

Community Detection in Networks: SDP relaxation and Computational Gaps

Yihong Wu

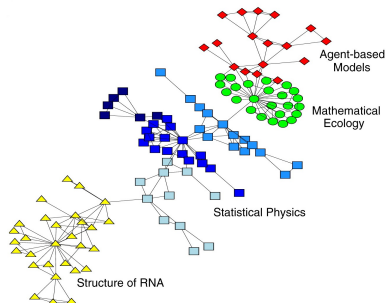
Department of ECE
University of Illinois at Urbana-Champaign
yihongwu@illinois.edu

Joint work with Bruce Hajek (Illinois) and Jiaming Xu (Wharton)

May 20, 2015

Community detection in networks

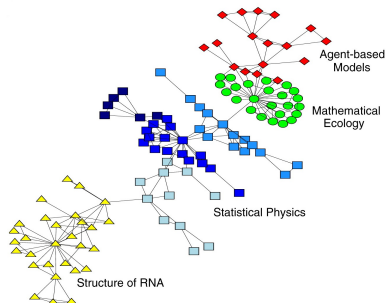
- Networks with community structures arise in many applications



Santa Fe Institute Collaboration network [Girvan-Newman '02]

Community detection in networks

- Networks with community structures arise in many applications

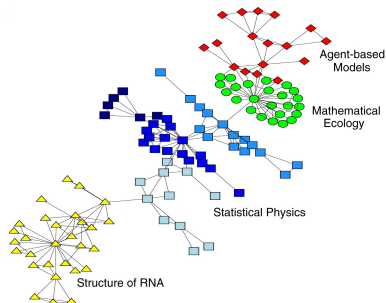


Santa Fe Institute Collaboration network [Girvan-Newman '02]

- Task: Discover underlying communities based on the network topology

Community detection in networks

- Networks with community structures arise in many applications



Santa Fe Institute Collaboration network [Girvan-Newman '02]

- Task: Discover underlying communities based on the network topology
- Applications: Friend or movie recommendation in online social networks

- The observed network is **sparse**
- Large solution space

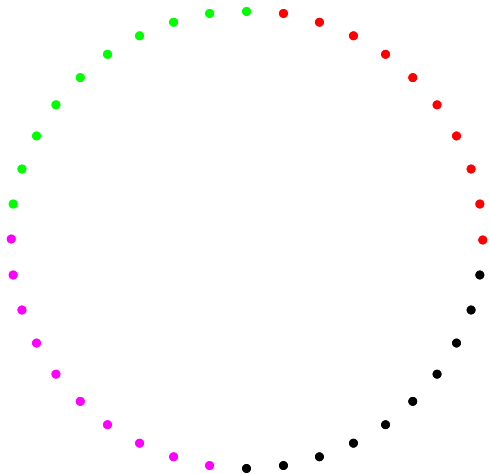
- The observed network is **sparse**
- Large solution space

Question

- Is there a computationally **efficient** and statistically **optimal** community detection algorithm?

Stochastic block model [Holland et al. '83]

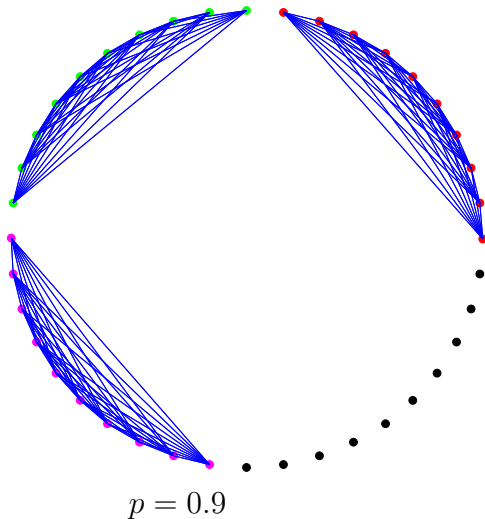
Planted partition model [Condon-Karp 01]



$$n = 40, K = 10, r = 3$$

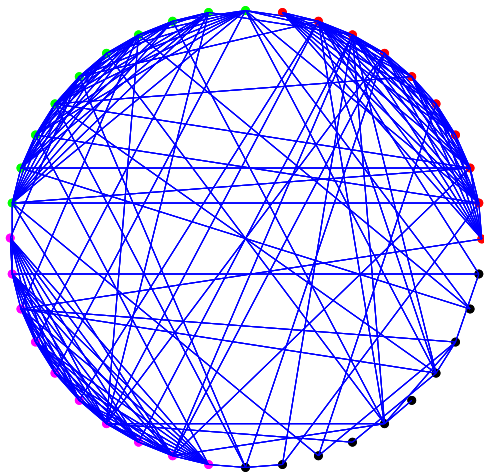
Stochastic block model [Holland et al. '83]

Planted partition model [Condon-Karp 01]



Stochastic block model [Holland et al. '83]

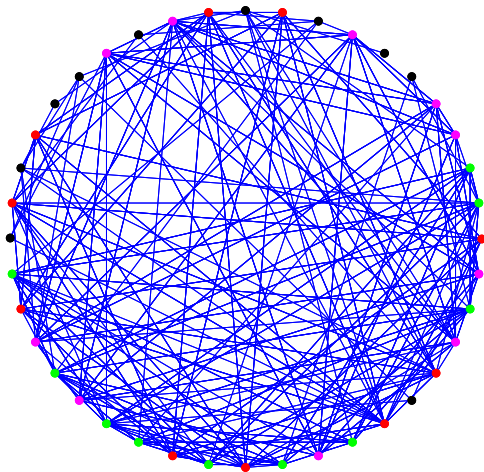
Planted partition model [Condon-Karp 01]



$$p = 0.9 \quad q = 0.1$$

Stochastic block model [Holland et al. '83]

Planted partition model [Condon-Karp 01]



$$p = 0.9 \quad q = 0.1$$

- True cluster: C^*
- Estimated cluster: \hat{C}
- Goal: **exact recovery** (strong consistency)

$$\mathbb{P}\{\hat{C} = C^*\} \xrightarrow{n \rightarrow \infty} 1$$

- True cluster: C^*
- Estimated cluster: \hat{C}
- Goal: **exact recovery** (strong consistency)

$$\mathbb{P}\{\hat{C} = C^*\} \xrightarrow{n \rightarrow \infty} 1$$

- Alternatives
 - ▶ almost exact recovery (weak consistency):
[Mossel-Neeman-Sly '13, Abbe-Sandon '15, Montanari '15]...
 - ▶ correlated recovery:
[Decelle-Krzakala-Moore-Zdeborova '11, Mossel-Neeman-Sly '12, Massoulié '13]...

