

Optimal rates for community estimation in the weighted stochastic block model

Po-Ling Loh

University of Cambridge
Department of Pure Mathematics and Mathematical Statistics

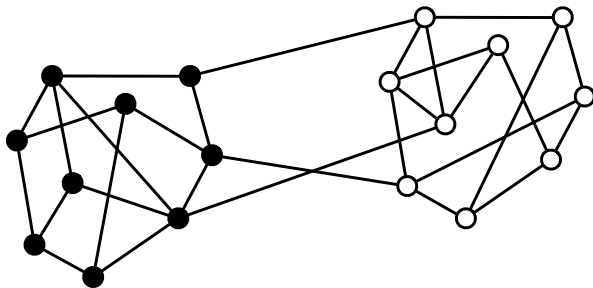
BIRS workshop on random graphs and statistical inference
August 10, 2021

Joint work with Min Xu (Rutgers) and Varun Jog (Cambridge)



Community recovery/estimation

- **Given:** Undirected graph $G = (V, E)$ on n nodes
- **Goal:** Partition nodes into communities based on relative connectivity



- Network structure may be assortative (homophilic), disassortative (heterophilic), etc.

Stochastic block model

- Probabilistic model introduced by Holland et al. '83 for generating community-structured data

Stochastic block model

- Probabilistic model introduced by Holland et al. '83 for generating community-structured data
- $K \geq 2$ communities
- For node i , let $\sigma(i) \in \{1, 2, \dots, K\}$ denote community assignment
- Presence of edge depends only on communities of incident vertices:

$$W_{i,j} \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p_{\sigma(i)\sigma(j)})$$

- What are good algorithms for community estimation?
- Can “optimal” accuracy be obtained with computationally tractable algorithms?

- For clustering algorithm $\hat{\sigma}$, define *misclassification error*

$$\ell(\hat{\sigma}(W), \sigma_0) := \min_{\tau \in \mathcal{S}_K} \frac{1}{n} d_H(\hat{\sigma}(W), \tau \circ \sigma_0)$$

- For clustering algorithm $\hat{\sigma}$, define *misclassification error*

$$\ell(\hat{\sigma}(W), \sigma_0) := \min_{\tau \in \mathcal{S}_K} \frac{1}{n} d_H(\hat{\sigma}(W), \tau \circ \sigma_0)$$

- Define *risk*

$$R(\hat{\sigma}, \sigma_0) := \mathbb{E} [\ell(\hat{\sigma}(W), \sigma_0)]$$

- For clustering algorithm $\hat{\sigma}$, define *misclassification error*

$$\ell(\hat{\sigma}(W), \sigma_0) := \min_{\tau \in \mathcal{S}_K} \frac{1}{n} d_H(\hat{\sigma}(W), \tau \circ \sigma_0)$$

- Define *risk*

$$R(\hat{\sigma}, \sigma_0) := \mathbb{E}[\ell(\hat{\sigma}(W), \sigma_0)]$$

- Characterize minimax risk

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, P \in \Theta} R(\hat{\sigma}, \sigma_0)$$

Sharp thresholds (Zhang & Zhou '15)

- Minimax bounds for misclassification error in (roughly) equal-sized communities:

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, P \in \Theta} R(\hat{\sigma}, \sigma_0) = \exp\left(- (1 + o(1)) \frac{nl_n}{K}\right)$$

where Θ is parameter space such that within-community edges have probability $\geq \frac{a}{n}$, between-community edges have probability $\leq \frac{b}{n}$, and

$$\begin{aligned} l_n &= -2 \log \left(\sqrt{\frac{a}{n}} \sqrt{\frac{b}{n}} + \sqrt{1 - \frac{a}{n}} \sqrt{1 - \frac{b}{n}} \right) \\ &= D_{1/2} \left(\text{Ber} \left(\frac{a}{n} \right) \parallel \text{Ber} \left(\frac{b}{n} \right) \right) \end{aligned}$$

is *Renyi divergence* of order $\frac{1}{2}$

- Implies weak consistency is governed by behavior of l_n : If $nl_n \rightarrow \infty$,

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, P \in \Theta} R(\hat{\sigma}, \sigma_0) \rightarrow 0$$

