

# Community Recovery on the Weighted Stochastic Block Model and Its Information-Theoretic Limits

Min Xu<sup>†</sup>  
minx@wharton.upenn.edu

Varun Jog<sup>‡</sup>  
vjog@wisc.edu

Po-Ling Loh<sup>‡\*</sup>  
loh@ece.wisc.edu

Department of Statistics<sup>†</sup>  
The Wharton School  
University of Pennsylvania  
Philadelphia, PA 19104

Departments of ECE<sup>‡</sup> & Statistics\*  
Grainger Institute of Engineering  
University of Wisconsin - Madison  
Madison, WI 53706

May 2017

## Abstract

Identifying communities in a network is an important problem in many fields, including social science, neuroscience, military intelligence, and genetic analysis. In the past decade, the Stochastic Block Model (SBM) has emerged as one of the most well-studied and well-understood statistical models for this problem. Yet, the SBM has an important limitation: it assumes that each network edge is drawn from a Bernoulli distribution. This is rather restrictive, since weighted edges are fairly ubiquitous in scientific applications, and disregarding edge weights naturally results in a loss of valuable information. In this paper, we study a weighted generalization of the SBM, where observations are collected in the form of a weighted adjacency matrix, and the weight of each edge is generated independently from a distribution determined by the community membership of its endpoints. We propose and analyze a novel algorithm for community estimation in the weighted SBM based on various subroutines involving transformation, discretization, spectral clustering, and appropriate refinements. We prove that our procedure is optimal in terms of its rate of convergence, and that the misclassification rate is characterized by the Renyi divergence between the distributions of within-community edges and between-community edges. In the regime where the edges are sparse, we also establish sharp thresholds for exact recovery of the communities. Our theoretical results substantially generalize previously established thresholds derived specifically for unweighted block models. Furthermore, our algorithm introduces a principled and computationally tractable method of incorporating edge weights to the analysis of network data.

## 1 Introduction

The recent explosion of interest in network data has created a need for new statistical methods for analyzing network datasets and interpreting results [8, 12, 17, 22]. One active area of research with diverse applications in many scientific fields pertains to community detection and estimation, where the information available consists of the presence or absence of edges between nodes in the graph, and the goal is to partition the nodes into disjoint groups based on their relative connectivity [9, 15, 19, 24, 25, 28].

A standard assumption in statistical modeling is that conditioned on the community labels of the nodes in the graph, edges are generated independently according to fixed distributions governing the connectivity of nodes within and between communities in the graph. This is the setting of the stochastic block model (SBM) [16, 29, 30]. In the homogeneous case, edges follow one distribution when both endpoints are in the same community, regardless of the community label; and edges follow a second distribution when the endpoints are in different communities. The majority of

existing literature on stochastic block models has focused on the case where no other information is available beyond the unweighted adjacency matrix, and much work in the information theory and statistics has focused on deriving thresholds for *exact* or *weak* recovery of community labels in terms of the underlying probability parameters and the size of the graph (e.g., [1–3, 13, 14, 18, 20, 21, 32]).

However, the pairwise connections in many real-world networks possess a natural weighting structure [7, 23]. For example, in social networks, information may be available quantifying the strength of a tie, such as the frequency of interactions between the individuals [27]; in cellular networks, information may be available quantifying the frequency of communication between users [6]; in gene co-expression networks, edges weights range from -1 to 1 and indicate the correlation between the expression levels of a gene pair; and in neural networks, edge weights may symbolize the level of neural activity between regions in the brain [26]. Of course, the connectivity data could be condensed into an adjacency matrix consisting of only zeros and ones, but this would result in a loss of valuable information that could be used to recover node communities.

In this paper, we analyze the “weighted” setting of the stochastic block model [4], where, after an edge is generated from a Bernoulli distribution, it is given an edge weight generated from one of two arbitrary densities  $p(x), q(x)$  depending on whether the edge is between-cluster or within-cluster. The weighted SBM presents a serious challenge in the design of algorithms because  $p(x), q(x)$  are unknown and must be estimated. Nonparametric estimation of a density is a difficult problem in its own right and it is made much harder in the weighted SBM because one does not know whether an edge weight is drawn from  $p(x)$  and  $q(x)$  without the latent cluster structure. There are various approaches to the weighted SBM. For example, Newman [23] assumes that the edge weights have discrete units and then converts a weighted graph into a multigraph; Aicher et al [4] assumes  $p(x), q(x)$  to be from a known exponential family and performs variational Bayesian inference. These approaches can be effective but they rely on strong assumptions to simplify the problems and nothing is known about their theoretical properties.

Our paper proposes a new discretization based approach that imposes weak assumptions and possesses strong guarantees. In the case of finitely-supported distributions, which correspond to a “labeled” or “colored” SBM, we demonstrate a method for choosing an initial label on which we apply a standard SBM estimation method to obtain an initial clustering. We then show how to use this initial rough clustering, together with the full set of edge labels, to obtain more accurate estimates of the true cluster assignments. In the case of continuous weight distributions, we propose a discretization strategy that will allow us to apply a recovery algorithm for the labeled case after appropriate preprocessing. Our method does not rely on prior knowledge of the densities  $p(x)$  and  $q(x)$  and does not rely on parametric assumptions.

Importantly, we show that the output of our algorithm is optimal, in the sense that under mild regularity assumptions on  $p(x)$  and  $q(x)$ , the misclustering error of our algorithm converges to zero at an optimal rate. Our analysis generalizes the results of Zhang and Zhou [32] and Gao et al [10], which show that the optimal rate of convergence of unweighted SBM is driven by the Renyi divergence of order  $1/2$  between two Bernoulli distributions, corresponding to the probability of generation for within-community and between-community edges. In fact, a similar phenomenon holds for the weighted SBM setting in our paper: the optimal error rate is also driven by a Renyi divergence of order  $1/2$  between two mixed distributions that capture both the divergence between the edge probabilities and the divergence between the edge weight densities  $p(x)$  and  $q(x)$ . Note that in order to achieve the optimal error rate, our discretization strategy must be chosen carefully when  $p(x)$  and  $q(x)$  are continuous distributions. Our proposed algorithm first transforms the distributions to be supported on  $[0, 1]$ , then bins the interval appropriately; in general, since  $p(x)$  and  $q(x)$  may vary with the size of the graph, the number of bins used will also need to grow slowly as the number of nodes increases. Our results has an interesting implication: although our algorithm is nonparametric, it is adaptive in the sense that it achieves the same optimal rate even if the edge

weight densities  $p(x), q(x)$  take on a parametric form such as Gaussian or Laplace. This is in contrast to most problems in statistics where nonparametric methods usually have slower rate of convergence than parametric methods in settings where a parametric form is known. This observation captures an important intuition behind our results, that on the weighted SBM, one do not need to estimate the densities well in order to cluster well.

The remainder of the paper is organized as follows: Section 2 introduces the mathematical framework of the weighted stochastic block model and defines the problems which we are trying to solve. Section 3 describes our proposed algorithm for finding communities on the weighted SBM. In Section 4 we provide the statements of our main results concerning the behavior of our algorithm in terms of misclassification error rates and exact recovery. Section 5 highlights the key technical components employed in the analysis of our algorithm. We close in Section 6 with further implications of our work and open questions related to our results.

## 2 Model and problem formulation

We begin with a formal definition of the weighted SBM and a description of our error metrics for clustering.

### 2.1 Model definition

Consider a network with  $n$  nodes and  $K \geq 2$  communities. In this paper, we suppose that the communities are approximately balanced; that is, there exists a *cluster-imbalance constant*  $\beta$  such that the cluster size  $n_k$  for each cluster  $k = 1, \dots, K$  satisfies  $\frac{\beta n}{K} \geq n_k \geq \frac{n}{\beta K}$ . For each node  $u$ , we let  $\sigma(u) \in \{1, 2, \dots, K\}$  denote the community assignment of the nodes.

**Definition 2.1.** (Homogeneous Stochastic Block Model) An edge random variable  $A_{uv}$  has the following distribution:

$$A_{uv} \sim \begin{cases} \text{Ber}(p) & \text{if } \sigma(u) = \sigma(v) \quad \text{and} \\ \text{Ber}(q) & \text{if } \sigma(u) \neq \sigma(v). \end{cases}$$

In the more general case of *heterogenous* SBM, we have a  $K \times K$  matrix  $P$  where each entry  $P_{ij} \in [0, 1]$ . The edge random variable is drawn from  $A_{uv} \sim \text{Ber}(P_{\sigma(u), \sigma(v)})$ . We focus on the homogeneous case in this paper but discuss how to extend our results to the heterogenous setting.

SBM gives a distribution over the set of all networks whose edges are binary. To adapt to networks with continuous edge weights, we generalize the homogenous SBM by adding a second step to the data generating process: an edge weight is sampled from a continuous distribution after it is generated.

**Definition 2.2.** (Weighted Homogeneous SBM) Let  $0 < P_0, Q_0 < 1$  and let  $p(x), q(x)$  be two densities. We first generate the edge presence indicator  $Z_{uv}$ :

$$Z_{uv} \sim \begin{cases} \text{Ber}(1 - P_0) & \text{if } \sigma(u) = \sigma(v) \quad \text{and} \\ \text{Ber}(1 - Q_0) & \text{if } \sigma(u) \neq \sigma(v). \end{cases}$$

The edge weight random variable is then:

$$A_{uv} \sim \begin{cases} 0 & \text{if } Z_{uv} = 0 \\ p(x) & \text{if } Z_{uv} = 1 \text{ and } \sigma(u) = \sigma(v) \quad \text{and} \\ q(x) & \text{if } Z_{uv} = 1 \text{ and } \sigma(u) \neq \sigma(v). \end{cases}$$

In this model, an edge is missing with probability either  $P_0$  or  $Q_0$  depending on whether the potential edge connects two nodes in the same cluster or in different clusters. If the edge is present,

then it is given an edge weight drawn from either the density  $p(x)$  or  $q(x)$ , depending again on the nature of the edge. If  $p(x)$  and  $q(x)$  are Dirac Delta mass at 1, then the weighted homogenous SBM reduces to homogeneous SBM with  $p = 1 - P_0$  and  $q = 1 - Q_0$ .

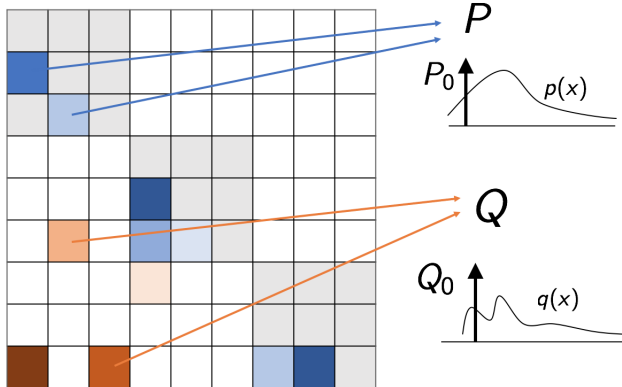


Figure 1: Weighted stochastic block model

The model defined in 2.2 is the focus of our method. However, it is useful to note that we can further generalize model 2.2 by allowing both weights and labels.

**Definition 2.3.** (Weighted and Labeled Homogenous SBM) Let  $P, Q$  be two general mixed distributions. The edge random variable  $A_{uv}$  is drawn as

$$A_{uv} \sim \begin{cases} P & \text{if } \sigma(u) = \sigma(v) \\ Q & \text{if } \sigma(u) \neq \sigma(v) \end{cases}$$

In the case where  $P, Q$  are mixed distributions with continuous part  $(1 - P_0)p(x)$  and  $(1 - Q_0)q(x)$  respectively and a discrete point mass of  $P_0, Q_0$  at zero respectively, then we get back the weighted SBM.

## 2.2 Community estimation

In this paper we aim to find a tractable community recovery algorithm whose misclustering error can be shown to converge to zero at an optimal rate.

### 2.2.1 Misclustering error rate

The goal of a community recovery algorithm is to take as input the adjacency matrix  $A$  and try to recover the community assignments. We evaluate a community recovery algorithm by looking at its mis-clustering error rate. To be precise, if  $\sigma_0$  is the true clustering and  $\hat{\sigma}$  is the clustering generated by a community recovery algorithm, then the misclustering error rate is the following loss function:

$$l(\hat{\sigma}, \sigma_0) \equiv \min_{\tau \in S_K} \frac{1}{n} \text{Hamming}(\hat{\sigma}, \tau \circ \sigma_0),$$

where  $\text{Hamming}(\cdot, \cdot)$  denotes the Hamming distance. We also write  $d(\cdot, \cdot)$  to denote  $\frac{1}{n}$  times the Hamming distance. In the definition of misclustering error rate, we minimize over the set of permutations  $\tau$  on  $K$  objects because clusterings are identifiable only up to a permutation of their labels. It is important to note that  $\hat{\sigma}$  is a random quantity both because the community recovery algorithm

may be stochastic and because the network  $A$  – the input to the algorithm – is random. Thus, we aim to bound  $l(\hat{\sigma}, \sigma_0)$  in probability.

Zhang and Zhou [32] and Gao et al [10] show that the minimax optimal rate of convergence for the unweighted stochastic block model is of the order  $\exp(-(1+o(1))\frac{nI_{\text{Ber}}}{K})$ .  $I_{\text{Ber}} = -2 \log \sqrt{P_0 Q_0} + \sqrt{(1-P_0)(1-Q_0)}$  is the Renyi divergence of order 1/2 between  $\text{Ber}(P_0)$  and  $\text{Ber}(Q_0)$ , where  $P_0, Q_0$  are the probabilities of absence for within-community and between-communities edges. Yun and Proutiere have also characterized, though they present the results differently, the optimal rate of convergence for the labeled stochastic block model. Our work extends these results to the weighted SBM and show that the optimal rate is again governed by a Renyi divergence.

Although Renyi divergence is of central importance in homogenous stochastic block model where the cluster sizes are approximately balanced, it is important to note that, in the case of cluster imbalance or in the case of *heterogenous* stochastic block model, Abbe and Sandon [3] and Yun and Proutiere [31] have shown that an information divergence that generalizes the Renyi is what drives the intrinsic difficulty of community recovery – a generalization that is referred to as the CH-divergence.

### 2.2.2 Other notions of recovery

A closely related problem is that of finding the exact recovery threshold. We say that the weighted stochastic block model has an exact recovery threshold if there is some function of the parameters  $\theta(P_0, Q_0, p(x), q(x), K, \beta, n)$  such that exact recovery is asymptotically almost always impossible if  $\theta < 1$  and almost always possible if  $\theta > 1$ . For the homogeneous unweighted stochastic block model, Abbe et al [2] have shown that, when  $\beta = 1, K = 2, 1 - P_0 = \frac{a \log n}{n}$ , and  $1 - Q_0 = \frac{b \log n}{n}$  (that is, the average degree is of order  $\log n$ ) for some constant  $a, b$ , then the exact threshold is  $\sqrt{a} - \sqrt{b}$ , that is, no exact recovery algorithm can succeed if  $\sqrt{a} - \sqrt{b} < 1$  and there exists a recovery algorithm that can succeed with probability tending to one if  $\sqrt{a} - \sqrt{b} > 1$ . This result was generalized by Zhang and Zhou [32] beyond the  $\log n$  degree setting where  $\frac{nI_{\text{Ber}}}{K \log n}$  was shown to be the threshold. Apart from exact recovery (also known as strong consistency) and weak recovery, a notion of partial recovery (also known as weak consistency) has also been considered [5, 21, 32]. This notion lies between the other two notions of recovery, and only requires the fraction of misclassified nodes to converge in probability to 0 as  $n$  becomes large. A very general result for the  $K = 2$  case, characterizing when exact and partial recovery are possible for the unweighted homogeneous stochastic block model, is provided in Mossel et al. [21].

## 3 Recovery algorithm

The weighted stochastic block model presents an extra layer of difficulty because the densities  $p(x)$  and  $q(x)$  are unknown: One consequence of not knowing  $p(x)$  and  $q(x)$  is that the MLE does not exist. To see this, first consider the MLE for the usual stochastic block model:

$$\hat{\sigma}_{MLE}^{SBM} = \arg \max_{\sigma} \sum_{\substack{(u,v) \in E \\ \sigma(u)=\sigma(v)}} \log \frac{p(1-q)}{q(1-p)}.$$

Since  $\log \frac{p(1-q)}{q(1-p)} > 0$ , the estimator  $\hat{\sigma}_{MLE}^{SBM}$  may be computed by searching for the clustering that maximizes the number of within-cluster edges. In contrast, likelihood maximization for the weighted SBM takes the form

$$\sup_{\sigma, p(x), q(x) \in \mathcal{P}} \sum_{\substack{(u,v) \in E \\ \sigma(u)=\sigma(v)}} \log \frac{p(A_{uv})(1-Q_0)}{q(A_{uv})(1-P_0)},$$

where  $\mathcal{P}$  is the set of all densities. The maximum does not exist because the maximizer of the likelihood does not exist for nonparametric density estimation. This remains true even if we restrict  $\mathcal{P}$  to be the set of all smooth densities with, say, bounded second derivatives. Our approach is therefore appreciably different and consists of combining the idea of discretization from nonparametric density estimation with clustering techniques for the unweighted stochastic block model.

### 3.1 Algorithm overview

We now outline the main components of our algorithm. The key ideas are to convert the edge weights into a finite set of labels by discretization, and then cluster nodes on the labeled network. We first provide a broad overview of our algorithm and then describe each step in detail. Given a weighted network represented as an adjacency matrix  $A$ , our estimation method has four steps. We summarize the flow of the algorithm below and also in Figure 3:

1. **Transformation & discretization.** We take as input a weighted matrix  $A$  and apply an invertible transformation function  $\Phi : \mathbb{R} \rightarrow [0, 1]$  to obtain a matrix  $\Phi(A)$  with weights between 0 and 1.

Next, we divide the  $[0, 1]$  interval into  $l = 1, \dots, L$  equally-spaced subintervals, which we call bins. We replace the real-valued weight entries  $\Phi(A)$  with a categorical label  $l \in \{1, \dots, L\}$ :  $[\Phi(A)]_{uv}$  is assigned label  $l$  if the value  $[\Phi(A)]_{uv}$  falls into bin  $l$ . We output a network with each edge assigned one of  $L$  possible colors. We continue to denote the adjacency matrix by  $A$ .

2. **Add noise.** For a fixed constant  $c > 0$ , let  $\delta = \frac{cL}{n}$ . We perform the following process on every edge of the labeled graph, independently of other edges: With probability  $1 - \delta$ , keep an edge as it is, and with probability  $\delta$ , erase the edge and replace it with an edge with label uniformly drawn from the set of  $L$  labels. Again, we continue to denote the modified adjacency matrix as  $A$ .
3. **Initialization Parts 1 & 2.** For each color  $l$ , we create a sub-network by including only edges of color  $l$ . For each sub-network, we perform spectral clustering. We output  $l^*$ , the color that induces the maximally separated spectral clustering.

Let  $A_{l^*}$  be the adjacency matrix for color  $l^*$ . For each  $u \in \{1, \dots, n\}$ , we perform spectral clustering on  $A_{l^*} \setminus \{u\}$ , which denotes the adjacency matrix with vertex  $u$  removed. We output  $n$  clusters  $\tilde{\sigma}_1, \dots, \tilde{\sigma}_n$ , where  $\tilde{\sigma}_u$  is a clustering on  $\{1, 2, \dots, n\} \setminus \{u\}$ , for  $1 \leq u \leq n$ .

4. **Refinement & consensus.** From each  $\tilde{\sigma}_u$ , we generate a clustering  $\hat{\sigma}_u$  on  $\{1, 2, \dots, n\}$  which retains the assignments specified by  $\tilde{\sigma}_u$  for  $\{1, 2, \dots, n\} \setminus \{u\}$ , and assigns  $\hat{\sigma}_u(u)$  by maximizing the likelihood taking into account only the neighborhood around  $u$ .

We then align the cluster assignments made in the previous step.

### 3.2 Transformation and discretization

These two steps are straightforward: In the transformation step, we apply an invertible CDF function  $\Phi : \mathbb{R} \rightarrow [0, 1]$  as the transformation function onto all the edge weights so that each entry of  $\Phi(A)$  lies in the interval  $[0, 1]$ . In the discretization step, we divide the interval  $[0, 1]$  into  $L$  equally-spaced bins of the form  $[a_l, b_l]$ , where  $a_1 = 0, b_L = 1$ , and  $b_l - a_l = \frac{1}{L}$ . An edge is assigned the label  $l$  if the weight of that edge lies in bin  $l$ .

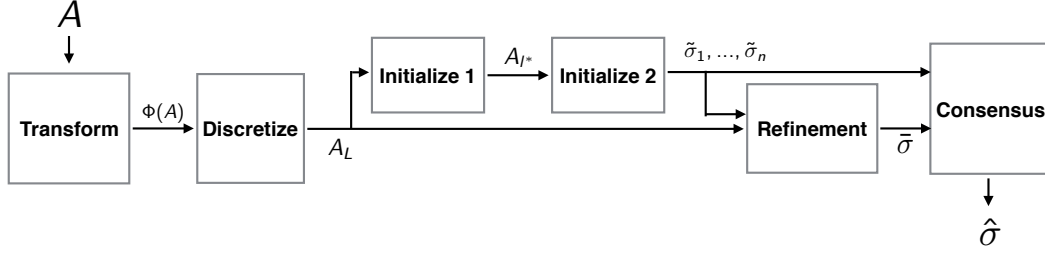


Figure 2: Add a box indicating the add noise step. Pipeline for the our proposed algorithm

---

**Algorithm 3.1** Transformation and Discretization

---

**Input:** A weighted network  $A$ , a positive integer  $L$ , and an invertible function  $\Phi : \mathbb{R} \rightarrow [0, 1]$ .

**Output:** A labeled network  $A$  with  $L$  labels

---

Divide  $[0, 1]$  into  $L$  bins, labeled  $Bin_1, \dots, Bin_L$ .

**for** every edge  $(u, v)$  **do**

    let  $l$  be the bin in which  $\Phi(A_{uv})$  falls.

    Give the edge  $(u, v)$  the label  $l$  in the labeled network  $A$

**end for**

Output  $A$

---

### 3.3 Add noise

This part of the algorithm is required for technical reasons. As detailed in the proof of Proposition 5.1 in Appendix A, deliberately forming a noisy version of the graph barely affects the separation between the distributions specifying within- and between-community edge labels, but has the desirable effect of ensuring that all edge labels occur with probability at least  $\frac{c}{n}$ . This property is crucial to our analysis in subsequent steps of the recovery algorithm.

In the description of the algorithm below, we treat the label 0 (i.e., an empty edge) as a separate label, so we have a network with a total of  $L + 1$  labels.

---

**Algorithm 3.2** Add noise

---

**Input:** A labeled network with  $L + 1$  labels and a constant  $c$

**Output:** A labeled network  $A$  with  $L + 1$  labels

---

**for** every edge  $(u, v)$  **do**

    With probability  $1 - \frac{c(L+1)}{n}$ , do nothing. With probability  $\frac{c(L+1)}{n}$  replace edge label with a label drawn uniformly at random from  $\{0, 1, 2, \dots, L\}$

**end for**

Output  $A$

---

### 3.4 Initialization

The initialization procedure takes as input a network with edges labeled  $\{1, \dots, L\}$ . The goal of the initialization procedure is to create a rough clustering  $\tilde{\sigma}$  that is suboptimal but still consistent.

As outlined in Algorithm 3.3, the rough clustering is based on a single label  $l^*$ , selected based on the maximum value of the estimated Renyi divergence between within-community and between-community distributions for the unweighted SBMs based on individual labels.

For technical reasons, we actually create  $n$  separate rough clusterings  $\{\tilde{\sigma}_u\}_{u=1,\dots,n}$ , where each  $\tilde{\sigma}_u : [n-1] \rightarrow [K]$  is a clustering of a network of  $n-1$  nodes where node  $u$  has been removed. The clusterings  $\{\tilde{\sigma}_u\}$  will later be combined into a single clustering algorithm.

---

**Algorithm 3.3** Initialization

---

**Input:** A labeled network  $A$  with  $L$  labels

**Output:** A set of clusterings  $\{\tilde{\sigma}_u\}_{u=1,\dots,n}$ , where  $\tilde{\sigma}_u$  is a clustering on  $\{1, 2, \dots, n\} \setminus \{u\}$

---

- 1: Separate  $A_L$  into  $L$  networks  $\{A_l\}_{l=1,\dots,L}$  where  $A_l$  contains only edges with label  $l$ .  $\triangleright$  Stage 1
  - 2: **for** each label  $l$  **do**
  - 3:   Compute  $\bar{d} = \frac{1}{n} \sum_{u=1}^n d_u$  as the average degree.
  - 4:   Perform spectral clustering with  $\tau = \bar{d}$  and  $\mu \geq C\beta$  to get  $\tilde{\sigma}_l$ , where  $C$  is an appropriately chosen large constant
  - 5:   estimate  $\hat{P}_l = \frac{\sum_{u \neq v : \tilde{\sigma}_l(u) = \tilde{\sigma}_l(v)} (A_l)_{uv}}{|\{u \neq v : \tilde{\sigma}_l(u) = \tilde{\sigma}_l(v)\}|}$  and  $\hat{Q}_l = \frac{\sum_{u \neq v : \tilde{\sigma}_l(u) \neq \tilde{\sigma}_l(v)} (A_l)_{uv}}{|\{u \neq v : \tilde{\sigma}_l(u) \neq \tilde{\sigma}_l(v)\}|}$ .
  - 6:    $\hat{I}_l \leftarrow \frac{(\hat{P}_l - \hat{Q}_l)^2}{\hat{P}_l \vee \hat{Q}_l}$
  - 7: **end for**
  - 8: Choose  $l^* = \arg \max_l \hat{I}_l$ . Let  $A_{l^*}$  be the network with only edges labeled  $l^*$
  - 9: **for** each node  $u$  **do**  $\triangleright$  Stage 2
  - 10:   Create network  $A_{l^*} \setminus \{u\}$  by removing node  $u$  from  $A_{l^*}$
  - 11:   Perform SPECTRAL CLUSTERING on  $A_{l^*} \setminus \{u\}$  to get  $\tilde{\sigma}_u$
  - 12: **end for**
  - 13: Output the set of clusterings  $\{\tilde{\sigma}_u\}_{u=1,\dots,n}$ .
- 

**Spectral clustering:** Note that Algorithm 3.3 involves several applications of SPECTRAL CLUSTERING. We describe the spectral clustering algorithm used as a subroutine in Algorithm 3.4 below:

---

**Algorithm 3.4** SPECTRAL CLUSTERING

---

**Input:** An unweighted network  $A$ , trim threshold  $\tau$ , number of communities  $K$ , tuning parameter  $\mu$

**Output:** A clustering  $\sigma$

---

- 1: For each node  $u$  whose degree  $d_u \geq \tau$ , set  $A_{uv} = 0$  to get  $T_\tau(A)$
  - 2: Let  $\hat{A}$  be the best rank- $K$  approximation to  $T_\tau(A)$  in spectral norm, formed by truncating the SVD
  - 3: For each node  $u$ , define the neighbor set  $N(u) = \{v : \|\hat{A}_u - \hat{A}_v\|_2^2 \leq \mu K^2 \frac{\bar{d}}{n}\}$
  - 4: Initialize  $S \leftarrow \emptyset$ . Select node  $u$  with the most neighbors and add  $u$  into  $S$  as  $S[1]$
  - 5: **for**  $i = 2, \dots, K$  **do**
  - 6:   Among all  $u$  such that  $|N(u)| \geq \frac{n}{\mu K}$ , select  $u^* = \arg \max_u \min_{v \in S} \|\hat{A}_u - \hat{A}_v\|_2$
  - 7:   Add  $u^*$  into  $S$  as  $S[i]$ .
  - 8: **end for**
  - 9: **for**  $u = 1, \dots, n$  **do**
  - 10:   Take  $\arg \min_i \|\hat{A}_u - \hat{A}_{S[i]}\|_2$  and assign  $\sigma(u) = i$
  - 11: **end for**
-



Importantly, note that we may always choose the parameter  $\mu$  sufficiently large such that Algorithm 3.4 generates a set  $S$  with  $|S| = K$ .

### 3.5 Refinement and consensus

Our refinement and consensus steps closely follow the method described by Gao et al. [10]. In the refinement step, we use the set of initial clusterings  $\{\tilde{\sigma}_u\}_{u=1,\dots,n}$  to generate a more accurate clustering for the labeled network by locally maximizing an approximate log-likelihood expression for each of the nodes  $u = 1, \dots, n$ . The consensus step then resolves a cluster label consistency problem arising after the refinement stage.

---

#### Algorithm 3.5 Refinement

---

**Input:** A labeled network  $A$  and a set of rough clusterings  $\{\tilde{\sigma}_u\}_{u=1,\dots,n}$ , where  $\tilde{\sigma}_u$  is a clustering on the set  $\{1, 2, \dots, n\} \setminus \{u\}$  for each  $u$

**Output:** a clustering  $\hat{\sigma}$  over the whole network

- 1: **for** each node  $u$  **do**
- 2:     Estimate  $\{\hat{P}_l, \hat{Q}_l\}_{l=0,\dots,L}$  from  $\tilde{\sigma}_u$
- 3:     Let  $\hat{\sigma}_u : [n] \rightarrow [K]$  where  $\hat{\sigma}_u(v) = \tilde{\sigma}_u(v)$  for all  $v \neq u$  and

$$\hat{\sigma}_u(u) = \arg \max_k \sum_{v : \tilde{\sigma}_u(v)=k, v \neq u} \sum_l \log \frac{\hat{P}_l}{\hat{Q}_l} \mathbf{1}(A_{uv} = l)$$

- 4: **end for**
- 5: Let  $\hat{\sigma}(1) = \hat{\sigma}_1(1)$  ▷ Consensus Stage
- 6: **for** each node  $u \neq 1$  **do**

$$\hat{\sigma}(u) = \arg \max_k |\{v : \hat{\sigma}_1(v) = k\} \cap \{v : \hat{\sigma}_u(v) = \hat{\sigma}_u(u)\}|$$

- 7: **end for**
  - 8: Output  $\hat{\sigma}$
- 

## 4 Analysis of misclustering error

For an unweighted stochastic block model, the key information quantity that governs the threshold behavior is  $I = -2 \log(\sqrt{pq} + \sqrt{(1-p)(1-q)})$ . This is the Renyi divergence of order  $\frac{1}{2}$  between the  $Ber(p)$  and  $Ber(q)$  distributions.

The Renyi divergence of order  $\frac{1}{2}$  is defined on pairs of general measures as

$$I = -2 \log \int \left( \frac{dP}{dQ} \right)^{1/2} dQ.$$

Interestingly, this generalized form of the Renyi divergence is also what governs both the rate of convergence of our proposed algorithm and the threshold behavior of the weighted stochastic block model. In the weighted stochastic block model setting where  $P$  and  $Q$  have continuous parts  $p(x)$  and  $q(x)$  and point masses at zero with probability  $P_0$  and  $Q_0$ , respectively, the Renyi divergence takes the form

$$I = -2 \log \left( \sqrt{P_0 Q_0} + \int \sqrt{(1-P_0)(1-Q_0)p(x)q(x)} dx \right)$$

When  $I \rightarrow 0$ , as is the case in our situation,  $I$  is asymptotically equal to the Hellinger distance:

$$\begin{aligned} I &= \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + \int (\sqrt{(1-P_0)p(x)} - \sqrt{(1-Q_0)q(x)})^2 dx \right\} (1 + o(1)) \\ &= \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + (\sqrt{1-P_0} - \sqrt{1-Q_0})^2 + \sqrt{(1-P_0)(1-Q_0)} \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right\} \\ &\quad \cdot (1 + o(1)) \end{aligned} \tag{4.1}$$

Equation (4.1) demonstrates that the Renyi divergence is driven by both the divergence between the edge probabilities  $1 - P_0$  and  $1 - Q_0$  and the divergence between the densities  $p(x)$  and  $q(x)$ . This is a novel feature of the weighted stochastic block model.

When  $p(x) = q(x)$ , the Renyi divergence  $I$  simplifies to the Renyi divergence of the unweighted SBM with edge probabilities  $P_0$  and  $Q_0$ . This is intuitive because if  $p(x) = q(x)$ , the values of the edge weights provide no additional information about the cluster structure. On the other hand, when  $P_0 = Q_0$ , the first two terms disappear and the Renyi divergence is driven only by the difference between the edge weight densities  $p(x)$  and  $q(x)$ . This is also intuitive because if  $P_0 = Q_0$ , the presence or absence of an edge offers no information about the cluster structure. On the other hand, smaller values of  $P_0$  (e.g., more dense graphs) correspond to larger values of  $I$ , since observing the values of more edge weights provides us with greater ability to distinguish clusters.

For the remainder of the paper, we denote  $H := \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$ . Note that  $H \leq 2$ . We will be interested in the regimes where  $H = o(1)$  and  $H = \Theta(1)$ , but  $(1 - P_0)(1 - Q_0) = o(1)$ . Both of these are important and nontrivial cases to consider: in the first case, the challenge is to distinguish two densities  $p(x)$  and  $q(x)$  that are becoming increasingly similar; in the second case, the challenge is to estimate the edge weight densities well when the graph may be very sparse. The algorithm we propose in Section 3 is suitable for both of these settings, but the theoretical analysis in each case is quite different.

## 4.1 Rate of convergence

Our analysis is asymptotic. We characterize the performance of the algorithm as  $n \rightarrow \infty$ . In our analysis, we treat  $p(x), q(x), P_0$ , and  $Q_0$  all as varying with  $n$ , but we leave the dependence on  $n$  implicit. All of our results will use the following assumption:

**Assumption A0:** There exist absolute constants  $c_0$  and  $C_0$  such that  $c_0 \leq \frac{1-P_0}{1-Q_0} \leq C_0$ . If  $P_0 \vee Q_0 > 0$ , we also assume  $c_0 \leq \frac{P_0}{Q_0} \leq C_0$ .

Assumption A0 states that the density of within-community edges has the same order as the density of between-community edges. This assumption is standard in the existing literature on unweighted stochastic block models.

**Definition 4.1.** Let  $S$  be either  $\mathbb{R}$  or  $\mathbb{R}^+$ . We say that  $\Phi : S \rightarrow [0, 1]$  is a *transformation function* if it is a differentiable bijection (hence a cumulative distribution function), and  $\phi = \Phi'$  satisfies (a)  $\int \phi(x)^s dx < \infty$  for any  $0 < s < 1$ , and (b)  $\left| \frac{\phi'(x)}{\phi(x)} \right|$  is bounded.

The transformation function  $\Phi$  induces a probability measure on  $S$ , and we let  $\Phi\{\cdot\}$  denote the  $\Phi$ -measure of a set. The additional regularity conditions stated below stipulate that  $p(x)$  and  $q(x)$  are smooth and the likelihood ratio  $\frac{p(x)}{q(x)}$  is well-behaved. Furthermore, the distribution  $\phi$  must be heavier-tailed than  $p(x)$  and  $q(x)$ . We state different sets of assumptions for the cases  $H = o(1)$  and  $H = \Theta(1)$ .

#### 4.1.1 The case $H = o(1)$

We state the required assumptions and then present our main result.

**Definition 4.2.** We call a nonnegative function  $f(x)$  is  $(c_{s1}, c_{s2}, C_s)$ -*bowl-shaped* if  $f(x)$  is nonincreasing for all  $x \leq c_{s1}$ , nondecreasing for all  $x \geq c_{s2}$ , and bounded by  $C_s$  for all  $x \in [c_{s1}, c_{s2}]$ .

**Regularity conditions:**

A1  $p(x), q(x) > 0$  on the interior of  $S$ ,  $\sup_x \{p(x) \vee q(x)\} < \infty$ , and  $\inf_{x \in S} \left\{ \frac{p(x) \vee q(x)}{\phi(x)} \right\} < \infty$ .

A2 There exists a subinterval  $R$  of  $S$  such that

- (a)  $\frac{1}{\rho} \leq \frac{p(x)}{q(x)} \leq \rho$  for  $x \in R$ , where  $\rho$  is an absolute constant, and
- (b)  $\Phi\{R^c\} = o(H)$ .

A3 Denoting  $\alpha^2 = \int_R q(x) \left( \frac{p(x)-q(x)}{q(x)} \right)^2 dx$  and  $\gamma(x) = \frac{q(x)-p(x)}{\alpha}$ , we have

$$\int_R q(x) \left| \frac{\gamma(x)}{q(x)} \right|^r dx < \infty,$$

for an absolute constant  $r > 4$ .

A4 There exists a function  $h(x)$  such that

- (a)  $h(x) \geq \max \left\{ \left| \frac{\gamma'(x)}{q(x)} \right|, \left| \frac{q'(x)}{q(x)} \right|, \left| \frac{\gamma(x)}{q(x)} \right| \right\}$ ,
- (b)  $h(x)$  is  $(c_{s1}, c_{s2}, C_s)$ -bowl-shaped for absolute constants  $c_{s1}, c_{s2}$ , and  $C_s$ , and
- (c) for an absolute constant  $t$  such that  $\frac{4}{r} < 2t < 1$ , we have

$$\int_R |h(x)|^{2t} \phi(x) dx < \infty.$$

A5  $(\log p)'(x), (\log q)'(x) \geq (\log \phi)'(x)$  for all  $x \leq c_{s1}$ , and  $(\log p)'(x), (\log q)'(x) \leq (\log \phi)'(x)$  for all  $x \geq c_{s2}$ .