Objectives of this talk

- **Statistical limit:** When is exact recovery possible (impossible)?
- **Computational limit:** When is exact recovery computationally easy (hard)?

- ① Linear community size: Sharp recovery via semidefinite programming
- ② Sublinear community size: Computational lower bounds

Two equal-sized communities

Model:

- n nodes partitioned into two communities of size $\frac{n}{2}$ ($\sigma_i = \pm 1$).
- $i \sim j$ independently w.p.
$$\begin{cases} p = \frac{a \log n}{n} & \sigma_i = \sigma_j \\ q = \frac{b \log n}{n} & \sigma_i \neq \sigma_j \end{cases}$$

Assuming $p > q$

- Maximum likelihood estimator (MLE)

$$\begin{aligned} \max_{\sigma} \quad & \langle A, \sigma \sigma^{\top} \rangle \\ \text{s.t.} \quad & \sigma_i \in \{\pm 1\}, \quad i \in [n] \\ & \sigma^{\top} \mathbf{1} = 0, \end{aligned}$$

Assuming $p > q$

- Maximum likelihood estimator (MLE)

$$\begin{aligned} \max_{\sigma} \langle A, \sigma \sigma^{\top} \rangle \\ \text{s.t. } \sigma_i \in \{\pm 1\}, \quad i \in [n] \\ \sigma^{\top} \mathbf{1} = 0, \end{aligned}$$

$\xleftrightarrow{\text{lift}}$

$$\begin{aligned} \max_Y \langle A, Y \rangle \\ \text{s.t. } \text{rank}(Y) = 1 \\ Y_{ii} = 1, \quad i \in [n] \\ \langle \mathbf{J}, Y \rangle = 0. \end{aligned}$$

where \mathbf{J} = all-one matrix

- Semidefinite programming (SDP) relaxation of MLE

$$\begin{aligned} \hat{Y}_{\text{SDP}} &= \arg \max_Y \langle A, Y \rangle \\ \text{s.t. } & Y \succeq 0 \\ & Y_{ii} = 1, \quad i \in [n] \\ & \langle \mathbf{J}, Y \rangle = 0. \end{aligned}$$

- similar SDP as in [Frieze-Jerrum '95] for MAX BISECTION
- average-case analysis on generative model (SBM)
- focus on **arg max** rather than approximating max
- goal: $\mathbb{P} \left\{ \hat{Y}_{\text{SDP}} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \right\} \rightarrow 1$

Theorem (Abbe-Bandeira-Hall '14, Mossel-Neeman-Sly '14)

- *If $(\sqrt{a} - \sqrt{b})^2 > 2$, recovery is achievable in polynomial-time.*
- *If $(\sqrt{a} - \sqrt{b})^2 < 2$, recovery is impossible.*

Optimal recovery via SDP

Theorem (Abbe-Bandeira-Hall '14, Mossel-Neeman-Sly '14)

- If $(\sqrt{a} - \sqrt{b})^2 > 2$, recovery is achievable in polynomial-time.
- If $(\sqrt{a} - \sqrt{b})^2 < 2$, recovery is impossible.

Theorem (Hajek-W.-Xu '14)

SDP achieves the optimal recovery threshold $(\sqrt{a} - \sqrt{b})^2 > 2$.

Theorem (Abbe-Bandeira-Hall '14, Mossel-Neeman-Sly '14)

- If $(\sqrt{a} - \sqrt{b})^2 > 2$, recovery is achievable in polynomial-time.
- If $(\sqrt{a} - \sqrt{b})^2 < 2$, recovery is impossible.

Theorem (Hajek-W.-Xu '14)

SDP achieves the optimal recovery threshold $(\sqrt{a} - \sqrt{b})^2 > 2$.

Remarks

- originally conjectured in [Abbe-Bandeira-Hall '14]
- independently proved by [Bandeira '15]
- $\mathbb{P} \left\{ \hat{Y}_{\text{SDP}} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \right\} = 1 - n^{-\Omega(1)}$

$$\begin{aligned} \max_Y & \langle A, Y \rangle \\ \text{s.t.} & Y \succeq 0 \\ & Y_{ii} = 1 \\ & \langle \mathbf{J}, Y \rangle = 0 \end{aligned}$$

$$\max_Y \langle A, Y \rangle$$

$$\text{s.t. } Y \succeq 0$$

$$Y_{ii} = 1$$

$$\langle \mathbf{J}, Y \rangle = 0$$

dual variables

$$S \succeq 0$$

$$D = \text{diag} \{d_i\}$$

$$\lambda \in \mathbb{R}$$

Lemma

$Y^* = \sigma^*(\sigma^*)^\top$ is unique solution if $\exists D, \lambda$ s.t. $S = \lambda \mathbf{J} + D - A$ satisfies

$$S\sigma = 0 \quad \text{and} \quad \lambda_2(S) > 0.$$

$$\max_Y \langle A, Y \rangle$$

dual variables

$$\text{s.t. } Y \succeq 0$$

$$S \succeq 0$$

$$Y_{ii} = 1$$

$$D = \text{diag} \{d_i\}$$

$$\langle \mathbf{J}, Y \rangle = 0$$

$$\lambda \in \mathbb{R}$$

Lemma

$Y^* = \sigma^*(\sigma^*)^\top$ is unique solution if $\exists D, \lambda$ s.t. $S = \lambda \mathbf{J} + D - A$ satisfies

$$S\sigma = 0 \quad \text{and} \quad \lambda_2(S) > 0.$$

$$\Rightarrow d_i = (\# \text{ of nbrs in own cluster}) - (\# \text{ of nbrs in other cluster})$$

$$= \begin{cases} e(i, C_1) - e(i, C_2) & i \in C_1 \\ e(i, C_2) - e(i, C_1) & i \in C_2 \end{cases}$$

