

Information Limits for Recovering a Hidden Community

Bruce Hajek ¹ Yihong Wu ² Jiaming Xu ³

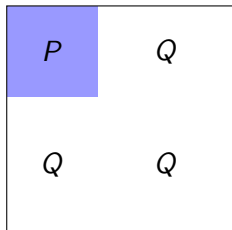
¹University of Illinois at Urbana-Champaign

²University of Illinois at Urbana-Champaign → Yale

³Simons Insitute, UC Berkeley → Purdue

July 12, 2016

Hidden community model [Deshpande-Montanari '13]



- Data: $n \times n$ symmetric matrix A with empty diagonal
- Community $C^* \subset [n]$ of size K uniform at random, such that

$$A_{ij} \sim \begin{cases} P & \text{both } i \text{ and } j \in C^* \\ Q & \text{otherwise} \end{cases}$$

- (K, P, Q) varies with n
- Goal is recovery of C^* from A (almost exactly or exactly)
- Fruitful venue for studying computational aspects of statistical problems

Planted dense subgraph

$$P = \text{Bern}(p), Q = \text{Bern}(q), \quad p > q$$

- A = adjacency matrix of $G(n, q)$ planted with $G(K, p)$
- [Alon et al '98, McSherry '01, Arias-Castro-Verzelen '14, Chen-Xu 14, Montanari '15, ...]

Planted dense subgraph

$$P = \text{Bern}(p), Q = \text{Bern}(q), \quad p > q$$

- A = adjacency matrix of $G(n, q)$ planted with $G(K, p)$
- [Alon et al '98, McSherry '01, Arias-Castro-Verzelen '14, Chen-Xu 14, Montanari '15, ...]

Submatrix localization

$$P = \mathcal{N}(0, \mu), Q = \mathcal{N}(0, 1), \quad \mu > 0$$

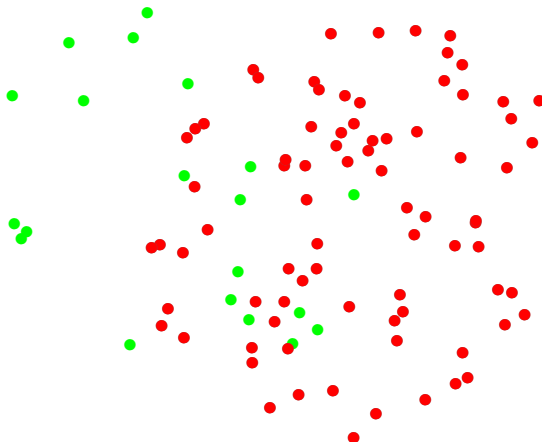
- $A = \begin{bmatrix} \mu & \\ & 0 \end{bmatrix} + \begin{bmatrix} \text{noise} \end{bmatrix}$
- [Shabalin et al '09, Butucea-Ingster '11, Kolar et al '11, Ma-W '13, Cai et al '15, ...]