The simplest setting for which the assumptions are satisfied is when  $p(x)$  and  $q(x)$  have compact support (so the transformation  $\Phi$  is not even necessary), bounded first derivatives, likelihood ratio  $\frac{p(x)}{q(x)}$  bounded away from 0 and infinity, and uniform convergence of  $|p(x) - q(x)| \rightarrow 0$ . However, this simple setting excludes many interesting cases, such as when  $p(x)$  and  $q(x)$  are Gaussian. To include such cases (cf. Section 4.1.4 below), we impose the more technical set of conditions.

We have the following result:

**Theorem 4.1.** Suppose  $\hat{\sigma}$  is the output of the algorithm in Section 3 with transformation  $\Phi$  and discretization level  $L$  chosen such that  $L \rightarrow \infty$ ,  $L = o(\frac{1}{H})$ , and  $L = o(nI)$ . Suppose  $P_0, Q_0$  satisfy Assumption A0 and  $p(x)$  and  $q(x)$  satisfy Assumptions A1–A5 with respect to  $\Phi$ . Also suppose  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = o(1)$ , the number of communities  $K$  is fixed, and  $I = o(1)$ . Then

$$\lim_{n \rightarrow \infty} P \left\{ l(\hat{\sigma}, \sigma_0) \leq \exp \left( -\frac{nI}{\beta K} (1 + o(1)) \right) \right\} \rightarrow 1.$$

The proof of Theorem 4.1 is outlined in Appendix F.1.

#### 4.1.2 The case $H = \Theta(1)$

We begin by stating the required assumptions.

**Regularity conditions:**

A1'  $p(x), q(x) > 0$  on the interior of  $S$ ,  $\sup_x \{p(x) \vee q(x)\} < \infty$ , and  $\inf_{x \in S} \left\{ \frac{p(x) \vee q(x)}{\phi(x)} \right\} < \infty$ .

A2' For an absolute constant  $r \geq 8$ ,  $\int \left| \log \frac{p(x)}{q(x)} \right|^r \phi(x) dx < \infty$ .

A3' There exists a function  $h(x)$  such that

- (a)  $h(x) \geq \max \left\{ \left| \frac{q'(x)}{q(x)} \right|, \left| \frac{p'(x)}{p(x)} \right| \right\}$ ,
- (b)  $h(x)$  is  $(c_{s1}, c_{s2}, C_s)$ -bowl-shaped for some absolute constants  $c_{s1}, c_{s2}$ , and  $C_s$ , and
- (c) For some constant  $t$  such that  $\frac{4}{r} < 2t < 1$ , we have

$$\sup_n \int |h_n(x)|^{2t} \phi(x) dx < \infty.$$

A4'  $(\log p)'(x), (\log q)'(x) \geq (\log \phi)'(x)$  for all  $x < c_{s1}$ , and  $(\log p)'(x), (\log q)'(x) \leq (\log \phi)'(x)$  for all  $x > c_{s2}$ .

These assumptions are similar in nature to Assumptions A1–A5, and one can show that the examples in Section 4.1.4 also satisfy Assumptions A1'–A4'. Our main result is the following:

**Theorem 4.2.** *Suppose  $\hat{\sigma}$  is the output of the algorithm in section 3 with transformation  $\Phi$  and discretization level  $L$  chosen such that  $L \rightarrow \infty$  and  $\frac{nI}{L \exp(L^{1/r})} \rightarrow \infty$ . Suppose that  $P_0$  and  $Q_0$  satisfy Assumption A0 and  $p(x)$  and  $q(x)$  satisfy Assumptions A1'–A4' with respect to  $\Phi$ . Also suppose  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \Theta(1)$ , the number of communities  $K$  is fixed, and  $I = o(1)$ . Then*

$$\lim_{n \rightarrow \infty} P \left\{ l(\hat{\sigma}, \sigma_0) \leq \exp \left( -\frac{nI}{\beta K} (1 + o(1)) \right) \right\} \rightarrow 1.$$

For the proof of Theorem 4.2, see Appendix F.2.

### 4.1.3 Additional discussion of assumptions

It is crucial to note that our algorithm does not require any prior knowledge about the form of  $p(x)$  and  $q(x)$ : The same algorithm and guarantees apply whether  $p(x)$  and  $q(x)$  are Gaussian, Laplace, or any other (possibly nonparametric) distributions, as long as they satisfy Assumptions A1–A5 with respect to the transformation function  $\Phi$ .

To aid the reader, we now provide a brief, non-technical interpretation of the regularity conditions described above.

#### Interpretation of Assumptions A1–A5:

- A1 Assumption A1 is simple; the second part states that  $\phi$  must have a tail at least as heavy as that of  $p(x)$  and  $q(x)$ .
- A2 In Assumption A2, we require that the likelihood ratio  $\frac{p(x)}{q(x)}$  be bounded away from 0 and  $\infty$  except on a region  $R^c \subseteq \mathbb{R}$ . Since  $H \rightarrow 0$ , we the densities  $p(x)$  and  $q(x)$  are becoming increasingly similar and  $R^c$  is shrinking. We require that the measure of  $R^c$ , with respect to  $\Phi$ , shrinks faster than  $H$ . This condition intuitively states that  $\left| \frac{p(x)}{q(x)} \right|$  and its reciprocal tend to infinity slowly with respect to  $x$ . For heavier-tailed functions  $\Phi$ , Assumption A2 is a stronger condition on  $\frac{p(x)}{q(x)}$ ; for lighter-tailed functions  $\Phi$ , Assumption A2 is a looser condition.

A3 In Assumption A3, note that  $H \rightarrow 0$  implies  $\alpha \rightarrow 0$ , as well, so  $\gamma(x) = \frac{p(x)-q(x)}{\alpha}$  is a function of constant order. The integrability condition on  $\gamma(x)$  states that  $|p(x) - q(x)|$  must converge to 0 almost uniformly for all  $x$  in the region  $R$ . (Having an  $L_\infty$ -bound on  $\gamma$  would imply uniform convergence.)

A4 Assumption A4 imposes smoothness on  $q(x)$  and  $\gamma(x)$ . The second part of Assumption A4 is a weak condition ensuring that  $h(x)$  does not oscillate with infinite frequency.

A5 Assumption A5 is again a condition stating that  $\phi$  has tails at least as heavy as  $p(x)$  and  $q(x)$ .

Note that an analogous interpretation may be used to describe to the conditions A1'-A4'.

An alternative way to interpret these assumptions is that, for a given transformation  $\Phi$ , there is a space  $\mathcal{P}_\Phi(\rho, c_{s1}, c_{s2}, C_s, r, t)$  of densities satisfying Assumptions A1–A5. We again note that  $\mathcal{P}_\Phi$  is actually a sequence of function spaces indexed by  $n$ , although the dependence is implicit in our notation. For a given  $\Phi$ , Assumptions A1–A5 impose a set of constraints on the densities  $p(x)$  and  $q(x)$ .

**Assumptions for parametric families:** When  $p(x)$  and  $q(x)$  belong to a parametric family, as in the examples discussed in Section 4.1.4, it is helpful to consider a simpler set of assumptions. Suppose  $p(x) = \exp(f_{\theta_1}(x))$  and  $q(x) = \exp(f_{\theta_0}(x))$ , where  $\{f_\theta(x)\}$  is a set of functions indexed by  $\theta \in \Theta \subseteq \mathbb{R}^{d_\Theta}$  and  $\Theta$  is compact.

Consider the following conditions:

B1  $\inf_x \{\log \phi(x) - f_\theta(x)\} > -\infty$  for all  $\theta \in \Theta$ .

B2 The Fisher information matrix  $G_\theta := \int (\nabla_\theta f_\theta(x))(\nabla_\theta f_\theta(x))^\top \exp(f_\theta(x)) dx$  is full-rank:

$$0 < c_{\min} < \inf_{\theta \in \Theta} \lambda_{\min}(G_\theta) \leq \sup_{\theta \in \Theta} \lambda_{\max}(G_\theta) < c_{\max} < \infty.$$

B3 There exist  $g_1(x) \geq \sup_{\theta \in \Theta} \|\nabla_\theta f_\theta(x)\|$  and  $g_{2,\theta}(x) \geq \max\{\|\nabla_\theta f'_\theta(x)\|, |f'_\theta(x)|\}$  such that  $g_1$  and  $g_{2,\theta}$  are  $(c_{s1}, c_{s2}, \widetilde{C}_s)$ -bowl-shaped, and

$$\int g_1(x)^r \phi(x) dx < \infty, \quad \text{and} \quad \sup_{\theta \in \Theta} \int g_{2,\theta}(x)^{4t} \phi(x) dx < \infty,$$

where  $t$  and  $r$  are constants satisfying  $\frac{8}{r} \leq 4t < 1$ .

B4  $\inf_{\theta \in \Theta} f'_\theta(x) \geq (\log \phi)'(x)$  for all  $x \leq c_{s1}$ , and  $\inf_{\theta \in \Theta} f'_\theta(x) \leq (\log \phi)'(x)$  for all  $x \geq c_{s2}$ .

We then have the following result, proved in Appendix G:

**Proposition 4.1.** *Suppose Assumptions B1–B5 hold.*

(a) *If  $\|\theta_1 - \theta_0\| \rightarrow 0$ , then  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \rightarrow 0$  and Assumptions A1–A5 are satisfied.*

(b) *If  $\|\theta_1 - \theta_0\| = \Theta(1)$ , then  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \Theta(1)$  and Assumptions A1'–A4' are satisfied.*

#### 4.1.4 Examples

Since Assumptions A1–A5 and A1'–A4' are rather technical, we illustrate them with concrete examples. Although we do not in general require  $p(x)$  and  $q(x)$  to belong to a parametric family, we

will discuss cases where  $p(x) = \exp(f_{\theta_1}(x))$  and  $q(x) = \exp(f_{\theta_0}(x))$ , where  $f_{\theta}(x)$  is a set of functions indexed by  $\theta \in \Theta \subset \mathbb{R}^{d_{\Theta}}$  and  $\Theta$  is compact. Generally speaking, when  $p(x)$  and  $q(x)$  have subexponential tails, we take

$$\phi(x) = \frac{e^{1-\sqrt{x+1}}}{4}, \quad \text{when } S = \mathbb{R}^+, \quad (4.2)$$

$$\phi(x) = \frac{e^{1-\sqrt{|x|+1}}}{8}, \quad \text{when } S = \mathbb{R}. \quad (4.3)$$

These functions  $\phi$  are similar to a generalized normal density, modified so that  $\left| \frac{\phi'(x)}{\phi(x)} \right|$  is bounded. It is easy to verify that  $\Phi(x) = \int_0^x \phi(t)dt$  (respectively,  $\Phi(x) = \int_{-\infty}^x \phi(t)dt$ ) is a valid transformation function.

**Example 4.1.** (Location-scale family over  $\mathbb{R}$ )

Let  $\exp(f(x))$  be a positive base density over  $\mathbb{R}$ . We can define a location-scale family parametrized by  $\mu$  and  $\sigma$  as  $\exp(f(\frac{x-\mu}{\sigma}))$ , and let

$$p(x) = \frac{1}{\sigma_1} \exp\left(f\left(\frac{x-\mu_1}{\sigma_1}\right)\right), \quad \text{and} \quad q(x) = \frac{1}{\sigma_0} \exp\left(f\left(\frac{x-\mu_0}{\sigma_0}\right)\right).$$

Define  $\theta = (\mu, \sigma)$  and  $\Theta = [-C_{\mu}, C_{\mu}] \times [\frac{1}{c_{\sigma}}, c_{\sigma}]$  for some constants  $C_{\mu}$  and  $c_{\sigma}$ , and let

$$f_{\theta} = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right).$$

Then  $p(x), q(x) \in \{f_{\theta}(x)\}_{\theta \in \Theta}$ .

If  $f(x)$  satisfies the conditions

- (a)  $|f^{(k)}(x)|$  is bounded for some  $k \geq 2$ , and
- (b) there exist absolute constants  $c$  and  $M$  such that  $f'(x) > M$  for  $x < -c$  and  $f'(x) < -M$  for  $x > c$ ,

then  $\{f_{\theta}\}_{\theta \in \Theta}$  satisfy Assumptions B1–B4 (hence,  $p(x)$  and  $q(x)$  satisfy Assumptions A1–A5 if  $H \rightarrow 0$ , and Assumptions A1' – A4' if  $H = \Theta(1)$ ), when  $\phi$  is chosen according to equation (4.3). Details are provided in Appendix G.2.

The above assumptions on  $f(x)$  are not strong, and are satisfied for **Gaussian location-scale families**, where the base density is the standard Gaussian density with

$$f(x) = -x^2 - \frac{1}{2} \log 2\pi,$$

and **Laplace location-scale families**, where the base density is the standard Laplace density with

$$f(x) = -|x| - \log 2.$$

**Example 4.2.** (Scale family over  $\mathbb{R}^+$ )

If the base density  $\exp(f(x))$  is supported and positive on  $\mathbb{R}^+$ , we can define a scale family parametrized by  $\sigma$  as  $\exp(f(\frac{x}{\sigma}))$ . Let

$$p(x) = \frac{1}{\sigma_1} \exp\left(f\left(\frac{x}{\sigma_1}\right)\right), \quad \text{and} \quad q(x) = \frac{1}{\sigma_0} \exp\left(f\left(\frac{x}{\sigma_0}\right)\right).$$

Define  $\theta = (\sigma)$  and  $\Theta = \left[\frac{1}{c_\sigma}, c_\sigma\right]$  for some constant  $c_\sigma$ , and let

$$f_\theta = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right).$$

Then  $p(x), q(x) \in \{f_\theta(x)\}_{\theta \in \Theta}$ .

Again, if  $f(x)$  satisfy conditions (a) and (b) from Example 4.1 above, then  $\{f_\theta\}_{\theta \in \Theta}$  satisfy Assumptions B1–B4 (hence  $p(x)$  and  $q(x)$  satisfy Assumptions A1–A5 if  $H \rightarrow 0$ , and Assumptions A1’–A4’ if  $H = \Theta(1)$ ), when  $\phi$  is chosen according to equation (4.2).

**Example 4.3.** (Gamma distribution)

Let

$$p(x) = \frac{\beta_1^{\alpha_1}}{\Gamma(\alpha_1)} x^{\alpha_1-1} e^{-\beta_1 x}, \quad \text{and} \quad q(x) = \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} x^{\alpha_0-1} e^{-\beta_0 x},$$

defined over  $\mathbb{R}^+$ , where  $\alpha_0, \alpha_1, \beta_0, \beta_1 > 0$ .

Let  $\theta = (\alpha, \beta)$  and  $\Theta = \left[\frac{1}{C}, C\right]^2$  for some constant  $C$ , and let

$$f_\theta(x) = (\alpha - 1) \log x - \beta x + \alpha \log \beta - \log \Gamma(\alpha).$$

Then  $p(x), q(x) \in \{\exp(f_\theta(x))\}_{\theta \in \Theta}$ . Furthermore, if we choose  $\phi$  as equation (4.2), then  $\{f_\theta\}_{\theta \in \Theta}$  satisfies Assumptions B1–B4.

## 4.2 Lower bound

We now derive a lower bound on the performance of **permutation-equivariant** clustering algorithms for the weighted stochastic block model.

**Definition 4.3.** For an  $n \times n$  matrix  $A$  and a permutation  $\pi \in S_n$ , define  $\pi A$  as an  $n \times n$  matrix such that  $[\pi A]_{uv} = A_{\pi^{-1}(u), \pi^{-1}(v)}$ . In other words,  $\pi A$  is the result of applying  $\pi$  to the rows and columns of  $A$ . Let  $\hat{\sigma}$  be a clustering algorithm; i.e.,  $\hat{\sigma}(A)$  is a clustering  $[n] \rightarrow [K]$  for any  $n \times n$  matrix  $A$ . We say that  $\hat{\sigma}$  is *permutation-equivariant* if, for any  $A$  and any  $\pi \in S_n$ , there exists  $\tau \in S_K$  such that

$$\hat{\sigma}(\pi A) = \tau \circ \hat{\sigma}(A) \circ \pi^{-1}.$$

In particular, if  $\hat{\sigma}$  is permutation-equivariant, we have

$$l(\hat{\sigma}(A) \circ \pi^{-1}, \hat{\sigma}(\pi A)) = d(\tau \circ \hat{\sigma}(A) \circ \pi^{-1}, \hat{\sigma}(\pi A)) = 0,$$

for any  $A$  and any  $\pi$ . Intuitively,  $\hat{\sigma}$  is permutation-equivariant if the output clustering does not depend on node labeling. Permutation-equivariance is a weak restriction—indeed, all existing algorithms for network clustering, including the one proposed in this paper, are permutation-equivariant.

**Theorem 4.3.** Suppose we have  $K$  clusters, of which at least one has size  $\frac{n}{\beta K}$  and at least one has size  $\frac{n}{\beta K} + 1$ , for some constant  $\beta \geq 1$ . Let  $\sigma_0$  denote the true clustering. Suppose  $I \rightarrow 0$  and  $P_0$  and  $Q_0$  satisfy Assumption A0, and let  $p(x)$  and  $q(x)$  be densities satisfying the following condition: If  $H = o(1)$ , then  $\sup_x \left| \log \frac{p(x)}{q(x)} \right| \leq C$  for some constant  $C$ ; and if  $H = \Theta(1)$ , then  $\int p(x) \left| \log \frac{p(x)}{q(x)} \right|^2 dx < \infty$  and  $\int q(x) \left| \log \frac{p(x)}{q(x)} \right|^2 dx < \infty$ . Then any permutation-equivariant algorithm  $\hat{\sigma}$  satisfies the following:

- (i) If  $\frac{nI}{K} \rightarrow \infty$ , then  $\mathbb{E}l(\hat{\sigma}, \sigma_0) \geq \exp(-(1 + o(1)) \frac{nI}{K})$ .
- (ii) If  $\frac{nI}{K} \rightarrow c < \infty$ , for some constant  $c$ , then  $\mathbb{E}l(\hat{\sigma}, \sigma_0) \geq c' > 0$ , for some constant  $c'$ .

The proof of Theorem 4.3 is provided in Appendix H. The proof employs the change-of-measure technique used by Yun and Proutiere to prove a similar lower bound on labeled stochastic block model [31]. We note that Theorem 4.3 applies to any  $p(x)$  and  $q(x)$  that satisfy the assumptions, rather than being a minimax lower bound involving a supremum over a function space.

**Remark 4.1.** It is interesting to observe that Theorem 4.3, in conjunction with Theorem 4.1, shows that one does not have to pay a price for making nonparametric assumptions: Our nonparametric method achieves the optimal rate of recovery even if the densities  $p(x)$  and  $q(x)$  take on a parametric form. This seemingly counterintuitive phenomenon arises because the cost of discretization is reflected in the  $o(1)$  term in the exponent and is thus of lower order.

## 5 Proof sketch: Recovery algorithm

A large portion of the Appendix is devoted to proving that our recovery algorithm succeeds and achieves the optimal error rates. We provide an outline of the proofs here.

We divide our argument into propositions that focus on successive stages of our algorithm. A birds-eye view of our method reveals that it consists of two major components: (1) convert a weighted network into a labeled network, and then (2) run a community recovery algorithm on the labeled network.

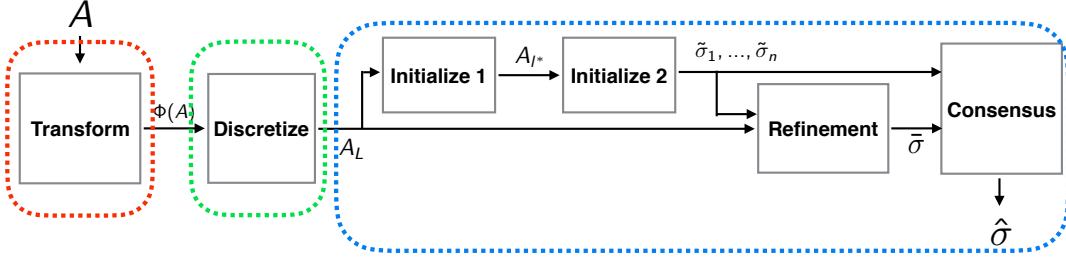


Figure 3: The add noise component also goes into the blue section. Analysis of the right-most blue region is in subsection 5.1, of the middle green region in subsection 5.2, and of the left-most red region in subsection 5.3

### 5.1 Analysis of community recovery on a labeled network

The workhorse behind our algorithm is a subroutine (right-most blue region in Figure 3) for recovering communities on a network where the edges have discrete labels  $l = 1, \dots, L$ . The following proposition characterizes the rate of convergence of the output of the subroutine, where within-community edges are assigned edge labels with probabilities  $\{P_l\}$ , and between-community edges are assigned edge labels according to  $\{Q_l\}$ .

**Proposition 5.1.** *Suppose the edge label probabilities satisfy  $\frac{1}{\rho_L} \leq \frac{P_l}{Q_l} \leq \rho_L$  for a sequence  $\rho_L = \Omega(1)$ . Define  $I_L = -2 \log \sum_{l=0}^L \sqrt{P_l Q_l}$  and suppose  $I_L \rightarrow 0$ . Suppose  $L = \Omega(1)$  and  $\frac{n I_L}{L \rho_L^4} \rightarrow \infty$ . Let  $\hat{\sigma}$  be the output of our algorithm. Then*

$$\lim_{n \rightarrow \infty} P \left( l(\hat{\sigma}, \sigma_0) \leq \exp \left( -\frac{n I_L}{\beta K} (1 + o(1)) \right) \right) \rightarrow 1.$$

Yun and Proutiere [31] have proposed an algorithm for the labeled SBM that achieves the same rate of convergence. Proposition 5.1 is more general than their result, however, in that we allow



both the number of labels  $L$  and the bound  $\rho_L$  on the ratio  $\frac{P_l}{Q_l}$  to diverge to infinity. This extension is critical in analyzing the weighted SBM, since to achieve consistency for continuous distributions, the discretization level  $L$  must increase with  $n$ .

## 5.2 Discretization of the Renyi divergence

The rate of convergence in Proposition 5.1 resembles the expressions in Theorems 4.1 and 4.2, except the Renyi divergence  $I$  is replaced by the discretized Renyi divergence  $I_L$ . Thus, we may derive Theorems 4.1 and 4.2 from Proposition 5.1 by showing that  $|I - I_L|$  is sufficiently small. It is easy to show that  $I_L \leq I$ , since discretization always leads to a loss of information. If  $p(x)$  and  $q(x)$  are sufficiently regular in that they may be well-approximated via discretization, one might expect that  $I_L$  to be not much smaller than  $I$ . Proposition 5.2 and 5.3 rigorizes this notion.

The following proposition is useful for proving Theorem 4.1:

**Proposition 5.2.** *Let  $p(z)$  and  $q(z)$  be two densities supported on  $[0, 1]$ . Suppose  $H = o(1)$ . Let  $L$  be a sequence such that  $L \rightarrow \infty$ . Suppose the following assumptions are satisfied:*

*C1  $p(z), q(z) > 0$  on  $(0, 1)$ , and  $\sup_z \{p(z) \vee q(z)\} < \infty$ .*

*C2 There exists a subinterval  $R \subseteq [0, 1]$  such that*

*(a)  $\frac{1}{\rho} \leq \left| \frac{p(z)}{q(z)} \right| \leq \rho$  for all  $z \in R$ , where  $\rho$  is an absolute constant, and*

*(b)  $\mu\{R^c\} = o(H)$ , where  $\mu$  is the Lebesgue measure.*

*C3 Let  $\alpha^2 = \int_R \frac{(p(z)-q(z))^2}{q(z)} dz$  and  $\gamma(z) = \frac{q(z)-p(z)}{\alpha}$ , and suppose  $\int_R q(z) \left| \frac{\gamma(z)}{q(z)} \right|^r dz < \infty$  for an absolute constant  $r \geq 4$ .*

*C4 There exists  $h(z)$  such that*

*(a)  $h(z) \geq \max \left\{ \left| \frac{\gamma'(z)}{q(z)} \right|, \left| \frac{q'(z)}{q(z)} \right| \right\}$ , and*

*(b)  $h(z)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped for absolute constants  $c'_{s1}, c'_{s2}$ , and  $C'_s$ , and*

*(c)  $\int_R |h(z)|^t dz < \infty$  for an absolute constant  $\frac{2}{r} < t < 1$ .*

*C5  $p'(z), q'(z) \geq 0$  for all  $z < c'_{s1}$ , and  $p'(z), q'(z) \leq 0$  for all  $z > c'_{s2}$ .*

*Suppose  $\frac{1}{c_0} \leq \frac{1-P_0}{1-Q_0} \leq c_0$ . Let  $\text{bin}_l = [a_l, b_l]$ , for  $l = 1, \dots, L$ , be a uniformly-spaced binning of  $[0, 1]$ , and let  $P_l = (1 - P_0) \int_{a_l}^{b_l} p(z) dz$  and  $Q_l = (1 - Q_0) \int_{a_l}^{b_l} q(z) dz$ . Suppose  $L \rightarrow \infty$  and  $L \leq \frac{2}{H}$ . Then*

$$\left| \frac{I - I_L}{I} \right| = o(1)$$

*and  $\frac{1}{4\rho c_0} \leq \frac{P_l}{Q_l} \leq 4\rho c_0$ , for all  $l$ .*

The following proposition is useful for proving Theorem 4.2:

**Proposition 5.3.** *Let  $p(z), q(z)$  be two densities supported on  $[0, 1]$ . Suppose  $H = \Theta(1)$ . Also suppose*

*C1'  $p(z), q(z) > 0$  on  $(0, 1)$ , and  $\sup_z \{p(z) \vee q(z)\} < \infty$ .*

*C2'  $\int \left| \log \frac{p(z)}{q(z)} \right| dz < \infty$ .*

C3' There exists  $h(z)$  such that

- (a)  $h(z) \geq \max \left\{ \left| \frac{p'(z)}{p(z)} \right|, \left| \frac{q'(z)}{q(z)} \right| \right\},$
- (b)  $h(z)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped, and
- (c)  $\int |h_n(z)|^t dz < \infty$  for some constant  $t$  such that  $\frac{2}{r} \leq t \leq 1$ .

C4'  $p'(z), q'(z) \geq 0$  for all  $z < c'_{s1}$ , and  $p'(z), q'(z) \leq 0$  for all  $z > c'_{s2}$ .

Let  $L$  be a sequence such that  $L \rightarrow \infty$ . Suppose  $\frac{1}{c_0} \leq \frac{1-P_0}{1-Q_0} \leq c_0$ . Let  $\text{bin}_l = [a_l, b_l]$ , for  $l = 1, \dots, L$  be a uniformly-spaced binning of  $[0, 1]$  and let  $P_l = (1 - P_0) \int_{a_l}^{b_l} p(z) dz$  and  $Q_l = (1 - Q_0) \int_{a_l}^{b_l} q(z) dz$ . Then

$$\left| \frac{I - I_L}{I} \right| = o(1)$$

and  $\frac{1}{4c_0} \exp(-L^{1/r}) \leq \frac{P_l}{Q_l} \leq 4c_0 \exp(L^{1/r})$ , for all  $l$ .

### 5.3 Analysis of the transformation function

Propositions 5.2 and 5.3 consider densities supported on  $[0, 1]$ . This is enough for us because once we transform the densities by an application of  $\Phi$ , the new densities are compactly supported and, importantly, the Renyi divergence  $I$  and the Hellinger divergence  $H$  are invariant with respect to the transformation  $\Phi$ .

To see this, let  $p(x)$  and  $q(x)$  denote densities over  $\mathbb{R}$ , and let  $p_\Phi(z)$  and  $q_\Phi(z)$  denote the transformed densities over  $[0, 1]$ . It is easy to see that  $p_\Phi(z) = \frac{p(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$  and  $q_\Phi(z) = \frac{q(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$ . Therefore, via the change of variables  $z = \Phi^{-1}(x)$ , we have the following relations:

$$\begin{aligned} \int_{\mathbb{R}} \sqrt{p(x)q(x)} dx &= \int_0^1 \sqrt{p_\Phi(z)q_\Phi(z)} dz, \\ \int_{\mathbb{R}} \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx &= \int_0^1 \left( \sqrt{p_\Phi(z)} - \sqrt{q_\Phi(z)} \right)^2 dz. \end{aligned}$$

Therefore, the divergences  $I$  and  $H$  between  $p(x)$  and  $q(x)$  are equal to the divergences between  $p_\Phi(z)$  and  $q_\Phi(z)$ .

To prove Theorems 4.1 and 4.2, it thus remains to show that if the densities  $p(x)$  and  $q(x)$  satisfy Assumptions A1–A5 (or A1'–A4'), the transformed densities  $p_\Phi(z)$  and  $q_\Phi(z)$  satisfy Assumptions C1–C5 (or C1'–C4') in Proposition 5.2 (or Proposition 5.3). This is done in Propositions F.1 and F.2.

## 6 Conclusion

We have provided a rate-optimal community estimation algorithm for the homogeneous weighted stochastic block model. In the setting where the average degree is of order  $\log n$  and the edge weight densities  $p(x)$  and  $q(x)$  are fixed, we have also characterized the exact recovery threshold. Our algorithm includes a preprocessing step consisting of transforming and discretizing the (possibly) continuous edge weights to obtain a simpler graph with edge weights supported on a finite discrete set. This approach may be useful for other network data analysis problems involving continuous distributions, where discrete versions of the problem are simpler to analyze.

Our paper is a first step toward understanding the weighted SBM under the same mathematical framework that has been so fruitful for the unweighted SBM. It is far from comprehensive, however, and many open questions remain. We describe a few here:

1. An important problem is to extend our analysis to the case of a *heterogenous* stochastic block model, where edge weight distributions depend on the exact community assignments of both endpoints. In such a setting, Abbe and Sandon [3] and Yun and Proutiere [31] have shown that a generalized information divergence—the CH divergence—governs the intrinsic difficulty of community recovery. We believe that a similar discretization-based approach should lead to analogous results in the case of a heterogeneous weighted SBM. The key challenge would be to show that discretization does not lose much information with respect to the CH-divergence.
2. Real-world networks often have nodes with very high degrees, which may adversely affect the accuracy of recovery methods for the stochastic block model. To solve this problem, degree-corrected SBMs [11, 33] have been proposed as an effective alternative to regular SBMs. It is straightforward to extend the concept of degree-correction to the weighted SBM, but it is unclear whether our discretization-based approach would be effective in obtaining optimal error rates.
3. It is easy to extend our results to the weighted *and* labeled SBMs if the number of labels is finite or assumed to be slowly growing. However, this excludes some interesting cases, including the setting where edge labels represent counts from a Poisson distribution. We suspect that in such a situation, it may be possible to combine low-probability labels in a clever way to obtain a discretization that is again amenable to our approach.

## References

- [1] E. Abbe, A. S. Bandeira, A. Bracher, and A. Singer. Decoding binary node labels from censored edge measurements: Phase transition and efficient recovery. *IEEE Transactions on Network Science and Engineering*, 1(1):10–22, 2014.
- [2] E. Abbe, A. S. Bandeira, and G. Hall. Exact recovery in the stochastic block model. *arXiv preprint arXiv:1405.3267*, 2014.
- [3] E. Abbe and C. Sandon. Community detection in general stochastic block models: Fundamental limits and efficient recovery algorithms. *arXiv preprint arXiv:1503.00609*, 2015.
- [4] C. Aicher, A. Z. Jacobs, and A. Clauset. Learning latent block structure in weighted networks. *Journal of Complex Networks*, page cnu026, 2014.
- [5] A. A. Amini, A. Chen, P. J. Bickel, E. Levina, et al. Pseudo-likelihood methods for community detection in large sparse networks. *The Annals of Statistics*, 41(4):2097–2122, 2013.
- [6] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008, 2008.
- [7] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang. Complex networks: Structure and dynamics. *Physics reports*, 424(4):175–308, 2006.
- [8] D. Easley and J. Kleinberg. *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, New York, NY, USA, 2010.
- [9] S. E. Fienberg, M. M. Meyer, and S. S. Wasserman. Statistical analysis of multiple sociometric relations. *Journal of the American Statistical Association*, 80(389):51–67, 1985.
- [10] C. Gao, Z. Ma, A. Y. Zhang, and H. H. Zhou. Achieving optimal misclassification proportion in stochastic block model. *arXiv preprint arXiv:1505.03772*, 2015.

- [11] C. Gao, Z. Ma, A. Y. Zhang, and H. H. Zhou. Community detection in degree-corrected block models. *arXiv preprint arXiv:1607.06993*, 2016.
- [12] A. Goldenberg, A. X. Zheng, S. E. Fienberg, and E. M. Airoldi. A survey of statistical network models. *Found. Trends Mach. Learn.*, 2(2):129–233, February 2010.
- [13] B. Hajek, Y. Wu, and J. Xu. Achieving exact cluster recovery threshold via semidefinite programming. *arXiv preprint arXiv:1412.6156*, 2014.
- [14] B. Hajek, Y. Wu, and J. Xu. Achieving exact cluster recovery threshold via semidefinite programming: Extensions. *arXiv preprint arXiv:1502.07738*, 2015.
- [15] E. Hartuv and R. Shamir. A clustering algorithm based on graph connectivity. *Information Processing Letters*, 76(4–6):175–181, 2000.
- [16] P. W. Holland, K. B. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5(2):109–137, 1983.
- [17] M. O. Jackson. *Social and Economic Networks*. Princeton University Press, 2010.
- [18] L. Massoulié. Community detection thresholds and the weak Ramanujan property. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, STOC ’14, pages 694–703. ACM, 2014.
- [19] F. McSherry. Spectral partitioning of random graphs. In *Foundations of Computer Science, 2001. Proceedings. 42nd IEEE Symposium on*, pages 529–537. IEEE, 2001.
- [20] E. Mossel, J. Neeman, and A. Sly. Stochastic Block Models and Reconstruction. *arXiv preprint arXiv:1202.1499*, 2012.
- [21] E. Mossel, J. Neeman, and A. Sly. Consistency thresholds for binary symmetric block models. *arXiv preprint arXiv:1407.1591*, 2014.
- [22] M. Newman, A.-L. Barabasi, and D. J. Watts. *The Structure and Dynamics of Networks: (Princeton Studies in Complexity)*. Princeton University Press, Princeton, NJ, USA, 2006.
- [23] M. E. J. Newman. Analysis of weighted networks. *Physical Review E*, 70(5):056131, 2004.
- [24] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004.
- [25] J. K. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multi-locus genotype data. *Genetics*, 155(2):945–959, 2000.
- [26] M. Rubinov and O. Sporns. Complex network measures of brain connectivity: Uses and interpretations. *NeuroImage*, 52(3):1059–1069, 2010. Computational Models of the Brain.
- [27] D.S. Sade. Sociometrics of Macaca mulatta: I. Linkages and cliques in grooming matrices. *Folia Primatologica*, 18(3–4):196–223, 1972.
- [28] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(8):888–905, August 2000.
- [29] S. Wasserman and C. Anderson. Stochastic a posteriori blockmodels: Construction and assessment. *Social Networks*, 9(1):1–36, 1987.

- [30] H. C. White, S. A. Boorman, and R. L. Breiger. Social structure from multiple networks: I. Blockmodels of roles and positions. *American Journal of Sociology*, 81(4):730–780, 1976.
- [31] S. Yun and A. Proutiere. Optimal cluster recovery in the labeled stochastic block model. In *Advances in Neural Information Processing Systems*, pages 965–973, 2016.
- [32] A. Y. Zhang and H. H. Zhou. Minimax rates of community detection in stochastic block model. *arXiv preprint arXiv:1507.05313*, 2015.
- [33] Y. Zhao, E. Levina, and J. Zhu. Consistency of community detection in networks under degree-corrected stochastic block models. *The Annals of Statistics*, 40(4):2266–2292, 2012.

## A Proof of Proposition 5.1

We structure the proof according to the flow of our algorithm. Since this proposition addresses the case of discrete labels, we do not need to consider the “transformation and discretization” step.

We begin by analyzing Algorithm 3.2, where we deliberately add noise to the graph by changing edge colors at random. Although adding random noise destroys information and increases the difficulty of community recovery, Lemma B.2 in Appendix B.5 shows that the process does not significantly affect the Renyi divergence  $I_L$ . Furthermore, the new probabilities of edge labels are at least  $\frac{c}{n}$ , which is important for our later analysis. To simplify notation, we continue to refer the new edge label probabilities as  $P_l$  and  $Q_l$  throughout the proof.

Next, our algorithm performs spectral clustering using only the edges with label  $l$ , and calculates  $\hat{I}_l := \frac{(\hat{P}_l - \hat{Q}_l)^2}{\hat{P}_l \vee \hat{Q}_l}$ , where  $\hat{P}_l$  and  $\hat{Q}_l$  are the estimated probabilities obtained by using the output of spectral clustering:

$$\hat{P}_l = \frac{\sum_{u \neq v : \tilde{\sigma}_l(u) = \tilde{\sigma}_l(v)} (A_l)_{uv}}{|u \neq v : \tilde{\sigma}_l(u) = \tilde{\sigma}_l(v)|}, \quad \text{and} \quad \hat{Q}_l = \frac{\sum_{u \neq v : \tilde{\sigma}_l(u) \neq \tilde{\sigma}_l(v)} (A_l)_{uv}}{|u \neq v : \tilde{\sigma}_l(u) \neq \tilde{\sigma}_l(v)|}.$$

We then select  $l^* \in \arg \max_{1 \leq l \leq L} \hat{I}_l$ . Note that if  $|\hat{P}_l - P_l|$  and  $|\hat{Q}_l - Q_l|$  are small,  $\hat{I}_l$  provides a measure of how “good” a color is for clustering: larger values of  $\hat{I}_l$  correspond to greater separation between  $P_l$  and  $Q_l$ . Naturally, the accuracy of the estimated edge probabilities  $\hat{P}_l$  and  $\hat{Q}_l$  depends on the accuracy of the spectral clustering step. Proposition B.1 below makes this statement rigorous. Before stating the proposition, we define the set of “good” labels as follows:

$$L_1 = \left\{ l : \frac{n(P_l - Q_l)^2}{P_l \vee Q_l} := \frac{\Delta_l^2}{P_l \vee Q_l} \geq 1 \right\}.$$

We bound the difference between the estimated and true probabilities for good and bad colors in Proposition B.1, the formal statement and proof of which are contained in Appendix B.1.

