Neeman (with Jess Banks, Cris Moore, Praneeth Netrapalli) Austin, May 9, 2016

Stochastic block model $\mathcal{G}(n, k, a, b)$

1. *n* nodes, *k* colors, about n/k nodes of each color 2. connect *u* to *v* with probability $\begin{cases} \frac{a}{n} & \text{if the same color} \\ \frac{b}{n} & \text{if different colors} \end{cases}$



Given the (uncolored) graph, recover the colors (up to permutation) better than a random guess.

Given the (uncolored) graph, recover the colors (up to permutation) better than a random guess.

Definition

Let $\sigma_v \in \{1, \ldots, k\}$ be the color of v. For another coloring τ ,

$$O_{lap}(\sigma,\tau) = \frac{\#\{v \in V : \sigma_v = -\tau_v\}}{n} - \frac{1}{k},$$

Given the (uncolored) graph, recover the colors (up to permutation) better than a random guess.

Definition

Let $\sigma_v \in \{1, \ldots, k\}$ be the color of v. For another coloring τ ,

$$O_{\text{lap}}(\sigma,\tau) = \max_{\pi} \frac{\#\{v \in V : \sigma_v = \pi(\tau_v)\}}{n} - \frac{1}{k},$$

where max is over all permutations π on $\{1, \ldots, k\}$.

Given the (uncolored) graph, recover the colors (up to permutation) better than a random guess.

Definition

Let $\sigma_v \in \{1, \ldots, k\}$ be the color of v. For another coloring τ ,

$$O_{lap}(\sigma,\tau) = \max_{\pi} \frac{\#\{v \in V : \sigma_v = \pi(\tau_v)\}}{n} - \frac{1}{k},$$

where max is over all permutations π on $\{1, \ldots, k\}$.

Definition

 $(G_n, \sigma_n) \sim \mathcal{G}(n, k, a, b)$ is detectable if there exists $\epsilon > 0$ and maps $A_n : \{\text{graphs}\} \rightarrow \{\text{labellings}\}$ such that

$$\liminf_{n\to\infty} \Pr(\mathsf{O}_{\mathsf{lap}}(\sigma_n, A_n(G_n)) > \epsilon) > \epsilon.$$

Otherwise it is undetectable.

Given the (uncolored) graph, did it come from $\mathcal{G}(n, k, a, b)$ or $\mathcal{G}(n, \frac{d}{n})$, where $d = \frac{a + (k-1)b}{k}$?

Given the (uncolored) graph, did it come from $\mathcal{G}(n, k, a, b)$ or $\mathcal{G}(n, \frac{d}{n})$, where $d = \frac{a + (k-1)b}{k}$?

Definition

Sequences \mathbb{P}_n and \mathbb{Q}_n of probability measures are

- contiguous if $\mathbb{P}_n(A_n) \to 0$ iff $\mathbb{Q}_n(A_n) \to 0$
- orthogonal if $\exists A_n$ with $\mathbb{P}_n(A_n) \to 0$ and $\mathbb{Q}_n(A_n) \to 1$.

Given the (uncolored) graph, did it come from $\mathcal{G}(n, k, a, b)$ or $\mathcal{G}(n, \frac{d}{n})$, where $d = \frac{a + (k-1)b}{k}$?

Definition

Sequences \mathbb{P}_n and \mathbb{Q}_n of probability measures are

- contiguous if $\mathbb{P}_n(A_n) \to 0$ iff $\mathbb{Q}_n(A_n) \to 0$
- orthogonal if $\exists A_n$ with $\mathbb{P}_n(A_n) \to 0$ and $\mathbb{Q}_n(A_n) \to 1$.

Say that $\mathcal{G}(n, k, a, b)$ is

- *distinguishable* if it is orthogonal to $\mathcal{G}(n, \frac{d}{n})$
- *indistinguishable* if it is contiguous with $\mathcal{G}(n, \frac{d}{n})$

Better parametrization

- $\frac{a}{n}$ = within-block edge probability
- $\frac{b}{n}$ = between-block edge probability
- $\cdot k =$ number of blocks

$$d = \frac{a + (k - 1)b}{k}$$
$$\lambda = \frac{a - b}{a + (k - 1)b}$$

Note $\lambda \in \left[-\frac{1}{k-1}, 1\right]$.