- Mean adj matrix: $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$

- Mean adj matrix: $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$
-

$$\begin{aligned}
 S &= \lambda\mathbf{J} - A + D \\
 &= \underbrace{\left(\lambda - \frac{p+q}{2}\right)\mathbf{J}} - \frac{p-q}{2}\sigma^*(\sigma^*)^\top + p\mathbf{I} + \underbrace{D - (A - \mathbb{E}[A])}
 \end{aligned}$$

- **Mean adj matrix:** $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$

-

$$\begin{aligned}
 S &= \lambda\mathbf{J} - A + D \\
 &= \underbrace{\left(\lambda - \frac{p+q}{2}\right)\mathbf{J}} - \frac{p-q}{2}\sigma^*(\sigma^*)^\top + p\mathbf{I} + \underbrace{D - (A - \mathbb{E}[A])}
 \end{aligned}$$

- $\lambda_2(S) = \inf_{x \perp \sigma^*, \|x\|_2=1} x^\top Sx > 0$ if $\min d_i \geq \|A - \mathbb{E}[A]\|$ and $\lambda \geq (p+q)/2$

- **Mean adj matrix:** $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$

-

$$\begin{aligned}
 S &= \lambda\mathbf{J} - A + D \\
 &= \underbrace{\left(\lambda - \frac{p+q}{2}\right)\mathbf{J}} - \frac{p-q}{2}\sigma^*(\sigma^*)^\top + p\mathbf{I} + \underbrace{D - (A - \mathbb{E}[A])}
 \end{aligned}$$

- $\lambda_2(S) = \inf_{x \perp \sigma^*, \|x\|_2=1} x^\top Sx > 0$ if $\min d_i \geq \|A - \mathbb{E}[A]\|$ and $\lambda \geq (p+q)/2$
- To finish the proof:

- **Mean adj matrix:** $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$

-

$$\begin{aligned} S &= \lambda\mathbf{J} - A + D \\ &= \underbrace{\left(\lambda - \frac{p+q}{2}\right)\mathbf{J}} - \frac{p-q}{2}\sigma^*(\sigma^*)^\top + p\mathbf{I} + \underbrace{D - (A - \mathbb{E}[A])} \end{aligned}$$

- $\lambda_2(S) = \inf_{x \perp \sigma^*, \|x\|_2=1} x^\top Sx > 0$ if $\min d_i \geq \|A - \mathbb{E}[A]\|$ and $\lambda \geq (p+q)/2$
- To finish the proof:
 - ① $\min d_i = \Omega_P(\log n)$ if $\sqrt{a} - \sqrt{b} > \sqrt{2}$

- **Mean adj matrix:** $\mathbb{E}[A] = \frac{p+q}{2}\mathbf{J} + \frac{p-q}{2}\sigma^*(\sigma^*)^\top - p\mathbf{I}$

-

$$\begin{aligned}
 S &= \lambda\mathbf{J} - A + D \\
 &= \underbrace{\left(\lambda - \frac{p+q}{2}\right)\mathbf{J}} - \frac{p-q}{2}\sigma^*(\sigma^*)^\top + p\mathbf{I} + \underbrace{D - (A - \mathbb{E}[A])}
 \end{aligned}$$

- $\lambda_2(S) = \inf_{x \perp \sigma^*, \|x\|_2=1} x^\top Sx > 0$ if $\min d_i \geq \|A - \mathbb{E}[A]\|$ and $\lambda \geq (p+q)/2$
- To finish the proof:
 - ① $\min d_i = \Omega_P(\log n)$ if $\sqrt{a} - \sqrt{b} > \sqrt{2}$
 - ② $\|A - \mathbb{E}[A]\| = O_P(\sqrt{\log n})$

① Necessity

$$\sqrt{a} - \sqrt{b} < \sqrt{2}$$

$\Rightarrow \min d_i < 0$ w.h.p.

$\Rightarrow \exists i : \#$ of nbrs in own cluster $<$ $\#$ of nbrs in other cluster

\Rightarrow MLE fails

① Necessity

$$\sqrt{a} - \sqrt{b} < \sqrt{2}$$

$\Rightarrow \min d_i < 0$ w.h.p.

$\Rightarrow \exists i : \#$ of nbrs in own cluster $<$ $\#$ of nbrs in other cluster

\Rightarrow MLE fails

② Proof of $\|A - \mathbb{E}[A]\| = O_P(\sqrt{\log n})$

- ▶ 2nd-order stochastic dominance argument [Tomozei-Massoulié '14] + result for iid matrix [Seginer '00]
- ▶ [Feige-Ofek '05]: $\mathcal{G}(n, \frac{C \log n}{n})$ for sufficiently large C
- ▶ [Bandeira-van Handel '14]: comparison argument

Multiple equal-sized communities

- $\{0, 1\}$ -cluster matrix:

$$Y^* = \sum_{k=1}^r \xi_k (\xi_k)^\top = \begin{array}{|c|c|c|} \hline 1 & & \\ \hline & 1 & 0 \\ \hline & & 1 \\ \hline 0 & & & 1 \\ \hline \end{array}$$

where ξ_k = indicator of the k^{th} cluster of size $K = n/r$

- $\{0, 1\}$ -cluster matrix:

$$Y^* = \sum_{k=1}^r \xi_k (\xi_k)^\top = \begin{array}{|c|c|c|} \hline 1 & & \\ \hline & 1 & 0 \\ \hline & & 1 \\ \hline 0 & & & 1 \\ \hline \end{array}$$

where ξ_k = indicator of the k^{th} cluster of size $K = n/r$

- SDP relaxation of MLE:

$$\max_Y \langle A, Y \rangle$$

$$\text{s.t. } Y \succeq 0$$

$$Y_{ii} = 1$$

$$Y_{ij} \geq 0$$

$$\sum_j Y_{ij} = K$$

Theorem ([Hajek-W.-Xu '15])

SDP achieves optimal threshold $(\sqrt{a} - \sqrt{b})^2 > r$.