**Proposition B.1.** *Suppose  $\sigma$  is a clustering with error rate at most  $\gamma$ ; i.e.,  $l(\sigma, \sigma_0) \leq \gamma$ , for sufficiently small  $\gamma \geq \frac{1}{n}$ . With probability at least  $1 - Ln^{-(3+\delta_p)}$ , for a small  $\delta_p > 0$ , the following hold:*

1. For  $l \in L_1$ , we have  $|\hat{P}_l - P_l| \leq \eta \Delta_l$  and  $|\hat{Q}_l - Q_l| \leq \eta \Delta_l$ .
2. For  $l \in L_1^c$ , we have  $|\hat{P}_l - P_l| \leq \eta \sqrt{\frac{P_l \vee Q_l}{n}}$  and  $|\hat{Q}_l - Q_l| \leq \eta \sqrt{\frac{P_l \vee Q_l}{n}}$ .

In both cases,  $\eta = C \sqrt{\gamma \log \frac{1}{\gamma}}$ , for an absolute constant  $C$ .

We now work toward obtaining a suitable initial clustering with small error rate  $\gamma$ . In Proposition B.2, we show that if the edge probabilities for a particular label are well-separated, the spectral clustering output of Algorithm 3.4 is reasonably accurate. We provide a rough statement of the proposition here, and refer to Appendix B.2 for the precise statement and proof.

**Proposition B.2.** *If  $P_l$  and  $Q_l$  satisfy  $C_1 \frac{(P_l \vee Q_l)}{n(P_l - Q_l)^2} \leq 1$  for an absolute constant  $C_1$ , the output  $\sigma^l$  of Algorithm 3.4 satisfies the inequality*

$$l(\sigma^l, \sigma_0) \leq C_2 \frac{(P_l \vee Q_l)}{n(P_l - Q_l)^2},$$

for a constant  $C_2$ , with probability at least  $1 - n^{-4}$ .

Thus, if we want to cluster to an arbitrary degree of accuracy, we need  $\frac{(P_l \vee Q_l)}{n(P_l - Q_l)^2} \rightarrow 0$  for at least one well-separated label  $l$ . We show that the label  $l^*$  selected in Algorithm 3.3 satisfies  $\frac{(P_{l^*} \vee Q_{l^*})}{n(P_{l^*} - Q_{l^*})^2} \rightarrow 0$  in the following proposition. A more detailed restatement and proof is contained in Appendix B.3.

**Proposition B.4.** *With probability at least  $1 - 2Ln^{-(3+\delta_p)}$ , we have  $\frac{n(P_{l^*} - Q_{l^*})^2}{\rho_L^4(P_{l^*} \vee Q_{l^*})} \rightarrow \infty$ .*

Now let  $E_1$  denote the high-probability event that the label  $l^*$  is chosen according to Proposition B.4. We perform spectral clustering  $n$  times, omitting one vertex and clustering on the remaining graph each time, and denote the resulting community assignments by  $\{\tilde{\sigma}_u\}_{1 \leq u \leq n}$ . Note that Proposition B.2, together with a simple union bound, implies that with probability at least  $1 - n^{-3}$ , all clusterings have error rate bounded by  $\gamma := C \frac{P_{l^*} \vee Q_{l^*}}{n(P_{l^*} - Q_{l^*})^2}$ , for some constant  $C$ . Denote this event by  $E_2$ . On  $E_1 \cap E_2$ , we then have  $\gamma \rho_L^4 \rightarrow 0$ . Thus, we may apply Proposition B.1 on each clustering  $\tilde{\sigma}_u$  to show that the conclusion holds simultaneously for all  $\tilde{\sigma}_u$ 's, with probability at least  $1 - Ln^{-(2+\delta_p)}$ . We denote this last event by  $E_3$ . Furthermore,  $\eta = \Theta\left(\sqrt{\gamma \log \frac{1}{\gamma}}\right)$ , so  $|\eta \rho_L| \rightarrow 0$ .

We now construct the clustering  $\hat{\sigma}_u$  by assigning vertex  $u$  to an appropriate community in  $\tilde{\sigma}_u$ , using the relation from Algorithm 3.5. In Proposition B.5, we show that with high probability, the assignment  $\hat{\sigma}_u(u)$  is “correct.” We provide a rough statement of the proposition here, and defer the exact statement and proof to Appendix B.4:

**Proposition B.5.** *Let  $\pi_u \in S_K$  be such that  $l(\sigma_0, \tilde{\sigma}_u) = d(\sigma_0, \pi_u(\tilde{\sigma}_u))$ . Conditioned on  $E_1 \cap E_2 \cap E_3$ , with probability at least  $1 - (K - 1) \exp\left(- (1 - o(1)) \frac{n}{\beta K} I_L\right)$ , we have  $\pi_u^{-1}(\sigma_0(u)) = \hat{\sigma}_u(u)$ .*

We briefly discuss the uniqueness of  $\pi_u$ . By construction, the error rate of  $\hat{\sigma}_u$  is at most  $\gamma + \frac{1}{n}$ , so  $l(\sigma_0, \hat{\sigma}_u) < \frac{1}{8\beta K}$  for sufficiently small  $\gamma$ . Now note that for any  $u$ , the minimum cluster size of the clustering  $\hat{\sigma}_u$  is at least  $\frac{n}{\beta K} - (n\gamma + 1) \geq \frac{n(1 - \beta K \gamma - \beta K/n)}{\beta K} \geq \frac{n}{2\beta K}$ , for small  $\gamma$ . A simple argument (cf. Lemma B.8) shows that  $\pi_u$  is the unique permutation to obtain such a small error rate.

Define  $\pi_1$  and  $\pi_u$  to be the permutations minimizing  $d(\sigma_0, \pi_1(\hat{\sigma}_1))$  and  $d(\sigma_0, \pi_u(\hat{\sigma}_u))$ , respectively, and let  $\xi_u$  denote the permutation minimizing  $d(\hat{\sigma}_1, \xi_u(\hat{\sigma}_u))$ . We know that  $d(\sigma_0, \pi_1(\hat{\sigma}_1)) < \frac{1}{8\beta K}$  and  $d(\sigma_0, \pi_u(\hat{\sigma}_u)) < \frac{1}{8\beta K}$ . Thus, the triangle inequality implies

$$d(\hat{\sigma}_1, \pi_1^{-1}(\pi_u(\hat{\sigma}_u))) = d(\pi_1(\hat{\sigma}_1), \pi_u(\hat{\sigma}_u)) \leq d(\sigma_0, \pi_1(\hat{\sigma}_1)) + d(\sigma_0, \pi_u(\hat{\sigma}_u)) < \frac{1}{4\beta K}.$$

Since the minimum cluster size of both  $\hat{\sigma}_1$  and  $\hat{\sigma}_u$  is  $\frac{n}{2\beta K}$ , Lemma B.8 implies that  $\xi_u = \pi_1^{-1} \circ \pi_u$ ; and by Lemma B.7, we also have  $\hat{\sigma}(u) = \xi_u(\hat{\sigma}_u(u))$ .

Restating Proposition B.5, we have

$$P\left(\hat{\sigma}_u(u) \neq \pi_u^{-1}(\sigma_0(u)) \mid E_1 \cap E_2 \cap E_3\right) \leq \exp\left(- (1 + o(1)) \frac{n I_L}{\beta K}\right).$$

Furthermore, the left-hand expression is equivalent to  $\xi_u^{-1}(\hat{\sigma}(u)) \neq \pi_u^{-1}(\sigma_0(u))$ , or  $\hat{\sigma}(u) \neq \xi_u \circ \pi_u^{-1}(\sigma_0(u)) = \pi_1^{-1}(\sigma_0(u))$ .

Altogether, we conclude that

$$\begin{aligned} P(\hat{\sigma}(u) \neq \pi_1^{-1}(\sigma_0(u))) &\leq \exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) + P(E_1^c) + P(E_2^c) + P(E_3^c) \\ &\leq \exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) + 2Ln^{-(3+\delta_p)} + n^{-3} + Ln^{-(2+\delta_p)} \\ &\leq \exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) + n^{-(2+\delta_p)}, \end{aligned}$$

where  $\eta' = o(1)$ .

Finally, suppose

$$\exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) \geq n^{-(1+\delta_p)}.$$

Defining  $\eta'' = \eta' + \beta\sqrt{\frac{K}{nI_L}} = o(1)$ , we have

$$\begin{aligned} P\left\{l(\hat{\sigma}, \sigma_0) > \exp\left(-(1-\eta'')\frac{nI_L}{\beta K}\right)\right\} &\leq \frac{\mathbb{E}l(\hat{\sigma}, \sigma_0)}{\exp\left(-(1-\eta'')\frac{nI_L}{\beta K}\right)} \\ &\leq \frac{1}{\exp\left(-(1-\eta'')\frac{nI_L}{\beta K}\right)} \frac{1}{n} \sum_{u=1}^n P(\hat{\sigma}(u) \neq \pi_1^{-1}(\sigma_0(u))) \\ &\leq \exp\left\{-(\eta'' - \eta')\frac{nI_L}{\beta K}\right\} + \frac{Cn^{-(2+\delta_p)}}{\exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right)} = o(1). \end{aligned}$$

On the other hand, if

$$\exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) \leq n^{-(1+\delta_p)},$$

we have

$$\begin{aligned} P\left\{l(\hat{\sigma}, \sigma_0) > \exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right)\right\} &\leq P(l(\hat{\sigma}, \sigma_0) > 0) \\ &\leq P(d(\hat{\sigma}, \pi^{-1}(\sigma_0)) > 0) \\ &\leq \sum_{u=1}^n P(\hat{\sigma}(u) \neq \pi_1^{-1}(\sigma_0(u))) \\ &\leq n \exp\left(-(1-\eta')\frac{nI_L}{\beta K}\right) + n^{-(1+\delta_p)} = o(1). \end{aligned}$$

This completes the proof of Proposition 5.1.

## B Supporting results for Proposition 5.1

We now provide proofs for the supporting results stated in Appendix A.

### B.1 Analysis of estimation error of $\hat{P}_l$ and $\hat{Q}_l$

**Proposition B.1.** *Let  $A$  be the adjacency matrix of a labeled network with true clustering assignment  $\sigma_0$ . Suppose  $\sigma$  is a random initial clustering satisfying  $l(\sigma, \sigma_0) \leq \gamma$ . Let  $\hat{P}_l = \frac{\sum_{u \neq v: \sigma(u)=\sigma(v)} \mathbf{1}(A_{uv}=l)}{|\{u \neq v: \sigma(u)=\sigma(v)\}|}$*

and  $\widehat{Q}_l = \frac{\sum_{u \neq v: \sigma(u) \neq \sigma(v)} \mathbf{1}(A_{uv}=l)}{|\{u \neq v: \sigma(u) = \sigma(v)\}|}$  be the MLE of  $P_l$  and  $Q_l$  based on  $\sigma$ . Let  $\delta_p$  be a positive, fixed, and arbitrarily small real number, and let  $c > 0$  be an absolute constant. Then with probability at least  $1 - Ln^{-(3+\delta_p)}$ , the following hold for all sufficiently small  $\gamma$ :

1. For all  $l$  such that  $P_l \vee Q_l \geq \frac{c}{n}$ , if  $\frac{n\Delta_l^2}{P_l \vee Q_l} \geq 1$ , then

$$|\widehat{P}_l - P_l| \leq \eta \Delta_l, \quad \text{and} \quad |\widehat{Q}_l - Q_l| \leq \eta \Delta_l.$$

2. For all  $l$  such that  $P_l \vee Q_l \geq \frac{c}{n}$ , if  $\frac{n\Delta_l^2}{P_l \vee Q_l} \leq 1$ , then

$$|\widehat{P}_l - P_l| \leq \eta \sqrt{\frac{P_l \vee Q_l}{n}}, \quad \text{and} \quad |\widehat{Q}_l - Q_l| \leq \eta \sqrt{\frac{P_l \vee Q_l}{n}}.$$

In both cases,  $\eta = C \sqrt{\gamma \log \frac{1}{\gamma}}$ , for an absolute constant  $C$ .

*Proof.* Our proof proceeds by showing that for any fixed assignment  $\sigma$  with error rate bounded by  $\gamma$ , the event described in Proposition B.1 holds with high probability. For a fixed assignment  $\sigma$ , we call a random graph a “bad graph for  $\sigma$ ” if the event does not hold. For each  $\sigma$ , we upper-bound the probability that a randomly chosen graph lies in the set of bad graphs for  $\sigma$ ; we then use a union bound over all choices of  $\sigma$  and all  $L$  colors to show that the probability of choosing a bad graph is bounded by  $Ln^{-(3+\delta_p)}$ .

We begin by bounding the bias of  $\widehat{P}_l$ . We have

$$\begin{aligned} \mathbb{E}(\widehat{P}_l) &= \frac{\sum_{u \neq v: \sigma(u) = \sigma(v)} \{\mathbf{1}(\sigma_0(u) = \sigma_0(v))P_l + \mathbf{1}(\sigma_0(u) \neq \sigma_0(v))Q_l\}}{|\{u \neq v: \sigma(u) = \sigma(v)\}|} \\ &= (1 - \lambda)P_l + \lambda Q_l = P_l + \lambda(Q_l - P_l), \end{aligned} \tag{B.1}$$

where  $\lambda := \frac{\sum_{u \neq v: \sigma(u) = \sigma(v)} \mathbf{1}(\sigma_0(u) \neq \sigma_0(v))}{|\{u \neq v: \sigma(u) = \sigma(v)\}|}$ . Thus,  $|\mathbb{E}(\widehat{P}_l) - P_l| \leq \lambda|Q_l - P_l|$ . Furthermore, if  $\widehat{n}_k$  denotes the number of vertices in cluster  $k$  according to  $\sigma$ , we have

$$\begin{aligned} \lambda &= \frac{\sum_{u \neq v: \sigma(u) = \sigma(v)} \mathbf{1}(\sigma_0(u) \neq \sigma_0(v))}{|\{u \neq v: \sigma(u) = \sigma(v)\}|} = \frac{\sum_k \sum_{u \neq v: \sigma(u) = \sigma(v) = k} \mathbf{1}(\sigma_0(u) \neq \sigma_0(v))}{\sum_k \widehat{n}_k(\widehat{n}_k - 1)} \\ &\leq \frac{\sum_k \sum_{u \neq v: \sigma(u) = \sigma(v) = k} \mathbf{1}(\neg(\sigma_0(u) = \sigma_0(v) = k))}{\sum_k \widehat{n}_k(\widehat{n}_k - 1)} \\ &\leq \frac{\sum_k \sum_{u \neq v: \sigma(u) = \sigma(v) = k} \{\mathbf{1}(\sigma_0(v) \neq k) + \mathbf{1}(\sigma_0(u) \neq k)\}}{\sum_k \widehat{n}_k(\widehat{n}_k - 1)}. \end{aligned}$$

Define  $\gamma_k = \frac{1}{n} \sum_{u: \sigma(u) = k} \mathbf{1}(\sigma_0(u) \neq k)$  to be the error rate within the estimated cluster  $k$ . Then  $\sum_k \gamma_k \leq \gamma$  and  $\sum_{u: \sigma(u) = k} \sum_{v: \sigma(v) = k} \mathbf{1}(\sigma_0(v) \neq k) = \gamma_k n \widehat{n}_k$ , implying that

$$\lambda \leq \frac{\sum_k 2\gamma_k n \widehat{n}_k}{\sum_k \widehat{n}_k(\widehat{n}_k - 1)} = \frac{n}{\sum_k \widehat{n}_k(\widehat{n}_k - 1)} \sum_k 2\gamma_k \widehat{n}_k \stackrel{(a)}{\leq} \frac{K}{n - K} \sum_k 2\gamma_k \widehat{n}_k \stackrel{(b)}{\leq} 4\gamma K,$$

where (a) uses the fact that

$$\sum_k \frac{\widehat{n}_k}{n} (\widehat{n}_k - 1) = n \sum_k \left( \frac{\widehat{n}_k}{n} \right)^2 - 1 \geq \frac{n}{K} - 1, \tag{B.2}$$

and (b) uses the assumption  $K < \frac{n}{2}$ . Altogether, we conclude that  $|\mathbb{E}(\widehat{P}_l) - P_l| \leq 4\gamma K \Delta_l$ . A similar calculation may be performed for  $\widehat{Q}_l$ , so

$$\max \{|\mathbb{E}(\widehat{P}_l) - P_l|, |\mathbb{E}(\widehat{Q}_l) - Q_l|\} \leq C_2 \gamma \Delta_l,$$



for a constant  $C_2 > 0$ . To simplify presentation, we define  $\eta_1 = C_2\gamma$ , so

$$\max \left\{ |\mathbb{E}(\widehat{P}_l) - P_l|, |\mathbb{E}(\widehat{Q}_l) - Q_l| \right\} \leq \eta_1 \Delta_l. \quad (\text{B.3})$$

We now turn to bounding  $|\widehat{P}_l - P_l|$  and  $|\widehat{Q}_l - Q_l|$ . Denoting  $\tilde{A}_{uv} = \mathbf{1}(A_{ij} = l)$  and using Bernstein's inequality, we have

$$P \left( \left| \sum_{u,v: \sigma(u)=\sigma(v)} (\tilde{A}_{uv} - \mathbb{E}\tilde{A}_{uv}) \right| > t \right) \leq 2 \exp \left( - \frac{t^2}{2 \sum_{u,v: \sigma(u)=\sigma(v)} \mathbb{E}\tilde{A}_{uv} + \frac{2}{3}t} \right).$$

By equation (B.1), we have

$$\sum_{u,v: \sigma(u)=\sigma(v)} \mathbb{E}\tilde{A}_{uv} = \sum_k \hat{n}_k(\hat{n}_k - 1) \mathbb{E}\widehat{P}_l \leq (P_l \vee Q_l) \sum_k \hat{n}_k(\hat{n}_k - 1),$$

implying that

$$P \left( \left| \sum_{u,v: \sigma(u)=\sigma(v)} (\tilde{A}_{uv} - \mathbb{E}\tilde{A}_{uv}) \right| > t \right) \leq 2 \exp \left( - \frac{t^2}{2(P_l \vee Q_l) \sum_k \hat{n}_k(\hat{n}_k - 1) + \frac{2}{3}t} \right).$$

Let

$$t^2 = 4 \left\{ \left( 2(P_l \vee Q_l) \sum_k \hat{n}_k(\hat{n}_k - 1) \right) \left( C_1 \gamma n \log \frac{1}{\gamma} + (3 + \delta_p) \log n \right) \right\} \\ \vee 4 \left\{ \left( C_1 \gamma n \log \frac{1}{\gamma} + (3 + \delta_p) \log n \right)^2 \right\},$$

for a constant  $C_1$  to be defined later. Let

$$A = 2(P_l \vee Q_l) \sum_k \hat{n}_k(\hat{n}_k - 1), \quad \text{and} \quad B = C_1 \gamma n \log \frac{1}{\gamma} + (3 + \delta_p) \log n.$$

We split into two cases:

1. Suppose  $A \geq B$ . Then  $t^2 = 4AB$ , and the probability term is at most

$$2 \exp \left( - \frac{4AB}{A + \frac{4}{3}\sqrt{AB}} \right) \leq 2 \exp \left( - \frac{4AB}{A + \frac{4}{3}A} \right) \leq 2 \exp(-B).$$

2. Suppose  $A \leq B$ . Then  $t^2 = 4B^2$ , and the probability term is at most

$$2 \exp \left( - \frac{4B^2}{A + \frac{4}{3}B} \right) \leq 2 \exp \left( - \frac{4B^2}{B + \frac{4}{3}B} \right) \leq 2 \exp(-B).$$

In either case, with probability at least  $1 - 2 \exp \left( - \left( C_1 \gamma n \log \frac{1}{\gamma} + (3 + \delta_p) \log n \right) \right)$ , we have

$$|\widehat{P}_l - \mathbb{E}(\widehat{P}_l)| = \frac{\sum_{u \neq v} \sigma(u)=\sigma(v) (\tilde{A}_{uv} - \mathbb{E}\tilde{A}_{uv})}{\sum_{u \neq v} \mathbf{1}(\sigma(u) = \sigma(v))} \leq \frac{t}{\sum_{u \neq v} \mathbf{1}(\sigma(u) = \sigma(v))}.$$

We now derive a more manageable upper bound for  $t$ . Using the notation from above, we have  $t^2 = \max(4AB, 4B^2) \leq 4(\sqrt{AB} + B)^2$ . Since  $\gamma \geq \frac{1}{n}$ , we have  $C_1 \gamma n \log \frac{1}{\gamma} + (3 + \delta_p) \log n \leq \tilde{C}_1 \gamma n \log \frac{1}{\gamma}$ , implying that

$$\begin{aligned} \frac{t}{\sum_{u \neq v} \mathbf{1}(\sigma(u) = \sigma(v))} &\leq 2 \frac{\sqrt{2(P_l \vee Q_l) \tilde{C}_1 \gamma n \log \frac{1}{\gamma}}}{\sqrt{\sum_k \hat{n}_k (\hat{n}_k - 1)}} + 2 \frac{\tilde{C}_1 \gamma n \log \frac{1}{\gamma}}{\sum_k \hat{n}_k (\hat{n}_k - 1)} \\ &\stackrel{(a)}{\leq} 2 \frac{\sqrt{2(P_l \vee Q_l) \tilde{C}_1 \gamma K \log \frac{1}{\gamma}}}{\sqrt{n - K}} + 2 \frac{\tilde{C}_1 \gamma K \log \frac{1}{\gamma}}{n - K} \\ &\stackrel{(b)}{\leq} 4 \sqrt{\frac{P_l \vee Q_l}{n}} \sqrt{\tilde{C}_1 K \gamma \log \frac{1}{\gamma}} + 4 \frac{\tilde{C}_1 K \gamma \log \frac{1}{\gamma}}{n}, \end{aligned}$$

where (a) uses inequality (B.2) and (b) uses the assumption  $n - K \geq \frac{n}{2}$ .

To further simplify the expression, note that  $P_l \vee Q_l \geq \frac{c}{n}$  implies  $\frac{1}{n} \leq \frac{1}{\sqrt{c}} \sqrt{\frac{P_l \vee Q_l}{n}}$ , so with probability at least  $1 - \exp(-C_1 \gamma n \log \frac{1}{\gamma} - (3 + \delta_p) \log n)$ , we have

$$|\hat{P}_l - \mathbb{E}(\hat{P}_l)| \leq \sqrt{\frac{P_l \vee Q_l}{n}} \left( C'_1 \sqrt{\gamma \log \frac{1}{\gamma}} + C'_2 \gamma \log \frac{1}{\gamma} \right), \quad (\text{B.4})$$

for suitable constants  $C'_1$  and  $C'_2$ . Using a similar calculation, we may show that there exist suitable constants  $C'_3$  and  $C'_4$  such that

$$|\hat{Q}_l - \mathbb{E}(\hat{Q}_l)| \leq \sqrt{\frac{P_l \vee Q_l}{n}} \left( C'_3 \sqrt{\gamma \log \frac{1}{\gamma}} + C'_4 \gamma \log \frac{1}{\gamma} \right).$$

When  $\gamma$  is sufficiently small, the first term dominates, so we may take choose the right-hand sides to be  $\eta_2 \sqrt{\frac{P_l \vee Q_l}{n}}$ , where  $\eta_2 = C_3 \sqrt{\gamma \log \frac{1}{\gamma}}$ .

Finally, note that there are at most  $\binom{n}{\gamma n} K^{\gamma n}$  possible  $\sigma$ 's satisfying the error bound. We have

$$\begin{aligned} \log \left( \binom{n}{\gamma n} K^{\gamma n} \right) &\leq \log \left( \frac{n^{\gamma n} e^{\gamma n}}{(\gamma n)^{\gamma n}} \frac{1}{\sqrt{2\pi \gamma n}} \right) + \gamma n \log K \\ &\leq \log \left( \frac{e^{\gamma n}}{\gamma^{\gamma n}} \right) - \frac{1}{2} \log 2\pi \gamma n + \gamma n \log K \\ &\leq \gamma n \log \frac{e}{\gamma} + \gamma n \log K \\ &= \gamma n \log \frac{Ke}{\gamma} \\ &\leq C_1 \gamma n \log \frac{1}{\gamma}, \end{aligned}$$

for a suitable constant  $C_1$ . Taking a union bound across all cluster assignments, we then conclude that the probability of inequality (B.4) holding simultaneously for all labels  $l$  is at least  $1 - Ln^{-(3+\delta_p)}$ .

Combining inequalities (B.3) and (B.4), we arrive at the bound

$$\max \left\{ |P_l - \hat{P}_l|, |Q_l - \hat{Q}_l| \right\} \leq \eta_1 \Delta_l + \eta_2 \sqrt{\frac{P_l \vee Q_l}{n}}. \quad (\text{B.5})$$

If  $\frac{n \Delta_l^2}{P_l \vee Q_l} \geq 1$ , we therefore have

$$\max \left\{ |P_l - \hat{P}_l|, |Q_l - \hat{Q}_l| \right\} \leq \eta_1 \Delta_l + \eta_2 \Delta_l = (\eta_1 + \eta_2) \Delta_l,$$

whereas if  $\frac{n\Delta_l^2}{P_l \vee Q_l} < 1$ , we have

$$\max \left\{ |P_l - \widehat{P}_l|, |Q_l - \widehat{Q}_l| \right\} \leq \eta_1 \sqrt{\frac{P_l \vee Q_l}{n}} + \eta_2 \sqrt{\frac{P_l \vee Q_l}{n}} = (\eta_1 + \eta_2) \sqrt{\frac{P_l \vee Q_l}{n}}.$$

Since  $\eta_2 = C_3 \sqrt{\gamma \log \frac{1}{\gamma}}$  dominates  $\eta_1 = C_2 \gamma$  for sufficiently small  $\gamma$ , the desired bounds follow.  $\square$

## B.2 Analysis of spectral clustering

**Proposition B.2.** *Suppose an unweighted adjacency matrix  $A$  is drawn from a homogeneous stochastic block model with probabilities  $p$  and  $q$  and cluster imbalance factor  $\beta$ . Suppose  $p, q \geq \frac{c}{n}$ . Then there exists a constant  $C$  such that if*

$$256\mu\beta C^2 K^3 \frac{(p \vee q)}{n(p-q)^2} \leq 1,$$

the output  $\sigma$  of Algorithm 3.4 with parameters  $\mu \geq 32C^2\beta$  and  $\tau = \bar{d}$  satisfies

$$l(\sigma, \sigma_0) \leq 64C^2\beta \frac{K^2(p \vee q)}{n(p-q)^2},$$

with probability at least  $1 - n^{-C'}$ , where  $C' > 4$ .

*Proof.* Note that the trim parameter  $\tau$  is a random variable, since the average degree  $\bar{d}$  is random. Since the community sizes are bounded by  $\frac{n}{\beta K}$ , we may find constants  $C_{d_1} < C_{d_2} = 1$ , depending only on  $K$  and  $\beta$ , such that

$$C_{d_1} n(p \vee q) \leq \mathbb{E}[\bar{d}] \leq C_{d_2} n(p \vee q).$$

Using Hoeffding's inequality, we conclude that with probability at least  $1 - \exp(-nC_{\bar{d}})$ , for some constant  $C_{\bar{d}}$ , we have

$$\frac{C_{d_1}}{2} n(p \vee q) \leq \bar{d} \leq 2C_{d_2} n(p \vee q). \quad (\text{B.6})$$

We now apply the following lemma:

**Lemma B.1.** (Lemma 5 of Gao et al. [10]) *Let  $P \in [0, 1]^{n \times n}$  be a symmetric matrix, and let  $p_{\max} := \max_{u \geq v} P_{uv}$ . Let  $A$  be an adjacency matrix such that  $A_{uu} = 0$  and  $A_{uv} \sim \text{Ber}(P_{uv})$  for  $u < v$ . For any  $C' > 0$  and  $0 < C_1 < C_2$ , there exists some  $C > 0$  such that*

$$\|T_\tau(A) - P\|_2 \leq C\sqrt{np_{\max} + 1}, \quad \forall \tau \in [C_1(np_{\max} + 1), C_2(np_{\max} + 1)],$$

with probability at least  $1 - n^{-C'}$ .

**Remark B.1.** Lemma 5 from Gao et al. [10] is stated slightly differently—for any  $C' > 0$ , there exist constants  $c$ ,  $C_1$ , and  $C_2$  such that the result holds with probability at least  $1 - n^{-C'}$ . However, our restatement follows immediately from slight modifications of the proof. Furthermore, note that the statement of Lemma B.1 refers to the output of spectral clustering with respect to a fixed trim parameter, but we will apply it in a setting where  $\tau$  is random.

Using Lemma B.1 with a fixed  $C' > 4$  and constants  $C_1 = \frac{C_{d_1}}{2}$  and  $C_2 = 2C_{d_2}$ , we conclude that there exists a constant  $C > 0$  such that

$$\|T_\tau(A) - P\|_2 \leq C\sqrt{n(p \vee q)}, \quad \forall \tau \in [C_1, C_2], \quad (\text{B.7})$$

with probability at least  $1 - n^{-C'}$ . In particular, this inequality will also hold for the random choice  $\tau = \bar{d}$ , by inequality (B.6) above. Furthermore, we may assume that  $C \geq 1$ , since inequality (B.7) also holds with  $C$  replaced by  $\max(1, C)$ . Thus,

$$\begin{aligned}\|\hat{A} - P\|_2 &\leq \|T_\tau(A) - P\|_2 + \|\hat{A} - T_\tau(A)\|_2 \\ &\stackrel{(a)}{\leq} 2\|T_\tau(A) - P\|_2 \\ &\leq 2C\sqrt{n(p \vee q)},\end{aligned}$$

where (a) follows because  $\hat{A}$  is the best rank- $K$  approximation of  $T_\tau(A)$  and  $\text{rank}(P) = K$ , so  $\|T_\tau(A) - \hat{A}\|_2 \leq \|T_\tau(A) - P\|_2$  by the Eckart-Young-Mirsky Theorem. This implies that

$$\sum_{u=1}^n \|\hat{A}_u - P_u\|_2^2 = \|\hat{A} - P\|_F^2 \leq K\|\hat{A} - P\|_2^2 \leq 4KC^2n(p \vee q).$$

We now denote the  $K$  distinct rows of  $P$  by  $\{\mathcal{Z}_i\}_{1 \leq i \leq K}$ , and for a vertex  $u$ , denote the row  $P_u$  by  $\mathcal{Z}(u)$ . Note that

$$\|\mathcal{Z}_i - \mathcal{Z}_j\|_2^2 \geq \frac{2n}{\beta K}(p - q)^2, \quad \forall i \neq j,$$

since each cluster contains at least  $\frac{n}{\beta K}$  vertices.

A vertex  $u$  is considered *valid* if  $\|\hat{A}_u - \mathcal{Z}_i\|_2^2 \leq \frac{1}{16} \frac{1}{\beta K}(p - q)^2 n$  for some  $\mathcal{Z}_i$ ; otherwise,  $u$  is *invalid*. Also define

$$\mathcal{Z}^*(u) := \arg \min_{\mathcal{Z}_i} \|\hat{A}_u - \mathcal{Z}_i\|_2^2,$$

so  $\mathcal{Z}^*(u)$  is the row of  $P$  closest to  $\hat{A}_u$ . Note that if  $u$  is valid, then  $\|\hat{A}_u - \mathcal{Z}^*(u)\|_2^2 \leq \frac{1}{16} \frac{1}{\beta K}(p - q)^2 n$ .

We show that the set  $S$  constructed in Algorithm 3.4 satisfies the following properties:

**Claim 1:**  $S$  contains only valid points.

**Claim 2:** For every pair of distinct nodes  $u, v \in S$ , we have  $\mathcal{Z}^*(u) \neq \mathcal{Z}^*(v)$ .

We first prove that the proposition follows from the claims. We denote the rows of  $\hat{A}$  corresponding to the members of  $S$  by  $\mathcal{S}_i$ , assigning indices so that  $\mathcal{S}_i$  is the surrogate for  $\mathcal{Z}_i$  (i.e., both  $\mathcal{S}_i$  and  $\mathcal{Z}_i$  are associated to a common vertex  $u$ ). In particular, note that

$$\|\mathcal{S}_i - \mathcal{Z}_i\|_2 \leq \frac{1}{16} \frac{1}{\beta K}(p - q)^2 n, \quad \forall 1 \leq i \leq K,$$

since  $S$  only contains valid points. Let  $\mathcal{S}(u)$  be the surrogate of  $\mathcal{Z}(u)$ , and denote  $\mathcal{S}^*(u) = \arg \min_{\mathcal{S}_i} \|\hat{A}_u - \mathcal{S}_i\|_2^2$ ; i.e., the member of  $S$  that is closest to  $\hat{A}_u$ . We say that a valid point  $u$  is *misclassified* if  $\mathcal{S}^*(u) \neq \mathcal{S}(u)$ . The number of mistakes we make is thus bounded by the number of invalid points plus the number of misclassified valid points. Note that if  $u$  is invalid, we have  $\|\hat{A}_u - P_u\|_2^2 \geq \frac{1}{16} \frac{1}{\beta K}(p - q)^2 n$ . We claim that the same inequality holds for any misclassified valid point  $u$ .

Consider such a point  $u$ . Since  $u$  is valid, there exists  $\mathcal{Z}_i$  such that

$$\|\hat{A}_u - \mathcal{Z}_i\|_2^2 \leq \frac{1}{16} \frac{1}{\beta K}(p - q)^2 n.$$

We claim that  $\mathcal{S}^*(u) = \mathcal{S}_i$ . For any  $j \neq i$ , we have

$$\|\hat{A}_u - \mathcal{S}_j\|_2 \geq \|\mathcal{Z}_i - \mathcal{Z}_j\|_2 - \|\mathcal{Z}_j - \mathcal{S}_j\|_2 - \|\mathcal{Z}_i - \hat{A}_u\|_2$$

$$\begin{aligned}
&\geq \sqrt{\frac{2}{\beta K}(p-q)^2n} - 2\sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n} \\
&> 2\sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n}.
\end{aligned}$$

Furthermore,

$$\|\hat{A}_u - \mathcal{S}_i\|_2 \leq \|\hat{A}_u - \mathcal{Z}_i\|_2 + \|\mathcal{S}_i - \mathcal{Z}_i\|_2 \leq 2\sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n}.$$

Thus, for any  $j \neq i$ , we have

$$\|\hat{A}_u - \mathcal{S}_j\|_2 > \|\hat{A}_u - \mathcal{S}_i\|_2,$$

implying that  $\mathcal{S}^*(u) = \mathcal{S}_i$ .

Since  $u$  is also misclassified, we have  $\mathcal{S}(u) \neq \mathcal{S}^*(u) = \mathcal{S}_i$ . Let  $\mathcal{S}(u) = \mathcal{S}_j$  and  $\mathcal{Z}(u) = \mathcal{Z}_j$ . We have the following sequence of inequalities:

$$\begin{aligned}
\|\hat{A}_u - \mathcal{Z}(u)\|_2 &= \|\hat{A}_u - \mathcal{Z}_j\|_2 \\
&\geq \|\hat{A}_u - \mathcal{S}_j\|_2 - \|\mathcal{S}_j - \mathcal{Z}_j\|_2 \\
&\geq 2\sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n} - \sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n} \\
&= \sqrt{\frac{1}{16}\frac{1}{\beta K}(p-q)^2n},
\end{aligned}$$

which is the bound we wanted to prove.

Finally, we conclude that the number of mistakes incurred by algorithm is bounded by

$$\frac{\sum_{u=1}^n \|\hat{A}_u - P_u\|_2^2}{\frac{1}{16\beta K}(p-q)^2n} \leq \frac{4KC^2n(p \vee q)}{\frac{1}{16\beta K}(p-q)^2n} \leq \frac{64\beta K^2C^2(p \vee q)}{(p-q)^2},$$

as wanted.

**Proof of Claim 1:** Recall the notation  $N(u) = \{v : \|\hat{A}_u - \hat{A}_v\|_2^2 \leq \mu K^2 \frac{\bar{d}}{n}\}$ . Furthermore, by a Chernoff bound, we have  $\bar{d} \leq 2(p \vee q)n$  with probability at least  $1 - \exp(-\bar{C}_d n)$ . We condition on this event so that if  $v \in N(u)$ , then  $\|\hat{A}_u - \hat{A}_v\|_2^2 \leq 2\mu K^2(p \vee q)$ . We prove the claim by showing that an invalid point  $u$  cannot have  $\frac{1}{\mu} \frac{n}{K}$  neighbors.