Phase diagram for k = 2



(Mossel/N/Sly, Massoulié)

Conjectured phase diagram for k = 20



(Decelle, Krzakala, Moore, Zdeborova)

What we know for k = 20



Theorem (Banks/Moore/N/Netrapalli)

$$d^{+} = \frac{2k \log k}{(1 + (k - 1)\lambda) \log(1 + (k - 1)\lambda) + (k - 1)(1 - \lambda) \log(1 - \lambda)}$$
$$d^{-} = \frac{2 \log(k - 1)}{\lambda^{2}(k - 1)}$$

- *d* > *d*⁺ implies detectability, distinguishability.
- $\cdot d < d^{-}$ implies undetectability, indistinguishability.

Theorem (Banks/Moore/N/Netrapalli)

$$d^{+} = \frac{2k \log k}{(1 + (k - 1)\lambda) \log(1 + (k - 1)\lambda) + (k - 1)(1 - \lambda) \log(1 - \lambda)}$$
$$d^{-} = \frac{2 \log(k - 1)}{\lambda^{2}(k - 1)}$$

- *d* > *d*⁺ *implies detectability, distinguishability.*
- $\cdot d < d^{-}$ implies undetectability, indistinguishability.

If k is large enough then there are λ such that $d^+ < \frac{1}{\lambda^2}$, giving the yellow region.

Theorem (Banks/Moore/N/Netrapalli)

$$d^{+} = \frac{2k \log k}{(1 + (k - 1)\lambda) \log(1 + (k - 1)\lambda) + (k - 1)(1 - \lambda) \log(1 - \lambda)}$$
$$d^{-} = \frac{2 \log(k - 1)}{\lambda^{2}(k - 1)}$$

- *d* > *d*⁺ implies detectability, distinguishability.
- *d* < *d*[−] implies undetectability, indistinguishability.

If k is large enough then there are λ such that $d^+ < \frac{1}{\lambda^2}$, giving the yellow region.

$$\lim_{k \to \infty} \frac{d^+}{d^-} = \frac{\mu^2}{(1+\mu)\log(1+\mu) - \mu} \text{ where } \mu = \frac{a-b}{d}.$$

If $\mu \approx \pm 1$ and $\lim_{k \to \infty} \frac{d^+}{d^-} \approx 1$ (planted coloring / giant)

The proofs



Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

For suitable *a*, *b*, *k*, w.h.p.

• *G*(*n*, *k*, *a*, *b*):

all good partitions are correlated with the truth.

Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

For suitable *a*, *b*, *k*, w.h.p.

• *G*(*n*, *k*, *a*, *b*):

all good partitions are correlated with the truth.

• $\mathcal{G}(n, \frac{d}{n})$:

there are no good partitions.

Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

For suitable *a*, *b*, *k*, w.h.p.

• *G*(*n*, *k*, *a*, *b*):

all good partitions are correlated with the truth.

• $\mathcal{G}(n, \frac{d}{n})$:

there are no good partitions.

Proof: concentration + union bound.

Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

For suitable *a*, *b*, *k*, w.h.p.

• *G*(*n*, *k*, *a*, *b*):

all good partitions are correlated with the truth.

• $\mathcal{G}(n, \frac{d}{n})$:

there are *no* good partitions.

Proof: concentration + union bound.

Distinguishing: check if there is a good partition.

Detecting: find a good partition.

Consider partitions of *G* into *k* equal parts. A partition is *good* if its average in-degree is $\approx \frac{a}{k}$ and its average out-degree is $\approx \frac{(k-1)b}{k}$.

For suitable *a*, *b*, *k*, w.h.p.

• *G*(*n*, *k*, *a*, *b*):

all good partitions are correlated with the truth.

• $\mathcal{G}(n, \frac{d}{n})$:

there are no good partitions.