Planted dense subgraph – graph view



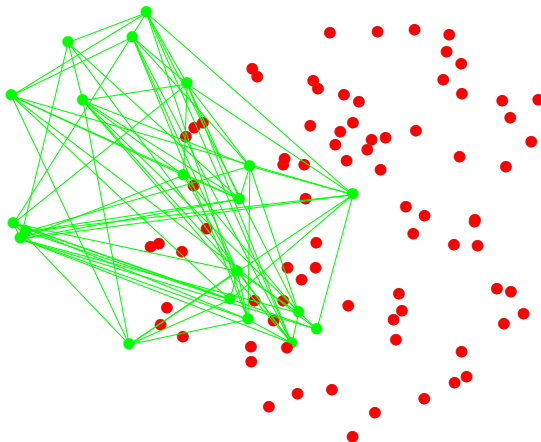
Planted dense subgraph – graph view

- 1 A community of K vertices are chosen randomly



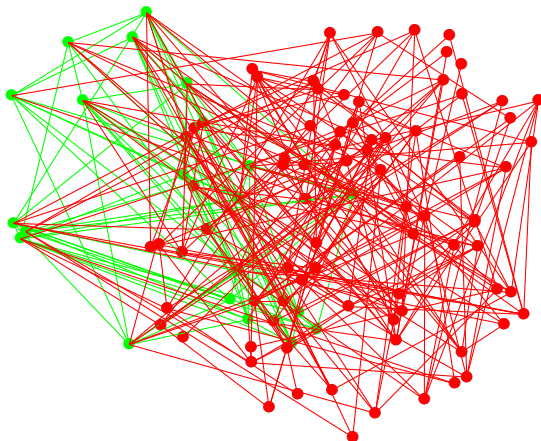
Planted dense subgraph – graph view

- ① A community of K vertices are chosen randomly
- ② For every pair of nodes in the community, add an edge w.p. p



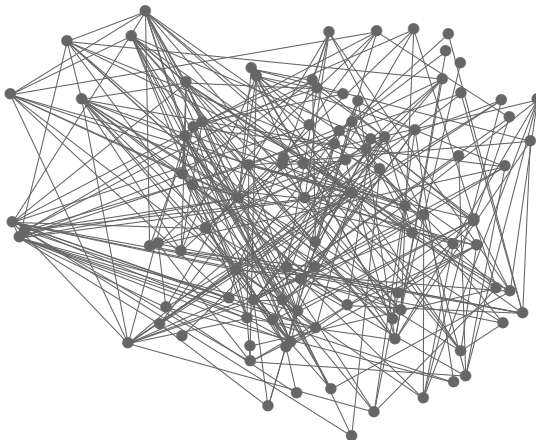
Planted dense subgraph – graph view

- ① A community of K vertices are chosen randomly
- ② For every pair of nodes in the community, add an edge w.p. p
- ③ For other pairs of nodes, add an edge w.p. q

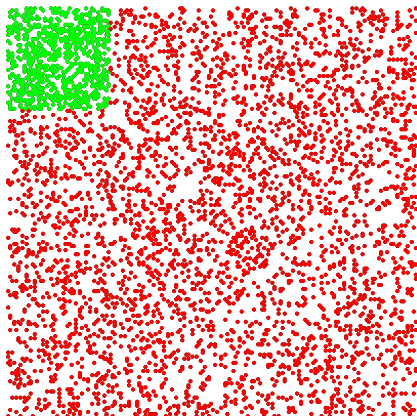


Planted dense subgraph – graph view

- ① A community of K vertices are chosen randomly
- ② For every pair of nodes in the community, add an edge w.p. p
- ③ For other pairs of nodes, add an edge w.p. q

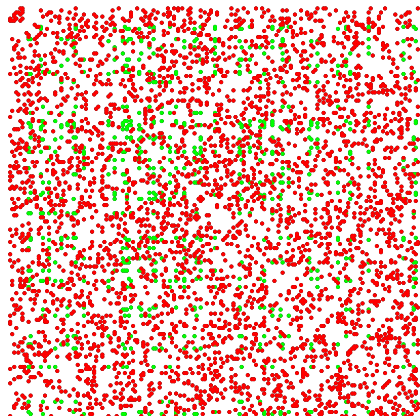


Planted dense subgraph – adjacency matrix view



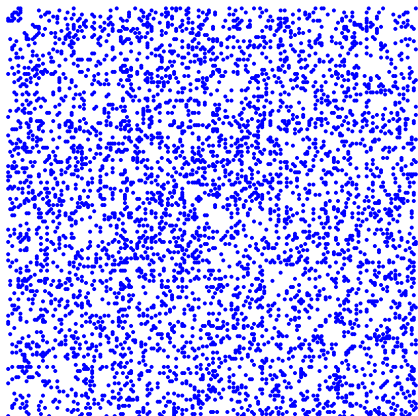
$$n = 200, K = 50, p = 0.3, q = 0.1$$

Planted dense subgraph – adjacency matrix view



$$n = 200, K = 50, p = 0.3, q = 0.1$$

Planted dense subgraph – adjacency matrix view



$$n = 200, K = 50, p = 0.3, q = 0.1$$

Assumption 1 $K \geq 2$, K/n is bounded away from one as $n \rightarrow \infty$.