Theorem ([Hajek-W.-Xu '15])

SDP achieves optimal threshold $(\sqrt{a} - \sqrt{b})^2 > r$.

Proof of correctness:

$$\begin{aligned} \max_Y \quad & \langle A, Y \rangle \\ \text{s.t.} \quad & Y \succeq 0 \quad S \succeq 0 \\ & Y_{ii} = 1 \quad d_i \\ & Y_{ij} \geq 0 \quad B \geq 0 \\ & \sum_j Y_{ij} = K \quad \lambda_i \end{aligned}$$

Construction of the dual witness

- For node $i \in C_k$,

$$\lambda_i = \frac{1}{K} \left(\max_{\ell \neq k} e(i, C_\ell) - Kq/2 + \sqrt{\log n}/2 \right)$$

$$d_i = e(i, C_k) - \max_{\ell \neq k} e(i, C_\ell) - \frac{1}{K} \sum_{j \in C_k} \max_{\ell \neq k} e(j, C_\ell) + Kq - \sqrt{\log n}$$

0			
	0		
		0	
			0

- $B =$ [table], where each [square] is rank-2, specified by

$$B_{C_k \times C_{k'}}(i, j) = \frac{1}{K} \left(\max_{\ell \neq k} e(i, C_{\ell'}) - e(i, C_{k'}) + \max_{\ell \neq k'} e(j, C_\ell) - e(j, C_k) \right. \\ \left. + \frac{e(C_k, C_{k'})}{K} - Kq + \sqrt{\log n} \right)$$

- $S = D - A - B + \lambda \mathbf{1}^\top + \mathbf{1} \lambda^\top$

$$\begin{aligned} & \max_Y \langle A, Y \rangle \\ & \text{s.t. } Y \succeq 0 \quad S \succeq 0 \\ & \quad Y_{ii} = 1 \quad d_i \\ & \quad Y_{ij} \geq 0 \quad B \geq 0 \\ & \quad \sum_j Y_{ij} = K \quad \lambda_i \end{aligned}$$

- $S\xi_k = 0$ for $k = 1, \dots, r$.
- $\lambda_{r+1}(S) > 0$ if $\min d_i \geq \|A - \mathbb{E}[A]\| = O_P(\sqrt{\log n})$
- $d_i = (\# \text{ of nbrs in own cluster}) - \text{maximal } (\# \text{ of nbrs in other clusters}) + O_P(\sqrt{\log n})$.
- Sharp threshold
 - ▶ $\sqrt{a} - \sqrt{b} > \sqrt{r} \Rightarrow \min d_i = \Omega(\log n) \Rightarrow$ **SDP succeeds**
 - ▶ $\sqrt{a} - \sqrt{b} < \sqrt{r} \Rightarrow \min d_i = -\Omega(\log n) \Rightarrow$ **MLE fails**

Unequal-sized clusters

Two unequal-sized clusters: known size

Two clusters of size K and $n - K$ ($K = \rho n$):

p	q
q	p

$$\begin{aligned}\hat{Y}_{\text{SDP}} &= \arg \max_Y \langle A, Y \rangle \\ \text{s.t. } & Y \succeq 0 \\ & Y_{ii} = 1, \quad i \in [n] \\ & \langle \mathbf{J}, Y \rangle = (2K - n)^2\end{aligned}$$

achieves optimal threshold $\eta(\rho, a, b) > 1$.

Two unequal-sized clusters: known size

Two clusters of size K and $n - K$ ($K = \rho n$):

p	q
q	p

$$\begin{aligned}\hat{Y}_{\text{SDP}} &= \arg \max_Y \langle A, Y \rangle \\ \text{s.t. } & Y \succeq 0 \\ & Y_{ii} = 1, \quad i \in [n] \\ & \langle \mathbf{J}, Y \rangle = (2K - n)^2\end{aligned}$$

achieves optimal threshold $\eta(\rho, a, b) > 1$.

Note: $\rho \mapsto \eta(\rho, a, b)$ is minimized at $\eta(1/2, a, b) = \frac{1}{2}(\sqrt{a} - \sqrt{b})^2 \Rightarrow$
“suggests” equal-sized case is the hardest for two communities

Two unequal-sized clusters: unknown size

Two clusters of size K and $n - K$ ($K = 0, 1, \dots, n$):

p	q
q	p

$$\begin{aligned}\widehat{Y}_{\text{SDP}} &= \arg \max_Y \langle A, Y \rangle - \lambda \langle \mathbf{J}, Y \rangle \\ \text{s.t. } & Y \succeq 0 \\ & Y_{ii} = 1, \quad i \in [n]\end{aligned}$$

with $\lambda = \frac{a-b}{\log a - \log b} \frac{\log n}{n}$ achieves optimal threshold $(\sqrt{a} - \sqrt{b})^2 > 2$.

Two unequal-sized clusters: unknown size

Two clusters of size K and $n - K$ ($K = 0, 1, \dots, n$):

p	q
q	p

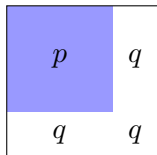
$$\begin{aligned}\hat{Y}_{\text{SDP}} &= \arg \max_Y \langle A, Y \rangle - \lambda \langle \mathbf{J}, Y \rangle \\ \text{s.t. } & Y \succeq 0 \\ & Y_{ii} = 1, \quad i \in [n]\end{aligned}$$

with $\lambda = \frac{a-b}{\log a - \log b} \frac{\log n}{n}$ achieves optimal threshold $(\sqrt{a} - \sqrt{b})^2 > 2$.

Note: If $K = \Omega(n)$, there exists a **data-driven** choice of λ .