Proof: concentration + union bound.

Distinguishing: check if there is a good partition.

Detecting: find a good partition.

Abbe/Sandon improved this for small *d* by taking the giant component and pruning trees.



If $\mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \to C < \infty$ then $\mathbb{Q}_n(A_n) \to 0 \Rightarrow \mathbb{P}_n(A_n) \to 0$.

If
$$\mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \to C < \infty$$
 then $\mathbb{Q}_n(A_n) \to 0 \Rightarrow \mathbb{P}_n(A_n) \to 0$
Set $\mathbb{P}_n = \mathcal{G}(n, k, a, b)$ and $\mathbb{Q}_n = \mathcal{G}(n, \frac{d}{n})$. Then
$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} = \frac{k^{-n} \sum_{\sigma} \prod_E \frac{a \text{ or } b}{n} \prod_{E^c} \left(1 - \frac{a \text{ or } b}{n} \right)}{\prod_E \frac{d}{n} \prod_{E^c} \left(1 - \frac{d}{n} \right)}$$

If
$$\mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C < \infty$$
 then $\mathbb{Q}_n(A_n) \to 0 \Rightarrow \mathbb{P}_n(A_n) \to 0$.
Set $\mathbb{P}_n = \mathcal{G}(n, k, a, b)$ and $\mathbb{Q}_n = \mathcal{G}(n, \frac{d}{n})$. Then
 $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} = \frac{k^{-n} \sum_{\sigma} \prod_E \frac{a \text{ or } b}{n} \prod_{E^c} \left(1 - \frac{a \text{ or } b}{n}\right)}{\prod_E \frac{d}{n} \prod_{E^c} \left(1 - \frac{d}{n}\right)}$
 $\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 = k^{-2n} \sum_{\sigma, \tau} \prod_E \frac{(a \text{ or } b)(a \text{ or } b)}{d^2} \prod_{E^c} \frac{(1 - \frac{a \text{ or } b}{n})(1 - \frac{a \text{ or } b}{n})}{(1 - \frac{d}{n})^2}$

If
$$\mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C < \infty$$
 then $\mathbb{Q}_n(A_n) \to 0 \Rightarrow \mathbb{P}_n(A_n) \to 0$.
Set $\mathbb{P}_n = \mathcal{G}(n, k, a, b)$ and $\mathbb{Q}_n = \mathcal{G}(n, \frac{d}{n})$. Then
 $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} = \frac{k^{-n} \sum_{\sigma} \prod_E \frac{a \text{ or } b}{n} \prod_{E^c} \left(1 - \frac{a \text{ or } b}{n}\right)}{\prod_E \frac{d}{n} \prod_{E^c} \left(1 - \frac{d}{n}\right)}$
 $\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 = k^{-2n} \sum_{\sigma, \tau} \prod_E \frac{(a \text{ or } b)(a \text{ or } b)}{d^2} \prod_{E^c} \frac{(1 - \frac{a \text{ or } b}{n})(1 - \frac{a \text{ or } b}{n})}{(1 - \frac{d}{n})^2}$

Under \mathbb{Q}_n , the events $(u, v) \in E$ are all independent, so can compute:

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 = C(1+o(1))\mathbb{E}\exp(X^T B X),$$

where X is a multinomial vector of length k^2 .

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T BZ)$$

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where $\psi(x) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$. Finite whenever $\lambda^2 d < 1$.

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where
$$\psi(\mathbf{x}) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$$
. Finite whenever $\lambda^2 d < 1$.

 $multinomials \leftrightarrow Gaussians$

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where
$$\psi(x) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$$
. Finite whenever $\lambda^2 d < 1$.

multinomials \leftrightarrow Gaussians $\Leftrightarrow \exp(X^T B X)$ uniformly integrable

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where $\psi(x) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$. Finite whenever $\lambda^2 d < 1$.

multinomials \leftrightarrow Gaussians $\Leftrightarrow \exp(X^T B X)$ uniformly integrable $\Leftrightarrow x^T B x - nH(x)$ maximized at $x = \mathbb{E} X$,

where H(x) is some kind of multivariate entropy.