By the definition of invalidity,  $\|\hat{A}_u - \mathcal{Z}_i\|_2^2 \geq \frac{1}{16\beta K}(p-q)^2n$ , for any  $\mathcal{Z}_i$ . Let  $v$  be a neighbor of  $u$ . By the triangle inequality, we then have

$$\begin{aligned}
\|\hat{A}_v - \mathcal{Z}(v)\|_2 &\geq \|\hat{A}_u - \mathcal{Z}(v)\|_2 - \|\hat{A}_u - \hat{A}_v\|_2 \\
&\geq \sqrt{\frac{1}{16\beta K}(p-q)^2n} - \sqrt{2\mu K^2(p \vee q)} \\
&\stackrel{(a)}{\geq} \sqrt{\frac{1}{16\beta K}(p-q)^2n} - \sqrt{\frac{1}{64\beta K}(p-q)^2n} = \sqrt{\frac{1}{64\beta K}(p-q)^2n},
\end{aligned}$$

where (a) follows from our assumption coupled with the choice of  $C \geq 1$ , which essentially states that

$$2\mu K^2(p \vee q) \leq \frac{1}{128\beta K}(p-q)^2n < \frac{1}{64\beta K}(p-q)^2n.$$

Thus, for every neighbor  $v$  of  $u$ , we must have  $\|\hat{A}_v - P_v\|_2^2 \geq \frac{1}{64\beta K}(p-q)^2n$ . The number of neighbors of  $u$  may be bounded by

$$\frac{\sum_{v=1}^n \|\hat{A}_v - P_v\|_2^2}{\frac{1}{64\beta K}(p-q)^2n} \leq \frac{4KC^2n(p \vee q)}{\frac{1}{64\beta K}(p-q)^2n} \leq \frac{256\beta K^2C^2(p \vee q)}{(p-q)^2}.$$

By assumption, this quantity is less than  $\frac{1}{\mu} \frac{n}{K}$ .

**Proof of Claim 2:** We first claim that in every cluster, at least half the points  $u$  satisfy  $\|\hat{A}_u - P_u\|_2^2 \leq \frac{1}{4}\mu K^2(p \vee q)$ . This is because the total error is bounded by  $\sum_{u=1}^n \|\hat{A}_u - P_u\|_2^2 \leq 4KC^2n(p \vee q)$ , so the total number of points that violate the condition is at most  $\frac{4KC^2n(p \vee q)}{\frac{1}{4}\mu K^2(p \vee q)} \leq \frac{n}{2\beta K}$ , using the assumption that  $\mu \geq 32C^2\beta$ .

For two points  $u$  and  $v$  in the same cluster satisfying  $\|\hat{A}_w - P_w\|_2^2 \leq \frac{1}{4}\mu K^2(p \vee q)$  for  $w \in \{u, v\}$ , we also have  $\|\hat{A}_u - \hat{A}_v\|_2^2 \leq \mu K^2(p \vee q)$ , by the triangle inequality. Thus, every cluster contains a point  $u$  such that  $N(u) \geq \frac{n}{2\beta K} \geq \frac{1}{\mu} \frac{n}{K}$ , since  $\mu \geq 32C^2\beta > 2\beta$  by our choice of  $C > 1$ .

Suppose that at iteration  $r$ , the set  $S$  consists of points  $s_1, \dots, s_r$ , where  $1 \leq r < K$ , and suppose for a contradiction that  $s_{r+1}$  is such that  $\mathcal{Z}(s_{r+1}) = \mathcal{Z}(s_i)$  for some  $1 \leq i \leq r$ . Since  $s_i$  and  $s_{r+1}$  are both valid points, the triangle inequality implies

$$\|\hat{A}_{s_{r+1}} - \hat{A}_{s_i}\| \leq \frac{1}{4\beta K}(p-q)^2n.$$

On the other hand, because  $S$  does not yet have cardinality  $K$ , some  $\mathcal{Z}_j$  must exist that does not have a surrogate in  $S$ . The cluster that corresponds to  $\mathcal{Z}_j$  must, by our neighborhood size analysis, contain a node  $u$  such that  $N(u) \geq \frac{1}{\mu} \frac{n}{K}$  and

$$\|\hat{A}_u - \mathcal{Z}_j\|_2^2 \leq \frac{1}{4}\mu K^2(p \vee q) \leq \frac{1}{16} \frac{1}{\beta K}(p-q)^2n,$$

where the second inequality follows by assumption. Since  $\mathcal{Z}_j \neq \mathcal{Z}(s_i)$  for any  $1 \leq i \leq r$ , we have  $\|\mathcal{Z}_j - \mathcal{Z}(s_i)\|_2^2 \geq 2\frac{1}{\beta K}(p-q)^2n$ , for all  $1 \leq i \leq r$ . By Claim 1, all  $s_i$ 's are valid, so

$$\|\hat{A}_{s_i} - \mathcal{Z}(s_i)\| \leq \frac{1}{16} \frac{1}{\beta K}(p-q)^2n.$$

Hence, by the triangle inequality, we have  $\|\hat{A}_u - \hat{A}_{s_i}\|_2^2 \geq \frac{1}{\beta K}(p-q)^2n$ , for all  $s_i \in S$ . This is a contradiction because  $u$  is further from every point in  $S$  than  $s_{r+1}$ , so our assumption that  $\mathcal{Z}(s_{r+1}) = \mathcal{Z}(s_i)$  must be incorrect.  $\square$

### B.3 Choosing the label $l^*$

First, we show that for sufficiently well-separated labels,  $\hat{I}_l$  is close to  $\frac{(P_l - Q_l)^2}{P_l \vee Q_l}$ . If the probabilities are not well-separated, we claim that  $\hat{I}_l$  is negligibly small.

**Proposition B.3.** Suppose  $\frac{1}{\rho_L} \leq \frac{P_l}{Q_l} \leq \rho_L$  for all  $l$ . Let  $\sigma^l$  be the output of spectral clustering based on  $\tilde{A}_{ij} = \mathbf{1}(A_{ij} = l)$ , and let  $\hat{P}_l$  and  $\hat{Q}_l$  be estimates of  $P_l$  and  $Q_l$  constructed from  $\sigma^l$ . There exist positive constants  $C_{test}, C_1, C_2, C$ , and  $\delta_p$  such that, with probability at least  $1 - Ln^{-3+\delta_p}$ , we have the following:

1. For all labels  $l$  satisfying  $P_l \vee Q_l > \frac{\epsilon}{n}$  and  $\Delta_l \geq \sqrt{C_{test}} \sqrt{\frac{P_l \vee Q_l}{n}}$ ,

$$C_1 \frac{|P_l - Q_l|}{\sqrt{P_l \vee Q_l}} \leq \frac{|\hat{P}_l - \hat{Q}_l|}{\sqrt{\hat{P}_l \vee \hat{Q}_l}} \leq C_2 \frac{|P_l - Q_l|}{\sqrt{P_l \vee Q_l}}. \quad (\text{B.8})$$

2. For all labels satisfying  $P_l \vee Q_l > \frac{c}{n}$  and  $\Delta_l < \sqrt{C_{test}} \sqrt{\frac{P_l \vee Q_l}{n}}$ ,

$$\frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \leq C \sqrt{\frac{1}{n}}. \quad (\text{B.9})$$

*Proof.* Recall from Proposition B.1 that given a clustering with error rate  $\gamma$ , and under the assumptions  $P_l \vee Q_l > \frac{c}{n}$  and  $\Delta_l^2 \geq \frac{P_l \vee Q_l}{n}$ , the estimated probabilities  $\widehat{P}_l$  and  $\widehat{Q}_l$  satisfy

$$|\widehat{P}_l - \widehat{P}| \leq \eta \Delta_l, \quad \text{and} \quad |\widehat{Q}_l - \widehat{Q}| \leq \eta \Delta_l,$$

with probability at least  $1 - n^{-(3+\delta_p)}$ . We first pick a value of  $\gamma$  such that  $\eta < \frac{1}{4}$ . We now ensure that the error rate obtained from Proposition B.2 matches our choice of  $\gamma$ . Recall that Proposition B.2 states that under the assumptions  $P_l \vee Q_l > \frac{c}{n}$  and

$$C_1 \frac{P_l \vee Q_l}{n(P_l - Q_l)^2} \leq 1,$$

we have

$$l(\sigma, \sigma_0) \leq C_2 \frac{P_l \vee Q_l}{n(P_l - Q_l)^2},$$

for appropriate constants  $C_1$  and  $C_2$ . In particular, for  $C_{test} \geq 1$  sufficiently large,

$$\begin{aligned} C_1 \frac{P_l \vee Q_l}{n(P_l - Q_l)^2} &\leq \frac{C_1}{C_{test}} < 1, \quad \text{and} \\ l(\sigma, \sigma_0) &\leq C_2 \frac{P_l \vee Q_l}{n(P_l - Q_l)^2} \leq \frac{C_2}{C_{test}} < \gamma, \end{aligned}$$

for all labels  $l$  such that  $\Delta_l \geq \sqrt{C_{test}} \sqrt{\frac{P_l \vee Q_l}{n}}$ . Next, note that

$$|\widehat{P}_l - \widehat{Q}_l| \leq |\widehat{P}_l - P_l| + |P_l - Q_l| + |\widehat{Q}_l - Q_l| \leq 2\eta \Delta_l + \Delta_l \leq \frac{3}{2} \Delta_l,$$

and

$$|\widehat{P}_l - \widehat{Q}_l| \geq |P_l - Q_l| - |\widehat{Q}_l - Q_l| - |\widehat{P}_l - P_l| \geq \Delta_l - 2\eta \Delta_l \geq \frac{1}{2} \Delta_l.$$

Furthermore,

$$\widehat{P}_l \vee \widehat{Q}_l \leq (P_l \vee Q_l) + \eta \Delta_l \leq (P_l \vee Q_l) + \eta(P_l \vee Q_l) \leq \frac{5}{4}(P_l \vee Q_l),$$

and

$$\widehat{P}_l \vee \widehat{Q}_l \geq (P_l \vee Q_l) - \eta \Delta_l \geq \frac{3}{4}(P_l \vee Q_l).$$

We conclude that

$$\frac{1}{\sqrt{5}} \frac{\Delta_l}{P_l \vee Q_l} \leq \frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \leq \frac{3}{\sqrt{3}} \frac{\Delta_l}{P_l \vee Q_l}.$$

Now suppose  $\Delta_l^2 < C_{test} \frac{P_l \vee Q_l}{n}$ . Note that this does not necessarily imply  $\Delta_l^2 \leq \frac{P_l \vee Q_l}{n}$ , since  $C_{test} \geq 1$ . However, we may still take the maximum of the bounds provided in Proposition B.1, so with probability at least  $1 - Ln^{-(3+\delta_p)}$ ,

$$|\widehat{P}_l - P_l| \leq \eta \left( \Delta_l \vee \sqrt{\frac{P_l \vee Q_l}{n}} \right) \leq \frac{\sqrt{C_{test}}}{4} \sqrt{\frac{P_l \vee Q_l}{n}},$$

using the choice  $\eta \leq \frac{1}{4}$ . An analogous bound holds for  $|\widehat{Q}_l - Q_l|$ . Hence,

$$\begin{aligned} |\widehat{P}_l - \widehat{Q}_l| &\leq \Delta_l + |\widehat{P}_l - P_l| + |\widehat{Q}_l - Q_l| \\ &\leq \Delta_l + \frac{\sqrt{C_{test}}}{2} \sqrt{\frac{P_l \vee Q_l}{n}} \\ &\leq \frac{3}{2} \sqrt{C_{test}} \sqrt{\frac{P_l \vee Q_l}{n}}. \end{aligned}$$

Now note that

$$\widehat{P}_l \vee \widehat{Q}_l \geq (P_l \vee Q_l) - \frac{\sqrt{C_{test}}}{4} \sqrt{\frac{P_l \vee Q_l}{n}} \geq C'(P_l \vee Q_l).$$

for some constant  $C'$ . It follows that

$$\frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \leq C \sqrt{\frac{1}{n}}.$$

□

We apply Proposition B.3 to conclude that Algorithm 3.3 succeeds in choosing a color  $l^*$  for which  $\frac{n(P_{l^*} - Q_{l^*})^2}{P_{l^*} \vee Q_{l^*}}$  is arbitrarily large:

**Proposition B.4.** *Suppose  $a_n := \frac{nI_L}{L\rho_L^4} \rightarrow \infty$ . For sufficiently large  $n$ , with probability at least  $1 - 2Ln^{-(3+\delta_p)}$ , we have  $\frac{n(P_{l^*} - Q_{l^*})^2}{(P_{l^*} \vee Q_{l^*})\rho_L^4} \geq Ca_n$ , for some constant  $C$ .*

*Proof.* Let  $C_{test}$  be the constant in Proposition B.3. By Lemma B.6, we know that  $I_L$  is of the same order as  $\sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l}$ , implying the existence of a label  $l$  such that  $\Delta_l \geq C_{test} \sqrt{\frac{P_l \vee Q_l}{n}}$  and  $\frac{n\Delta_l^2}{P_l \vee Q_l} \geq C \frac{nI_L}{L} = Ca_n \rho_L^4$ , for a constant  $C$ . Suppose the event of Proposition B.3 holds, which happens with probability at least  $1 - Ln^{-(3+\delta_p)}$ .

**Step 1.** We claim that  $l^*$  satisfies  $\Delta_{l^*} \geq C_{test} \sqrt{\frac{P_{l^*} \vee Q_{l^*}}{n}}$ . Let  $l$  be a label such that  $\frac{n\Delta_l^2}{P_l \vee Q_l} \geq Ca_n \rho_L^4$ , and suppose the claim is false. By Proposition B.3 and the maximality of  $l^*$ , we have

$$\frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \stackrel{(a)}{\leq} \frac{|\widehat{P}_{l^*} - \widehat{Q}_{l^*}|}{\sqrt{\widehat{P}_{l^*} \vee \widehat{Q}_{l^*}}} \leq C \sqrt{\frac{1}{n}},$$

Proposition B.3 also implies that

$$\frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \geq C' \frac{|P_l - Q_l|}{\sqrt{P_l \vee Q_l}} \geq C'' \sqrt{\frac{a_n \rho_L^4}{n}}.$$

However, this is a contradiction, since  $a_n \rightarrow \infty$  and  $\rho_L \geq 1$ .

**Step 2:** Again, let  $l$  be a label such that  $\frac{n\Delta_l^2}{P_l \vee Q_l} \geq Ca_n \rho_L^4$ . By Proposition B.3, we then have

$$\frac{|P_{l^*} - Q_{l^*}|}{\sqrt{P_{l^*} \vee Q_{l^*}}} \geq C \frac{|\widehat{P}_{l^*} - \widehat{Q}_{l^*}|}{\sqrt{\widehat{P}_{l^*} \vee \widehat{Q}_{l^*}}} \geq C \frac{|\widehat{P}_l - \widehat{Q}_l|}{\sqrt{\widehat{P}_l \vee \widehat{Q}_l}} \geq C' \frac{|P_l - Q_l|}{\sqrt{P_l \vee Q_l}} \geq C'' \sqrt{\frac{a_n \rho_L^4}{n}},$$

implying the desired result. □



## B.4 Analysis of error probability for a single node

**Proposition B.5.** *Let  $u$  be an arbitrary fixed node, and let  $\tilde{\sigma}_u$  be the output of Algorithm 3.3. Suppose  $\pi_u \in S_K$  satisfies*

$$l(\sigma_0, \tilde{\sigma}_u) = d(\sigma_0, \pi_u(\tilde{\sigma}_u)),$$

where both  $l$  and  $d$  are taken with respect to the set  $\{1, 2, \dots, n\} \setminus \{u\}$ . Conditioned on the events that the error rate  $\gamma$  of  $\tilde{\sigma}_u$  satisfies  $\gamma \rho_L^4 \rightarrow 0$ , and also the event that the result of Proposition B.1 holds for a sequence  $\eta$  satisfying  $\eta \rho_L^2 \rightarrow 0$ , we have

$$\pi_u^{-1}(\sigma_0(u)) = \arg \max_k \sum_{v: \tilde{\sigma}_u(v)=k} \sum_l \log \frac{\hat{P}_l}{\hat{Q}_l} \mathbf{1}(A_{uv} = l),$$

with probability at least  $1 - (K-1) \exp\left(-(1-o(1))\frac{n}{\beta K} I_L\right)$ .

*Proof.* Throughout the proof, we assume that  $n$  is large enough so  $\frac{1}{2} \sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 \leq \frac{1}{2}$ . Suppose without loss of generality that  $\sigma_0(u) = 1$ . We misclassify  $u$  into community  $k$  if

$$\sum_{v: \tilde{\sigma}_u(v)=k} \sum_l \log \frac{\hat{P}_l}{\hat{Q}_l} \mathbf{1}(A_{uv} = l) \geq \sum_{v: \tilde{\sigma}_u(v)=1} \sum_l \log \frac{\hat{P}_l}{\hat{Q}_l} \mathbf{1}(A_{uv} = l),$$

or equivalently,

$$\sum_{v: \tilde{\sigma}_u(v)=k} \bar{A}_{uv} - \sum_{v: \tilde{\sigma}_u(v)=1} \bar{A}_{uv} \geq 0, \quad (\text{B.10})$$

where  $\bar{A}_{uv} \equiv \sum_l \log \frac{\hat{P}_l}{\hat{Q}_l} \mathbf{1}(A_{uv} = l)$ . Note that the edges from  $u$  are independent of the clustering  $\tilde{\sigma}_u$ , since this clustering was obtained by running the algorithm with vertex  $u$  excluded.

Define  $m_1 = |\{v : \tilde{\sigma}_u(v) = 1\}|$  and  $m_k = |\{v : \tilde{\sigma}_u(v) = k\}|$ , and let  $m'_1 = \{v : \tilde{\sigma}_u(v) = 1, \sigma_0(v) = 1\}$  be the points correctly clustered by  $\sigma_u$ . Let  $m'_k = \{v : \tilde{\sigma}_u(v) \neq 1, \sigma_0(v) = k\}$  denote the points correctly classified by  $\tilde{\sigma}_u$  in community  $k$ . With these definitions, the probability of the bad event in equation (B.10) is the probability of the event

$$\left( \sum_{i=1}^{m'_k} \tilde{Y}_i + \sum_{i=1}^{m_k - m'_k} \tilde{X}_i \right) - \left( \sum_{i=1}^{m'_1} \tilde{X}_i + \sum_{i=1}^{m_1 - m'_1} \tilde{Y}_i \right) \geq 0,$$

where  $\tilde{X}_i = \log \frac{\hat{P}_l}{\hat{Q}_l}$  with probability  $P_l$  and  $\tilde{Y}_i = \log \frac{\hat{P}_l}{\hat{Q}_l}$  with probability  $Q_l$ . (For simplicity, we abuse notation and write  $\tilde{Y}_i$  and  $\tilde{X}_i$  in both bracketed terms. These random variables are not the same, but are independent and identical copies.) This is equal to the probability of the event

$$\exp \left( t \left( \sum_{i=1}^{m'_k} \tilde{Y}_i + \sum_{i=1}^{m_k - m'_k} \tilde{X}_i - \sum_{i=1}^{m'_1} \tilde{X}_i - \sum_{i=1}^{m_1 - m'_1} \tilde{Y}_i \right) \right) \geq 1.$$

We further bound this probability as follows:

$$P \left( \exp \left( t \left( \sum_{i=1}^{m'_k} \tilde{Y}_i + \sum_{i=1}^{m_k - m'_k} \tilde{X}_i - \sum_{i=1}^{m'_1} \tilde{X}_i - \sum_{i=1}^{m_1 - m'_1} \tilde{Y}_i \right) \right) \geq 1 \right)$$

$$\begin{aligned}
&\leq \mathbb{E} \left[ \exp \left( t \left( \sum_{i=1}^{m'_k} \widetilde{Y}_i + \sum_{i=1}^{m_k-m'_k} \widetilde{X}_i - \sum_{i=1}^{m'_1} \widetilde{X}_i - \sum_{i=1}^{m_1-m'_1} \widetilde{Y}_i \right) \right) \right] \\
&= \mathbb{E}[\exp(t\widetilde{Y}_i)]^{m'_k} \mathbb{E}[\exp(t\widetilde{X}_i)]^{m_k-m'_k} \mathbb{E}[\exp(-t\widetilde{X}_i)]^{m'_1} \mathbb{E}[\exp(-t\widetilde{Y}_i)]^{m_1-m'_1} \\
&= \left( \sum_l e^{t \log \frac{\widehat{P}_l}{\widehat{Q}_l} Q_l} \right)^{m'_k} \left( \sum_l e^{t \log \frac{\widehat{P}_l}{\widehat{Q}_l} P_l} \right)^{m_k-m'_k} \left( \sum_l e^{-t \log \frac{\widehat{P}_l}{\widehat{Q}_l} P_l} \right)^{m'_1} \left( \sum_l e^{-t \log \frac{\widehat{P}_l}{\widehat{Q}_l} Q_l} \right)^{m_1-m'_1}.
\end{aligned}$$

We will set  $t = \frac{1}{2}$ , in which case

$$\begin{aligned}
&\left( \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l \right)^{m'_k} \left( \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} P_l \right)^{m_k-m'_k} \left( \sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} Q_l \right)^{m_1-m'_1} \left( \sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l \right)^{m'_1} \\
&= \left( \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} P_l}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l} \right)^{m_k-m'_k} \left( \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} Q_l}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l} \right)^{m_1-m'_1} \quad (\text{B.11})
\end{aligned}$$

$$\left( \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l \right)^{m_k} \left( \sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l \right)^{m_1}. \quad (\text{B.12})$$

We bound terms (B.11) and (B.12) separately. Loosely speaking, we will show that term (B.11) is bounded in magnitude by  $\exp(o(I_L) \frac{n}{K})$ , and term (B.12) is bounded by  $\exp(-\frac{n}{\beta K}(1+o(1))I_L)$ .

**Bound for term (B.11).** We derive a number of separate lemmas bounding various intermediate terms in the computation. In particular, we use the bounds from Lemmas B.5, B.4, and B.6 in the following sequence of inequalities:

$$\begin{aligned}
\left| 1 - \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} P_l}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l} \right| &= \left| \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} (P_l - Q_l)}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l} \right| \\
&\stackrel{(a)}{\leq} \frac{8}{\sum_l \sqrt{P_l Q_l}} \left| \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} (P_l - Q_l) \right| \\
&\stackrel{(b)}{\leq} 16 \left| \sum_l \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) (P_l - Q_l) \right| \\
&\leq 16 \left| \sum_{l \in L_1} \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) (P_l - Q_l) \right| + 16 \sum_{l \notin L_1} \left| \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right| |P_l - Q_l| \\
&\stackrel{(c)}{\leq} 16 \sum_{l \in L_1} \frac{\Delta_l^2}{Q_l} (1 + \eta') + \sum_{i \notin L_1} 32 \rho_L \frac{\Delta_l}{\sqrt{n(P_l \vee Q_l)}} \\
&\leq 16 \rho_L \sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l} (1 + \eta') + \sum_{l \notin L_1} 32 \rho_L \frac{\Delta_l}{\sqrt{n(P_l \vee Q_l)}} \\
&\stackrel{(d)}{\leq} C I_L \rho_L (1 + \eta') + C' \rho_L \frac{L}{n} \stackrel{(e)}{\leq} C \rho_L I_L.
\end{aligned}$$

In (a), we have used Lemma B.5. In (b), we have used the fact that  $\sum_l \sqrt{P_l Q_l} \rightarrow 1$ , so this sum exceeds  $\frac{1}{2}$  when  $n$  is sufficiently large. In (c), we have employed Lemma B.4, which appropriately bounds the term  $\left(\sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1\right)$ . Here,  $\eta' = o(1)$ . Inequality (d) follows from Lemma B.6. Finally, inequality (e) follows from the assumption that  $\frac{I_L n}{L \rho_L^4} \rightarrow \infty$  (note that  $\rho_L \geq 1$ ) and by appropriately redefining  $\eta'$ . Identical analysis shows that

$$\left| 1 - \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} Q_l}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l} \right| \leq C \rho_L I_L.$$

Finally, note that  $|x| \leq \exp(|1 - x|)$ , so term (B.11) may be bounded as

$$\begin{aligned} \left( \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} P_l}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l} \right)^{m_k - m'_k} \left( \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} Q_l}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l} \right)^{m_1 - m'_1} &\leq \exp(C \rho_L I_L (m_k - m'_k + m_1 - m'_1)) \\ &\leq \exp(C I_L \rho_L \gamma n). \end{aligned}$$

Since  $\gamma \rho_L = o(1)$ , we conclude that term (B.11) is bounded by  $\exp\left(\frac{n}{K} o(I_L)\right)$ , as desired.

**Bound for term (B.12).** Let  $\widehat{I} = -\log\left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right) \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)$ . With this definition, we have

$$\left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right)^{m_k} \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)^{m_1} = \exp(-\widehat{I})^{\frac{m_k + m_1}{2}} \left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right)^{\frac{m_k - m_1}{2}} \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)^{\frac{m_1 - m_k}{2}}.$$

We claim that the following statements are true:

1.  $m_1, m_k \geq \frac{n}{\beta K} (1 - \beta K \gamma)$ .
2.  $\widehat{I} \geq I_L (1 + o(1))$ .
3.  $\left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right)^{\frac{m_k - m_1}{2}} \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)^{\frac{m_1 - m_k}{2}} = \exp\left(\frac{n}{K} o(I_L)\right)$ .

Let us first assume these statements are true, and bound term (B.12). We have

$$\begin{aligned} \exp(-\widehat{I})^{\frac{m_1 + m_k}{2}} \left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right)^{\frac{m_k - m_1}{2}} \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)^{\frac{m_1 - m_k}{2}} \\ \leq \exp\left(-I_L (1 + o(1)) \frac{n}{\beta K} \cdot (1 - \beta K \gamma) + \frac{n}{K} o(I_L)\right) \leq \exp\left(-(1 + o(1)) \frac{n}{\beta K} I_L\right), \end{aligned}$$

where the last inequality holds because  $\gamma = o(1)$ . It remains to prove the three claims.

**Claim 1:** This is straightforward. The labeling  $\tilde{\sigma}_u$  has at most  $\gamma n$  errors, so  $m_1 \geq m'_1 \geq \frac{n}{\beta K} - \gamma n$ . A similar argument works for  $m_k$ .

**Claim 2:** We show that the estimation error of  $\widehat{P}_l, \widehat{Q}_l$  does not make  $\widehat{I}$  too small. We begin by writing

$$\widehat{I} - I_L = -\log \frac{\left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right) \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right)}{\left(\sum_l \sqrt{P_l Q_l}\right)^2}.$$

Let us first consider the numerator:

$$\begin{aligned} & \left(\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l\right) \left(\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l\right) \\ &= \left(\sum_l \sqrt{P_l Q_l} \sqrt{\frac{\widehat{P}_l}{P_l} \frac{Q_l}{\widehat{Q}_l}}\right) \left(\sum_l \sqrt{P_l Q_l} \sqrt{\frac{P_l}{\widehat{P}_l} \frac{\widehat{Q}_l}{Q_l}}\right) \\ &= \sum_l P_l Q_l + 2 \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} + \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left(\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2\right) \\ &= \left(\sum_l \sqrt{P_l Q_l}\right)^2 + \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left(\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2\right), \end{aligned}$$

where we define

$$T_{l,l'} = \frac{\widehat{P}_l Q_l P_{l'} \widehat{Q}_{l'}}{P_l \widehat{Q}_l \widehat{P}_{l'} Q_{l'}}.$$

Furthermore,

$$\begin{aligned} \widehat{I} - I_L &= -\log \left( 1 + \frac{\sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left(\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2\right)}{\left(\sum_l \sqrt{P_l Q_l}\right)^2} \right) \\ &\geq -\log \left( 1 + 4 \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left(\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2\right) \right) \quad (\text{assuming } \sum_l \sqrt{P_l Q_l} \geq 1/2) \\ &\geq -4 \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left(\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2\right). \end{aligned} \tag{B.13}$$

We now bound  $|T_{l,l'} - 1|$ :

$$\begin{aligned} |T_{l,l'} - 1| &= \left| \frac{\widehat{P}_l Q_l P_{l'} \widehat{Q}_{l'}}{P_l \widehat{Q}_l \widehat{P}_{l'} Q_{l'}} - 1 \right| \\ &= \left| \left(1 - \frac{P_l - \widehat{P}_l}{P_l}\right) \left(1 - \frac{\widehat{Q}_l - Q_l}{\widehat{Q}_l}\right) \left(1 - \frac{\widehat{P}_{l'} - P_{l'}}{\widehat{P}_{l'}}\right) \left(1 - \frac{Q_{l'} - \widehat{Q}_{l'}}{Q_{l'}}\right) - 1 \right| \\ &\stackrel{(a)}{\leq} 2 \left( \frac{|P_l - \widehat{P}_l|}{P_l} + \frac{|\widehat{Q}_l - Q_l|}{\widehat{Q}_l} + \frac{|\widehat{P}_{l'} - P_{l'}|}{\widehat{P}_{l'}} + \frac{|Q_{l'} - \widehat{Q}_{l'}|}{Q_{l'}} \right) \\ &\stackrel{(b)}{\leq} 4 \left( \frac{|P_l - \widehat{P}_l|}{P_l} + \frac{|\widehat{Q}_l - Q_l|}{Q_l} + \frac{|\widehat{P}_{l'} - P_{l'}|}{P_{l'}} + \frac{|Q_{l'} - \widehat{Q}_{l'}|}{Q_{l'}} \right), \end{aligned}$$

where (a) and (b) follow from Lemma B.3. Since we only work with pairs  $(l, l')$  such that  $l' > l$ , we may choose any ordering we like. Thus, suppose the  $l$ 's are ordered in decreasing order of

$\frac{|\widehat{P}_l - P_l|}{P_l} + \frac{|\widehat{Q}_l - Q_l|}{Q_l}$ . For all pairs  $l < l'$ , we then have

$$|T_{l,l'} - 1| \leq 8 \left( \frac{|\widehat{P}_l - P_l|}{P_l} + \frac{|\widehat{Q}_l - Q_l|}{Q_l} \right).$$

By Proposition B.1, we have

$$\frac{|P_l - \widehat{P}_l|}{P_l} + \frac{|\widehat{Q}_l - Q_l|}{Q_l} \leq \eta \Delta_l \left( \frac{1}{P_l} + \frac{1}{Q_l} \right) \leq \frac{\eta \Delta_l}{P_l \vee Q_l} \cdot 2\rho_L \leq \eta' \frac{\Delta_l}{P_l \vee Q_l},$$

for any  $l \in L_1$ . For  $l \notin L_1$ , we have

$$\frac{|P_l - \widehat{P}_l|}{P_l} \leq \eta \sqrt{\frac{P_l \vee Q_l}{nP_l^2}} = \eta \frac{P_l \vee Q_l}{P_l} \sqrt{\frac{1}{n(P_l \vee Q_l)}} \leq \eta \rho_L \sqrt{\frac{1}{n(P_l \vee Q_l)}} \leq \eta' \sqrt{\frac{1}{n(P_l \vee Q_l)}},$$

and similarly for the  $\frac{|\widehat{Q}_l - Q_l|}{Q_l}$  term. Plugging these bounds into the previous derivation, we obtain

$$|T_{l,l'} - 1| \leq \begin{cases} \eta' \frac{\Delta_l}{P_l \vee Q_l}, & \text{for } l \in L_1, \\ \eta' \frac{1}{\sqrt{n(P_l \vee Q_l)}}, & \text{for } l \notin L_1. \end{cases}$$

We now use the Taylor approximation of  $\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2$  around  $T_{l,l'} = 1$ :

$$\sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2 = \frac{1}{4}(T_{l,l'} - 1)^2 + O(T_{l,l'} - 1)^3.$$

Continuing the bound (B.13), we then obtain

$$\begin{aligned} \widehat{I} - I_L &\geq -4 \sum_{l < l'} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left( \sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2 \right) \\ &\geq -4 \sum_{l \in L_1} \sum_{l' > l} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left( \sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2 \right) \\ &\quad - 4 \sum_{l \notin L_1} \sum_{l' > l} \sqrt{P_l Q_l P_{l'} Q_{l'}} \left( \sqrt{T_{l,l'}} + \frac{1}{\sqrt{T_{l,l'}}} - 2 \right) \\ &\geq - \sum_{l \in L_1} \sum_{l' > l} \sqrt{P_l Q_l P_{l'} Q_{l'}} \eta' \left( \frac{\Delta_l}{P_l \vee Q_l} \right)^2 - \sum_{l \notin L_1} \sum_{l' > l} \sqrt{P_l Q_l P_{l'} Q_{l'}} \eta' \frac{1}{n(P_l \vee Q_l)} \\ &\geq -\eta' \left( \sum_{l \in L_1} \frac{\Delta_l^2 \sqrt{P_l Q_l}}{(P_l \vee Q_l)^2} \right) \left( \sum_{l'} \sqrt{P_{l'} Q_{l'}} \right) - \eta' \left( \sum_{l \notin L_1} \frac{\sqrt{P_l Q_l}}{n(P_l \vee Q_l)} \right) \left( \sum_{l'} \sqrt{P_{l'} Q_{l'}} \right) \\ &\geq -\eta' \left( \sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l} \right) \left( \sum_{l'} \sqrt{P_{l'} Q_{l'}} \right) - \eta' \left( \sum_{l \notin L_1} \frac{1}{n} \right) \left( \sum_{l'} \sqrt{P_{l'} Q_{l'}} \right) \\ &\stackrel{(a)}{=} -o(I_L), \end{aligned}$$

where (a) follows from the fact that  $\sum_{l'} \sqrt{P_{l'} Q_{l'}} \leq 1$ , the statement  $\sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l} = \Theta(I_L)$  from Lemma B.6, and our assumption that  $\sum_{l \notin L_1} \frac{1}{n} \leq \frac{L}{n} = o(I_L)$ . This proves the claim.

**Claim 3.** We rewrite the term in claim 3 as follows:

$$\begin{aligned}
& \left( \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l \right)^{\frac{m_k - m_1}{2}} \left( \sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l \right)^{\frac{m_1 - m_k}{2}} \\
&= \left( \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l \right)^{\frac{m_k - m_1}{2}} \left( \sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l \right)^{\frac{m_1 - m_k}{2}} \left( \frac{\sum_l \sqrt{\widehat{P}_l \widehat{Q}_l}}{\sum_l \sqrt{\widehat{P}_l \widehat{Q}_l}} \right)^{\frac{m_1 - m_k}{2}} \\
&= \left( \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} \widehat{Q}_l} \right)^{\frac{m_k - m_1}{2}} \left( \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} P_l}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} \widehat{P}_l} \right)^{\frac{m_1 - m_k}{2}}.
\end{aligned}$$

Assume  $m_k \geq m_1$ . The reverse case may be analyzed in an identical manner. We may rewrite the term as

$$\left( 1 + \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} (Q_l - \widehat{Q}_l)}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} \widehat{Q}_l} \right)^{\frac{m_k - m_1}{2}} \left( 1 + \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} (\widehat{P}_l - P_l)}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} \widehat{P}_l} \right)^{\frac{m_k - m_1}{2}}.$$

Note that  $\sum_l \sqrt{\widehat{P}_l \widehat{Q}_l} \rightarrow 1$ , so Lemma B.5 implies that the denominators are  $\Theta(1)$ . We bound the numerator as follows:

$$\begin{aligned}
\left| \sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} (Q_l - \widehat{Q}_l) \right| &= \left| \sum_l \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) (Q_l - \widehat{Q}_l) \right| \\
&\leq \left| \sum_{l \in L_1} \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) (Q_l - \widehat{Q}_l) \right| + \left| \sum_{l \notin L_1} \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) (Q_l - \widehat{Q}_l) \right| \\
&\stackrel{(a)}{\leq} \left| \sum_{l \in L_1} \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) \eta \Delta_l \right| + \left| \sum_{l \notin L_1} \left( \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right) \eta \sqrt{\frac{P_l \vee Q_l}{n}} \right| \\
&\stackrel{(b)}{\leq} \sum_{l \in L_1} \eta \frac{\Delta_l^2}{\widehat{Q}_l} + \sum_{l \notin L_1} \eta \rho_L \frac{1}{n} \\
&\leq \eta \rho_L \sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l} + \sum_{l \notin L_1} \eta \rho_L \frac{1}{n} \\
&\stackrel{(c)}{\leq} \eta' I_L + \eta' \frac{L}{n} \\
&\stackrel{(d)}{\leq} \eta' I_L.
\end{aligned}$$

In the above sequence of inequalities, step (a) follows from Proposition B.1, step (b) follows from Lemma B.4, step (c) follows from Lemma B.6 and the assumption  $\eta \rho_L \rightarrow 0$ , and step (d) follows from our assumption  $\frac{L}{n} = o(I_L)$ . Thus, we obtain

$$\left( 1 + \frac{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} (Q_l - \widehat{Q}_l)}{\sum_l \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} \widehat{Q}_l} \right)^{\frac{m_k - m_1}{2}} \left( 1 + \frac{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} (\widehat{P}_l - P_l)}{\sum_l \sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} \widehat{P}_l} \right)^{\frac{m_k - m_1}{2}}$$

$$\leq \exp((m_k - m_1) \log(1 + o(I_L))) \leq \exp\left(\frac{n}{K} o(I_L)\right),$$

proving the claim.