Assumption 1 $K \geq 2$, K/n is bounded away from one as $n \rightarrow \infty$.

Assumption 2 There exists a constant C such that for all n ,

$$\psi_Q''(\lambda) \leq C \min\{D(P\|Q), D(Q\|P)\}, \quad \forall \lambda \in [-1, 1].$$

where $\psi_Q(\lambda) = \log \mathbb{E}_Q[\exp(\lambda L)]$ with $L \triangleq \frac{dP}{dQ}$

Assumption 1 $K \geq 2$, K/n is bounded away from one as $n \rightarrow \infty$.

Assumption 2 There exists a constant C such that for all n ,

$$\psi_Q''(\lambda) \leq C \min\{D(P\|Q), D(Q\|P)\}, \quad \forall \lambda \in [-1, 1].$$

where $\psi_Q(\lambda) = \log \mathbb{E}_Q[\exp(\lambda L)]$ with $L \triangleq \frac{dP}{dQ}$

Notes:

- $\psi_Q''(\lambda) = \text{var}_{Q_\lambda}(L)$, where $dQ_\lambda = \exp(\lambda L - \psi_Q(\lambda))dQ$,
- Assumption 2 holds if $L \triangleq \frac{dP}{dQ}$ is bounded. e.g. in Bernoulli case, if $\frac{p}{q}$ and $\frac{1-p}{1-q}$ are bounded away from 0 and ∞ .
- Assumption 2 holds if $P = \mathcal{N}(\mu, 1)$ and $Q = \mathcal{N}(0, 1)$ (**more generally, exp. families** under mild conditions): Then $L(x) = \mu(x - \frac{\mu}{2})$, $D(P\|Q) = D(Q\|P) = \mu^2/2$ and $\psi_Q''(\lambda) \equiv \mu^2$.

Main result for weak (aka almost exact) recovery

Theorem (Weak recovery)

If

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} > 2, \quad (1)$$

then

$$\mathbb{P}\{|\hat{C}_{\text{ML}} \triangle C^*| \leq 2K\epsilon\} \geq 1 - e^{-\Omega(K/\epsilon)},$$

where $\epsilon = 1/\sqrt{KD(P\|Q)}$.

If there exists \hat{C} such that $\mathbb{E}[|\hat{C} \triangle C^|] = o(K)$, then*

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} \geq 2. \quad (2)$$

Main result for weak (aka almost exact) recovery

Theorem (Weak recovery)

If

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} > 2, \quad (1)$$

then

$$\mathbb{P}\{|\hat{C}_{\text{ML}} \triangle C^*| \leq 2K\epsilon\} \geq 1 - e^{-\Omega(K/\epsilon)},$$

where $\epsilon = 1/\sqrt{KD(P\|Q)}$. [union bound]

If there exists \hat{C} such that $\mathbb{E}[|\hat{C} \triangle C^*|] = o(K)$, then

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} \geq 2. \quad (2)$$

Main result for weak (aka almost exact) recovery

Theorem (Weak recovery)

If

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} > 2, \quad (1)$$

then

$$\mathbb{P}\{|\hat{C}_{\text{ML}} \triangle C^*| \leq 2K\epsilon\} \geq 1 - e^{-\Omega(K/\epsilon)},$$

where $\epsilon = 1/\sqrt{KD(P\|Q)}$. [union bound]

If there exists \hat{C} such that $\mathbb{E}[|\hat{C} \triangle C^*|] = o(K)$, then

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} \geq 2. \quad (2)$$

[genie argument]

Main result for weak (aka almost exact) recovery

Theorem (Weak recovery)

If

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} > 2, \quad (1)$$

then

$$\mathbb{P}\{|\hat{C}_{\text{ML}} \triangle C^*| \leq 2K\epsilon\} \geq 1 - e^{-\Omega(K/\epsilon)},$$

where $\epsilon = 1/\sqrt{KD(P\|Q)}$. [union bound]