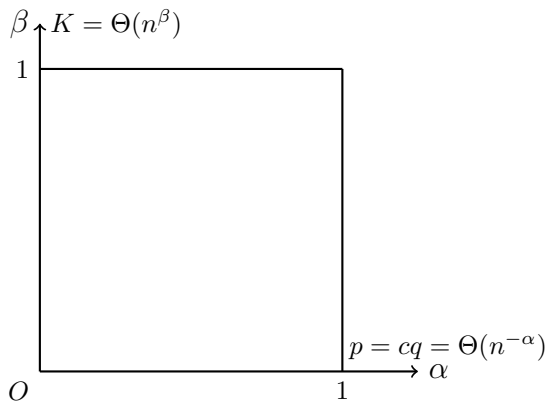
- **Binary censored block model:** $\mathcal{G}(n, \frac{a \log n}{n})$ observe edge label flipped w.p. ϵ
 - ▶ SDP achieves sharp threshold $a (\sqrt{1-\epsilon} - \sqrt{\epsilon})^2 > 1$
 - ▶ Closes the gap in [Abbe-Bandeira-Bracher-Singer '14]
- **General SBM:**
 - ▶ Optimality of SDP relaxation remains open (but within a factor of 4)
 - ▶ Sharp threshold is found in [Abbe-Sandon '15] via a two-stage procedure.

Detecting a single cluster

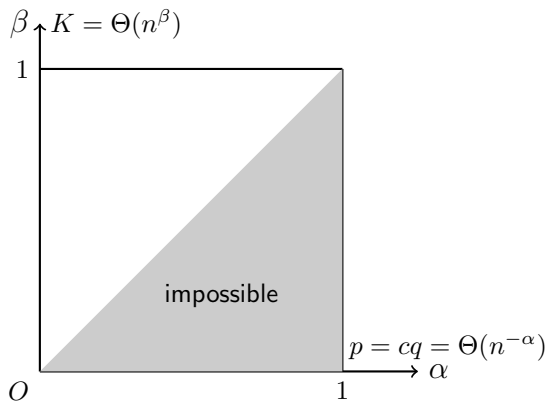


- One cluster of size K plus $n - K$ outliers
- Connectivity p within cluster and q otherwise
- Also known as **Planted Dense Subgraph** model
- Linear community size: $K = \rho n$ and SDP achieves sharp threshold
- Next focus on $K = \Theta(n^\beta)$.

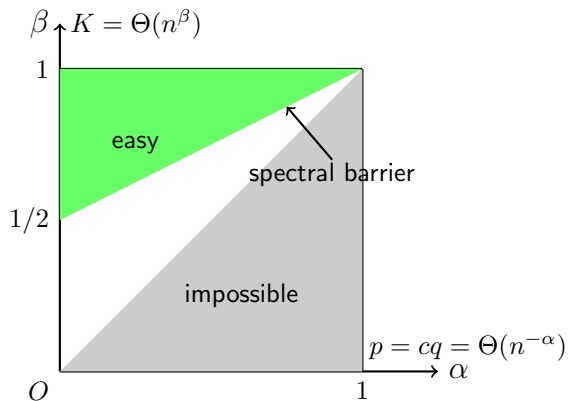
Conjecture on computational limit



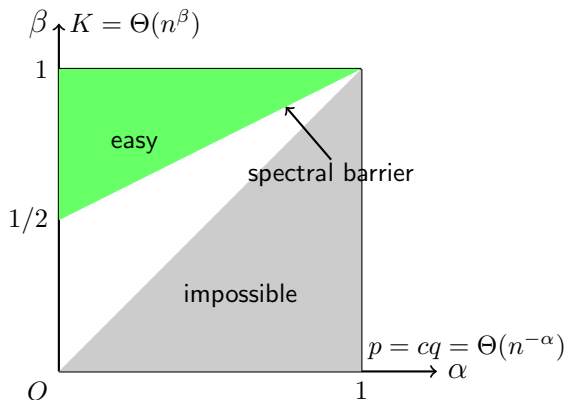
Conjecture on computational limit



Conjecture on computational limit



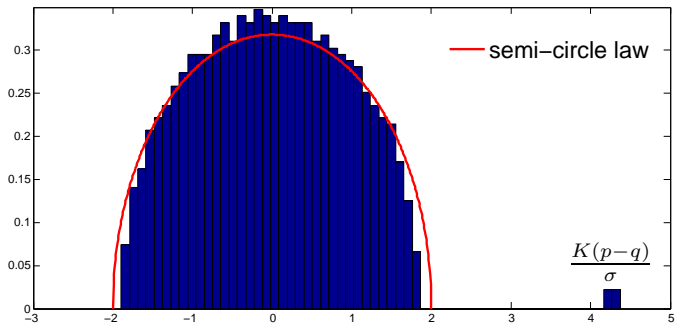
Conjecture on computational limit



Conjecture [Chen-Xu '14]: no polynomial-time algorithm succeeds beyond the spectral barrier [Nadakuditi-Newman '12]

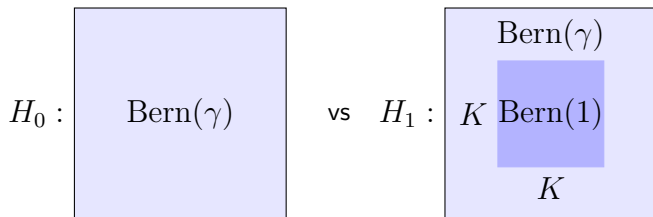
$$A = \begin{array}{c} K \\ K \end{array} \begin{array}{|c|} \hline p \\ \hline \end{array} \begin{array}{c} \\ q \end{array} + A - \mathbb{E}[A]$$

$$A = \begin{matrix} & & K \\ K & \begin{matrix} p \\ \square \end{matrix} & \\ & & q \end{matrix} + A - \mathbb{E}[A]$$



Eigenvalue distribution of $\frac{A - q\mathbf{1}\mathbf{1}^\top}{\sigma}$ for $\sigma = \sqrt{q(1-q)n}$