Sharp thresholds (Zhang & Zhou '15)

- Implies weak consistency is governed by behavior of I_n : If $nI_n \rightarrow \infty$,

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, P \in \Theta} R(\hat{\sigma}, \sigma_0) \rightarrow 0$$

- **Rough intuition:** $I_n = D_{1/2} \left(\text{Ber} \left(\frac{a}{n} \right) \parallel \text{Ber} \left(\frac{b}{n} \right) \right)$ quantifies *distinguishability* of within-community and between-community edges

Sharp thresholds (Zhang & Zhou '15)

- Implies weak consistency is governed by behavior of I_n : If $nI_n \rightarrow \infty$,

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, P \in \Theta} R(\hat{\sigma}, \sigma_0) \rightarrow 0$$

- **Rough intuition:** $I_n = D_{1/2} \left(\text{Ber} \left(\frac{a}{n} \right) \parallel \text{Ber} \left(\frac{b}{n} \right) \right)$ quantifies *distinguishability* of within-community and between-community edges
- **Question:** What if edge distributions are not Bernoulli?

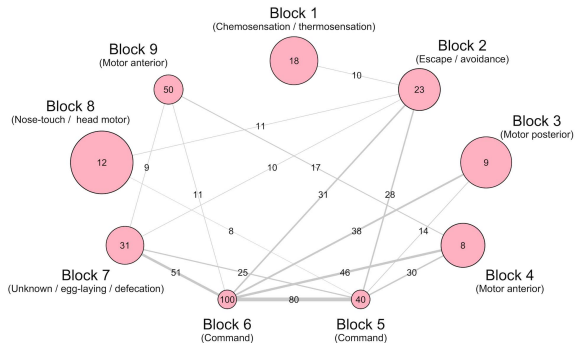
Motivating examples

- Weighted/labeled graphs occur naturally in social networks, airline networks, neural networks, etc.



C. elegans nervous system:

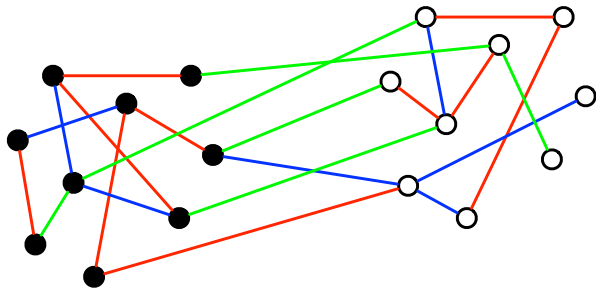
- Chemical synapses
- Gap junctions
- Neuromuscular junctions



“Stochastic blockmodeling of the modules
and core of the *Caenorhabditis elegans* connectome”
Pavlovic et al. (2014)

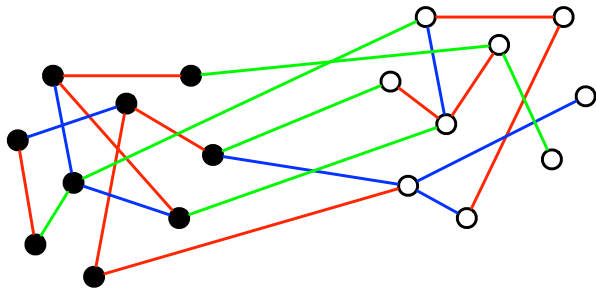
Weighted stochastic block model (discrete case)

- Entries of adjacency matrix drawn from general distributions



Weighted stochastic block model (discrete case)

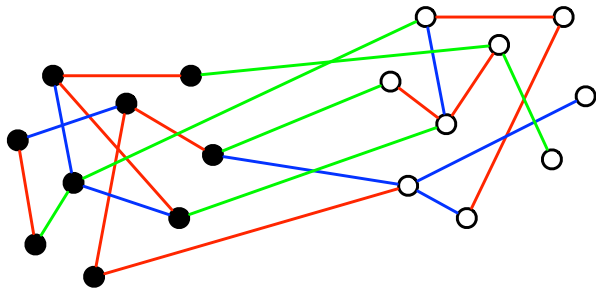
- Entries of adjacency matrix drawn from general distributions