If there exists \hat{C} such that $\mathbb{E}[|\hat{C} \triangle C^*|] = o(K)$, then

$$K \cdot D(P\|Q) \rightarrow \infty \quad \text{and} \quad \liminf_{n \rightarrow \infty} \frac{(K-1)D(P\|Q)}{\log \frac{n}{K}} \geq 2. \quad (2)$$

[genie argument] [mutual information / rate distortion]

Main result for exact recovery

Theorem (Exact recovery)

If *weak recovery suff. condition (1)* and the following hold:

$$\liminf_{n \rightarrow \infty} \frac{KE_Q\left(\frac{1}{K} \log \frac{n}{K}\right)}{\log n} > 1. \quad (3)$$

then the maximum likelihood estimator satisfies $P\{\hat{C}_{ML} = C^*\} \rightarrow 1$.

If there exists an estimator \hat{C} such that $P\{\hat{C} = C^*\} \rightarrow 1$, then *weak recovery nece. condition (2)* and the following hold:

$$\liminf_{n \rightarrow \infty} \frac{KE_Q\left(\frac{1}{K} \log \frac{n}{K}\right)}{\log n} \geq 1. \quad (4)$$

Note: $E_Q(\theta) = \psi_Q^*(\theta) \triangleq \sup_{\lambda \in \mathbb{R}} \lambda\theta - \psi_Q(\lambda)$ (**Legendre transform of the log moment generating function**)

VOTING WITH SUCCESSIVE WITHHOLDING ALGORITHM

- ① Input: $n \in \mathbb{N}$, $K > 0$, distributions P, Q ; observed matrix A ; $\delta \in (0, 1)$ with $1/\delta, n\delta \in \mathbb{N}$.
- ② ([Partition](#)): Partition $[n]$ into $1/\delta$ subsets S_k of size $n\delta$.
- ③ ([Approximate Recovery](#)) For each $k = 1, \dots, 1/\delta$, produce $\hat{C}_k \approx C^* \cap ([n] \setminus S_k)$.
- ④ ([Cleanup](#)) For each $k = 1, \dots, 1/\delta$ compute $r_i = \sum_{j \in \hat{C}_k} L_{ij}$ for all $i \in S_k$ and return \check{C} , the set of K indices in $[n]$ with the largest values of r_i .

Lemma (Voting Converse for Exact Recovery)

Assume that $K \rightarrow \infty$ and $\limsup K/n < 1$. *Let L_i denote i.i.d. copies of $\log \frac{dP}{dQ}$.* If there exists an estimator \hat{C} such that $P\{\hat{C} = C^*\} \rightarrow 1$, then for any $K_o \rightarrow \infty$ such that $K_o = o(K)$, there exists a threshold θ_n

$$P \left[\sum_{i=1}^{K-K_o} L_i \leq (K-1)\theta_n - (K_o-1)D(P\|Q) - 6\sigma \right] \leq \frac{2}{K_o}$$
$$Q \left[\sum_{i=1}^{K-1} L_i \geq (K-1)\theta_n \right] \leq \frac{1}{n-K},$$

where $\sigma^2 = K_o \text{var}_P(L_1)$.

Let ALGX=semi-definite programming relaxation or belief propagation
For both planted dense subgraph (Bernoulli) and submatrix localization (Gaussian)

- $K = \omega(\frac{n}{\log n})$: ALGX attains the info-theoretic limit with sharp constants
- $K = \Theta(\frac{n}{\log n})$: ALGX is order-wise optimal, but strictly suboptimal by a constant factor
- $K = o(\frac{n}{\log n})$ and $K \rightarrow \infty$: ALGX is order-wise suboptimal

Can the computational gap for exact recovery be bridged by any polynomial time algorithm? (SoS hardness result or reduction to planted clique would offer further evidence for “no” answer.)

Can the computational gap for exact recovery be bridged by any polynomial time algorithm? (SoS hardness result or reduction to planted clique would offer further evidence for “no” answer.)

Thank you!