**Combining bounds for terms (B.11) and (B.12):** Multiplying the bounds for terms (B.11) and (B.12) shows that the probability of misclassifying  $u$  into some cluster  $k \neq 1$  is at most  $\exp\left((1 + o(1)) \frac{n I_L}{\beta K}\right)$ . Taking a union bound over all clusters  $k \neq 1$  completes the proof.  $\square$

## B.5 Additional lemmas for Proposition 5.1

**Lemma B.2.** *Let  $L, P_l, Q_l, \rho_L$ , and  $I_L$  satisfy the assumptions in Proposition 5.1. Define the new probabilities of edge labels as follows:*

$$P'_l := P_l(1 - \delta) + \frac{\delta}{L+1}, \quad \text{and} \quad Q'_l := Q_l(1 - \delta) + \frac{\delta}{L+1},$$

for all  $0 \leq l \leq L$ , where  $\delta = \frac{c(L+1)}{n}$ . Let  $I'_L$  denote the Renyi divergence between  $P'_l$  and  $Q'_l$ . Then for all sufficiently large  $n$ , we have  $P'_l, Q'_l > \frac{c}{n}$  for all  $0 \leq l \leq L$ , and

$$I'_L = I_L(1 + o(1)).$$

*Proof.* Clearly,  $P'_l, Q'_l > \frac{c}{n}$ . For the second part of the lemma, we begin by writing

$$I_L = -2 \log \sum_{l=0}^L \sqrt{P_l Q_l} = -2 \log \left( 1 - \frac{1}{2} \sum_{l=0}^L (\sqrt{P_l} - \sqrt{Q_l})^2 \right) = \left( \sum_{l=0}^L (\sqrt{P_l} - \sqrt{Q_l})^2 \right) (1 + o(1)).$$

Similarly, we have  $I'_L = \left( \sum_{l=0}^L (\sqrt{P'_l} - \sqrt{Q'_l})^2 \right) (1 + o(1))$ , so it is enough to show that

$$\left( \sum_{l=0}^L (\sqrt{P_l} - \sqrt{Q_l})^2 \right) = \left( \sum_{l=0}^L (\sqrt{P'_l} - \sqrt{Q'_l})^2 \right) (1 + o(1)).$$

We consider two cases:  $\rho_L = \omega(1)$  and  $\rho_L = \Theta(1)$ . If  $\rho_L = \omega(1)$ , we choose  $a = \frac{n I_L}{\rho_L (L+1)}$ . If  $\rho_L = \Theta(1)$ , we choose  $a = o\left(\frac{n I_L}{L+1}\right)$  such that  $a \rightarrow \infty$ . Note that in both cases, we have  $\frac{a}{\rho_L} \rightarrow \infty$  and  $\frac{a(L+1)}{n} = o(I_L)$ . We now break the set of labels into two groups, where  $G_1$  contains all labels satisfying  $P_l \vee Q_l \leq \frac{a}{n}$ , and  $G_2 = G_1^c$ .

Let  $\Delta_l := |P_l - Q_l|$  and  $\Delta'_l := |P'_l - Q'_l|$ . For labels in  $G_1$ , we have  $\Delta_l \leq \frac{a}{n}$ . Thus,

$$(\sqrt{P_l} - \sqrt{Q_l})^2 = \frac{\Delta_l^2}{(\sqrt{P_l} + \sqrt{Q_l})^2} \leq \Delta_l \leq \frac{a}{n},$$

and

$$(\sqrt{P'_l} - \sqrt{Q'_l})^2 = \frac{(\Delta'_l)^2}{(\sqrt{P'_l} + \sqrt{Q'_l})^2} \leq \Delta'_l = (1 - \delta) \Delta_l \leq (1 - \delta) \frac{a}{n}.$$

Therefore,

$$\left| \sum_{l \in G_1} (\sqrt{P_l} - \sqrt{Q_l})^2 - \sum_{l \in G_1} (\sqrt{P'_l} - \sqrt{Q'_l})^2 \right| \leq \frac{a(L+1)}{n} = o(I_L).$$

For labels in  $G_2$ , we may write

$$\left| \sum_{l \in G_2} (\sqrt{P_l} - \sqrt{Q_l})^2 - \sum_{l \in G_2} (\sqrt{P'_l} - \sqrt{Q'_l})^2 \right| = \sum_{l \in G_2} \frac{\Delta_l^2}{(\sqrt{P_l} + \sqrt{Q_l})^2} \left| 1 - (1 - \delta)^2 \frac{(\sqrt{P_l} + \sqrt{Q_l})^2}{(\sqrt{P'_l} + \sqrt{Q'_l})^2} \right|.$$

We analyze the term inside the absolute value as follows:

$$\frac{(\sqrt{P_l} + \sqrt{Q_l})^2}{(\sqrt{P'_l} + \sqrt{Q'_l})^2} = \left( \frac{\sqrt{P_l} + \sqrt{Q_l}}{\sqrt{P_l} \sqrt{1 - \delta + \frac{c}{nP_l}} + \sqrt{Q_l} \sqrt{1 - \delta + \frac{c}{nQ_l}}} \right)^2 = \left( \frac{\sqrt{P_l} + \sqrt{Q_l}}{\sqrt{P_l} \sqrt{1 - \delta + \frac{c}{nP_l}} + \sqrt{Q_l} \sqrt{1 - \delta + \frac{c}{nQ_l}}} \right)^2.$$

Since  $P_l \vee Q_l > \frac{a}{n}$ , we have

$$\frac{1}{n(P_l \vee Q_l)} < \frac{1}{a} = o(1), \quad \text{so} \quad \frac{1}{nP_l} \leq \frac{\rho_L}{n(P_l \vee Q_l)} < \frac{\rho_L}{a} = o(1).$$

Furthermore, since  $\delta = o(1)$ , we have

$$\left( \frac{\sqrt{P_l} + \sqrt{Q_l}}{\sqrt{P_l} \sqrt{1 - \delta + \frac{c}{nP_l}} + \sqrt{Q_l} \sqrt{1 - \delta + \frac{c}{nQ_l}}} \right)^2 = \left( \frac{\sqrt{P_l} + \sqrt{Q_l}}{\sqrt{P_l}(1 + o(1)) + \sqrt{Q_l}(1 + o(1))} \right)^2 = 1 + o(1).$$

Hence, we may conclude that

$$\sum_{l \in G_2} \frac{\Delta_l^2}{(\sqrt{P_l} + \sqrt{Q_l})^2} \left| 1 - (1 - \delta)^2 \frac{(\sqrt{P_l} + \sqrt{Q_l})^2}{(\sqrt{P'_l} + \sqrt{Q'_l})^2} \right| = \sum_{l \in G_2} \frac{\Delta_l^2}{(\sqrt{P_l} + \sqrt{Q_l})^2} \cdot o(1) = o(I_L).$$

Combining the results for labels in  $G_1$  and  $G_2$ , we conclude that  $I'_L = (1 + o(1))I_L$ .  $\square$

We often use the bound  $\frac{1}{2}P \leq \widehat{P}_l \leq 2P_l$ . The following lemma justifies this:

**Lemma B.3.** *Let  $l$  be any label and suppose  $\frac{1}{\rho_L} \leq \frac{P_l}{Q_l} \leq \rho_L$  where  $\rho_L > 1$ ; also suppose  $P_l, Q_l \geq \frac{c}{n}$ . Conditioned on the event that the conclusion of Proposition B.1 holds with a sequence  $\eta$  such that  $\eta\rho_L^2 \rightarrow 0$ , we have*

$$\max_l \frac{|\widehat{P}_l - P_l|}{P_l} \rightarrow 0, \quad \text{and} \quad \max_l \frac{|\widehat{Q}_l - Q_l|}{Q_l} \rightarrow 0.$$

*In particular, for sufficiently small  $\eta$ , we have  $\frac{1}{2}P_l \leq \widehat{P}_l \leq 2P_l$ , and likewise for  $Q_l$ .*

*Proof.* We prove the statement first for  $P_l$ ; the same argument applies to  $Q_l$ . By Proposition B.1, we have that either  $|\widehat{P}_l - P_l| \leq \eta\Delta_l$  or  $|\widehat{P}_l - P_l| \leq \eta\sqrt{\frac{P_l \vee Q_l}{n}}$ . First suppose  $|\widehat{P}_l - P_l| \leq \eta\Delta_l$ . Then

$$\frac{|\widehat{P}_l - P_l|}{P_l} \leq \eta \frac{\Delta_l}{P_l} \leq \eta(1 + \rho_L) \rightarrow 0.$$

If instead  $|\widehat{P}_l - P_l| \leq \eta\sqrt{\frac{P_l \vee Q_l}{n}}$ , we have

$$\frac{|\widehat{P}_l - P_l|}{P_l} \leq \eta\sqrt{\frac{\rho_L}{P_l n}} \leq \eta\sqrt{\rho_L} \sqrt{\frac{1}{c}} \rightarrow 0,$$

where we use the fact that  $\frac{P_l \vee Q_l}{P_l}$  is at most  $\rho_L$ .  $\square$



**Lemma B.4.** Suppose  $\frac{1}{\rho_L} \leq \frac{P_l}{Q_l} \leq \rho_L$  and  $P_l, Q_l \geq \frac{c}{n}$  for all  $l$ , where  $\rho_L > 1$ . Conditioned on the event that the conclusion of Proposition B.1 holds for a sequence  $\eta$  such that  $\eta\rho_L^2 \rightarrow 0$ :

1. For all  $l$  satisfying  $n\frac{\Delta_l^2}{P_l \vee Q_l} \geq 1$ , we have

$$\left| \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right| \leq \left| \frac{P_l - Q_l}{Q_l} \right| (1 + \eta'),$$

where  $\eta' \rightarrow 0$  and  $\eta'$  does not depend on the color  $l$ .

2. For all  $l$  satisfying  $n\frac{\Delta_l^2}{P_l \vee Q_l} < 1$ , we have

$$\left| \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right| \leq 2\rho_L \frac{1}{\sqrt{n(P_l \vee Q_l)}}.$$

By symmetry, the same bounds also hold for  $\sqrt{\frac{\widehat{Q}_l}{\widehat{P}_l}} - 1$ .

*Proof.* First suppose  $l$  satisfies  $\frac{n\Delta_l^2}{P_l \vee Q_l} \geq 1$ . By Lemma B.3, we have  $\frac{\widehat{Q}_l - Q_l}{Q_l} = \eta'$ , where  $\eta' \rightarrow 0$ . In the following derivation, we use  $\eta'$  to denote a sequence such that  $\eta' = o(1)$ ; the actual value of  $\eta'$  may change from instance to instance. We use  $\eta$  to denote a sequence where  $\eta\rho_L = o(1)$ . We have

$$\begin{aligned} \frac{\widehat{P}_l}{\widehat{Q}_l} - 1 &= \frac{\widehat{P}_l - P_l + P_l}{\widehat{Q}_l - Q_l + Q_l} - 1 = \frac{\frac{\widehat{P}_l - P_l}{Q_l} + \frac{P_l}{Q_l}}{\frac{\widehat{Q}_l - Q_l}{Q_l} + 1} - 1 \\ &= \left( \frac{P_l}{Q_l} + \frac{\widehat{P}_l - P_l}{Q_l} \right) \left( 1 - \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') \right) - 1 \\ &= \frac{P_l}{Q_l} + \frac{\widehat{P}_l - P_l}{Q_l} - \frac{P_l}{Q_l} \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') - \frac{\widehat{P}_l - P_l}{Q_l} \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') - 1 \\ &\stackrel{(a)}{=} \frac{P_l}{Q_l} + \frac{\eta\Delta_l}{Q_l} + \rho_L \frac{\eta\Delta_l}{Q_l} (1 + \eta') + \frac{\eta\Delta_l}{Q_l} \eta' - 1 \\ &= \frac{P_l}{Q_l} + \eta \frac{\Delta_l}{Q_l} + \eta\rho_L \frac{\Delta_l}{Q_l} - 1 \\ &= \frac{P_l - Q_l}{Q_l} (1 + \eta'). \end{aligned}$$

In (a), we have used the fact that  $|\widehat{P}_l - P_l| \leq \eta\Delta_l$  and  $|\widehat{Q}_l - Q_l| \leq \eta\Delta_l$  by Proposition B.1. Applying the inequality  $|\sqrt{x} - 1| \leq |x - 1|$  then completes the proof of the first case.

The proof of the second case is almost identical. Suppose  $l$  satisfies  $\frac{n\Delta_l^2}{P_l \vee Q_l} < 1$ . Then

$$|\widehat{P}_l - P_l| = \eta \sqrt{\frac{P_l \vee Q_l}{n}}, \quad \text{and} \quad |\widehat{Q}_l - Q_l| = \eta \sqrt{\frac{P_l \vee Q_l}{n}}.$$

By Lemma B.3, we have  $\frac{\widehat{Q}_l - Q_l}{Q_l} = \eta'$ , where  $\eta' = o(1)$ . Hence,

$$\frac{\widehat{P}_l}{\widehat{Q}_l} - 1 = \left( \frac{P_l}{Q_l} + \frac{\widehat{P}_l - P_l}{Q_l} \right) \left( 1 - \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') \right) - 1$$

$$\begin{aligned}
&= \frac{P_l}{Q_l} + \frac{\widehat{P}_l - P_l}{Q_l} - \frac{P_l}{Q_l} \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') - \frac{\widehat{P}_l - P_l}{Q_l} \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') - 1 \\
&= \frac{P_l}{Q_l} + \eta \sqrt{\frac{P_l \vee Q_l}{nQ_l^2}} + \rho_L \eta \sqrt{\frac{P_l \vee Q_l}{nQ_l^2}} + \eta \sqrt{\frac{P_l \vee Q_l}{nQ_l^2}} - 1 \\
&= \frac{P_l}{Q_l} + \rho_L \eta \sqrt{\frac{(P_l \vee Q_l)^2}{nQ_l^2(P_l \vee Q_l)}} - 1 \\
&= \frac{P_l - Q_l}{Q_l} + \eta \rho_L^2 \sqrt{\frac{1}{n(P_l \vee Q_l)}} \\
&= \frac{P_l - Q_l}{Q_l} + \eta' \rho_L \sqrt{\frac{1}{n(P_l \vee Q_l)}}.
\end{aligned}$$

Using the inequality  $|\sqrt{1+x} - 1| \leq x$  for  $x \geq 0$ , we conclude that

$$\begin{aligned}
\left| \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} - 1 \right| &= \left| \sqrt{1 + \frac{P_l - Q_l}{Q_l} + \eta' \rho_L \frac{1}{\sqrt{n(P_l \vee Q_l)}}} - 1 \right| \\
&\leq \left| \frac{P_l - Q_l}{Q_l} + \eta' \rho_L \frac{1}{\sqrt{n(P_l \vee Q_l)}} \right| \\
&\stackrel{(a)}{\leq} 2\rho_L \frac{1}{\sqrt{n(P_l \vee Q_l)}},
\end{aligned}$$

where (a) follows because we have assumed that  $\frac{n\Delta_l^2}{P_l \vee Q_l} < 1$ , implying

$$\left| \frac{P_l - Q_l}{Q_l} \right| \leq \sqrt{\frac{P_l \vee Q_l}{nQ_l^2}} = \sqrt{\frac{(P_l \vee Q_l)^2}{nQ_l^2(P_l \vee Q_l)}} \leq \rho_L \frac{1}{\sqrt{n(P_l \vee Q_l)}}.$$

□

The following lemma provides a bound for  $\sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l$ :

**Lemma B.5.** *Suppose*

$$\frac{|\widehat{Q}_l - Q_l|}{Q_l} = \eta', \quad \text{and} \quad \frac{|\widehat{P}_l - P_l|}{P_l} = \eta',$$

where  $\eta' = o(1)$ . For all sufficiently small  $\eta$ , we have

$$\frac{1}{8} \sqrt{P_l Q_l} \leq \frac{1}{2} \sqrt{\widehat{P}_l \widehat{Q}_l} \leq \sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l.$$

*Proof.* We have the sequence of equalities

$$\begin{aligned}
\sqrt{\frac{\widehat{P}_l}{\widehat{Q}_l}} Q_l &= \sqrt{\widehat{P}_l \widehat{Q}_l} \frac{Q_l}{\widehat{Q}_l} = \sqrt{\widehat{P}_l \widehat{Q}_l} \frac{1}{\frac{\widehat{Q}_l - Q_l}{Q_l} + 1} \\
&= \sqrt{\widehat{P}_l \widehat{Q}_l} \left( 1 - \frac{\widehat{Q}_l - Q_l}{Q_l} (1 + \eta') \right) = \sqrt{\widehat{P}_l \widehat{Q}_l} (1 - \eta),
\end{aligned}$$

where the penultimate equality uses the fact that  $\frac{\widehat{Q}_l - Q_l}{Q_l} = \eta' \rightarrow 0$ . Taking  $\eta$  sufficiently small yields the upper bound.

For sufficiently small  $\eta$ , we also have  $\widehat{P}_l \geq \frac{1}{2} P_l$  and  $\widehat{Q}_l \geq \frac{1}{2} Q_l$ , yielding the lower bound. □

**Lemma B.6.** Define  $L_1 = \{l : \frac{n\Delta_l^2}{P_l \vee Q_l} \geq C_{test}^2\}$ . Then

$$C_1 \sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l} \leq I_L \leq C_2 \sum_{l \in L_1} \frac{\Delta_l^2}{P_l \vee Q_l}, \quad (\text{B.14})$$

for some constants  $C_1$  and  $C_2$ .

*Proof.* Throughout the proof, let  $\eta'$  denote a sequence converging to 0, and let  $C$  denote a  $\Theta(1)$  sequence that may change from line to line. First observe that

$$I_L = -2 \log \sum_l \sqrt{P_l Q_l} = -2 \log \left( 1 - \frac{\sum_l (\sqrt{P_l} - \sqrt{Q_l})^2}{2} \right).$$

Using the fact that  $I_L \rightarrow 0$  as  $n \rightarrow \infty$ , so  $\sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 \rightarrow 0$ , we have the following bound for sufficiently large  $n$ :

$$\frac{1}{2} \sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 \leq I_L \leq 2 \sum_l (\sqrt{P_l} - \sqrt{Q_l})^2.$$

Since  $\sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 = \sum_l \frac{\Delta_l^2}{(\sqrt{P_l} + \sqrt{Q_l})^2}$ , there exist constants  $\tilde{C}_1$  and  $\tilde{C}_2$  such that

$$\tilde{C}_1 \sum_l \frac{\Delta_l^2}{P_l \vee Q_l} \leq I_L \leq \tilde{C}_2 \sum_l \frac{\Delta_l^2}{P_l \vee Q_l}.$$

Now note that

$$\sum_{l \in L_1^c} \frac{\Delta_l^2}{P_l \vee Q_l} \leq \frac{LC_{test}^2}{n},$$

and by our assumptions,  $\frac{L}{n} = o(I_L)$ . Hence, the sum over labels in  $L_1$  must be  $\Theta(I_L)$ , which is equivalent to inequality (B.14).  $\square$

We state Lemma 4 from Gao et al. [10], which analyzes the consensus step of the algorithm:

**Lemma B.7.** Let  $\sigma$  and  $\sigma'$  be two clusters such that, for some constant  $C \geq 1$ , the minimum cluster size is at least  $\frac{n}{Ck}$ . Define the map  $\xi : [k] \rightarrow [k]$  according to

$$\xi(k) = \arg \max_{k'} |\{v : \sigma(v) = k\} \cap \{v : \sigma'(v) = k'\}|.$$

If  $\min_{\pi \in S_k} l(\pi(\sigma), \sigma') < \frac{1}{Ck}$ , we have  $\xi \in S_k$  and  $l(\xi(\sigma), \sigma') = \min_{\pi \in S_k} l(\pi(\sigma), \sigma')$ .

We include a simple additional lemma:

**Lemma B.8.** Let  $\sigma, \sigma' : [n] \rightarrow [K]$  be two clusterings where the minimum cluster size of  $\sigma$  is  $T$ . Let  $\pi, \xi \in S_K$  be such that  $d(\pi(\sigma), \sigma') < \frac{T}{2n}$  and  $d(\xi(\sigma), \sigma') < \frac{T}{2n}$ . Then  $\pi = \xi$ .

*Proof.* Suppose the contrary, and choose any  $k$  such that  $\pi(k) \neq \xi(k)$ . We then have

$$|\{\sigma(u) = k\} \cap \{\sigma'(u) \neq \pi(k)\}| < n \cdot d(\pi(\sigma), \sigma') < \frac{T}{2},$$

implying that  $|\{\sigma(u) = k\} \cap \{\sigma'(u) = \pi(k)\}| > \frac{T}{2}$ . But then

$$n \cdot d(\xi(\sigma), \sigma') \geq |\{\sigma(u) = k\} \cap \{\sigma'(u) \neq \xi(k)\}| \geq |\{\sigma(u) = k\} \cap \{\sigma'(u) = \pi(k)\}| \geq \frac{T}{2},$$

a contradiction.  $\square$

## C Proof of Proposition 5.2

*Proof.* We use the notation

$$\begin{aligned}\tilde{P}_l &:= (1 - P_0) \int_{a_l}^{b_l} p(x) dx := (1 - P_0) P_l, \\ \tilde{Q}_l &:= (1 - Q_0) \int_{a_l}^{b_l} q(x) dx := (1 - Q_0) Q_l.\end{aligned}$$

We first show that the likelihood ratio  $\frac{\tilde{P}_l}{\tilde{Q}_l} = \frac{1-P_0}{1-Q_0} \frac{P_l}{Q_l}$  satisfies the claimed bounds. Consider an  $l$  such that  $\text{bin}_l \cap R^c = \emptyset$ . For all  $x \in \text{bin}_l$ , we have  $\frac{1}{\rho} \leq \frac{p(x)}{q(x)} \leq \rho$  by Assumption C2. It follows that  $\frac{P_l}{Q_l} = \frac{\int_{\text{bin}_l} p(x) dx}{\int_{\text{bin}_l} q(x) dx} \leq \rho$ . The lower bound is derived in the same manner. Since  $\frac{1-P_0}{1-Q_0} \in [\frac{1}{c_0}, c_0]$ , we conclude that  $\frac{1}{\rho c_0} \leq \frac{P_l}{Q_l} \leq \rho c_0$ .

Now suppose  $\text{bin}_l \cap R^c \neq \emptyset$ . Since  $R$  is an interval and that  $\mu\{R^c\} = o(H)$ , and since  $L \leq \frac{2}{H}$ , we conclude that  $\mu\{R^c\} < \frac{1}{2L}$  for all sufficiently large  $n$ . Thus, only bins  $[0, \frac{1}{L}]$  and  $[1 - \frac{1}{L}, 1]$  can satisfy  $\text{bin}_l \cap R^c \neq \emptyset$ . Note that Assumption C5 implies both  $p(x)$  and  $q(x)$  are increasing in  $[0, \frac{1}{L}]$  and decreasing in  $[1 - \frac{1}{L}, 1]$ . Define  $P'_l = \int_{\text{bin}_l \cap R} p(x) dx$  and  $Q'_l = \int_{\text{bin}_l \cap R} q(x) dx$ , and define  $P''_l = \int_{\text{bin}_l \cap R^c} p(x) dx$  and  $Q''_l = \int_{\text{bin}_l \cap R^c} q(x) dx$ . Then  $P_l = P'_l + P''_l$  and  $Q_l = Q'_l + Q''_l$ . Note that  $\frac{1}{\rho} \leq \frac{P'_l}{Q'_l} \leq \rho$  by the same argument as before. Furthermore, using the monotonic properties of  $p(x)$  and  $q(x)$  in the relevant intervals, we have

$$P'_l \geq \min_{x \in \text{bin}_l \cap R} \frac{p(x)}{2L} \geq \max_{x \in \text{bin}_l \cap R^c} \frac{p(x)}{2L} \geq P''_l,$$

where the first inequality follows because  $\mu(R^c) \leq \frac{1}{2L}$ , and the second inequality follows from Assumption C5. Similarly,  $Q'_l \geq Q''_l$ . Thus,

$$\frac{P_l}{Q_l} \leq \frac{2P'_l}{Q'_l} \leq 2\rho, \quad \text{and} \quad \frac{P_l}{Q_l} \geq \frac{P'_l}{2Q'_l} \geq \frac{1}{2\rho}.$$

Using the bound on  $\frac{1-P_0}{1-Q_0}$  completes the proof.

We now proceed with bounding  $|I - I_L|$ . Using the simple relation between Renyi divergence and Hellinger distance detailed in Lemma I.1, we have

$$\begin{aligned}I &= (1 + o(1)) \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + \int \left( \sqrt{(1 - P_0)p(x)} - \sqrt{(1 - Q_0)q(x)} \right)^2 dx \right\} \\ &= (1 + o(1)) \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + (\sqrt{1 - P_0} - \sqrt{1 - Q_0})^2 \right. \\ &\quad \left. + \sqrt{(1 - P_0)(1 - Q_0)} \int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx \right\}.\end{aligned}$$

Likewise,

$$\begin{aligned}I_L &= (1 + o(1)) \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + \sum_{l=1}^L (\sqrt{\tilde{P}_l} - \sqrt{\tilde{Q}_l})^2 dx \right\} \\ &= (1 + o(1)) \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + (\sqrt{1 - P_0} - \sqrt{1 - Q_0})^2 + \sqrt{(1 - P_0)(1 - Q_0)} \sum_{l=1}^L (\sqrt{P_l} - \sqrt{Q_l})^2 \right\}.\end{aligned}$$

The key step in completing our proof is the following proposition, proved in Appendix D.1:

**Proposition C.1.** *Under Assumptions C1–C5, we have*

$$\left| \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx - \sum_{l=1}^L (\sqrt{P_l} - \sqrt{Q_l})^2 \right| = o \left( \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right).$$

The claimed result follows from Proposition C.1 by noticing that

$$\begin{aligned} I_L &= (1 + o(1)) \left\{ (\sqrt{P_0} - \sqrt{Q_0})^2 + (\sqrt{1 - P_0} - \sqrt{1 - Q_0})^2 \right. \\ &\quad \left. + \sqrt{(1 - P_0)(1 - Q_0)}(1 + o(1)) \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right\} \\ &= (1 + o(1))I. \end{aligned}$$

The proof of Proposition C.1 contains a number of subparts, which we briefly outline below. Since  $p(x)$  and  $q(x)$  are easier to handle on the interval  $R$ , we initially only concern ourselves with comparing

$$H_R := \int_R (\sqrt{p(x)} - \sqrt{q(x)})^2 dx, \quad \text{and} \quad H_L^R := \sum_{l=1}^L (\sqrt{P_l'} - \sqrt{Q_l'})^2.$$

We first notice that  $\{\text{bin}_l\} \cap R$  constitute an approximately-uniform binning of  $R$ ; i.e., there exist constants  $c_{\text{bin}}$  and  $C_{\text{bin}}$  such that  $\frac{c_{\text{bin}}}{L} \leq |\text{bin}_l \cap R| \leq \frac{C_{\text{bin}}}{L}$ . This is reasoned as follows: Since  $R$  is an interval, we know that  $\text{bin}_l \cap R$  is an interval, as well. The inequality  $|\text{bin}_l \cap R^c| \leq \mu\{R^c\} \leq \frac{1}{2L}$  then implies  $\frac{1}{2L} \leq |\text{bin}_l \cap R| \leq \frac{1}{L}$ .

In a series of lemmas, we show that the approximately-uniform binning of  $R$  leads to several useful bounds on  $H^R$  and  $H_L^R$ . In particular, Lemma D.1 shows that as long as  $L$  grows, we have  $d_L \rightarrow \frac{1}{4}$ , where  $d_L := \sum_l Q_l' \left( \frac{1}{2} \frac{\gamma_l'}{Q_l'} \right)$  and  $\gamma_l' = Q_l' - P_l'$ . In Lemma D.2, we show that  $H^R = \frac{\alpha^2}{4}(1 + \eta)$ , where  $\eta = \Theta(\alpha)$ . Similarly, Lemma D.3 establishes that  $H_L^R = d_L \alpha^2(1 + \eta_L)$ . We combine the results of Lemmas D.1, D.2, and D.3 in Lemma D.4, to show that  $|H^R - H_L^R| = o(H^R)$ . The last step is to bound the difference between the sums and integrals over  $R$  and the entire real line.  $\square$

## D Appendix for Proposition 5.2

### D.1 Proof of Proposition C.1

Let  $a_L$  be an  $o(1)$  sequence such that  $\mu(R^c) \leq a_L H$ . We divide the set of bins into three subsets:

$$\begin{aligned} L_1 &= \{l : \text{bin}_l \cap R^c = \emptyset\}, \\ L_2 &= \{l : \text{bin}_l \cap R^c \neq \emptyset, P_l \vee Q_l \geq 2Ca_L H\}, \\ L_3 &= \{l : \text{bin}_l \cap R^c \neq \emptyset, P_l \vee Q_l \leq 2Ca_L H\}. \end{aligned}$$

For each bin  $l$ , define  $P_l' = \int_{\text{bin}_l \cap R} p(x) dx$  and  $P_l'' = \int_{\text{bin}_l \cap R^c} p(x) dx$ , and likewise define  $Q_l'$  and  $Q_l''$ . We now proceed in two steps:

**Step 1:** We first claim that for all  $l \in L_2$ ,

$$\left| (\sqrt{P_l} - \sqrt{Q_l})^2 - (\sqrt{P_l'} - \sqrt{Q_l'})^2 \right| \leq a_L H.$$

Since  $\mu(R^c) \leq a_L H$ , we have  $P_l'' = \int_{\text{bin}_l \cap R^c} p(x) dx \leq C a_L H$ , and likewise for  $Q_l''$ . Then

$$\begin{aligned} (\sqrt{P_l} - \sqrt{Q_l})^2 - (\sqrt{P_l'} - \sqrt{Q_l'})^2 &= P_l + Q_l - P_l' - Q_l' - 2\sqrt{P_l Q_l} + 2\sqrt{P_l' Q_l'} \\ &\stackrel{(a)}{\leq} P_l'' + Q_l'' - 2\sqrt{P_l'' Q_l''} \\ &\leq P_l'' + Q_l'' \\ &\leq 2C a_L H. \end{aligned}$$

Here, inequality (a) holds by the following reasoning: By the AM-GM inequality, we have  $2\sqrt{P_l' Q_l' P_l'' Q_l''} \leq P_l' Q_l'' + P_l'' Q_l'$ . Thus,

$$P_l' Q_l' + P_l'' Q_l'' + 2\sqrt{P_l' Q_l' P_l'' Q_l''} \leq (P_l' + P_l'')(Q_l' + Q_l'') = P_l Q_l.$$

Taking square roots, we conclude that  $\sqrt{P_l' Q_l'} + \sqrt{P_l'' Q_l''} \leq \sqrt{P_l Q_l}$ , yielding (a).

On the other hand, we have

$$\begin{aligned} \sqrt{P_l Q_l} - \sqrt{P_l' Q_l'} &= \frac{P_l Q_l - P_l' Q_l'}{\sqrt{P_l Q_l} + \sqrt{P_l' Q_l'}} \\ &= \frac{P_l' Q_l'' + P_l'' Q_l' + P_l'' Q_l''}{\sqrt{P_l Q_l} + \sqrt{P_l' Q_l'}} \\ &\leq \frac{P_l' Q_l'' + P_l'' Q_l' + P_l'' Q_l''}{2\sqrt{P_l' Q_l'}} \\ &\leq Q_l'' \frac{P_l'}{2\sqrt{P_l' Q_l'}} + P_l'' \frac{Q_l'}{2\sqrt{P_l' Q_l'}} + Q_l'' \frac{P_l''}{2\sqrt{P_l' Q_l'}}. \end{aligned}$$

Note that because  $P_l'$  and  $Q_l'$  are defined on  $R$ , we have

$$\left| \frac{P_l'}{Q_l'} \right| = \left| \int_{\text{bin}_l \cap R} \frac{p(x)}{Q_l'} dx \right| \leq \int_{\text{bin}_l \cap R} \left| \frac{p(x)}{q(x)} \right| \frac{q(x)}{Q_l'} dx \leq \rho.$$

Thus,  $\sqrt{\frac{P_l'}{Q_l'}} \vee \sqrt{\frac{Q_l'}{P_l'}} \leq \sqrt{\rho}$ . This bounds the terms  $Q_l'' \frac{P_l'}{2\sqrt{P_l' Q_l'}} + P_l'' \frac{Q_l'}{2\sqrt{P_l' Q_l'}} \leq \sqrt{\rho}(Q_l'' + P_l'')$ .

We still need to bound the last term  $\frac{Q_l'' P_l''}{2\sqrt{P_l' Q_l'}}$ . Since  $l \in L_2$ , either  $P_l \geq 2C a_L H$  or  $Q_l \geq 2C a_L H$ . Suppose the former inequality holds; the latter case may be handled in an identical manner. Since  $P_l'' \leq C a_L H$  and  $P_l \geq 2C a_L H$ , we have  $P_l'' \leq P_l'$ , so

$$\frac{Q_l'' P_l''}{2\sqrt{P_l' Q_l'}} \leq Q_l'' \frac{P_l'}{2\sqrt{P_l' Q_l'}} \leq \sqrt{\rho} Q_l''.$$

Putting everything together, we have

$$\sqrt{P_l Q_l} - \sqrt{P_l' Q_l'} \leq 2\sqrt{\rho}(Q_l'' + P_l'').$$

Thus,

$$\begin{aligned} (\sqrt{P_l} - \sqrt{Q_l})^2 - (\sqrt{P_l'} - \sqrt{Q_l'})^2 &= P_l + Q_l - P_l' - Q_l' - 2\sqrt{P_l Q_l} + 2\sqrt{P_l' Q_l'} \\ &\geq P_l'' + Q_l'' - 4\sqrt{\rho}(Q_l'' + P_l'') \\ &\geq -(4\sqrt{\rho} - 1)(P_l'' + Q_l'') \\ &\geq -(4\sqrt{\rho} - 1) \cdot 2C a_L H. \end{aligned}$$

Combining these two bounds yields

$$\left| (\sqrt{P_l} - \sqrt{Q_l})^2 - (\sqrt{P'_l} - \sqrt{Q'_l})^2 \right| \leq C_{C,\rho} a_L H,$$

for an appropriate constant  $C_{C,\rho}$ . This completes step 1.

**Step 2:** In step 2, we verify that  $\{\text{bin}_l\}_{l \in L_1} \cup \{\text{bin}_l \cap R\}_{l \in L_2} \cup \{\text{bin}_l \cap R\}_{l \in L_3}$  constitutes a valid approximately-uniform binning of  $R$ . First, since  $R$  is an interval, it is easy to see that  $\text{bin}_l \cap R$  is also an interval. Second, we have  $|\text{bin}_l \cap R^c| \leq \mu\{R^c\} \leq a_L H$ . Since  $\frac{1}{H} \leq L$  by assumption, we have  $\mu\{R^c\} \leq \frac{a_L}{L}$ , so there exists a constant  $C_{\text{bin}}$  such that  $\frac{C_{\text{bin}}}{L} \leq |\text{bin}_l \cap R| \leq \frac{1}{L}$ .

**Step 3:** We now turn to main step of the proof. We may bound  $|H - H_L|$  as

$$\begin{aligned} & \left| \sum_{l=1}^L (\sqrt{P_l} - \sqrt{Q_l})^2 - \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| \\ &= \left| \sum_{l \in L_1} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_2} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_3} (\sqrt{P_l} - \sqrt{Q_l})^2 - \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| \\ &\stackrel{(a)}{\leq} \left| \sum_{l \in L_1} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_2} (\sqrt{P_l} - \sqrt{Q_l})^2 - \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| + 8C_{a_L} H \\ &\stackrel{(b)}{\leq} \left| \sum_{l \in L_1} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_2} (\sqrt{P'_l} - \sqrt{Q'_l})^2 - \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| + C_{C,\rho} a_L H \\ &\stackrel{(c)}{\leq} \left| \sum_{l \in L_1} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_2} (\sqrt{P'_l} - \sqrt{Q'_l})^2 + \sum_{l \in L_3} (\sqrt{P'_l} - \sqrt{Q'_l})^2 - \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| + C_{C,\rho} a_L H \\ &\stackrel{(d)}{\leq} \left| \sum_{l \in L_1} (\sqrt{P_l} - \sqrt{Q_l})^2 + \sum_{l \in L_2} (\sqrt{P'_l} - \sqrt{Q'_l})^2 + \sum_{l \in L_3} (\sqrt{P'_l} - \sqrt{Q'_l})^2 - \int_R (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right| + C_{C,\rho} a_L H \\ &\stackrel{(e)}{\leq} C_{C,\rho} a_L H, \end{aligned}$$

where (a) follows because  $P_l \vee Q_l \leq 2C_{a_L} H$  for all  $l \in L_3$ , and  $|L_3| \leq 2$ ; (b) follows from step 1 and the fact that  $|L_2| \leq 2$ ; (c) follows because  $P'_l \leq P_l$ , so  $\sum_{l \in L_3} (\sqrt{P'_l} - \sqrt{Q'_l})^2 \leq 2C_{a_L} H$ ; (d) follows because  $\int_{R^c} (\sqrt{p(x)} - \sqrt{q(x)})^2 \leq C\mu\{R^c\} = C_{a_L} H$ ; and (e) follows by Lemma D.4. Since  $a_L \rightarrow 0$ , the conclusion follows.

## D.2 Lemmas for Proposition 5.2

**Lemma D.1.** Let  $d_L = \sum_l Q'_l \left( \frac{\gamma'_l}{Q_l} \right)^2$ . Suppose Assumptions C1–C5 hold. Then  $\lim_{L \rightarrow \infty} d_L = \frac{1}{4}$ .