- p_n, q_n denote within-community, between-community distributions

Weighted stochastic block model (discrete case)

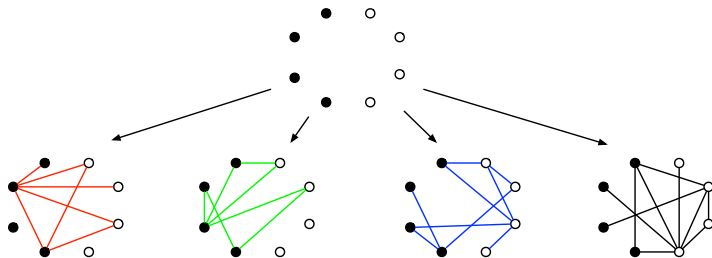
- Entries of adjacency matrix drawn from general distributions



- p_n, q_n denote within-community, between-community distributions
- Edge labels may contain valuable information for community recovery

Example: Multilayer networks

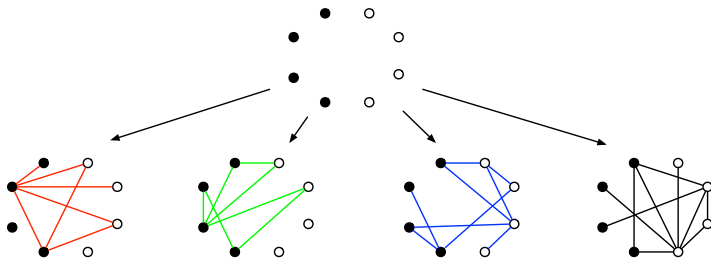
- Nodes in graph give rise to m different edge sets/“views”



- **Goal:** Combine all views to recover underlying community structure

Example: Multilayer networks

- Nodes in graph give rise to m different edge sets/“views”



- **Goal:** Combine all views to recover underlying community structure
- Can be viewed as single weight distribution taking 2^m possible values

- MLE formulation has good theoretical properties, but computationally infeasible

Algorithms for unweighted SBMs

- MLE formulation has good theoretical properties, but computationally infeasible
- Popular approach based on **spectral clustering**: W is close to $\mathbb{E}[W]$, which (viewed columnwise) has well-separated cluster structure

Algorithms for unweighted SBMs

- MLE formulation has good theoretical properties, but computationally infeasible
- Popular approach based on **spectral clustering**: W is close to $\mathbb{E}[W]$, which (viewed columnwise) has well-separated cluster structure
 - Performing spectral clustering on W yields weak recovery under appropriate problem scaling (Lei & Rinaldo '15)

Algorithms for unweighted SBMs

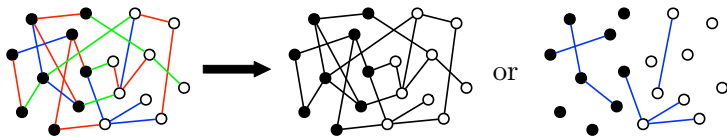
- MLE formulation has good theoretical properties, but computationally infeasible
- Popular approach based on **spectral clustering**: W is close to $\mathbb{E}[W]$, which (viewed columnwise) has well-separated cluster structure
 - Performing spectral clustering on W yields weak recovery under appropriate problem scaling (Lei & Rinaldo '15)
 - To achieve *optimal misclassification error*, require additional refinement steps involving local MLE calculations (Gao et al. '15)

Algorithms for unweighted SBMs

- MLE formulation has good theoretical properties, but computationally infeasible
- Popular approach based on **spectral clustering**: W is close to $\mathbb{E}[W]$, which (viewed columnwise) has well-separated cluster structure
 - Performing spectral clustering on W yields weak recovery under appropriate problem scaling (Lei & Rinaldo '15)
 - To achieve *optimal misclassification error*, require additional refinement steps involving local MLE calculations (Gao et al. '15)
- Naively applying spectral clustering to weighted adjacency matrix *cannot* be optimal, since numerical labels may be arbitrary

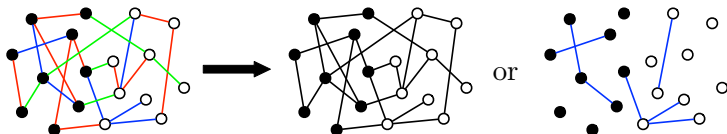
Reducing to unweighted case

- **Idea:** Map weighted graph to unweighted graph, then perform recovery algorithm for unweighted SBM