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where $\psi(x) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$. Finite whenever $\lambda^2 d < 1$.

multinomials \leftrightarrow Gaussians $\Leftrightarrow \exp(X^T B X)$ uniformly integrable $\Leftrightarrow x^T B x - nH(x)$ maximized at $x = \mathbb{E} X$,

where H(x) is some kind of multivariate entropy.

Achlioptas-Naor: sufficient condition for the maximum to be at $x = \mathbb{E}X$.

Replacing multinomials with Gaussians,

$$\mathbb{E}_{\mathbb{Q}_n}\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}\right)^2 \to C\mathbb{E}\exp(Z^T B Z) = \psi(\lambda^2 d)^{k-1}$$

where $\psi(x) = \frac{e^{-x/2-x^2/4}}{\sqrt{1-x}}$. Finite whenever $\lambda^2 d < 1$.

multinomials \leftrightarrow Gaussians $\Leftrightarrow \exp(X^T B X)$ uniformly integrable $\Leftrightarrow x^T B x - nH(x)$ maximized at $x = \mathbb{E} X$,

where H(x) is some kind of multivariate entropy.

Achlioptas-Naor: sufficient condition for the maximum to be at $x = \mathbb{E}X$. (They were studying planted colorings.)

For the other direction $(\mathbb{P}_n(A_n) \to 0 \Rightarrow \mathbb{Q}_n(A_n) \to 0)$, want to show $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}$ bounded away from zero.

For the other direction $(\mathbb{P}_n(A_n) \to 0 \Rightarrow \mathbb{Q}_n(A_n) \to 0)$, want to show $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}$ bounded away from zero.

Small subgraph conditioning (Robinson/Wormald): $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}$ is essentially a function of the number of short cycles; it converges to an explicit limiting random variable that is never zero. For the other direction $(\mathbb{P}_n(A_n) \to 0 \Rightarrow \mathbb{Q}_n(A_n) \to 0)$, want to show $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}$ bounded away from zero.

Small subgraph conditioning (Robinson/Wormald): $\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}$ is essentially a function of the number of short cycles; it converges to an explicit limiting random variable that is never zero.

Main thing to check: convergence of second moment.

Suffices to show that the distribution of *G* is not much affected by conditioning on σ_u , σ_v .

$$d_{TV}(\mathbb{P}_{1,n},\mathbb{P}_{2,n})\to 0$$

$$d_{TV}(\mathbb{P}_{1,n},\mathbb{P}_{2,n}) \to 0$$
$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left| \frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right| \to 0$$

$$d_{TV}(\mathbb{P}_{1,n},\mathbb{P}_{2,n}) \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left| \frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right| \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right)^2 \to 0$$

$$d_{TV}(\mathbb{P}_{1,n}, \mathbb{P}_{2,n}) \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left| \frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right| \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right)^2 \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \frac{d\mathbb{P}_{i,n}}{d\mathbb{Q}_n} \frac{d\mathbb{P}_{j,n}}{d\mathbb{Q}_n} \to \psi(\lambda^2 d)^{k-1}$$

$$d_{TV}(\mathbb{P}_{1,n}, \mathbb{P}_{2,n}) \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left| \frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right| \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \left(\frac{d\mathbb{P}_{1,n}}{d\mathbb{Q}_n} - \frac{d\mathbb{P}_{2,n}}{d\mathbb{Q}_n} \right)^2 \to 0$$

$$\Leftrightarrow \mathbb{E}_{\mathbb{Q}_n} \frac{d\mathbb{P}_{i,n}}{d\mathbb{Q}_n} \frac{d\mathbb{P}_{j,n}}{d\mathbb{Q}_n} \to \psi(\lambda^2 d)^{k-1}.$$

Similar to previous second moment computation.

Indistinguishability and undetectability follow from an explicit second moment calculation. Use Achlioptas-Naor to estimate the set of parameters where the second moment is finite.





Thank you!