*Proof.* Let  $h(x)$  be as defined in Assumption C3; in particular,  $|h(x)| \geq \left| \frac{\gamma'(x)}{q(x)} \right| \vee \left| \frac{q'(x)}{q(x)} \right|$ . For a parameter  $0 < \tau < 1$  to be chosen later, we call  $\text{bin}_l$  *good* if

$$\sup_{x \in \text{bin}_l} |h(x)| \leq L^\tau.$$

We first argue that the proportion of bad bins converges to 0 as  $L \rightarrow \infty$ . Since  $h(x)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped, the set  $\{x : |h(x)| \geq L^\tau\}$  is a union of at most two intervals, for all  $L \geq C_s'^{1/\tau}$ . Using

the notation  $B_l = b_l - a_l$ , we have

$$\sum_{l \in \{l : |h(x)| \geq L^\tau\}} B_l \leq \mu(\{x : |h(x)| \geq L^\tau\}) + \frac{4C_{\text{bin}}}{L} \stackrel{(a)}{\leq} \frac{C}{L^{\tau t}} + \frac{4C_{\text{bin}}}{L} \stackrel{(b)}{\leq} \frac{C}{L^{\tau t}},$$

where (a) follows because  $\int_R |h(x)|^t dx < \infty$  by Assumption C4; and (b) follows because  $t \leq 1$  by Assumption C4, and the fact that  $\tau t < 1$  by choice. In particular, the number of bad bins may be bounded as follows:

$$\#\{l : |h(x)| \geq L^\tau\} \leq \frac{CL^{-\tau t}L}{C_{\text{bin}}} \leq CL^{1-\tau t},$$

where we redefine the constant  $C$  suitably. For a bad bin  $l$ , we may bound  $Q'_l \left(\frac{\gamma'_l}{Q'_l}\right)^2$  as follows:

$$\begin{aligned} Q'_l \left(\frac{\gamma'_l}{Q'_l}\right)^2 &= Q'_l \left(\frac{1}{Q'_l} \int_{\text{bin}_l} \gamma(x) dx\right)^2 \\ &= Q'_l \left(\int_{\text{bin}_l} \frac{\gamma(x)}{q(x)} \frac{q(x)}{Q'_l} dx\right)^2 \\ &\stackrel{(a)}{\leq} Q'_l \int_{\text{bin}_l} \frac{q(x)}{Q'_l} \left(\frac{\gamma(x)}{q(x)}\right)^2 dx \\ &\leq \int_{\text{bin}_l} q(x) \left(\frac{\gamma(x)}{q(x)}\right)^2 dx \\ &\stackrel{(b)}{\leq} \left(\int_{\text{bin}_l} q(x) \left|\frac{\gamma(x)}{q(x)}\right|^r dx\right)^{2/r} \left(\int_{\text{bin}_l} q(x) dx\right)^{(r-2)/r} \\ &\stackrel{(c)}{\leq} C(C_{\text{bin}})^{(r-2)/r} L^{-(r-2)/r} = CL^{-(r-2)/r}. \end{aligned}$$

Here, (a) follows from Jensen's inequality, (b) follows from Hölder's inequality, and (c) follows because  $\int_R q(x) \left|\frac{\gamma(x)}{q(x)}\right|^r dx < \infty$  by Assumption C3 and the fact that  $\int_{\text{bin}_l} q(x) dx \leq \frac{CC_{\text{bin}}}{L}$ .

We now have

$$\begin{aligned} d_L &= \sum_{l=1}^L Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 = \sum_{l \text{ good}} Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 + \sum_{l \text{ bad}} Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 \\ &\leq \sum_{l \text{ good}} Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 + CL^{-(r-2)/r} |\{l : l \text{ bad}\}| \\ &\leq \sum_{l \text{ good}} Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 + CL^{1-\tau t - \frac{(r-2)}{r}} = \sum_{l \text{ good}} Q'_l \left(\frac{1}{2} \frac{\gamma'_l}{Q'_l}\right)^2 + CL^{\frac{2}{r} - \tau t}. \end{aligned}$$

For each good bin  $l$ , define  $x_l = \arg \max_{x \in \text{bin}_l} |q(x)|$ . The maximum is attainable since  $q$  is continuous and bounded. Furthermore,

$$\begin{aligned} Q'_l &= \int_{\text{bin}_l} q(x) dx = \int_{a_l}^{b_l} q(x) dx \\ &= \int_{a_l}^{b_l} q(x_l) + q'(c_x)(x - x_l) dx \quad (\text{for some } c_x \in [a_l, b_l]) \end{aligned}$$



$$= B_l q(x_l) + \int_{a_l}^{b_l} q'(c_x)(x - x_l) dx = B_l q(x_l) + B_l^2 \xi_l,$$

where we define  $\xi_l := \frac{1}{B_l^2} \int_{a_l}^{b_l} q'(c_x)(x - x_l) dx$ . We also have

$$\begin{aligned} B_l \left| \frac{\xi_l}{q(x_l)} \right| &\leq \frac{1}{B_l} \int_{a_l}^{b_l} \left| \frac{q'(c_x)}{q(x_l)} \right| |x - x_l| dx \stackrel{(a)}{\leq} \frac{1}{B_l} \int_{a_l}^{b_l} \left| \frac{q'(c_x)}{q(c_x)} \right| |x - x_l| dx \\ &\stackrel{(b)}{\leq} \frac{1}{B_l} \int_{a_l}^{b_l} L^\tau |x - x_l| dx \leq C_{\text{bin}} L^{\tau-1}, \end{aligned}$$

where (a) follows because  $q(c_x) \leq q(x_l)$ , and (b) follows because  $l$  is a good bin, so  $\left| \frac{q'(c_x)}{q(c_x)} \right| \leq L^\tau$ . The last inequality follows because  $B_l \leq \frac{C_{\text{bin}}}{L}$ . We may perform a similar analysis on  $\gamma$ :

$$\gamma'_l = \int_{\text{bin}_l} \gamma(x) dx = \int_{a_l}^{b_l} \gamma(x_l) + \gamma'(c_x)(x - x_l) dx = B_l \gamma(x_l) + B_l^2 \xi'_l,$$

where  $\xi'_l := \frac{1}{B_l^2} \int_{a_l}^{b_l} \gamma'(c_x)(x - x_l) dx$ . It is straightforward to verify that  $B_l \left| \frac{\xi'_l}{q(x_l)} \right| \leq \frac{1}{2} C_{\text{bin}} L^{\tau-1}$ . For any bin  $l$ , we also have

$$Q'_l = \int_{\text{bin}_l} q(x) dx \leq C B_l,$$

where  $C$  is the bound on  $p(x) \vee q(x)$ . Now we look at a single  $Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2$  term for a good bin  $l$ :

$$\begin{aligned} Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 &= \frac{\gamma_l'^2}{Q_l'^2} = \frac{(B_l \gamma(x_l) + B_l^2 \xi'_l)^2}{B_l q(x_l) + B_l^2 \xi_l} \\ &= B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} + B_l \frac{\xi'_l}{q(x_l)} \right)^2 \left( \frac{1}{1 + B_l \frac{\xi_l}{q(x_l)}} \right) \\ &= B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} + B_l \frac{\xi'_l}{q(x_l)} \right)^2 \left( 1 - B_l \frac{\xi_l}{q(x_l)} + \eta_l (B_l \frac{\xi_l}{q(x_l)})^2 \right). \end{aligned}$$

To arrive at the last equality, we assume that  $L \geq C_{\text{bin}}^{1/(1-\tau)}$ . Then  $\left| B_l \frac{\xi'_l}{q(x_l)} \right| \leq \frac{1}{2}$ , so we may take a Taylor approximation. Here,  $\eta_l$  is a constant satisfying  $|\eta_l| \leq 16$ . Expanding the right-hand side, we have

$$\begin{aligned} Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 &= \left( B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 + 2 B_l q(x_l) \frac{\gamma(x_l)}{q(x_l)} B_l \frac{\xi'_l}{q(x_l)} + B_l q(x_l) \left( B_l \frac{\xi'_l}{q(x_l)} \right)^2 \right) \\ &\quad \cdot \left( 1 - B_l \frac{\xi_l}{q(x_l)} + \eta_l (B_l \frac{\xi_l}{q(x_l)})^2 \right). \end{aligned}$$

Again, note that  $\left| B_l \frac{\xi'_l}{q(x_l)} \right| \leq \frac{C_{\text{bin}}}{2} L^{\tau-1}$  and  $\left| B_l \frac{\xi_l}{q(x_l)} \right| \leq \frac{C_{\text{bin}}}{2} L^{\tau-1}$ . Suppose  $L \geq (2C_{\text{bin}})^{1/(1-\tau)}$ , so  $\frac{C_{\text{bin}}}{2} L^{\tau-1} \leq \frac{1}{4}$ . Then

$$\left| B_l \frac{\xi_l}{q(x_l)} \right| + \left| \eta_l (B_l \frac{\xi_l}{q(x_l)})^2 \right| \leq C_{\text{bin}} L^{\tau-1}, \quad \text{and} \quad \left| 1 - B_l \frac{\xi_l}{q(x_l)} + \eta_l (B_l \frac{\xi_l}{q(x_l)})^2 \right| \leq 2.$$

We now bound

$$\begin{aligned}
& \left| Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 - B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 \right| \\
& \leq B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 C_{\text{bin}} L^{\tau-1} + 2B_l q(x_l) \frac{\gamma(x_l)}{q(x_l)} C_{\text{bin}} L^{\tau-1} + B_l q(x_l) C_{\text{bin}}^2 L^{2(\tau-1)}.
\end{aligned}$$

The third term is bounded by  $C_1 L^{2\tau-3}$ , for a suitable constant  $C_1$ . To bound the second term, we split into two cases:

**Case 1:**  $\left| \frac{\gamma(x_l)}{q(x_l)} \right| \geq 1$ . Then  $q(x) \left| \frac{\gamma(x_l)}{q(x_l)} \right| \leq q(x) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2$ .

**Case 2:**  $\left| \frac{\gamma(x_l)}{q(x_l)} \right| \leq 1$ . The second term is bounded by  $2B_l C C_{\text{bin}} L^{\tau-1} \leq C_2 L^{\tau-2}$ , for some constant  $C_2$ .

In either case, we have

$$\left| Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 - B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 \right| \leq C_3 B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 L^{\tau-1} + C_4 L^{\tau-2}.$$

Define  $d_R = \sum_{l \text{ good}} B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2$ . Then

$$\begin{aligned}
|d_L - d_R| &= \left| \sum_l Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 - \sum_{l \text{ good}} B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 \right| \\
&\leq \sum_{l \text{ good}} \left| Q'_l \left( \frac{\gamma'_l}{Q'_l} \right)^2 - B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 \right| + C_{M,M',K,C} L^{\frac{2}{r}-\tau t} \\
&\leq C_3 d_R L^{\tau-1} + L \cdot C_4 L^{\tau-2} + C_5 L^{\frac{2}{r}-\tau t} \\
&\leq C_3 d_R L^{\tau-1} + C_4 L^{\tau-1} + C_5 L^{\frac{2}{r}-\tau t} \\
&\leq C_3 d_R L^{\frac{2-\tau t}{r(1+t)}} + C_6 L^{\frac{2-\tau t}{r(1+t)}},
\end{aligned}$$

where we have made the choice  $\tau = \frac{2+r}{r(1+t)}$  in the last inequality to balance  $L^{\tau-1}$  and  $L^{2/r-\tau t}$ . Notice that  $0 < \tau < 1$  by Assumption C4, since  $rt > 2$ . Furthermore,  $|d_L - d_R| = o(d_R) + o(1)$ .

In a similar manner, we bound  $|d_R - d|$ . We use the same definition of good and bad bins as before, and obtain

$$\begin{aligned}
d &= \int_R q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx = \sum_{l=1}^L \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx \\
&= \sum_{l \text{ good}} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx + \sum_{l \text{ bad}} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx \\
&\leq \sum_{l \text{ good}} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx + |\{l : l \text{ bad}\}| C L^{-\frac{2}{r}} \\
&\leq \sum_{l \text{ good}} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx + C L^{\frac{2}{r}-\tau t}.
\end{aligned}$$

The bound on the second term follows from the previous analysis. For the first term, note that for all  $x \in \text{bin}_l$ , we have

$$q(x) = q(x_l) + q'(c_l)(x - x_l), \quad \text{and} \quad \gamma(x) = \gamma(x_l) + \gamma'(c'_l)(x - x_l),$$

where  $c_l, c'_l \in \text{bin}_l$  depend implicitly on  $x$ . For  $\text{bin}_l$ , we have

$$\begin{aligned} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 &= \int_{\text{bin}_l} \frac{(\gamma(x_l) + \gamma'(c'_l)(x - x_l))^2}{q(x_l) + q'(c_l)(x - x_l)} dx \\ &= \int_{\text{bin}_l} q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} + \frac{\gamma'(c'_l)}{q(x_l)}(x - x_l) \right)^2 \left( \frac{1}{1 + \frac{q'(c_l)}{q(x_l)}(x - x_l)} \right) dx. \end{aligned}$$

Denote

$$T_1 = \frac{q'(c_l)}{q(x_l)}(x - x_l), \quad \text{and} \quad T_2 = \frac{\gamma'(c'_l)}{q(x_l)}(x - x_l).$$

Observe that  $|x - x_l| \leq B_l$  and

$$\left| \frac{\gamma'(c'_l)}{q(x_l)} \right| \leq \left| \frac{\gamma'(c'_l)}{q(c'_l)} \right| \leq L^\tau.$$

Similarly,  $\left| \frac{q'(c_l)}{q(x_l)} \right| \leq L^\tau$ . Hence,  $|T_1|, |T_2| \leq C_{\text{bin}} L^{\tau-1}$ . Now suppose  $C_{\text{bin}} L^{\tau-1} \leq \frac{1}{2}$ , which is satisfied if  $L \geq (2C_{\text{bin}})^{\frac{1}{1-\tau}}$ . We obtain

$$\begin{aligned} \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx &= \int_{\text{bin}_l} q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} + T_2 \right)^2 \left( \frac{1}{1 + T_1} \right) dx \\ &= \int_{\text{bin}_l} \left( q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 + q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right) T_2 + q(x_l) T_2^2 \right) (1 - T_1 + \eta T_1^2) dx, \end{aligned}$$

where  $\eta$  is some function of  $x$  satisfying  $|\eta| \leq 16$ . Thus,

$$\begin{aligned} \left| \int_{\text{bin}_l} q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 dx - B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 \right| \\ \leq B_l q(x_l) \left( \frac{\gamma(x_l)}{q(x_l)} \right)^2 C_{\text{bin}} L^{\tau-1} + B_l q(x_l) \frac{\gamma(x_l)}{q(x_l)} C_{\text{bin}} L^{\tau-1} + B_l q(x_l) C_{\text{bin}}^2 L^{2(\tau-1)}. \end{aligned}$$

The same analysis used to bound  $|d_L - d_R|$  implies that  $|d - d_R| = o(d_R) + o(1)$ . Since  $d = \frac{1}{4}$ , we have  $d_R \rightarrow \frac{1}{4}$ , which in turn implies  $d_L \rightarrow \frac{1}{4}$ . This completes the proof.  $\square$

**Lemma D.2.** *Let*

$$\begin{aligned} H^R &= \int_R (\sqrt{p(x)} - \sqrt{q(x)})^2 dx, \\ \delta(x) &= q(x) - p(x), \\ \alpha^2 &= \int_R q(x) \left( \frac{\delta(x)}{q(x)} \right)^2 dx, \\ \gamma(x) &= \frac{\delta(x)}{\alpha}. \end{aligned}$$

Suppose Assumptions C1–C5 hold. Then  $H^R = \frac{\alpha^2}{4}(1 + \eta)$ , where  $|\eta| \leq C(\alpha + \alpha^2)$  for some constant  $C$ . In particular, if  $H^R \rightarrow 0$ , then  $\alpha \rightarrow 0$  and  $\eta \rightarrow 0$ .

*Proof.* We write

$$H^R = \int_R (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \int_R (\sqrt{q(x)} - \sqrt{q(x) - \delta(x)})^2 dx = \int_R q(x) \left( 1 - \sqrt{1 - \frac{\delta(x)}{q(x)}} \right)^2 dx.$$

By convention, let  $\frac{\delta(x)}{q(x)} = 0$  whenever  $q(x) = p(x) = 0$ . Thus, we may define  $\xi(x) = 1 - \frac{1}{2} \frac{\delta(x)}{q(x)} - \sqrt{1 - \frac{\delta(x)}{q(x)}}$  for  $x \in [0, 1]$  and rewrite

$$\begin{aligned} H^R &= \int_R q(x) \left( 1 - \left( 1 - \frac{1}{2} \frac{\delta(x)}{q(x)} + \xi(x) \right) \right)^2 dx \\ &= \int_R q(x) \left( \frac{1}{2} \frac{\delta(x)}{q(x)} + \xi(x) \right)^2 dx \\ &= \int_R q(x) \left( \frac{1}{2} \frac{\delta(x)}{q(x)} \right)^2 (1 + \xi_2(x))^2 dx, \end{aligned}$$

where  $\xi_2(x) = \frac{2\xi(x)}{\delta(x)/q(x)}$  if  $\delta(x) \neq 0$ , and  $\xi_2(x) = 0$  if  $\delta(x) = 0$ . Thus,

$$\int_R \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx = (1 + \eta) \frac{\alpha^2}{4},$$

where

$$\eta = \frac{\int_R q(x) \left( \frac{1}{2} \frac{\delta(x)}{q(x)} \right)^2 (\xi_2(x)^2 + 2\xi_2(x)) dx}{\alpha^2/4} = \int_R q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 (\xi_2(x)^2 + 2\xi_2(x)) dx.$$

By Lemma L.3, we have  $\xi_2(x) \leq 2 \left| \frac{\delta(x)}{q(x)} \right|$ , implying that

$$\begin{aligned} |\eta| &\leq \int_R q(x) \left( \frac{\gamma(x)}{q(x)} \right)^2 \left( 4 \left| \frac{\delta(x)}{q(x)} \right|^2 + 4 \left| \frac{\delta(x)}{q(x)} \right| \right) dx \\ &= 4\alpha^2 \int_R q(x) \left| \frac{\gamma(x)}{q(x)} \right|^4 dx + 4\alpha \int_R q(x) \left| \frac{\gamma(x)}{q(x)} \right|^3 dx \\ &\leq C(\alpha^2 + \alpha), \end{aligned}$$

using the finiteness of integrals in Assumption C3. □

**Lemma D.3.** *Let*

$$\begin{aligned} H_L^R &= \sum_{l=1}^L \left( \sqrt{P'_l} - \sqrt{Q'_l} \right)^2, \\ \delta(x) &= q(x) - p(x), \\ \alpha^2 &= \int_R q(x) \left( \frac{\delta(x)}{q(x)} \right)^2 dx, \\ \gamma(x) &= \frac{\delta(x)}{\alpha} dx. \end{aligned}$$

Suppose that Assumptions C1–C5 hold. Then  $H_L^R = d_L(1 + \eta_L)$ , where  $d_L = \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^2 dx$ ,  $\gamma'_l = \frac{Q'_l - P'_l}{\alpha}$  and  $\sup_L |\eta_L| \leq C(\alpha + \alpha^2)$ , for some constant  $C$ .

*Proof.* Let  $\delta_l = Q'_l - P'_l$ . We have

$$H_L^R = \sum_{l=1}^L (\sqrt{P'_l} - \sqrt{Q'_l})^2 = \sum_{l=1}^L Q'_l \left( 1 - \sqrt{\frac{P'_l}{Q'_l}} \right)^2$$

$$= \sum_{l=1}^L Q'_l \left( 1 - \sqrt{1 - \frac{\delta_l}{Q'_l}} \right)^2 = \sum_{l=1}^L Q'_l \left( 1 - \left( 1 - \frac{1}{2} \frac{\delta_l}{Q'_l} - \xi_l \right) \right)^2,$$

where by convention, we define  $\frac{\delta_l}{Q'_l} = 0$  when  $Q'_l, P'_l = 0$ , and we use the shorthand  $\xi_l = 1 - \frac{1}{2} \frac{\delta_l}{Q'_l} - \sqrt{1 - \frac{\delta_l}{Q'_l}}$ . Hence,

$$H_L^R = \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} + \xi_l \right)^2 = \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2 (1 + \xi_{2l})^2,$$

where  $\xi_{2l} = 0$  if  $\frac{\delta_l}{Q'_l} = 0$ , and  $\xi_{2l} = 2\xi_l \frac{Q'_l}{\delta_l}$  otherwise. Then

$$H_L^R = (1 + \eta_L) \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2,$$

where  $\eta_L = \frac{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2 (2\xi_{2l} + \xi_{2l}^2)}{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2}$ . By Lemma [L.3](#), we have  $|\xi_{2l}| \leq 2 \left| \frac{\delta_l}{Q'_l} \right|$ . Therefore,

$$\begin{aligned} |\eta_L| &= \left| \frac{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2 (2\xi_{2l} - \xi_{2l}^2)}{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\delta_l}{Q'_l} \right)^2} \right| \leq \frac{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^2 (2|\xi_{2l}| + \xi_{2l}^2)}{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^2} \\ &\leq 4 \frac{\alpha \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^3 + \alpha^2 \sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^4}{\sum_{l=1}^L Q'_l \left( \frac{1}{2} \frac{\gamma'_l}{Q'_l} \right)^2}. \end{aligned}$$

The denominator tends to  $\frac{1}{4}$  by Lemma [D.1](#) and may be bounded by  $1/(2C')$  for large enough  $L$ . To bound the numerator, note that for a single  $l$ , we have

$$\int_{a_l}^{b_l} \frac{q(x)}{Q'_l} \left| \frac{\gamma(x)}{q(x)} \right|^3 dx \geq \left| \int_{\text{bin}_l} \frac{q(x)}{Q'_l} \frac{\gamma(x)}{q(x)} dx \right|^3 = \left| \frac{\gamma'_l}{Q'_l} \right|^3.$$

Therefore,

$$\begin{aligned} \sum_{l=1}^L Q'_l \left| \frac{\gamma'_l}{Q'_l} \right|^3 &\leq \int_R q(x) \left| \frac{\gamma(x)}{q(x)} \right|^3 \leq M, \text{ and} \\ \sum_{l=1}^L Q'_l \left| \frac{\gamma'_l}{Q'_l} \right|^4 &\leq \int_R q(x) \left| \frac{\gamma(x)}{q(x)} \right|^4 \leq M, \end{aligned}$$

implying that  $|\eta_L| \leq (2\alpha + \alpha^2)2C'M$ .  $\square$

**Lemma D.4.** Suppose Assumptions A1–A4 hold. For any sequences  $L_n, \alpha_n \rightarrow \infty$ , we have  $H_L^R = H^R(1 + o(1))$ ; i.e.,

$$\left| \frac{\sum_l (\sqrt{P'_l} - \sqrt{Q_l})^2}{\int_R (\sqrt{p(x)} - \sqrt{q(x)})^2 dx} - 1 \right| \rightarrow 0.$$

*Proof.* By Lemmas [D.2](#) and [D.3](#), we have

$$|H_L^R - H^R| = \left| d_L \alpha^2 (1 + \eta_L) - \frac{\alpha^2}{4} (1 + \eta) \right|,$$

implying that

$$\left| \frac{H_L^R}{H^R} - 1 \right| = \left| 4d_L \frac{(1 + \eta_L)}{1 + \eta} - 1 \right|,$$

where  $|\eta|, |\eta_L| \leq C_1(\alpha + \alpha^2)$  for all  $L$ . Thus,

$$\lim_{\alpha_n \rightarrow 0} \sup_L \left| \frac{1 + \eta_L}{1 + \eta} - 1 \right| = 0.$$

Furthermore, by Lemma D.1, we have  $|4d_L - 1| \rightarrow 0$ , uniformly for all  $\alpha$ . Thus,  $\left| \frac{H_L^R}{H^R} - 1 \right| \rightarrow 0$ , completing the proof.  $\square$

## E Proof of Proposition 5.3

First, we prove the bounds on  $\frac{\tilde{P}_l}{Q_l}$ . Define

$$R = \left\{ x \in [0, 1] : \left| \log \frac{p(x)}{q(x)} \right| \leq C(2L)^{1/r} \right\},$$

where  $C = \left( \int \left| \log \frac{p(x)}{q(x)} \right|^r dx \right)^{1/r}$  is a constant. Since  $\int \left| \log \frac{p(x)}{q(x)} \right|^r dx < \infty$ , Markov's Inequality implies  $\mu\{R^c\} \leq \frac{1}{2L}$ .

The remainder of the proof follows the argument used to prove Proposition 5.2, except for the final step, where we need to show that

$$\left| \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx - \sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 \right| = o\left( \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \right) = o(1).$$

We establish this fact in the following proposition:

**Proposition E.1.** *Let Assumptions C1' and C3' be satisfied. Let  $\text{bin}_l = [a_l, b_l]$  be a uniform binning of  $[0, 1]$ , for  $l = 1, \dots, L$ , and let  $P_l = \int_{\text{bin}_l} p(x) dx$  and  $Q_l = \int_{\text{bin}_l} q(x) dx$ . Then*

$$\left| \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx - \sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 \right| \rightarrow 0.$$

*Proof.* We use a similar argument to the proof of Proposition D.1. First observe that

$$\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \int p(x) dx + \int q(x) dx - 2 \int \sqrt{p(x)q(x)} dx = 2 - 2 \int \sqrt{p(x)q(x)}$$

and

$$\sum_l (\sqrt{P_l} - \sqrt{Q_l})^2 = \sum_l P_l + \sum_l Q_l - 2 \sum_l \sqrt{P_l Q_l}.$$

Thus, we only need to show that

$$\left| \int \sqrt{p(x)q(x)} dx - \sum_l \sqrt{P_l Q_l} \right| \rightarrow 0.$$

We have  $|h(x)| \geq \left| \frac{p'(x)}{p(x)} \right| \vee \left| \frac{q'(x)}{q(x)} \right|$ . Let  $0 < \tau < 1$ . We call  $\text{bin}_l$  *good* if

$$\sup_{x \in \text{bin}_l} |h(x)| \leq L^\tau.$$

We now argue that the proportion of bad bins converges to 0 as  $L \rightarrow \infty$ : Since  $h(x)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped, the set  $\{x : |h_n(x)| \geq L^\tau\}$  is a union of at most two intervals, for all  $L \geq C_s'^{1/\tau}$ . Hence,

$$\begin{aligned} \sum_{l \in \{l : \sup_{x \in \text{bin}_l} |h(x)| \geq L^\tau\}} B_l &\leq \mu \left( \left\{ x : \sup_{x \in \text{bin}_l} |h(x)| \geq L^\tau \right\} \right) + 4C_{\text{bin}} L^{-1} \\ &\stackrel{(a)}{\leq} CL^{-\tau t} + 4C_{\text{bin}} L^{-1} \stackrel{(b)}{\leq} CL^{-\tau t}, \end{aligned}$$

where (a) follows because  $\int_R |h(x)|^t dx < \infty$  by Assumption C3'; and (b) follows because  $t \leq 1$ , so  $\tau t < 1$  and the first term dominates. We now bound the number of bad bins:

$$\#\{l : |h(x)| \geq L^\tau\} \leq \frac{CL^{-\tau t} L}{C_{\text{bin}}} \leq CL^{1-\tau t}.$$

For a bad bin, we have  $P_l, Q_l \leq \frac{CC_{\text{bin}}}{L}$  and  $\int_{\text{bin}_l} (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \leq \frac{2CC_{\text{bin}}}{L}$ .

We now consider a good bin  $l$ . Let  $x_l$  be  $\arg \max_{x \in \text{bin}_l} p(x)$ . The argmax is attainable since  $p$  is continuous and bounded. We have

$$P_l = \int_{a_l}^{b_l} p(x) dx = \int_{a_l}^{b_l} p(x_l) + p'(c_x)(x - x_l) dx = B_l p(x_l) + B_l^2 \xi_l,$$

where  $\xi_l = \frac{1}{B_l^2} \int_{a_l}^{b_l} p'(c_x)(x - x_l) dx$ . Furthermore,

$$\begin{aligned} B_l \left| \frac{\xi_l}{p(x_l)} \right| &\leq \frac{1}{B_l} \int_{a_l}^{b_l} \left| \frac{p'(c_x)}{p(x_l)} \right| |x - x_l| dx \leq \frac{1}{B_l} \int_{a_l}^{b_l} \left| \frac{p'(c_x)}{p(c_x)} \right| |x - x_l| dx \\ &\leq \frac{1}{B_l} \int_{a_l}^{b_l} L^\tau |x - x_l| dx \leq C_{\text{bin}} L^{\tau-1}. \end{aligned}$$

Likewise, define  $x'_l = \arg \max_{x \in \text{bin}_l} q(x)$ . We have  $Q_l = B_l q(x'_l) + B_l^2 \xi'_l$ , where

$$\xi'_l := \frac{1}{B_l} \int_{a_l}^{b_l} q'(c_x)(x - x'_l) dx.$$

We can also bound  $B_l \left| \frac{\xi'_l}{q(x'_l)} \right| \leq C_{\text{bin}} L^{\tau-1}$ . Thus,

$$\begin{aligned} \sqrt{P_l Q_l} &= \sqrt{(B_l p(x_l) + B_l^2 \xi_l)(B_l q(x'_l) + B_l^2 \xi'_l)} \\ &= \sqrt{p(x_l) q(x'_l)} \sqrt{(B_l + B_l^2 \frac{\xi_l}{p(x_l)})(B_l + B_l^2 \frac{\xi'_l}{q(x'_l)})} \\ &= \sqrt{p(x_l) q(x'_l)} B_l \sqrt{(1 + B_l \frac{\xi_l}{p(x_l)})(1 + B_l \frac{\xi'_l}{q(x'_l)})}. \end{aligned}$$

By our bounds on  $B_l \frac{\xi_l}{p(x_l)}$  and  $B_l \frac{\xi'_l}{q(x'_l)}$ , we can bound the nuisance term as

$$\sqrt{(1 + B_l \frac{\xi_l}{p(x_l)})(1 + B_l \frac{\xi'_l}{q(x'_l)})} \leq \sqrt{1 + C_{\text{bin}} L^{\tau-1}(1 + o(1))} \leq 1 + \frac{1}{2} L^{\tau-1}(1 + o(1)).$$

It is clear that  $B_l \sqrt{p(x_l) q(x'_l)} \leq B_l C$ . Therefore,

$$\left| \sqrt{P_l Q_l} - \sqrt{p(x_l) q(x'_l)} B_l \right| \leq B_l C L^{\tau-1}(1 + o(1)), \quad (\text{E.1})$$

and likewise,

$$\begin{aligned}
\int_{a_l}^{b_l} \sqrt{p(x)q(x)} dx &= \int_{a_l}^{b_l} \sqrt{p(x)q(x)} dx \\
&= \int_{a_l}^{b_l} \sqrt{(p(x_l) + p'(c_x)(x - x_l))(q(x'_l) + q'(c'_x)(x - x'_l))} dx \\
&= \int_{a_l}^{b_l} \sqrt{p(x_l)q(x'_l)} \left( \sqrt{1 + (x - x_l) \frac{p'(c_x)}{p(x_l)} + (x - x'_l) \frac{q'(c'_x)}{q(x'_l)} + (x - x_l)(x - x'_l) \frac{p'(c_x)}{p(x_l)} \frac{q'(c'_x)}{q(x'_l)}} \right) dx.
\end{aligned}$$

Since

$$\begin{aligned}
\left| (x - x_l) \frac{p'(c_x)}{p(x_l)} \right| &\leq B_l \left| \frac{p'(c_x)}{p(c_x)} \right| \leq L^{\tau-1}, \\
\left| (x - x_l) \frac{q'(c'_x)}{q(x'_l)} \right| &\leq B_l \left| \frac{q'(c'_x)}{q(c'_x)} \right| \leq L^{\tau-1},
\end{aligned}$$

we may bound the nuisance term as follows:

$$\begin{aligned}
\sqrt{1 + (x - x_l) \frac{p'(c_x)}{p(x_l)} + (x - x'_l) \frac{q'(c'_x)}{q(x'_l)} + (x - x_l)(x - x'_l) \frac{p'(c_x)}{p(x_l)} \frac{q'(c'_x)}{q(x'_l)}} &\leq \sqrt{1 + C_{\text{bin}} L^{\tau-1} (1 + o(1))} \\
&\leq 1 + \frac{1}{2} C_{\text{bin}} L^{\tau-1} (1 + o(1)).
\end{aligned}$$

The term  $B_l \sqrt{p(x_l)q(x'_l)}$  is bounded by  $B_l C$ . Hence,

$$\left| \int_{a_l}^{b_l} \sqrt{p(x)q(x)} dx - B_l \sqrt{p(x_l)q(x'_l)} \right| \leq B_l C C_{\text{bin}} L^{\tau-1} \quad (\text{E.2})$$

By combining inequalities (E.1) and (E.2), we have

$$\left| \sqrt{P_l Q_l} - \int_{a_l}^{b_l} \sqrt{p(x)q(x)} dx \right| \leq B_l C C_{\text{bin}} L^{\tau-1}.$$

Hence,

$$\begin{aligned}
\left| \sum_l \sqrt{P_l Q_l} - \int \sqrt{p(x)q(x)} dx \right| &\leq \sum_{l: l \text{ bad}} B_l C + \sum_{l: l \text{ good}} \left| \sqrt{P_l Q_l} - \int_{a_l}^{b_l} \sqrt{p(x)q(x)} dx \right| \\
&\leq C L^{-\tau t} + \sum_{l: l \text{ good}} B_l C C_{\text{bin}} L^{\tau-1} \\
&\leq C L^{-\tau t} + C C_{\text{bin}} L^{\tau-1}.
\end{aligned}$$

Setting  $\tau = \frac{1}{1+t}$ , we obtain

$$\left| \sum_l \sqrt{P_l Q_l} - \int \sqrt{p(x)q(x)} dx \right| \rightarrow 0,$$

completing the proof.  $\square$

## F Proofs of Theorems 4.1 and Theorem 4.2

We now outline the proofs of Theorems 4.1 and 4.2, with proofs of supporting propositions in the succeeding subsections.



### F.1 Main argument: Proof of Theorem 4.1

By the argument outlined in Section 5.3, the divergences  $I$  and  $H$  do not change after transforming the densities  $p(x)$  and  $q(x)$  according to  $\Phi$ . Proposition F.1 shows that under Assumptions A1–A5, Assumptions C1–C5 are also satisfied.

Furthermore, our assumption that  $L = o(\frac{1}{H})$  implies  $L \leq \frac{2}{H}$  for sufficiently large  $L$ . Hence, Proposition 5.2 applies, and we may conclude that after transformation and discretization, the label probabilities satisfy  $\frac{1}{2c_0\rho} \leq \frac{P_l}{Q_l} \leq 2c_0\rho$ , for all  $l$ . Using the assumption  $L = o(nI)$  and the fact that  $I_L = I(1 + o(1))$  from Proposition 5.2, we also have  $L = o(nI_L)$ , so we may use Proposition 5.1 (with  $\rho_L = 2c_0\rho$ ) to obtain

$$\lim_{n \rightarrow \infty} P\left(l(\hat{\sigma}, \sigma_0) \leq \exp\left(-\frac{nI_L}{\beta K}(1 + o(1))\right)\right) \rightarrow 1.$$

The theorem then follows from the fact that  $I_L = I(1 + o(1))$ .

### F.2 Main argument: Proof of Theorem 4.2

The proof parallels the argument for Theorem 4.1 outlined above. Proposition F.2 establishes that Assumptions A1'–A4' imply Assumptions C1'–C4'. Hence, Proposition 5.3 applies, and we may conclude that after transformation and discretization, the label probabilities satisfy

$$\frac{1}{2c_0 \exp(L^{1/r})} \leq \frac{P_l}{Q_l} \leq 2c_0 \exp(L^{1/r}),$$

for all  $l$ , and  $I_L = I(1 + o(1))$ . Therefore, we may again apply Proposition 5.1 (with  $\rho_L = 2c_0 \exp(L^{1/r})$ ) to obtain

$$\lim_{n \rightarrow \infty} P\left(l(\hat{\sigma}, \sigma_0) \leq \exp\left(-\frac{nI_L}{\beta K}(1 + o(1))\right)\right) \rightarrow 1.$$

The theorem follows from the fact that  $I_L = I(1 + o(1))$ .

### F.3 Transformation Analysis

**Proposition F.1.** *Let  $p(x)$  and  $q(x)$  be densities over  $S$ , where  $S = \mathbb{R}$  or  $S = \mathbb{R}^+$ , and let  $\Phi : S \rightarrow [0, 1]$  be a CDF such that  $\phi = \Phi'$  is positive and continuous. Suppose Assumptions A1–A5 hold.*

*The following conditions are satisfied for  $p_\Phi(z) = \frac{p(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$  and  $q_\Phi(z) = \frac{q(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$ :*

*C1  $p_\Phi(z), q_\Phi(z) > 0$  on  $(0, 1)$ , and  $\sup_z \{p_\Phi(z) \vee q_\Phi(z)\} < \infty$ .*

*C2 There exists a subinterval  $R \subseteq [0, 1]$  such that*

*(a) for all  $z \in R$ ,  $\frac{1}{\rho} \leq \left| \frac{p_\Phi(z)}{q_\Phi(z)} \right| \leq \rho$ , where  $\rho$  is an absolute constant, and*

*(b)  $\mu\{R^c\} = o(H)$ , where  $\mu$  is the Lebesgue measure.*

*C3 Let  $\alpha^2 = \int_R \frac{(p_\Phi(z) - q_\Phi(z))^2}{q_\Phi(z)} dz$  and  $\gamma_\Phi(z) = \frac{q_\Phi(z) - p_\Phi(z)}{\alpha}$ . Then  $\int_R q_\Phi(z) \left| \frac{\gamma_\Phi(z)}{q_\Phi(z)} \right|^r dz < \infty$ , for an absolute constant  $r \geq 4$ .*

*C4 There exists  $h_\Phi(z)$  such that*

*(a)  $h_\Phi(z) \geq \max \left\{ \left| \frac{\gamma'_\Phi(z)}{q_\Phi(z)} \right|, \left| \frac{q'_\Phi(z)}{q_\Phi(z)} \right| \right\}$ ,*

*(b)  $h_\Phi(z)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped, for absolute constants  $c'_{s1}, c'_{s2}$ , and  $C'_s$ , and*

(c)  $\int_R |h_\Phi(z)|^t dz < \infty$  for an absolute constant  $\frac{2}{r} < t < 1$ .