Reducing to unweighted case

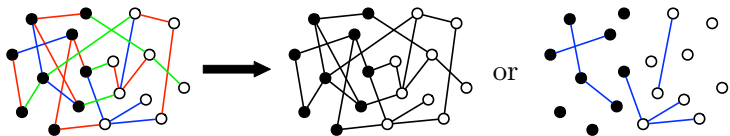
- **Idea:** Map weighted graph to unweighted graph, then perform recovery algorithm for unweighted SBM



- **Questions:** Could this possibly lead to optimal estimation rates...? How to perform mapping?

Reducing to unweighted case

- **Idea:** Map weighted graph to unweighted graph, then perform recovery algorithm for unweighted SBM



- **Questions:** Could this possibly lead to optimal estimation rates...? How to perform mapping?
- Modifying Zhang & Zhou '15 analysis yields lower bound

$$\inf_{\hat{\sigma}} \sup_{\sigma_0, (p_n, q_n) \in (\mathcal{P}, \mathcal{Q})} R(\hat{\sigma}, \sigma_0) \geq \exp \left(-(1 + o(1)) \frac{nl_n}{K} \right),$$

but this rate cannot be achieved simply by dropping edge weights

Characterization via Hellinger distance

- In setting where $I_n \rightarrow 0$,

$$\begin{aligned} I_n &= -2 \log \left(\sum_{\ell=0}^L \sqrt{p_n(\ell)q_n(\ell)} \right) \\ &= \left(\sum_{\ell=0}^L \left(\sqrt{p_n(\ell)} - \sqrt{q_n(\ell)} \right)^2 \right) (1 + o(1)) \end{aligned}$$

- Latter expression known as *Hellinger distance* between p_n and q_n

- **Key insight:** If $nl_n \rightarrow \infty$ (regime where weak recovery is possible), $\exists \ell^* \in \{1, \dots, L\}$ such that

$$n \left(\sqrt{p_n(\ell^*)} - \sqrt{q_n(\ell^*)} \right)^2 \rightarrow \infty,$$

so community recovery based on ℓ^* alone achieves weak recovery

- **Key insight:** If $nl_n \rightarrow \infty$ (regime where weak recovery is possible), $\exists \ell^* \in \{1, \dots, L\}$ such that

$$n \left(\sqrt{p_n(\ell^*)} - \sqrt{q_n(\ell^*)} \right)^2 \rightarrow \infty,$$

so community recovery based on ℓ^* alone achieves weak recovery

- **Problem:** How to identify ℓ^* based on observing graph?

- **Key insight:** If $nl_n \rightarrow \infty$ (regime where weak recovery is possible), $\exists \ell^* \in \{1, \dots, L\}$ such that

$$n \left(\sqrt{p_n(\ell^*)} - \sqrt{q_n(\ell^*)} \right)^2 \rightarrow \infty,$$

so community recovery based on ℓ^* alone achieves weak recovery

- **Problem:** How to identify ℓ^* based on observing graph?
- How to refine initial clustering based on ℓ^* to obtain optimal error rate?

- For each $1 \leq \ell \leq L$, perform spectral clustering on unweighted adjacency matrix A_ℓ to obtain community assignments $\hat{\sigma}_\ell$

- For each $1 \leq \ell \leq L$, perform spectral clustering on unweighted adjacency matrix A_ℓ to obtain community assignments $\hat{\sigma}_\ell$
- Estimate within- and between-community edge probabilities:

$$\hat{P}_\ell = \frac{\sum_{u \neq v: \hat{\sigma}_\ell(u) = \hat{\sigma}_\ell(v)} (A_\ell)_{uv}}{|\{u \neq v : \hat{\sigma}_\ell(u) = \hat{\sigma}_\ell(v)\}|}, \quad \hat{Q}_\ell = \frac{\sum_{u \neq v: \hat{\sigma}_\ell(u) \neq \hat{\sigma}_\ell(v)} (A_\ell)_{uv}}{|\{u \neq v : \hat{\sigma}_\ell(u) \neq \hat{\sigma}_\ell(v)\}|}$$