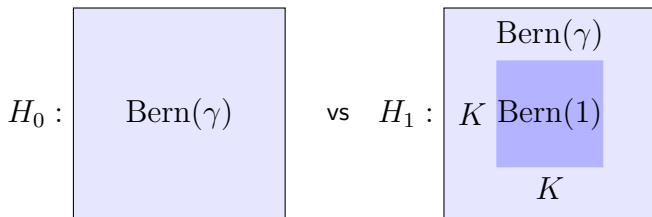
Planted clique hardness hypothesis



Intermediate regime: $\log n \ll K \ll \sqrt{n}$, $\gamma = \Theta(1)$

- detection is possible but believed to have high computational complexity: [Alon et al. '11] [Feldman et al. '13] [Deshpande-Montanari '15] [Meka-Potechin-Wigderson '15]

Planted clique hardness hypothesis

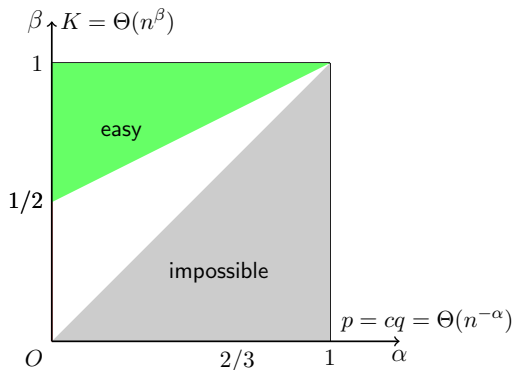


Intermediate regime: $\log n \ll K \ll \sqrt{n}$, $\gamma = \Theta(1)$

- detection is possible but believed to have high computational complexity: [Alon et al. '11] [Feldman et al. '13] [Deshpande-Montanari '15] [Meka-Potechin-Wigderson '15]
- various hardness results assuming Planted Clique hardness
 - ▶ detecting **sparse principal component** [Berthet-Rigollet '13]: $\gamma = \frac{1}{2}$
 - ▶ detecting **sparse submatrix** [Ma-W. '13, Cai-Liang-Rakhlin '15]: $\gamma = \frac{1}{2}$
 - ▶ cryptography [Applebaum-Barak-Wigderson '10]: $\gamma = 2^{-\log^{0.99} n}$

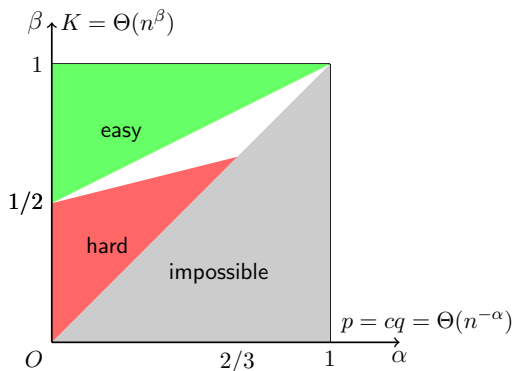
Hard regime for recovering a single cluster

Assuming Planted Clique hardness for **any constant** $\gamma > 0$



Hard regime for recovering a single cluster

Assuming Planted Clique hardness for **any constant** $\gamma > 0$



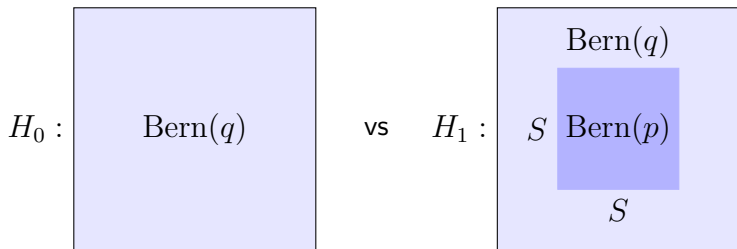
Recovering a single cluster in the red regime is at least as hard as detecting a clique of size $K = o(\sqrt{n})$

Proof step 1: Recovery is harder than *detection*

Recovery versus Detection [Arias-Castro-Verzelen '14] :

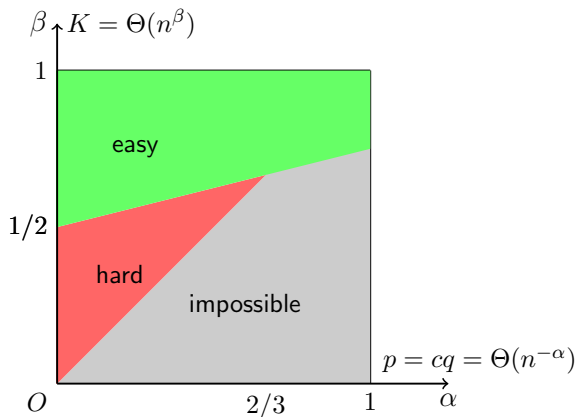
Proof step 1: Recovery is harder than *detection*

Recovery versus Detection [Arias-Castro-Verzelen '14] :



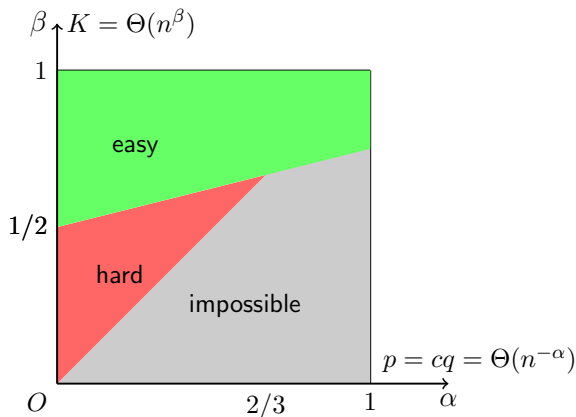
Each node is included in S with probability $\frac{K}{n}$