C5  $p'_\Phi(z), q'_\Phi(z) \geq 0$  and for all  $z < c'_{s1}$ , and  $p'_\Phi(z), q'_\Phi(z) \leq 0$  for all  $z > c'_{s2}$ .

*Proof.* **C1** follows from A1 and the condition that  $\phi$  is positive and continuous.

To prove **C2**, assume that A2 is true, and let  $R$  be a subinterval of  $\mathbb{R}$  such that  $\frac{1}{\rho} \leq \frac{p(x)}{q(x)} \leq \rho$ . Define  $R_\Phi = \{z \in [0, 1] : \Phi^{-1}(z) \in R\}$ , so  $\mathbf{1}_{R_\Phi}(z) = \mathbf{1}_R(\Phi^{-1}(z))$ . Then  $R_\Phi$  is clearly an interval, and  $\Phi\{R^c\} = \mu\{R_\Phi^c\}$ .

**C3** follows from A3 via a change of variables.

It remains to prove **C4** and **C5**. We first prove **C5**. Note that

$$p'_\Phi(z) = \frac{p'(\Phi^{-1}(z)) - p(\Phi^{-1}(z)) \frac{\phi'(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}}{\phi(\Phi^{-1}(z))^2}.$$

Therefore,  $p'_\Phi(z) \geq 0$  if and only if  $p'(x) \geq p(x) \frac{\phi'(x)}{\phi(x)}$ , and likewise for  $q'_\Phi(z)$ .

Moving onto **C4**, we first construct  $h(z)$ . For ease of presentation, let  $x = \Phi^{-1}(z)$ . We then have

$$\frac{q'_\Phi(z)}{q_\Phi(z)} = \frac{q'(x)}{q(x)} \frac{1}{\phi(x)} - \frac{\phi'(x)}{\phi(x)} \frac{1}{\phi(x)},$$

implying that

$$\left| \frac{q'_\Phi(z)}{q_\Phi(z)} \right| \leq \left| \frac{q'(x)}{q(x)} \frac{1}{\phi(x)} \right| + \left| \frac{\phi'(x)}{\phi(x)} \frac{1}{\phi(x)} \right| \lesssim (h(x) + 1) \frac{1}{\phi(x)},$$

where the last inequality follows because  $\left| \frac{\phi'(x)}{\phi(x)} \right|$  is bounded. Furthermore,

$$\begin{aligned} \frac{\gamma'_\Phi(z)}{q_\Phi(z)} &= \frac{1}{\alpha} \frac{p'(x) - p(x) \frac{\phi'(x)}{\phi(x)} - q'(x) + q(x) \frac{\phi'(x)}{\phi(x)}}{q(x)\phi(x)} \\ &= \left( \frac{1}{\alpha} \frac{p'(x) - q'(x)}{q(x)} - \frac{1}{\alpha} \frac{p(x) - q(x)}{q(x)} \frac{\phi'(x)}{\phi(x)} \right) \frac{1}{\phi(x)}, \end{aligned}$$

so

$$\begin{aligned} \left| \frac{\gamma'_\Phi(z)}{q_\Phi(z)} \right| &\leq \left| \frac{1}{\alpha} \frac{p'(x) - q'(x)}{q(x)} \right| \frac{1}{\phi(x)} + \left| \frac{1}{\alpha} \frac{p(x) - q(x)}{q(x)} \right| \left| \frac{\phi'(x)}{\phi(x)} \right| \frac{1}{\phi(x)} \\ &= \left| \frac{\gamma'(x)}{q(x)} \right| \frac{1}{\phi(x)} + \left| \frac{\gamma(x)}{q(x)} \right| \left| \frac{\phi'(x)}{\phi(x)} \right| \frac{1}{\phi(x)} \\ &\stackrel{(a)}{\lesssim} h(x) \frac{1}{\phi(x)}, \end{aligned}$$

where (a) follows because  $\left| \frac{\phi'(x)}{\phi(x)} \right|$  is bounded. We want to take  $h_\Phi(z) \simeq (h(x) + 1) \frac{1}{\phi(x)}$ , but we use a modified upper bound to ensure that  $h_\Phi(z)$  is bowl-shaped. Let  $\psi(x) = \max \left\{ \frac{1}{\phi(c_{s1})}, \frac{1}{\phi(c_{s2})}, \frac{1}{\phi(x)} \right\}$ . We then take

$$h_\Phi(z) \simeq h(x)\psi(x) = h(\Phi^{-1}(z))\psi(\Phi^{-1}(z)).$$

Note that  $(h(x) + 1)$  is  $(c_{s1}, c_{s2}, C_s + 1)$ -bowl-shaped, and  $\phi$  is unimodal, so  $\frac{1}{\phi(x)}$  is quasi-convex. Hence,  $\psi(x)$  is quasi-convex and has a mode lying in  $[c_{s1}, c_{s2}]$ . Therefore,  $(h(x) + 1)\psi(x)$  is  $(c_{s1}, c_{s2}, C'_s)$ -bowl-shaped, where  $C'_s \simeq (C_s + 1) \left( \frac{1}{\phi(c_{s1})} \vee \frac{1}{\phi(c_{s2})} \right)$ . This shows that  $h_\Phi(z)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped for  $c'_{s1} = \Phi(c_{s1})$  and  $c'_{s2} = \Phi(c_{s2})$ .

Finally, we need to verify the integrability conditions:

$$\begin{aligned}
\int |h_\Phi(z)|^t dz &\simeq \int (h(\Phi^{-1}(z)) + 1)^t \psi(\Phi^{-1}(z))^t dz \\
&\stackrel{(a)}{=} \int_S (h(x) + 1)^t \psi(x)^t \phi(x) dx \\
&\leq \left\{ \int_S (h(x) + 1)^{2t} \phi(x) dx \right\}^{1/2} \left\{ \int_S \psi(x)^{2t} \phi(x) dx \right\}^{1/2},
\end{aligned}$$

where (a) follows from a change of variables. To bound the first term, note that

$$\begin{aligned}
\int_S (h(x) + 1)^{2t} \phi(x) dx &\leq \int_S h(x)^{2t} \phi(x) dx + \int_S \phi(x) dx \\
&\leq \int_S h(x)^{2t} \phi(x) dx + 1.
\end{aligned}$$

The first inequality follows since  $2t < 1$  by assumption. Note that  $\int_S h(x)^{2t} \phi(x) dx < \infty$ .

We now bound the second term:

$$\int_S \psi(x)^{2t} \phi(x) dx \leq \int_S \phi(c_{s1})^{-2t} \phi(x) dx + \int_S \phi(c_{s2})^{-2t} \phi(x) dx + \int_S \phi(x)^{-2t} \phi(x) dx.$$

The first two terms are constants. The last term is  $\int_S \phi(x)^{1-2t} dx$ , which is finite because  $1 - 2t > 0$  and  $\phi$  is a valid transformation function.  $\square$

**Proposition F.2.** *Suppose Assumptions A1'–A4' hold. The following conditions are satisfied for  $p_\Phi(z) = \frac{p(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$  and  $q_\Phi(z) = \frac{q(\Phi^{-1}(z))}{\phi(\Phi^{-1}(z))}$ :*

C1'  $p_\Phi(z), q_\Phi(z) > 0$  on  $(0, 1)$ , and  $\sup_z \{p_\Phi(z) \vee q_\Phi(z)\} < \infty$ .

C2' For some  $r > 2$ ,  $\int \left| \log \frac{p_\Phi(z)}{q_\Phi(z)} \right|^r dz < \infty$ .

C3' There exists  $h_\Phi(z)$  such that

- (a)  $h_\Phi(z) \geq \max \left\{ \left| \frac{p'_\Phi(z)}{p_\Phi(z)} \right|, \left| \frac{q'_\Phi(z)}{q_\Phi(z)} \right| \right\}$ ,
- (b)  $h_\Phi(z)$  is  $(c'_{s1}, c'_{s2}, C'_s)$ -bowl-shaped, and
- (c)  $\int_R |h_{\Phi,n}(z)|^t dz < \infty$ , for some constant  $t$  such that  $\frac{2}{r} \leq t \leq 1$ .

C4' We have that  $p'_\Phi(z), q'_\Phi(z) \geq 0$  for all  $z < c'_{s1}$ , and  $p'_\Phi(z), q'_\Phi(z) \leq 0$  for all  $z > c'_{s2}$ .

*Proof.* The proof is identical to that of Proposition F.1, so we omit the details.  $\square$

## G Proof of Proposition 4.1

First suppose  $\|\theta_1 - \theta_0\| \rightarrow 0$ . In Lemma G.2, we show that  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx \rightarrow 0$ . Assumptions A1 and A5 follow directly from Assumptions B1 and B5, respectively.

We now prove A2. Let  $\rho$  be a constant and define

$$R = \left\{ x : g_1(x) \leq \frac{\log \rho}{\|\theta_1 - \theta_0\|} \right\}, \tag{G.1}$$

where  $g_1(x)$  is the upper bound on  $\sup_{\theta \in \Theta} \|\nabla_{\theta} f_{\theta}(x)\|$  defined in Assumption B3. Since  $g_1(x)$  is bowl-shaped, we conclude that  $R$  is an interval if  $\log \rho \geq C_s \text{diam}(\Theta)$ . Note that

$$\log \frac{p(x)}{q(x)} = f_{\theta_1}(x) - f_{\theta_0}(x) = (\theta_1 - \theta_0)^{\top} \nabla_{\theta} f_{\bar{\theta}}(x).$$

This implies

$$\left| \log \frac{p(x)}{q(x)} \right| \leq \|\theta_1 - \theta_0\| \|\nabla_{\theta} f_{\bar{\theta}}(x)\| \leq \|\theta_1 - \theta_0\| \sup_{\theta} \|\nabla_{\theta} f_{\theta}(x)\| \leq \|\theta_1 - \theta_0\| \cdot g_1(x),$$

where  $\bar{\theta}$  is a convex combination of  $\theta_0, \theta_1$ . Therefore, for all  $x \in R$ , we have  $\frac{1}{\rho} \leq \frac{p(x)}{q(x)} \leq \rho$ .

Since we know from Assumption B3 that  $\int |g_1(x)|^r \phi(x) dx < \infty$ , Markov's inequality gives

$$\Phi(R^c) = \Phi \left\{ x : g_1(x) > \frac{\log \rho}{\|\theta_1 - \theta_0\|} \right\} \leq C \frac{\|\theta_1 - \theta_0\|^r}{(\log \rho)^r} \stackrel{(a)}{=} \Theta(H^{r/2}) = o(H),$$

where (a) follows from Lemma G.2. The last equality follows from the assumption that  $r > 2$ . This proves A2.

Now we move on to **A3**. By Lemma G.1, we have

$$\begin{aligned} \frac{1}{\alpha} \left| \frac{p(x) - q(x)}{q(x)} \right| &\lesssim \frac{1}{\|\theta_1 - \theta_0\|} \left| \frac{p(x)}{q(x)} - 1 \right| \\ &= \frac{1}{\|\theta_1 - \theta_0\|} |\exp(f_{\theta_1}(x) - f_{\theta_0}(x)) - 1| \\ &\stackrel{(a)}{=} \frac{1}{\|\theta_1 - \theta_0\|} |(\theta_1 - \theta_0)^{\top} \nabla_{\theta} f_{\bar{\theta}}(x)| \exp(f_{\bar{\theta}}(x) - f_{\theta_0}(x)) \\ &\leq \|\nabla_{\theta} f_{\bar{\theta}}(x)\| \exp(f_{\bar{\theta}}(x) - f_{\theta_0}(x)) \\ &\stackrel{(b)}{=} \|\nabla_{\theta} f_{\bar{\theta}}(x)\| \exp((\bar{\theta} - \theta_1)^{\top} \nabla_{\theta} f_{\bar{\theta}}(x)) \\ &\leq \|\nabla_{\theta} f_{\bar{\theta}}(x)\| \exp(\|\theta_1 - \theta_0\| \|\nabla_{\theta} f_{\bar{\theta}}(x)\|), \end{aligned}$$

where in (a),  $\bar{\theta}$  is a convex combination of  $\theta_1$  and  $\theta_0$ , and in (b),  $\tilde{\theta}$  is a convex combination of  $\bar{\theta}$  and  $\theta_0$ . Assumption B3 implies that both  $\|\nabla_{\theta} f_{\bar{\theta}}(x)\|$  and  $\|\nabla_{\theta} f_{\tilde{\theta}}(x)\|$  are upper-bounded by  $g_1(x)$ , so

$$\frac{1}{\alpha} \left| \frac{p(x) - q(x)}{q(x)} \right| \lesssim g_1(x) \exp(\|\theta_0 - \theta_1\| g_1(x)). \quad (\text{G.2})$$

Therefore,

$$\begin{aligned} \int_R \left( \frac{1}{\alpha} \left| \frac{p(x)}{q(x)} - 1 \right| \right)^r q(x) dx &\lesssim \int_R g_1(x)^r \exp(r \|\theta_0 - \theta_1\| g_1(x)) q(x) dx \\ &\stackrel{(a)}{\leq} \int_R g_1(x)^r \rho^r q(x) dx \\ &\stackrel{(b)}{\lesssim} \int_R g_1(x)^r \phi(x) dx, \end{aligned}$$

where (a) follows from the definition of  $R$  and (b) follows because  $\frac{q(x)}{\phi(x)}$  is bounded. This proves A3.

To prove **A4**, we first construct  $h(x)$ . By equation G.2, we have

$$\left| \frac{\gamma(x)}{q(x)} \right| \lesssim g_1(x) \exp(\|\theta_0 - \theta_1\| g_1(x)).$$

By Assumption B3, we also have  $\left| \frac{q'(x)}{q(x)} \right| = |f'_{\theta_0}(x)| \leq g_{2,\theta_0}(x)$ , and

$$\begin{aligned} \left| \frac{\gamma'(x)}{q(x)} \right| &= \left| \frac{1}{\alpha} \frac{p'(x) - q'(x)}{q(x)} \right| \\ &\lesssim \frac{1}{\|\theta_0 - \theta_1\|} \left| \frac{p'(x) - q'(x)}{q(x)} \right| \\ &= \frac{1}{\|\theta_0 - \theta_1\|} \left| f'_{\theta_1} \frac{p(x)}{q(x)} - f'_{\theta_0}(x) \right| \\ &= \frac{1}{\|\theta_0 - \theta_1\|} \left| (f'_{\theta_1}(x) - f'_{\theta_0}(x)) \frac{p(x)}{q(x)} + f'_{\theta_0} \left( \frac{p(x)}{q(x)} - 1 \right) \right| \\ &\leq \|\nabla_{\theta} f'_{\bar{\theta}}(x)\| \frac{p(x)}{q(x)} + \frac{1}{\|\theta_1 - \theta_0\|} \left| \frac{p(x)}{q(x)} - 1 \right| |f'_{\theta_0}(x)|, \end{aligned}$$

where  $\bar{\theta}$  is a convex combination of  $\theta_0$  and  $\theta_1$ .

Using Assumption B3 and inequality (G.2), we have

$$\left| \frac{\gamma'(x)}{q(x)} \right| \lesssim g_{2,\bar{\theta}}(x) \exp\left(\|\theta_0 - \theta_1\|g_1(x)\right) + g_1(x) \exp\left(\|\theta_0 - \theta_1\|g_1(x)\right)g_{2,\theta_0}(x).$$

Hence, we may choose

$$\begin{aligned} h(x) &\simeq g_{2,\bar{\theta}}(x) \exp\left(\|\theta_0 - \theta_1\|g_1(x)\right) + g_1(x) \exp\left(\|\theta_0 - \theta_1\|g_1(x)\right)g_{2,\theta_0}(x) \\ &\quad + g_1(x) \exp\left(\|\theta_0 - \theta_1\|g_1(x)\right). \end{aligned}$$

Since all the component functions are  $(c_{s1}, c_{s2}, \widetilde{C}_s)$  bowl-shaped,  $h(x)$  is  $(c_{s1}, c_{s2}, C_s)$ -bowl-shaped, where  $C_s = 3\widetilde{C}_s^2 \exp(\text{diam}(\Theta)\widetilde{C}_s)$ . Furthermore,

$$\begin{aligned} \int_R h(x)^{2t} \phi(x) dx &\stackrel{(a)}{\lesssim} \int_R \left( g_{2,\bar{\theta}}^{2t} \rho^{2t} + g_1(x)^{2t} \rho^{2t} g_{2,\theta_0}^{2t} + g_1(x)^{2t} \rho^{2t} \right) \phi(x) dx \\ &\lesssim \int_R g_{2,\bar{\theta}}(x)^{2t} \phi(x) dx + \int_R g_1(x)^{2t} g_{2,\theta_0}^{2t} \phi(x) dx + \int_R g_1(x)^{2t} \phi(x) dx, \end{aligned}$$

where (a) follows because on  $R$ , we have  $\|\theta_0 - \theta_1\|g_1(x) \leq \log \rho$ . By Assumption B3, the first and third terms are finite, uniformly over all  $\theta_1, \theta_0 \in \Theta$ . It is straightforward to show that the second term is also finite, by an application of the Cauchy-Schwartz inequality.

Now suppose  $\|\theta_1 - \theta_0\| = \Theta(1)$ . Lemma G.2 implies  $\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \Theta(1)$ . Assumptions **A1'** and **A4'** follow directly from Assumptions B1 and B5.

To prove **A2'**, note that from a previous derivation, we have

$$\left| \log \frac{p(x)}{q(x)} \right| \leq \|\theta_1 - \theta_0\| \sup_{\theta} \|\nabla_{\theta} f_{\theta}(x)\| \leq \|\theta_1 - \theta_0\|g_1(x).$$

Since  $\|\theta_1 - \theta_0\| \leq \text{diam}(\Theta)$ , which is a constant, and  $\int g_1(x)^r \phi(x) dx < \infty$  by Assumption B3, we obtain **A2'**.

To prove **A3'**, note that

$$\frac{q'(x)}{q(x)} = f'_{\theta_0}(x), \quad \text{and} \quad \frac{p'(x)}{p(x)} = f'_{\theta_1}(x).$$

Therefore, the choice  $h(x) = g_{2,\theta_0}(x) + g_{2,\theta_1}(x)$  upper-bounds  $\left| \frac{q'(x)}{q(x)} \right|$  and  $\left| \frac{p'(x)}{p(x)} \right|$ , by Assumption B3. Furthermore,  $h(x)$  is  $(c_{s1}, c_{s2}, C_s)$ -bowl-shaped, where  $C_s = 2\widetilde{C}_s$ .

To prove the last integrability condition, note that

$$\int h(x)^{2t} \phi(x) dx \leq \int g_{2,\theta_0}(x)^{2t} \phi(x) dx + \int g_{2,\theta_1}(x)^{2t} \phi(x) dx.$$

Hence,

$$\int h(x)^{2t} \phi(x) dx \leq 2 \sup_{\theta \in \Theta} \int g_{2,\theta}(x)^{2t} \phi(x) dx.$$

## G.1 Supporting lemmas

**Lemma G.1.** *Under Assumptions B1–B5, we have  $\alpha \asymp \|\theta_1 - \theta_0\|$ .*

*Proof.* We write

$$\begin{aligned} \alpha^2 &= \int_R \left( \frac{p(x)}{q(x)} - 1 \right)^2 q(x) dx \\ &= \int_R \left| \exp(f_{\theta_1}(x) - f_{\theta_0}(x)) - 1 \right|^2 q(x) dx \\ &= \int_R \left( (\theta_1 - \theta_0)^\top \nabla_\theta f_{\bar{\theta}}(x) \exp(f_{\bar{\theta}}(x) - f_{\theta_0}(x)) \right)^2 q(x) dx. \end{aligned}$$

First we show an upper bound:

$$\begin{aligned} \alpha^2 &\leq \int_R \|\theta_1 - \theta_0\|^2 \|\nabla_\theta f_{\bar{\theta}}(x)\|^2 \exp(f_{\bar{\theta}}(x) - f_{\theta_0}(x)) \exp(f_{\bar{\theta}}(x)) dx \\ &\leq \int_R \|\theta_1 - \theta_0\|^2 \|\nabla_\theta f_{\bar{\theta}}(x)\|^2 \exp\left(\|\theta_1 - \theta_0\| \|\nabla_\theta f_{\bar{\theta}}(x)\|\right) \exp(f_{\bar{\theta}}(x)) dx. \end{aligned}$$

On  $R$ , we have  $\|\theta_1 - \theta_0\| \sup_\theta \|\nabla_\theta f_\theta(x)\| \leq \log \rho$ . Hence,

$$\begin{aligned} \alpha^2 &\leq \|\theta_1 - \theta_0\|^2 \int_R \|\nabla_\theta f_{\bar{\theta}}(x)\|^2 e^{\log \rho} \exp(f_{\bar{\theta}}(x)) dx \\ &\leq \|\theta_1 - \theta_0\|^2 \rho \int_{-\infty}^{\infty} \|\nabla_\theta f_{\bar{\theta}}(x)\|^2 \exp(f_{\bar{\theta}}(x)) dx \\ &\stackrel{(a)}{\lesssim} \|\theta_1 - \theta_0\|^2, \end{aligned}$$

where (a) follows from Assumptions B1 and B4. We now establish a lower bound:

$$\begin{aligned} \alpha^2 &\geq \int_R \left( (\theta_1 - \theta_0)^\top \nabla_\theta f_{\bar{\theta}}(x) \right)^2 \exp\left(-|f_{\bar{\theta}}(x) - f_{\theta_0}(x)|\right) \exp(f_{\bar{\theta}}(x)) dx \\ &\geq \int_R \left( (\theta_1 - \theta_0)^\top \nabla_\theta f_{\bar{\theta}}(x) \right)^2 \exp\left(-\|\theta_1 - \theta_0\| \|\nabla_\theta f_{\bar{\theta}}(x)\|\right) \exp(f_{\bar{\theta}}(x)) dx \\ &\stackrel{(a)}{\geq} \frac{1}{\rho} (\theta_1 - \theta_0)^\top \left( \int_R (\nabla_\theta f_{\bar{\theta}}(x)) (\nabla_\theta f_{\bar{\theta}}(x))^\top \exp(f_{\bar{\theta}}(x)) dx \right) (\theta_1 - \theta_0), \end{aligned}$$

where (a) follows from Assumption B3. Define

$$\widetilde{G}_{\bar{\theta}} = \int_R (\nabla_\theta f_{\bar{\theta}}(x)) (\nabla_\theta f_{\bar{\theta}}(x))^\top \exp(f_{\bar{\theta}}(x)) dx.$$

As  $\rho$  increases,  $R \rightarrow S$ . Therefore, there exists an absolute constant such that for all  $\rho$  greater than or equal to this constant, we have  $\lambda_{\min}(\widetilde{G}_{\bar{\theta}}) > \frac{1}{2} \lambda_{\min}(G_{\bar{\theta}}) > 0$ . Hence,  $\alpha^2 \gtrsim \|\theta_1 - \theta_0\|^2$ .  $\square$

**Lemma G.2.** *The Hellinger distance satisfies the bound*

$$\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = c \|\theta_0 - \theta_1\|_2^2,$$

where  $c_{\min} \leq c \leq \frac{1}{4} c_{\max} d_{\Theta}$ .

*Proof.* Expanding the left-hand side, we have

$$\begin{aligned} \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx &= \int q(x) \left( \sqrt{\frac{p(x)}{q(x)}} - 1 \right)^2 dx \\ &= \int q(x) \left( \exp\left(\frac{f_{\theta_1}(x)}{2} - \frac{f_{\theta_0}(x)}{2}\right) - 1 \right)^2 dx. \end{aligned}$$

Let  $h(\theta) = \exp\left(\frac{f_{\theta}(x)}{2} - \frac{f_{\theta_0}(x)}{2}\right)$ . It is clear that  $h(\theta_0) = 1$  and that we wish to bound  $h(\theta_1) - h(\theta_0)$ . We bound this as follows:

$$\begin{aligned} |h(\theta_1) - h(\theta_0)| &= |(\theta_1 - \theta_0)^\top \nabla_{\theta} h(\bar{\theta})| \\ &= \left| \frac{1}{2} (\theta_1 - \theta_0)^\top \nabla_{\theta} f_{\bar{\theta}}(x) \exp\left(\frac{f_{\bar{\theta}}(x)}{2} - \frac{f_{\theta_0}(x)}{2}\right) \right| \\ &\leq \frac{1}{2} \|\theta_1 - \theta_0\| \|\nabla_{\theta} f_{\bar{\theta}}(x)\| \exp\left(\frac{f_{\bar{\theta}}(x)}{2} - \frac{f_{\theta_0}(x)}{2}\right), \end{aligned}$$

where  $\bar{\theta} \in \Theta$  is some convex combination of  $\theta_1, \theta_0$ . Thus, we have

$$\begin{aligned} \int q(x) \left( \exp\left(\frac{f_{\theta_1}(x)}{2} - \frac{f_{\theta_0}(x)}{2}\right) - 1 \right)^2 dx &\leq \int q(x) \frac{1}{4} \|\theta_1 - \theta_0\|^2 \|\nabla_{\theta} f_{\bar{\theta}}(x)\|^2 \exp(f_{\bar{\theta}}(x) - f_{\theta_0}(x)) dx \\ &= \frac{1}{4} \|\theta_1 - \theta_0\|^2 \int \|\nabla_{\theta} f_{\bar{\theta}}(x)\|^2 \exp(f_{\bar{\theta}}(x)) dx \\ &\leq \frac{1}{4} \|\theta_1 - \theta_0\|^2 \text{tr}(G_{\bar{\theta}}) \\ &\leq \frac{1}{4} \|\theta_1 - \theta_0\|^2 c_{\max} d_{\Theta}, \end{aligned}$$

where  $\Theta \subseteq \mathbb{R}^{d_{\Theta}}$ . Furthermore,

$$\begin{aligned} \int q(x) \left( \left( \frac{f_{\theta_1}(x)}{2} - \frac{f_{\theta_0}(x)}{2} \right) - 1 \right)^2 dx &= \int \left( (\theta_1 - \theta_0)^\top \nabla_{\theta} f_{\bar{\theta}}(x) \right)^2 \exp(f_{\bar{\theta}}(x)) dx \\ &= (\theta_1 - \theta_0)^\top G_{\bar{\theta}} (\theta_1 - \theta_0) \\ &\geq c_{\min} \|\theta_1 - \theta_0\|^2. \end{aligned}$$

□

## G.2 Proof of examples

**Proposition G.1.** *Let  $\exp(f(x))$  be a positive density over  $\mathbb{R}$ , where*

- (a)  $|f^{(k)}(x)|$  is bounded for some  $k \geq 2$ , and
- (b) there exist constants  $c$  and  $M$  such that  $f'(x) > M$  for  $x < -c$  and  $f'(x) < -M$  for  $x > c$ .

Let  $\theta = (\mu, \sigma)$  and  $\Theta = [-C_\mu, C_\mu] \times [\frac{1}{c_\sigma}, c_\sigma]$ , for some absolute constants  $C_\mu$  and  $c_\sigma$ , and let

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

Then  $\{f_\theta(x)\}_{\theta \in \Theta}$  satisfies Assumptions B1–B4 with respect to  $\phi$  defined in equation (4.3).

*Proof.* Before we prove the claims, let us derive some useful properties of  $f$ .

First, for any  $x > c$ , we have

$$f(x) = \int_0^x f'(t)dt = \int_0^c f'(t)dt + \int_c^x f'(t)dt \lesssim 1 - \int_c^x Mdt \lesssim 1 - x.$$

Similarly, for any  $x < -c$ , we have  $f(x) \lesssim 1 + x$ . Therefore,  $f(x) \lesssim 1 - |x|$ .

Likewise, we have

$$f\left(\frac{x-\mu}{\sigma}\right) \lesssim 1 - \left|\frac{x-\mu}{\sigma}\right| \lesssim 1 - \left|\frac{x}{\sigma}\right| + \frac{\mu}{\sigma} \stackrel{(a)}{\lesssim} 1 - |x|,$$

where (a) follows because  $\sigma \geq \frac{1}{c_\sigma}$  and  $|\mu| \leq C_\mu$ , for some absolute constants  $c_\sigma$  and  $C_\mu$ . Thus, the density  $\exp f\left(\frac{x-\mu}{\sigma}\right)$  is sub-exponential.

Since  $f^{(k)}(x)$  is bounded, L'Hopital's rule implies  $|f'(x)| \lesssim |x|^{k-1} + 1$  and  $|f''(x)| \lesssim |x|^{k-2} + 1$ . Furthermore,

$$f'\left(\frac{x-\mu}{\sigma}\right) \lesssim \left|\frac{x-\mu}{\sigma}\right|^{k-1} + 1 \stackrel{(a)}{\lesssim} \left|\frac{x}{\sigma}\right|^{k-1} + \left|\frac{\mu}{\sigma}\right|^{k-1} + 1 \stackrel{(b)}{\lesssim} |x|^{k-1} + 1, \quad (\text{G.3})$$

where (a) follows because  $k$  is a constant and (b) follows because  $|\mu| \leq C_\mu$  and  $\sigma \geq \frac{1}{c_\sigma}$ , by assumption.

Now we prove the first claim **B1**. We have

$$\begin{aligned} \log \phi(x) - f_\theta(x) &= \log \frac{e}{8} - \sqrt{|x|+1} - f\left(\frac{x-\mu}{\sigma}\right) - \log \sigma \\ &\geq -\sqrt{|x|+1} - f\left(\frac{x-\mu}{\sigma}\right) - \log \frac{1}{c_\sigma} + \log \frac{e}{8} \\ &\geq -\sqrt{|x|+1} - C(1-|x|) - \log \frac{1}{c_\sigma} + \log \frac{e}{8} > -\infty. \end{aligned}$$

Moving on to **B2**, we have

$$\nabla f_\theta(x) = \begin{bmatrix} -\frac{1}{\sigma} f'\left(\frac{x-\mu}{\sigma}\right) \\ -\left(\frac{x-\mu}{\sigma^2}\right) f'\left(\frac{x-\mu}{\sigma}\right) - \frac{1}{\sigma} \end{bmatrix} = -\frac{1}{\sigma} f'\left(\frac{x-\mu}{\sigma}\right) \begin{bmatrix} 1 \\ \frac{x-\mu}{\sigma} + 1 \end{bmatrix}.$$

To show that  $\lambda_{\max}(G_\theta) < \infty$ , it is sufficient to show that

$$\begin{aligned} \int \frac{1}{\sigma} f'\left(\frac{x-\mu}{\sigma}\right)^2 \left(\frac{x-\mu}{\sigma} + 1\right) \exp f\left(\frac{x-\mu}{\sigma}\right) dx &< \infty, \quad \text{and} \\ \int \frac{1}{\sigma} f'\left(\frac{x-\mu}{\sigma}\right)^2 \left(\frac{x-\mu}{\sigma} + 1\right)^2 \exp f\left(\frac{x-\mu}{\sigma}\right) dx &< \infty. \end{aligned}$$

Since  $|f'\left(\frac{x-\mu}{\sigma}\right)| \lesssim \left|\frac{x-\mu}{\sigma}\right|^{k-1} + 1$  and  $\exp f\left(\frac{x-\mu}{\sigma}\right)$  is sub-exponential with all moments finite, we conclude that both integrals converge.

To show that  $\lambda_{\min}(G_\theta) > 0$ , we need to show that  $\det(G_\theta) > 0$ . Let  $g(x) = \frac{1}{\sigma} f'\left(\frac{x-\mu}{\sigma}\right)^2 \exp f\left(\frac{x-\mu}{\sigma}\right)$ , and note that  $g$  is positive and integrable. The integral of  $g$  is not 0, since  $|f'(x)| \geq M$  for all  $|x| > c$ . Thus,  $g$  may be normalized to a density  $\bar{g}$ .



Showing that  $\det(G_\theta) > 0$  is equivalent to showing that

$$\int g(x)dx \int \left(\frac{x-\mu}{\sigma} + 1\right)^2 g(x)dx > \left(\int \left(\frac{x-\mu}{\sigma} + 1\right) g(x)dx\right)^2,$$

which is equivalent to showing that

$$\mathbb{E}_{\bar{g}} \left[ \left( \frac{X-\mu}{\sigma} + 1 \right)^2 \right] - \left( \mathbb{E}_{\bar{g}} \left[ \frac{X-\mu}{\sigma} + 1 \right] \right)^2 > 0,$$

or  $\text{Var}_{\bar{g}}(X) > 0$ . This follows because  $g(x) \neq 0$ .

To verify **B3**, note that

$$\begin{aligned} \|\nabla f_\theta\| &= \left| \frac{1}{\sigma} f' \left( \frac{x-\mu}{\sigma} \right) \right| \sqrt{1 + ((x-\mu)/\sigma + 1)^2} \\ &\lesssim (1 + |x|^{k-1})(1 + |(x-\mu)/\sigma|) \\ &\lesssim 1 + |x|^k. \end{aligned}$$

Thus, we set  $g_1(x) = C(1 + |x|^k)$  for some absolute constant  $C$ . Note that  $g_1(x)$  is clearly bowl-shaped and  $\int g_1(x)^r \phi(x)dx$  is finite, since all moments of  $\phi$  are finite. To construct  $g_{2,\theta}$ , note that

$$f'_\theta(x) = \frac{1}{\sigma} f' \left( \frac{x-\mu}{\sigma} \right), \quad \text{and} \quad \nabla f'_\theta(x) = \begin{bmatrix} -\frac{1}{\sigma^2} f'' \left( \frac{x-\mu}{\sigma} \right) & -\frac{1}{\sigma^2} f' \left( \frac{x-\mu}{\sigma} \right) - \frac{x-\mu}{\sigma^3} f'' \left( \frac{x-\mu}{\sigma} \right) \end{bmatrix}.$$

Therefore,  $|f'_\theta(x)| \lesssim 1 + |x|^{k-1}$ , and

$$\begin{aligned} \|\nabla f'_\theta(x)\| &\leq \frac{1}{\sigma^2} \left| f'' \left( \frac{x-\mu}{\sigma} \right) \right| + \frac{1}{\sigma^2} \left| f' \left( \frac{x-\mu}{\sigma} \right) \right| + \frac{1}{\sigma^2} \left| f'' \left( \frac{x-\mu}{\sigma} \right) \right| \left| \frac{x-\mu}{\sigma} \right| \\ &\stackrel{(a)}{\lesssim} (1 + |x|^{k-2}) + (1 + |x|^{k-1}) + (1 + |x|^{k-2})(1 + |x|) \\ &\lesssim 1 + |x|^{k-1}, \end{aligned}$$

where (a) follows because  $|f''(x)| \lesssim 1 + |x|^{k-2}$ . Thus, we may take  $g_{2,\theta}(x) = C(1 + |x|^{k-1})$ . This is clearly bowl-shaped and integrable, as well.

Finally, we prove **B4**. We have

$$(\log \phi)'(x) = \begin{cases} \frac{1}{2} \frac{1}{\sqrt{1-x}}, & \text{if } x < 0 \\ -\frac{1}{2} \frac{1}{\sqrt{1+x}}, & \text{if } x > 0. \end{cases}$$

In particular,  $(\log \phi)'(x) \rightarrow 0$  as  $|x| \rightarrow \infty$ .

We also know that  $f'(x) \geq M$  for all  $x \leq -c$ , and  $f'(x) \leq -M$  for all  $x \geq c$ . If  $x \leq -\frac{c}{c_\sigma} - C_\mu$ , then  $\frac{x-\mu}{\sigma} \leq -c$  and

$$f'_\theta(x) = \frac{1}{\sigma} f' \left( \frac{x-\mu}{\sigma} \right) \geq \frac{M}{c_\sigma}.$$

If  $x \geq \frac{c}{c_\sigma} + C_\mu$ , then  $\frac{x-\mu}{\sigma} \geq c$  and

$$f'_\theta(x) = \frac{1}{\sigma} f' \left( \frac{x-\mu}{\sigma} \right) \leq -\frac{M}{c_\sigma}.$$

Thus, there exist  $c_{s1} < 0$  and  $c_{s2} > 0$  such that B4 holds. □

**Proposition G.2.** Let  $\exp(f(x))$  be a positive density over  $\mathbb{R}^+$ , where

(a)  $|f^{(k)}(x)|$  is bounded for some  $k \geq 2$ , and

(b) there exist constants  $c$  and  $M$  such that  $f'(x) < -M$  for  $x > c$ .

Let  $\theta = \sigma$  and  $\Theta = [\frac{1}{c_\sigma}, c_\sigma]$  for some absolute constant  $c_\sigma$ , and let

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

Then  $\{f_\theta(x)\}_{\theta \in \Theta}$  satisfies Assumptions B1–B4 with respect to  $\phi$  defined in equation (4.2).