- For each $1 \leq \ell \leq L$, perform spectral clustering on unweighted adjacency matrix A_ℓ to obtain community assignments $\hat{\sigma}_\ell$
- Estimate within- and between-community edge probabilities:

$$\hat{P}_\ell = \frac{\sum_{u \neq v: \hat{\sigma}_\ell(u) = \hat{\sigma}_\ell(v)} (A_\ell)_{uv}}{|\{u \neq v : \hat{\sigma}_\ell(u) = \hat{\sigma}_\ell(v)\}|}, \quad \hat{Q}_\ell = \frac{\sum_{u \neq v: \hat{\sigma}_\ell(u) \neq \hat{\sigma}_\ell(v)} (A_\ell)_{uv}}{|\{u \neq v : \hat{\sigma}_\ell(u) \neq \hat{\sigma}_\ell(v)\}|}$$

- Select

$$\ell^* := \arg \max_{\ell} \left\{ \frac{(\hat{P}_\ell - \hat{Q}_\ell)^2}{\hat{P}_\ell \vee \hat{Q}_\ell} \right\}$$

to be color achieving “maximal separation”

Achieving optimal rate

- However, recovery algorithm $\hat{\sigma}_{\ell^*}$ yields *suboptimal* error rate

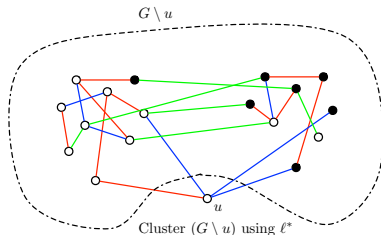
Achieving optimal rate

- However, recovery algorithm $\hat{\sigma}_{\ell^*}$ yields *suboptimal* error rate
- Additional steps (based on Gao et al. '15) can bootstrap $\hat{\sigma}_{\ell^*}$ to optimal recovery rates:

Achieving optimal rate

- However, recovery algorithm $\hat{\sigma}_{\ell^*}$ yields *suboptimal* error rate
- Additional steps (based on Gao et al. '15) can bootstrap $\hat{\sigma}_{\ell^*}$ to optimal recovery rates:
 - 1 For each node u ,
 - Cluster $G \setminus u$ via spectral method with color ℓ^* to obtain $\hat{\sigma}_u$
 - Use local MLE to estimate assignment of u :

$$\hat{\sigma}_u(u) = \arg \max_k \sum_{v: \hat{\sigma}_u(v)=k, v \neq u} \sum_{\ell=0}^L \log \frac{\hat{P}_{\ell}^*}{\hat{Q}_{\ell}^*}(A_{\ell})_{uv}$$



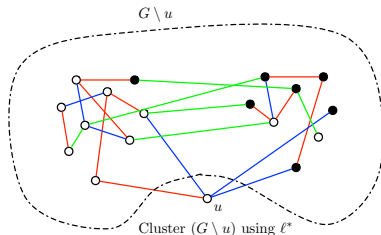
Achieving optimal rate

- However, recovery algorithm $\hat{\sigma}_{\ell^*}$ yields *suboptimal* error rate
- Additional steps (based on Gao et al. '15) can bootstrap $\hat{\sigma}_{\ell^*}$ to optimal recovery rates:

1 For each node u ,

- Cluster $G \setminus u$ via spectral method with color ℓ^* to obtain $\hat{\sigma}_u$
- Use local MLE to estimate assignment of u :

$$\hat{\sigma}_u(u) = \arg \max_k \sum_{v: \hat{\sigma}_u(v)=k, v \neq u} \sum_{\ell=0}^L \log \frac{\hat{P}_\ell^*}{\hat{Q}_\ell^*}(A_\ell)_{uv}$$



2 Combine assignments $\{\hat{\sigma}_u\}_{u \in V}$ to obtain overall clustering

Theorem: Optimal recovery (discrete)

Theorem

Suppose $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$. Then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\ell(\hat{\sigma}(W), \sigma_0) \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right) \right) = 1.$$

- Conditions $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$ imply weight distributions converge together, but not *too* quickly

Theorem: Optimal recovery (discrete)

Theorem

Suppose $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$. Then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\ell(\hat{\sigma}(W), \sigma_0) \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right) \right) = 1.$$

- Conditions $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$ imply weight distributions converge together, but not *too* quickly
- Matches lower bound stated above: Misclassification error governed by Renyi divergence

- Also desirable to perform community recovery in settings where edge weights are continuous
- **Examples:** Gaussian distributions, exponential distributions, etc.