Proof step 2: Hardness for *detecting* a single cluster

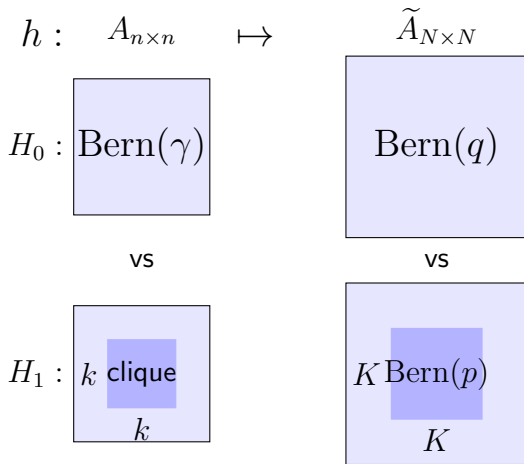


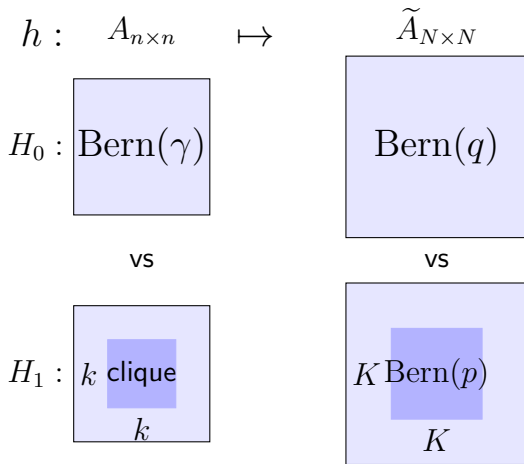
- Detecting a single cluster in the red regime is at least as hard as detecting a clique of size $K = o(\sqrt{n})$

Proof step 2: Hardness for *detecting* a single cluster



- Detecting a single cluster in the red regime is at least as hard as detecting a clique of size $K = o(\sqrt{n})$
- Reduced from **Planted Clique detection** in polynomial time

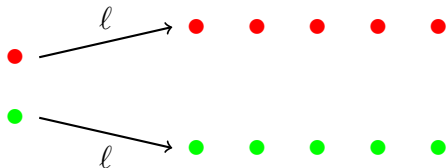




$h : A \mapsto \tilde{A}$ is **agnostic** to the clique and can be computed in P-time

Given an integer ℓ , two probability distributions P, Q on $\{0, 1, \dots, \ell^2\}$

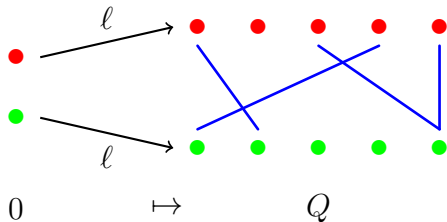
Split each node
into ℓ new nodes
 $N = n\ell, K = k\ell$



Given an integer ℓ , two probability distributions P, Q on $\{0, 1, \dots, \ell^2\}$

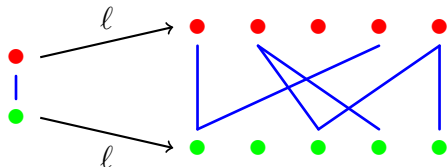
Split each node
into ℓ new nodes
 $N = n\ell, K = k\ell$

Assign edges with
distributions P, Q

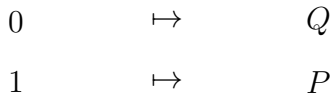


Given an integer ℓ , two probability distributions P, Q on $\{0, 1, \dots, \ell^2\}$

Split each node
into ℓ new nodes
 $N = n\ell, K = k\ell$

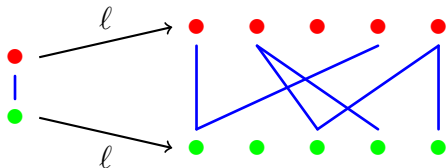


Assign edges with
distributions P, Q



Given an integer ℓ , two probability distributions P, Q on $\{0, 1, \dots, \ell^2\}$

Split each node
into ℓ new nodes
 $N = n\ell, K = k\ell$



Assign edges with
distributions P, Q

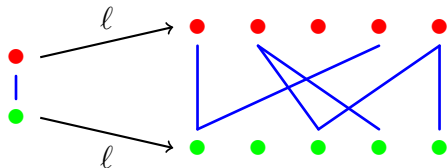
0	\mapsto	Q
1	\mapsto	P

H_0 : Bern(γ) $(1 - \gamma)Q + \gamma P$

H_1 : Bern(1) (in-clique) P (in-cluster)

Given an integer ℓ , two probability distributions P, Q on $\{0, 1, \dots, \ell^2\}$

Split each node
into ℓ new nodes
 $N = n\ell, K = k\ell$



Assign edges with
distributions P, Q

0 \mapsto Q
1 \mapsto P

H_0 : $\text{Bern}(\gamma)$ $(1 - \gamma)Q + \gamma P$

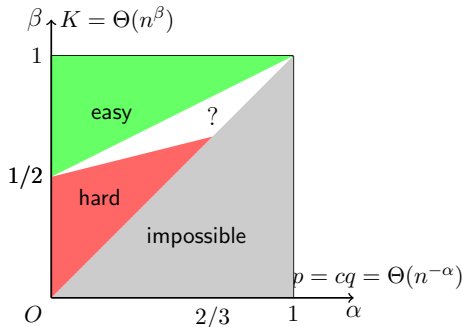
H_1 : $\text{Bern}(1)$ (in-clique) P (in-cluster)

How to choose P, Q ?