- Also desirable to perform community recovery in settings where edge weights are continuous
- **Examples:** Gaussian distributions, exponential distributions, etc.
- Previous work relatively scarce, depends on *parametric* assumptions

- Also desirable to perform community recovery in settings where edge weights are continuous
- **Examples:** Gaussian distributions, exponential distributions, etc.
- Previous work relatively scarce, depends on *parametric* assumptions
- **Our contribution:** Nonparametric, computationally feasible recovery algorithm achieving optimal error rates

Mathematical formulation

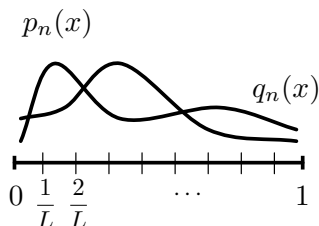
- Suppose weight distribution is **mixture** of point mass at 0 with probability $P_{0,n}$, continuous distribution with density $p_n(x)$ with probability $1 - P_{0,n}$ (similarly for $Q_{0,n}$ and $q_n(x)$)

Mathematical formulation

- Suppose weight distribution is **mixture** of point mass at 0 with probability $P_{0,n}$, continuous distribution with density $p_n(x)$ with probability $1 - P_{0,n}$ (similarly for $Q_{0,n}$ and $q_n(x)$)
- In order to obtain optimal error rate, **discretize** $p_n(x)$ and $q_n(x)$ more and more finely and reduce to previous algorithm

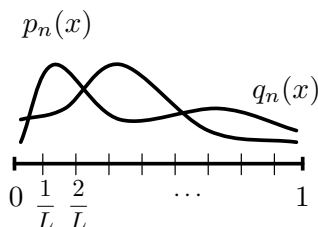
Mathematical formulation

- Suppose weight distribution is **mixture** of point mass at 0 with probability $P_{0,n}$, continuous distribution with density $p_n(x)$ with probability $1 - P_{0,n}$ (similarly for $Q_{0,n}$ and $q_n(x)$)
- In order to obtain optimal error rate, **discretize** $p_n(x)$ and $q_n(x)$ more and more finely and reduce to previous algorithm
 - If $p_n(x)$ and $q_n(x)$ have bounded support $S = [0, 1]$, partition into L_n equal-length intervals



Mathematical formulation

- Suppose weight distribution is **mixture** of point mass at 0 with probability $P_{0,n}$, continuous distribution with density $p_n(x)$ with probability $1 - P_{0,n}$ (similarly for $Q_{0,n}$ and $q_n(x)$)
- In order to obtain optimal error rate, **discretize** $p_n(x)$ and $q_n(x)$ more and more finely and reduce to previous algorithm
 - If $p_n(x)$ and $q_n(x)$ have bounded support $S = [0, 1]$, partition into L_n equal-length intervals



- Otherwise, first apply transformation $\Phi : S \rightarrow [0, 1]$ and then partition into L_n equal-length intervals

Theorem

Suppose $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$ and p_n and q_n satisfy appropriate regularity conditions with respect to Φ . If $L_n \rightarrow \infty$ is chosen such that

$\frac{nl_n}{L_n \exp(L_n^{1/r})} \rightarrow \infty$, then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\ell(\hat{\sigma}(W), \sigma_0) \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right) \right) = 1.$$

Theorem

Suppose $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$ and p_n and q_n satisfy appropriate regularity conditions with respect to Φ . If $L_n \rightarrow \infty$ is chosen such that

$\frac{nl_n}{L_n \exp(L_n^{1/r})} \rightarrow \infty$, then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\ell(\hat{\sigma}(W), \sigma_0) \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right) \right) = 1.$$

- Note that if $nl_n \rightarrow \infty$, appropriate discretization level L_n always exists

Theorem

Suppose $l_n \rightarrow 0$ and $nl_n \rightarrow \infty$ and p_n and q_n satisfy appropriate regularity conditions with respect to Φ . If $L_n \rightarrow \infty$ is chosen such that $\frac{nl_n}{L_n \exp(L_n^{1/r})} \rightarrow \infty$, then

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\ell(\hat{\sigma}(W), \sigma_0) \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right) \right) = 1.$$

- Note that if $nl_n \rightarrow \infty$, appropriate discretization level L_n always exists
- Under additional condition $\limsup_n \frac{nl_n}{\log n} \leq 1$, can obtain

$$\mathbb{E} [\ell(\hat{\sigma}(W), \sigma_0)] \leq \exp \left(-\frac{nl_n}{K} (1 + o(1)) \right),$$

agreeing with lower bound on risk

Examples

- When $S = [0, 1]$, simplest case satisfying conditions is when $p_n, q_n, |p'_n|$, and $|q'_n|$ are uniformly bounded away from 0 and ∞

Examples

- When $S = [0, 1]$, simplest case satisfying conditions is when $p_n, q_n, |p'_n|$, and $|q'_n|$ are uniformly bounded away from 0 and ∞
- When $S = \mathbb{R}$, consider *location-scale family*:

$$f_{\mu, \sigma}(x) = f\left(\frac{x - \mu}{\sigma}\right) - \log \sigma,$$

with parameter space $\mu \in [-C_\mu, C_\mu]$ and $\sigma \in \left[\frac{1}{c_\sigma}, c_\sigma\right]$

Examples

- When $S = [0, 1]$, simplest case satisfying conditions is when $p_n, q_n, |p'_n|$, and $|q'_n|$ are uniformly bounded away from 0 and ∞
- When $S = \mathbb{R}$, consider *location-scale family*:

$$f_{\mu,\sigma}(x) = f\left(\frac{x - \mu}{\sigma}\right) - \log \sigma,$$

with parameter space $\mu \in [-C_\mu, C_\mu]$ and $\sigma \in \left[\frac{1}{c_\sigma}, c_\sigma\right]$

- Densities given by

$$p_n(x) = \exp(f_{\mu_{1,n}, \sigma_{1,n}}(x)), \quad q_n(x) = \exp(f_{\mu_{2,n}, \sigma_{2,n}}(x))$$

Examples

- When $S = [0, 1]$, simplest case satisfying conditions is when $p_n, q_n, |p'_n|$, and $|q'_n|$ are uniformly bounded away from 0 and ∞
- When $S = \mathbb{R}$, consider *location-scale family*:

$$f_{\mu,\sigma}(x) = f\left(\frac{x - \mu}{\sigma}\right) - \log \sigma,$$

with parameter space $\mu \in [-C_\mu, C_\mu]$ and $\sigma \in \left[\frac{1}{c_\sigma}, c_\sigma\right]$

- Densities given by

$$p_n(x) = \exp(f_{\mu_{1,n}, \sigma_{1,n}}(x)), \quad q_n(x) = \exp(f_{\mu_{2,n}, \sigma_{2,n}}(x))$$

- Requirements:
 - $|f^{(k)}(x)|$ bounded for some $k \geq 2$,
 - $\exists c, M$ such that $f'(x) > M$ for $x < -c$ and $f'(x) < -M$ for $x > c$
(satisfied for Gaussian and Laplace families)

- Misclassification error rate for weighted SBMs sharply characterized by Renyi divergence

$$I_n = -2 \log \int \sqrt{p_n(x)q_n(x)} dx$$

- Misclassification error rate for weighted SBMs sharply characterized by Renyi divergence

$$I_n = -2 \log \int \sqrt{p_n(x)q_n(x)} dx$$

- Computationally feasible algorithm for optimal recovery:
 - 1 (If necessary) discretize continuous weight distribution
 - 2 Cluster according to each individual color, find ℓ^* with greatest separation
 - 3 Refine weakly consistent clustering based on ℓ^* to obtain optimal error rate

- Xu, Jog & Loh (2020). Optimal rates for community estimation in the weighted stochastic block model. *Annals of Statistics*.

Thank you!