- Matching H_0 : $(1 - \gamma)Q + \gamma P = \text{Binom}(\ell^2, q)$
- Matching H_1 approximately: $P \approx \text{Binom}(\ell^2, p)$ in total variation
- Main effort: the law of the resulting graph is close to SBM in total variation

- Versatility of SDP as a simple, general purpose, computationally feasible methodology for community detection
- Construction of dual witness **lacks a general recipe**

Concluding remarks



References

- B. Hajek, Y. W. & J. Xu (2014). *Computational lower bounds for community detection on random graphs*. [arXiv:1406.6625](https://arxiv.org/abs/1406.6625) (COLT '15)
- B. Hajek, Y. W. & J. Xu (2014). *Achieving exact cluster recovery threshold via semidefinite programming*. [arXiv:1412.6156](https://arxiv.org/abs/1412.6156)
- B. Hajek, Y. W. & J. Xu (2015). *Achieving exact cluster recovery threshold via semidefinite programming: Extensions*. [arXiv:1502.07738](https://arxiv.org/abs/1502.07738)

Formal statement of hardness of detecting a cluster

γ : edge probability in Planted Clique

Theorem

Assume Planted Clique Hypothesis holds for all $0 < \gamma \leq 1/2$. Let $\alpha > 0$ and $0 < \beta < 1$ be such that

$$\alpha < \beta < \frac{1}{2} + \frac{\alpha}{4}.$$

Then there exists a sequence $\{(N_\ell, K_\ell, q_\ell)\}_{\ell \in \mathbb{N}}$ satisfying $\lim_{\ell \rightarrow \infty} \frac{-\log q_\ell}{\log N_\ell} = \alpha$ and $\lim_{\ell \rightarrow \infty} \frac{\log K_\ell}{\log N_\ell} = \beta$ such that for any sequence of randomized polynomial-time tests ϕ_ℓ for the PDS($N_\ell, K_\ell, 2q_\ell, q_\ell$) problem, the Type-I+II error probability is lower bounded by 1.

Proof ideas: Reduce **from** Planted Clique in polynomial-time
Map approximately:

- $\mathcal{G}(n, \gamma) \mapsto \mathcal{G}(N, q)$
- $\mathcal{G}(n, k, \gamma, 1) \mapsto \mathcal{G}(N, K, q, p)$

Bound the total variation distance

Lemma

Let $\ell, n \in \mathbb{N}$, $k \in [n]$ and $\gamma \in (0, \frac{1}{2}]$. Let $N = \ell n$, $K = k\ell$, $p = 2q$ and $m_0 = \lfloor \log_2(1/\gamma) \rfloor$. Assume that $16q\ell^2 \leq 1$ and $k \geq 6e\ell$. If $G \sim \mathcal{G}(n, \gamma)$, then $\tilde{G} \sim \mathcal{G}(N, q)$. If $G \sim \mathcal{G}(n, k, 1, \gamma)$, then

$$d_{\text{TV}}(P_{\tilde{G}}, \mathcal{G}(N, K, p, q)) \lesssim e^{-K} + ke^{-\ell} + k^2(q\ell^2)^{m_0+1} + \sqrt{e q \ell^2 - 1}$$

Bound the total variation distance

Lemma

Let $\ell, n \in \mathbb{N}$, $k \in [n]$ and $\gamma \in (0, \frac{1}{2}]$. Let $N = \ell n$, $K = k\ell$, $p = 2q$ and $m_0 = \lfloor \log_2(1/\gamma) \rfloor$. Assume that $16q\ell^2 \leq 1$ and $k \geq 6e\ell$. If $G \sim \mathcal{G}(n, \gamma)$, then $\tilde{G} \sim \mathcal{G}(N, q)$. If $G \sim \mathcal{G}(n, k, 1, \gamma)$, then

$$d_{\text{TV}}(P_{\tilde{G}}, \mathcal{G}(N, K, p, q)) \lesssim e^{-K} + ke^{-\ell} + k^2(q\ell^2)^{m_0+1} + \sqrt{e q \ell^2 - 1}$$

Proof ideas: $d_{\text{TV}}(P, Q) \leq \frac{1}{2} \sqrt{\chi^2(P, Q)}$ and use **negative associations** [Dubhashi-Ranjan '98] to get rid of dependency in calculating the χ^2 distance.

Bound the total variation distance

Lemma

Let $\ell, n \in \mathbb{N}$, $k \in [n]$ and $\gamma \in (0, \frac{1}{2}]$. Let $N = \ell n$, $K = k\ell$, $p = 2q$ and $m_0 = \lfloor \log_2(1/\gamma) \rfloor$. Assume that $16q\ell^2 \leq 1$ and $k \geq 6\ell$. If $G \sim \mathcal{G}(n, \gamma)$, then $\tilde{G} \sim \mathcal{G}(N, q)$. If $G \sim \mathcal{G}(n, k, 1, \gamma)$, then

$$d_{\text{TV}}(P_{\tilde{G}}, \mathcal{G}(N, K, p, q)) \lesssim e^{-K} + ke^{-\ell} + k^2(q\ell^2)^{m_0+1} + \sqrt{e^{q\ell^2} - 1}$$

Proof ideas: $d_{\text{TV}}(P, Q) \leq \frac{1}{2}\sqrt{\chi^2(P, Q)}$ and use **negative associations** [Dubhashi-Ranjan '98] to get rid of dependency in calculating the χ^2 distance.

Apply the Lemma by choosing $q = \ell^{-2-\delta}$ so that $q\ell^2 \rightarrow 0$: $N = \ell^{\frac{2+\delta}{\alpha}}$, $K = \ell^{\frac{(2+\delta)\beta}{\alpha}}$, $n = \ell^{\frac{2+\delta}{\alpha}-1}$, $k = \ell^{\frac{(2+\delta)\beta}{\alpha}-1}$. Easy to check that

$$\alpha < \beta < \frac{1}{2} - \delta + \frac{\alpha(1+2\delta)}{4+2\delta} \Rightarrow \frac{\log k}{\log n} \leq \frac{1}{2} - \delta$$

Theorem

Let A denote a symmetric and zero-diagonal random matrix, where the entries $\{A_{ij} : i < j\}$ are independent and $[0, 1]$ -valued. Assume that $\mathbb{E}[A_{ij}] \leq p$, where $c_0 \log n/n \leq p \leq 1 - c_1$ for arbitrary constants $c_0 > 0$ and $c_1 > 0$. Then for any $c > 0$, there exists $c' > 0$ such that for any $n \geq 1$,

$$\mathbb{P} \left\{ \|A - \mathbb{E}[A]\|_2 \leq c' \sqrt{np} \right\} \geq 1 - n^{-c}.$$