The proof is almost identical to that of proposition G.1.

**Proposition G.3.** Let  $\theta = (\alpha, \beta)$  and  $\Theta = [\frac{1}{c}, c]^2$  for some constant  $c$ , and let

$$f_\theta = (\alpha - 1) \log x - \beta x + \alpha \log \beta - \log \Gamma(\alpha).$$

Then  $\{f_\theta(x)\}_{\theta \in \Theta}$  satisfies Assumptions B1–B4 with respect to  $\phi$  defined in equation (4.2).

*Proof.* We first prove **B1**. We have  $\log \phi(x) = \log \frac{e}{4} - \sqrt{x+1}$ , so

$$\log \phi(x) - f_\theta(x) = -\sqrt{x+1} - (\alpha - 1) \log x + \beta x + \log \frac{e}{4} - \alpha \log \beta - \log \Gamma(\alpha) > -\infty.$$

To prove **B2**, note that  $G_\theta = \int (H_\theta f_\theta(x)) \exp f_\theta(x) dx$ , where  $H_\theta$  is the Hessian operator. Hence,

$$\nabla f_\theta(x) = \begin{bmatrix} \log x + \log \beta - d_\alpha \log \Gamma(\alpha) \\ -x + \frac{\alpha}{\beta} \end{bmatrix},$$

implying that

$$H_\theta f_\theta(x) = \begin{bmatrix} d_\alpha^2 \log \Gamma(\alpha) & \frac{1}{\beta} \\ \frac{1}{\beta} & -\frac{\alpha}{\beta^2} \end{bmatrix}.$$

Therefore,  $G_\theta = H_\theta f_\theta(x)$  which is clearly full-rank.

To prove **B3**, we write

$$\begin{aligned} \|\nabla f_\theta(x)\| &\leq |\log x| + x + |\log \beta| + |d_\alpha \log \Gamma(\alpha)| + \frac{\alpha}{\beta} \\ &\stackrel{(a)}{\leq} |\log x| + x + C, \end{aligned}$$

where (a) follows because  $c \geq \beta, \alpha \geq \frac{1}{c}$ . Therefore, we take  $g_1(x) = |\log x| + x + C$ . Then

$$\begin{aligned} \int g_1(x)^r \phi(x) dx &= \int_0^\infty (|\log x| + x + C)^r \phi(x) dx \\ &\lesssim \int_0^\infty |\log x|^r \phi(x) + \int_0^\infty x^r \phi(x) dx. \end{aligned}$$

Observe that both terms are finite for our choice of  $\phi$ . Furthermore,

$$f'_\theta(x) = \frac{(\alpha - 1)}{x} - \beta, \quad \text{and} \quad \nabla f'_\theta(x) = \begin{bmatrix} \frac{1}{x} \\ -1 \end{bmatrix},$$

implying that  $|f'_\theta(x)| \lesssim 1 + x^{-1}$  and  $\|\nabla f'_\theta(x)\| \lesssim 1 + x^{-1}$ . Thus,  $g_{2,\theta} = C(1 + x^{-1})$  satisfies

$$\int_0^\infty g_{2,\theta}(x)^{4t} \phi(x) dx \lesssim 1 + \int_0^\infty x^{-4t} \phi(x) dx.$$

Since  $4t < 1$  by assumption, the integral converges.

**B4** readily follows because  $\beta \geq \frac{1}{c} > 0$ . □

## H Appendix for Theorem 4.3

We begin by defining some notation. Let  $\hat{\sigma}$  denote a clustering algorithm and  $A$  denote a weighted network such that  $\hat{\sigma}(A) : [n] \rightarrow [K]$  is the clustering obtained by  $\hat{\sigma}$  based on the input  $A$ . Let

$$S_K[\hat{\sigma}(A), \sigma_0] := \arg \min_{\rho \in S_K} d_H(\rho \circ \hat{\sigma}(A), \sigma_0),$$

where  $d_H(\cdot, \cdot)$  denotes the Hamming distance, and define

$$\mathcal{E}[\hat{\sigma}(A), \sigma_0] := \left\{ v : (\rho \circ \hat{\sigma}(A))(v) \neq \sigma_0(v), \text{ for some } \rho \in S_K[\hat{\sigma}(A), \sigma_0] \right\}. \quad (\text{H.1})$$

When  $S_K[\hat{\sigma}(A), \sigma_0]$  is a singleton, the set  $\mathcal{E}[\hat{\sigma}(A), \sigma_0]$  contains all nodes misclustered by  $\hat{\sigma}(A)$  in relation to  $\sigma_0$ . When  $S_K[\hat{\sigma}(A), \sigma_0]$  contains multiple elements, we continue to call  $\mathcal{E}[\hat{\sigma}(A), \sigma_0]$  the set of *misclustered* nodes.

### H.1 Proof of Theorem 4.3

Throughout this proof, let  $C$  denote a  $\Theta(1)$  sequence whose value may change from instance to instance. Let

$$\tilde{l}(\hat{\sigma}(A), \sigma_0) = \frac{1}{n} \sum_{v=1}^n \mathbb{1}\{v \in \mathcal{E}[\hat{\sigma}(A), \sigma_0]\},$$

where  $\mathcal{E}[\hat{\sigma}(A), \sigma_0]$  is defined in equation (H.1). In particular, note that if  $|S_K[\hat{\sigma}(A), \sigma_0]| = 1$ , we have  $\tilde{l} = l$ . We have the following claims:

**Claim 1:** If  $\frac{nI}{K} \rightarrow \infty$ , then  $\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0) \geq C \exp\left(-(1 + o(1))\frac{nI}{\beta K}\right)$ .

**Claim 2:** If  $\frac{nI}{K} \leq c < \infty$ , then  $\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0) \geq c' > 0$ , for some constants  $c$  and  $c'$ .

We first prove that the theorem follows from the claims. If  $P\left(l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K}\right) \geq \frac{1}{2}\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0)$ , we have

$$\mathbb{E}l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K} P\left(l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K}\right) \geq \frac{1}{4\beta K} \mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0).$$

On the other hand, if  $P\left(l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K}\right) < \frac{1}{2}\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0)$ , we have

$$\begin{aligned} \mathbb{E}l(\hat{\sigma}(A), \sigma_0) &\geq \mathbb{E}\left[l(\hat{\sigma}(A), \sigma_0) \mid l(\hat{\sigma}(A), \sigma_0) < \frac{1}{2\beta K}\right] P\left(l(\hat{\sigma}(A), \sigma_0) < \frac{1}{2\beta K}\right) \\ &\stackrel{(a)}{=} \mathbb{E}\left[\tilde{l}(\hat{\sigma}(A), \sigma_0) \mid l(\hat{\sigma}(A), \sigma_0) < \frac{1}{2\beta K}\right] P\left(l(\hat{\sigma}(A), \sigma_0) < \frac{1}{2\beta K}\right) \\ &= \mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0) - \mathbb{E}\left[\tilde{l}(\hat{\sigma}(A), \sigma_0) \mid l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K}\right] P\left(l(\hat{\sigma}(A), \sigma_0) \geq \frac{1}{2\beta K}\right) \\ &\geq \mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0) - \frac{1}{2}\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0) \\ &= \frac{1}{2}\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0), \end{aligned}$$

where (a) holds by invoking Lemma B.8. Thus, any lower bound on  $\mathbb{E}\tilde{l}(\hat{\sigma}(A), \sigma_0)$  translates into a lower bound on  $\mathbb{E}l(\hat{\sigma}(A), \sigma_0)$  scaled by a suitable constant, implying the desired result.

We now focus on proving the claims. Without loss of generality, suppose cluster 1 has size  $\frac{n}{\beta K} + 1$  and cluster 2 has size  $\frac{n}{\beta K}$ . Also suppose nodes 1 and 2 are such that  $\sigma_0(1) = 1$  and  $\sigma_0(2) = 2$ . Let  $C_i = \{u : \sigma_0(u) = i\}$  denote the  $i^{\text{th}}$  cluster.

Let  $\sigma_0^1 := \sigma_0$ , and let  $\sigma_0^2 : [n] \rightarrow [K]$  be the cluster assignment satisfying  $\sigma_0^2(v) = \sigma_0(v)$  for all  $v \neq 1$ , and  $\sigma_0^2(1) = 2$ . Let  $\sigma^*$  be a random cluster assignment, where

$$\sigma^* = \begin{cases} \sigma_0^1, & \text{with probability } \frac{1}{2}, \\ \sigma_0^2, & \text{with probability } \frac{1}{2}. \end{cases}$$

Let  $\Phi$  denote the probability measure on  $(\sigma^*, A, \hat{\sigma}(A))$ , defined by

$$P_\Phi(\sigma^*, A, \hat{\sigma}(A)) = P(\sigma^*)P_{SBM}(A | \sigma^*)P_{alg}(\hat{\sigma}(A) | A),$$

where  $P_{SBM}(A | \sigma^*)$  is the measure on the weighted graph defined by the weighted SBM treating  $\sigma^*$  as the true cluster assignment, and  $P_{alg}$  represents any randomness in the clustering algorithm. Let  $\Psi$  denote an alternative probability measure defined by

$$P_\Psi(\sigma^*, A, \hat{\sigma}(A)) = P(\sigma^*)P_\Psi(A | \sigma^*)P_{alg}(\hat{\sigma}(A) | A),$$

where  $P_\Psi(A | \sigma^*)$  is defined as follows:

1. If  $u, v \neq 1$ , then  $A_{uv}$  is distributed just as in  $P_{SBM}(A | \sigma^*)$ .
2. If  $v = 1$  and  $u \notin C_1 \cup C_2$ , then  $A_{uv}$  is distributed just as in  $P_{SBM}(A | \sigma^*)$ .
3. If  $v = 1$  and  $u \in C_1 \cup C_2$ , then  $A_{uv}$  is distributed as  $Y^*$ , where  $Y^*$  is the distribution that minimizes  $D$  in Lemma H.2; i.e.,  $Y_0^* \propto (P_0 Q_0)^{1/2}$  and  $(1 - Y_0^*)y^*(x) \propto \sqrt{(1 - P_0)p(x)(1 - Q_0)q(x)}$ .

Note that  $P_\Psi(A | \sigma^*) = P_\Psi(A)$  actually does not depend on whether  $\sigma^* = \sigma_0^1$  or  $\sigma_0^2$ .

Furthermore, we have

$$\begin{aligned} \mathcal{Q} &:= \log \frac{dP_\Psi}{dP_\Phi} = \log \frac{dP_{SBM}(A | \sigma^*)}{dP_\Psi(A | \sigma^*)} \\ &= \sum_{u \in C_{\sigma^*(1)}} \log \frac{Y(A_{u,1})}{P(A_{u,1})} + \sum_{u \in C_1 \cup C_2 \setminus C_{\sigma^*(1)}} \log \frac{Y(A_{u,1})}{Q(A_{u,1})}, \end{aligned}$$

where we use the notation  $P(A_{u,1}) = P_0$  if  $A_{u,1} = 0$  and  $P(A_{u,1}) = (1 - P_0)p(A_{u,1})$  if  $A_{u,1} \neq 0$ , and similarly for  $Y$ . Let

$$E = \left\{ 1 \notin \mathcal{E}[\hat{\sigma}(A), \sigma^*] \text{ and } \tilde{l}(\hat{\sigma}(A), \sigma^*) \leq \frac{1}{4\beta K} \right\}.$$

For an arbitrary function  $f(n)$  to be defined later, we may write

$$P_\Psi(\mathcal{Q} \leq f(n)) = P_\Psi(\mathcal{Q} \leq f(n), \neg E) + P_\Psi(\mathcal{Q} \leq f(n), E). \quad (\text{H.2})$$

We bound the first term as follows:

$$\begin{aligned} P_\Psi(\mathcal{Q} \leq f(n), \neg E) &= \int_{\mathcal{Q} \leq f(n), \neg E} dP_\Psi = \int_{\mathcal{Q} \leq f(n), \neg E} \exp(\mathcal{Q}) dP_\Phi \\ &\leq \exp(f(n)) P_\Phi(\mathcal{Q} \leq f(n), \neg E) \\ &\leq \exp(f(n)) P_\Phi(\neg E) \\ &\leq \exp(f(n)) \left( P_\Phi(1 \in \mathcal{E}[\hat{\sigma}(A), \sigma^*]) + P_\Phi \left( \tilde{l}(\hat{\sigma}(A), \sigma^*) \geq \frac{1}{4\beta K} \right) \right). \end{aligned}$$

Furthermore,

$$\mathbb{E}_\Phi \tilde{l}(\hat{\sigma}(A), \sigma^*) = \frac{1}{n} \sum_{v=1}^n P_\Phi(v \in \mathcal{E}[\hat{\sigma}(A), \sigma^*])$$

$$\begin{aligned}
&\geq \frac{1}{n} \sum_{v \in C_{\sigma^*(1)}} P_{\Phi}(v \in \mathcal{E}[\widehat{\sigma}(A), \sigma^*]) \\
&\stackrel{(a)}{=} \frac{|C_{\sigma^*(1)}|}{n} P_{\Phi}(1 \in \mathcal{E}[\widehat{\sigma}(A), \sigma^*]) \\
&\geq \frac{1}{\beta K} P_{\Phi}(1 \in \mathcal{E}[\widehat{\sigma}(A), \sigma^*]),
\end{aligned}$$

where (a) follows from Corollary H.1, and

$$\begin{aligned}
\mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*) &\geq \mathbb{E}_{\Phi} \left[ \widetilde{l}(\widehat{\sigma}(A), \sigma^*) \mid \widetilde{l}(\widehat{\sigma}(A), \sigma^*) \geq \frac{1}{4\beta K} \right] P_{\Phi} \left( \widetilde{l}(\widehat{\sigma}(A), \sigma^*) \geq \frac{1}{4\beta K} \right) \\
&\geq \frac{1}{4\beta K} P_{\Phi} \left( \widetilde{l}(\widehat{\sigma}(A), \sigma^*) \geq \frac{1}{4\beta K} \right).
\end{aligned}$$

so we have the bound

$$P_{\Psi}(\mathcal{Q} \leq f(n), \neg E) \leq \exp(f(n)) \cdot 5\beta K \cdot \mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*).$$

We now turn to the second term in equation (H.2). We have

$$\begin{aligned}
P_{\Psi}(E) &= \frac{1}{2} P_{\Psi} \left( 1 \notin \mathcal{E}[\widehat{\sigma}(A), \sigma_0^1] \text{ and } \widetilde{l}(\widehat{\sigma}(A), \sigma_0^1) \leq \frac{1}{4\beta K} \right) \\
&\quad + \frac{1}{2} P_{\Psi} \left( 1 \notin \mathcal{E}[\widehat{\sigma}(A), \sigma_0^2] \text{ and } \widetilde{l}(\widehat{\sigma}(A), \sigma_0^2) \leq \frac{1}{4\beta K} \right). \quad (\text{H.3})
\end{aligned}$$

If  $l(\widehat{\sigma}(A), \sigma_0^1) \leq \widetilde{l}(\widehat{\sigma}(A), \sigma_0^1) \leq \frac{1}{4\beta K}$ , Lemma B.8 implies  $S_K[\widehat{\sigma}(A), \sigma_0^1]$  contains only one element, which we denote by  $\rho$ . Since  $d_H(\sigma_0^1, \sigma_0^2) = 1$ , we have  $\frac{1}{n} d_H(\rho \circ \widehat{\sigma}(A), \sigma_0^2) \leq \frac{1}{4\beta K} + \frac{1}{n} \leq \frac{1}{2\beta K}$ , so applying Lemma B.8 again, we conclude that  $\rho \in S_K[\widehat{\sigma}(A), \sigma_0^2]$ , as well. However,  $(\rho \circ \widehat{\sigma}(A))(1)$  cannot be equal to both  $\sigma_0^1(1) = 1$  and  $\sigma_0^2(1) = 2$ . Hence, we cannot simultaneously have  $1 \notin \mathcal{E}[\widehat{\sigma}(A), \sigma_0^1]$  and  $1 \notin \mathcal{E}[\widehat{\sigma}(A), \sigma_0^2]$ . In particular, the two events in equation (H.3) are disjoint, so

$$P_{\Psi}(\mathcal{Q} \leq f(n), E) \leq P_{\Psi}(E) \leq \frac{1}{2}.$$

Plugging back into equation (H.2), we conclude that

$$P_{\Psi}(\mathcal{Q} \leq f(n)) \leq \exp(f(n)) \cdot 5\beta K \cdot \mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*) + \frac{1}{2},$$

so setting  $f(n) = \log \frac{1}{20\beta K \mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*)}$ , we have

$$P_{\Psi} \left( \mathcal{Q} \leq \log \frac{1}{20\beta K \mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*)} \right) \leq \frac{3}{4}.$$

By Chebyshev's inequality, we also have

$$P_{\Psi} \left( \mathcal{Q} \leq \mathbb{E}_{\Psi} \mathcal{Q} + \sqrt{5V_{\Psi}(\mathcal{Q})} \right) \geq 4/5,$$

where  $V_{\Psi}(\mathcal{Q})$  is the variance of  $\mathcal{Q}$  under  $\Psi$ . Hence,  $\log \frac{1}{20\beta K \mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma_0^1)} \leq \mathbb{E}_{\Psi} \mathcal{Q} + \sqrt{5V_{\Psi}(\mathcal{Q})}$ , or equivalently,

$$\mathbb{E}_{\Phi} \widetilde{l}(\widehat{\sigma}(A), \sigma^*) \geq \frac{1}{20\beta K} \exp \left( -(\mathbb{E}_{\Psi} \mathcal{Q} + \sqrt{5V_{\Psi}(\mathcal{Q})}) \right).$$

We now compute  $\mathbb{E}_\Psi \mathcal{Q}$  and  $V_\Psi(\mathcal{Q})$ . Note that

$$\mathbb{E}_\Psi \mathcal{Q} = \frac{1}{2} \mathbb{E}_\Psi[\mathcal{Q} | \sigma^* = \sigma_0^1] + \frac{1}{2} \mathbb{E}_\Psi[\mathcal{Q} | \sigma^* = \sigma_0^2].$$

Furthermore, by Lemma H.2, we have

$$\begin{aligned} \mathbb{E}_\Psi[\mathcal{Q} | \sigma^* = \sigma_0^1] &= \mathbb{E}_\Psi \left[ \sum_{u: u \neq 1, \sigma_0^1(u)=1} \log \frac{Y(A_{u,1})}{P(A_{u,1})} + \sum_{u: \sigma_0^1(u)=2} \log \frac{Y(A_{u,1})}{Q(A_{u,1})} \right] \\ &= \frac{n}{\beta K} \int \log \frac{dY}{dP} dY + \frac{n}{\beta K} \int \log \frac{dY}{dQ} dY \\ &= \frac{n}{\beta K} 2D = \frac{n}{\beta K} I. \end{aligned}$$

Similarly, we have  $\mathbb{E}_\Psi[\mathcal{Q} | \sigma^* = \sigma_0^2] = \frac{nI}{\beta K}$ , so  $\mathbb{E}_\Psi \mathcal{Q} = \frac{nI}{\beta K}$ . We show in Lemma H.3 that the following bound holds for the variance:

$$\sqrt{5V_\Psi(\mathcal{Q})} \leq C \sqrt{\frac{nI}{\beta K}}.$$

Now note that if  $\frac{nI}{\beta K} \rightarrow \infty$ , we have  $\sqrt{\frac{nI}{\beta K}} = o\left(\frac{nI}{\beta K}\right)$ , so  $\sqrt{5V_\Psi(\mathcal{Q})} = o\left(\frac{nI}{\beta K}\right)$ . Therefore,

$$\mathbb{E}_\Phi \tilde{l}(\hat{\sigma}(A), \sigma^*) \geq C \exp\left(-\left(1 + o(1)\right) \frac{nI}{\beta K}\right).$$

If instead  $\frac{nI}{\beta K} \rightarrow c < \infty$ , then  $\mathbb{E}_\Psi \mathcal{Q} = c(1 + o(1))$  and  $\sqrt{5V_\Psi(\mathcal{Q})} \leq C(1 + o(1))$ , so

$$\mathbb{E}_\Phi \tilde{l}(\hat{\sigma}(A), \sigma^*) \geq c' > 0,$$

for some constant  $c'$ .

Now define two measures  $P_1, P_2$  on  $(A, \hat{\sigma}(A))$ , as follows:

$$\begin{aligned} P_1(A, \hat{\sigma}(A)) &= P_{SBM}(A | \sigma_0^1) P_{alg}(\hat{\sigma}(A) | A), \\ P_2(A, \hat{\sigma}(A)) &= P_{SBM}(A | \sigma_0^2) P_{alg}(\hat{\sigma}(A) | A). \end{aligned}$$

Note that  $\mathbb{E}_\Phi[\tilde{l}(\hat{\sigma}, \sigma^*) | \sigma^* = \sigma_0^1] = \mathbb{E}_1[\tilde{l}(\hat{\sigma}, \sigma_0^1)]$  and  $\mathbb{E}_\Phi[\tilde{l}(\hat{\sigma}, \sigma^*) | \sigma^* = \sigma_0^2] = \mathbb{E}_2[\tilde{l}(\hat{\sigma}, \sigma_0^2)]$ , where  $\mathbb{E}_1$  and  $\mathbb{E}_2$  are expectations taken with respect to  $P_1$  and  $P_2$ , respectively. We claim that  $\mathbb{E}_1 \tilde{l}(\hat{\sigma}, \sigma_0^1) = \mathbb{E}_2 \tilde{l}(\hat{\sigma}, \sigma_0^2)$ , in which case  $\mathbb{E}_\Phi \tilde{l}(\hat{\sigma}, \sigma^*) = \mathbb{E}_1 \tilde{l}(\hat{\sigma}, \sigma_0^1) = \mathbb{E} \tilde{l}(\hat{\sigma}, \sigma_0)$  and the claims follow.

Define a permutation  $\pi \in S_n$  that swaps  $\{2, \dots, \frac{n}{\beta K} + 1\}$  with  $\{\frac{n}{\beta K} + 2, \dots, 2\frac{n}{\beta K} + 1\}$  and satisfies  $\pi(u) = u$  for  $u = 1$  and  $u \geq 2\frac{n}{\beta K} + 2$ . Clearly,  $\sigma_0^2 = \tau \circ \sigma_0^1 \circ \pi^{-1}$ , where  $\tau \in S_K$  swaps cluster labels 1 and 2. Now let  $A$  be fixed and let  $\rho \in S_K$  be arbitrary. We have

$$\begin{aligned} d_H(\rho \circ \hat{\sigma}(A), \sigma_0^1) &= d_H(\rho \circ \hat{\sigma}(A), \tau^{-1} \circ \sigma_0^2 \circ \pi) \\ &= d_H(\rho \circ \hat{\sigma}(A) \circ \pi^{-1}, \tau^{-1} \circ \sigma_0^2) \\ &= d_H(\rho \circ \xi^{-1} \circ \hat{\sigma}(\pi A), \tau^{-1} \circ \sigma_0^2) \\ &= d_H(\tau \circ \rho \circ \xi^{-1} \circ \hat{\sigma}(\pi A), \sigma_0^2). \end{aligned}$$

Therefore,  $\rho \mapsto \tau \circ \rho \circ \xi^{-1}$  is a bijection between  $S_K[\hat{\sigma}(A), \sigma_0^1]$  and  $S_K[\hat{\sigma}(\pi A), \sigma_0^2]$ . Furthermore, if  $v$  is a node such that  $(\rho \circ \hat{\sigma}(A))(v) \neq \sigma_0^1(v)$ , we equivalently have

$$(\rho \circ \hat{\sigma}(A))(v) \neq (\tau^{-1} \circ \sigma_0^2 \circ \pi)(v) \iff (\rho \circ \hat{\sigma}(A) \circ \pi^{-1})(u) \neq (\tau^{-1} \circ \sigma_0^2)(u), \quad \text{where } \pi(v) = u,$$

so  $(\tau \circ \rho \circ \xi^{-1} \circ \widehat{\sigma}(\pi A))(u) \neq \sigma_0^2(u)$ . Thus,  $v \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0^1]$  if and only if  $\pi(v) \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0^2]$ . Finally, we conclude that

$$\begin{aligned} \mathbb{E}_1 \widetilde{l}(\widehat{\sigma}(A), \sigma_0^1) &= \frac{1}{n} \sum_{v=1}^n P_1(v \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0^1]) \\ &= \frac{1}{n} \sum_{v=1}^n P_1(\pi(v) \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0^2]) \\ &\stackrel{(a)}{=} \frac{1}{n} \sum_{v=1}^n P_2(\pi(v) \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0^2]) \\ &= \mathbb{E}_2 \widetilde{l}(\widehat{\sigma}(A), \sigma_0^2), \end{aligned}$$

where (a) follows because  $[\pi A]_{ij} = A_{\pi^{-1}(i), \pi^{-1}(j)}$ , implying that if  $A$  is distributed according to  $P_{SBM}(A | \sigma_0^1)$ , then  $\pi A$  is distributed according to  $P_{SBM}(A | \sigma_0^2)$ . This concludes the proof.

## H.2 Properties of permutation-equivariant estimators

Permutation-equivariant estimators possess symmetry properties. The following lemma formalizes one symmetry property useful for the proof of Theorem 4.3:

**Lemma H.1.** *Let the true clustering  $\sigma_0$  be arbitrary. Suppose the weight matrix  $A$  is drawn from an arbitrary probability measure and  $\widehat{\sigma}$  is any permutation-equivariant estimator. Let  $u$  and  $v$  be two nodes such that there exists  $\pi \in S_n$  satisfying*

- (1)  $\pi(u) = v$ ,
- (2)  $\pi$  is measure-preserving; i.e.,  $A \stackrel{d}{=} \pi A$ , and
- (3)  $\pi$  preserves the true clustering; i.e., there exists  $\tau \in S_K$  such that  $\tau \circ \sigma_0 \circ \pi^{-1} = \sigma_0$ .

Then

$$P(u \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]) = P(v \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]).$$

*Proof.* Since  $\widehat{\sigma}(A) \stackrel{d}{=} \widehat{\sigma}(\pi A)$ , we have

$$P(v \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]) = P(v \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0]).$$

We claim that  $u \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]$  if and only if  $v \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0]$ , implying the desired result:

$$P(u \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]) = P(v \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0]) = P(v \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]).$$

Consider a fixed matrix  $A$ , and let  $\tau \in S_K$  satisfy  $\tau \circ \sigma_0 \circ \pi^{-1} = \sigma_0$ . Let  $\xi \in S_K$  be the permutation such that  $\widehat{\sigma}(\pi A) = \xi \circ \widehat{\sigma}(A) \circ \pi^{-1}$ . For any  $\rho \in S_K$ , we have

$$\begin{aligned} d_H(\rho \circ \widehat{\sigma}(A), \sigma_0) &= d_H(\tau \circ \rho \circ \xi^{-1} \circ \xi \circ \widehat{\sigma}(A) \circ \pi^{-1}, \tau \circ \sigma_0 \circ \pi^{-1}) \\ &= d_H(\tau \circ \rho \circ \xi^{-1} \circ \widehat{\sigma}(\pi A), \sigma_0). \end{aligned}$$

Therefore,  $\rho \in S_K[\widehat{\sigma}(A), \sigma_0]$  if and only if  $\tau \circ \rho \circ \xi^{-1} \in S_K[\widehat{\sigma}(\pi A), \sigma_0]$ . In particular, if  $v \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0]$ , we have  $\tau \circ \rho \circ \xi^{-1} \circ \widehat{\sigma}(\pi A)(v) \neq \sigma_0(v)$  for some  $\rho \in S_K[\widehat{\sigma}(A), \sigma_0]$ . Then

$$\begin{aligned} \widehat{\sigma}(A)(u) &= \widehat{\sigma}(A)(\pi^{-1}(v)) = \xi^{-1} \circ \xi \circ \widehat{\sigma}(A) \circ \pi^{-1}(v) = \xi^{-1} \circ \widehat{\sigma}(\pi A)(v) \\ &\neq \rho^{-1} \circ \tau^{-1} \circ \sigma_0(v) = \rho^{-1} \circ \tau^{-1} \circ \sigma_0(\pi(u)) = \rho^{-1}(\sigma_0(u)). \end{aligned}$$

Thus,  $u \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]$ . Similar reasoning shows that if  $u \in \mathcal{E}[\widehat{\sigma}(A), \sigma_0]$ , then  $v \in \mathcal{E}[\widehat{\sigma}(\pi A), \sigma_0]$ .  $\square$

**Corollary H.1.** *Let the true clustering  $\sigma_0$  be arbitrary. Suppose the weight matrix  $A$  is drawn from an arbitrary probability measure and  $\hat{\sigma}$  is any permutation-equivariant estimator. Let  $u$  and  $v$  be two nodes lying in equally-sized clusters. Then*

$$P(u \in \mathcal{E}[\hat{\sigma}(A), \sigma_0]) = P(v \in \mathcal{E}[\hat{\sigma}(A), \sigma_0]).$$

*Proof.* By Lemma H.1, it suffices to construct a permutation  $\pi \in S_n$  satisfying conditions (1)–(3).

First suppose  $u$  and  $v$  lie in the same cluster. It is easy to see that the conditions are satisfied when  $\pi$  is the permutation that transposes  $u$  and  $v$  and  $\tau$  is the identity.

If  $u$  and  $v$  lie in different clusters, suppose without loss of generality that  $u$  is in cluster 1 and  $v$  is in cluster 2, where clusters 1 and 2 have the same size. Let  $\pi$  be the permutation that exchanges all nodes in cluster 1 with all nodes in cluster 2. The conditions are satisfied when  $\tau$  is the permutation that transposes cluster labels 1 and 2.  $\square$

### H.3 Properties of Renyi divergence

We first state a lemma that provides an alternative characterization of the Renyi divergence:

**Lemma H.2.** *Let  $P$  and  $Q$  be two probability measures on  $\mathbb{R}$  that are absolutely continuous with respect to each other, with respective point masses  $P_0$  and  $Q_0$  at zero. The Renyi divergence satisfies  $I = 2D$ , where*

$$D := \inf_{Y \in \mathcal{P}} \max \left\{ \int \log \frac{dY}{dP} dY, \int \log \frac{dY}{dQ} dY \right\},$$

and  $\mathcal{P}$  denotes the set of probability measures absolutely continuous with respect to both  $P$  and  $Q$ .

*Proof.* First, note that  $D$  is finite by making the choice  $Y = P$ . We claim that

$$D = \inf_{Y \in \mathcal{P}} \left\{ \int \log \frac{dY}{dP} dY : \int \log \frac{dP}{dQ} dY = 0 \right\}. \quad (\text{H.4})$$

This is because for any  $Y \in \mathcal{P}$  such that  $\int \log \frac{dP}{dQ} dY \neq 0$ , we have  $\int \log \frac{dY}{dP} dY \neq \int \log \frac{dY}{dQ} dY$ . Suppose without loss of generality that the first quantity is larger. Then it is possible to take  $\tilde{Y} = (1 - \epsilon)Y + \epsilon P$  for  $\epsilon$  small enough such that  $\max \left\{ \int \log \frac{d\tilde{Y}}{dP} d\tilde{Y}, \int \log \frac{d\tilde{Y}}{dQ} d\tilde{Y} \right\}$  strictly decreases, so the infimum in the definition of  $D$  could not have been achieved.

Since the new formulation (H.4) is convex in  $Y$ , we may solve to obtain the optimal  $Y^* \in \mathcal{P}$ , defined by  $Y_0^* = \frac{P_0^{1/2} Q_0^{1/2}}{Z}$  and  $(1 - Y_0^*)y^*(x) = \frac{((1 - P_0)p(x))^{1/2}((1 - Q_0)q(x))^{1/2}}{Z}$ . The quantity  $Z$  is the normalization term:  $Z = P_0^{1/2} Q_0^{1/2} + \int \sqrt{(1 - P_0)p(x)(1 - Q_0)q(x)} dx$ . Then

$$\begin{aligned} \int \log \frac{dY^*}{dP} dY^* &= \log \frac{1}{Z} \left\{ \left( \frac{Q_0}{P_0} \right)^{1/2} Y_0^* + \int \left( \frac{(1 - P_0)p(x)}{(1 - Q_0)q(x)} \right)^{1/2} (1 - Y_0^*)y^*(x) dx \right\} \\ &= \log \frac{dP}{dQ} dY^* - \log Z = -\log Z. \end{aligned}$$

It is straightforward to verify that  $-2\log Z = I$ .  $\square$

### H.4 Bounding the variance

**Lemma H.3.** *For a suitable constant  $C$ , we have  $\sqrt{5V_\Psi(\mathcal{Q})} \leq C \sqrt{\frac{nI}{\beta K}}$ .*



*Proof.* We begin with the decomposition

$$\begin{aligned} V_\Psi(\mathcal{Q}) &= V(\mathbb{E}_\Psi[\mathcal{Q} | \sigma^*]) + E[V_\Psi(\mathcal{Q} | \sigma^*)] \\ &= E[V_\Psi(\mathcal{Q} | \sigma^*)] \\ &= \frac{1}{2}V_\Psi(\mathcal{Q} | \sigma^* = \sigma_0^1) + \frac{1}{2}V_\Psi(\mathcal{Q} | \sigma^* = \sigma_0^2). \end{aligned}$$

Then

$$\begin{aligned} V_\Psi(\mathcal{Q} | \sigma^* = \sigma_0^1) &= \sum_{u: u \neq 1, \sigma_0^1(u)=1} V_\Psi\left(\log \frac{Y(A_{v_1 u})}{P(A_{v_1 u})}\right) + \sum_{u: \sigma_0^1(u)=2} V_\Psi\left(\log \frac{Y(A_{v_1 u})}{Q(A_{v_1 u})}\right) \\ &\leq \frac{n}{\beta K} \mathbb{E}_\Psi\left(\log \frac{Y(A_{v_1 u})}{P(A_{v_1 u})}\right)^2 + \frac{n}{\beta K} \mathbb{E}_\Psi\left(\log \frac{Y(A_{v_1 u})}{Q(A_{v_1 u})}\right)^2. \end{aligned}$$

We will show that  $\mathbb{E}_\Psi\left(\log \frac{Y(A_{v_1 u})}{P(A_{v_1 u})}\right)^2$  is bounded by  $CI$ , so  $\sqrt{5V_\Psi(\mathcal{Q})} \leq C\sqrt{\frac{nI}{\beta K}}$ . We have

$$\begin{aligned} \mathbb{E}_\Psi\left(\log \frac{Y(A_{uv^*})}{P(A_{uv^*})}\right)^2 &= \int \left(\log \frac{dY}{dP}\right)^2 dY \\ &= Y_0 \log^2 \frac{Y_0}{P_0} + (1 - Y_0) \int y(x) \log^2 \frac{(1 - Y_0)y(x)}{(1 - P_0)p(x)} dx. \end{aligned} \quad (\text{H.5})$$

To bound the first term, we write

$$\begin{aligned} \left|\log \frac{Y_0}{P_0}\right| &= \left|\frac{1}{2} \log \frac{Q_0}{P_0} - \log Z\right| \\ &\leq \frac{1}{2} \left|\log \left(1 - \frac{P_0 - Q_0}{P_0}\right)\right| + \frac{I}{2} \\ &\stackrel{(a)}{\leq} \frac{1}{2} \left|\frac{P_0 - Q_0}{P_0}\right| + \left|\frac{P_0 - Q_0}{P_0}\right|^2 C + \frac{I}{2} \\ &\stackrel{(b)}{\leq} C \left|\frac{P_0 - Q_0}{P_0}\right| + CI, \end{aligned}$$

where (a) follows from Lemma 1.2 and the fact that  $\frac{Q_0}{P_0}$  is bounded, and (b) follows from the fact that  $\left|\frac{P_0 - Q_0}{P_0}\right| = \left|1 - \frac{Q_0}{P_0}\right| \leq 1 + \left|\frac{Q_0}{P_0}\right|$  is bounded. Therefore,

$$\begin{aligned} Y_0 \log^2 \frac{Y_0}{P_0} &\leq Y_0 \left(C \left|\frac{P_0 - Q_0}{P_0}\right| + CI\right)^2 \\ &\stackrel{(a)}{\leq} Y_0 \frac{|P_0 - Q_0|^2}{P_0^2} C + Y_0 I^2 C \\ &\stackrel{(b)}{\leq} \frac{|P_0 - Q_0|^2}{P_0} C + I^2 C, \end{aligned}$$

where (a) follows because  $(x+y)^2 \leq 2x^2 + 2y^2$ , and (b) follows because  $Y_0 = \frac{\sqrt{P_0 Q_0}}{Z} = (1+o(1))CP_0$ . Since  $I = o(1)$ , we have  $I^2 \leq I$ . Also,

$$I \geq (1+o(1))(\sqrt{P_0} - \sqrt{Q_0})^2 = (1+o(1)) \frac{(P_0 - Q_0)^2}{(\sqrt{P_0} + \sqrt{Q_0})^2} = C(1+o(1)) \frac{(P_0 - Q_0)^2}{P_0},$$

from which we conclude that  $Y_0 \log^2 \frac{Y_0}{P_0} \leq CI$ .

Now we turn our attention to the second term in equation (H.5). We have

$$\begin{aligned} \left| \log \frac{(1-Y_0)y(x)}{(1-P_0)p(x)} \right| &\leq \frac{1}{2} \left| \log \frac{(1-Q_0)q(x)}{(1-P_0)p(x)} - \log Z \right| \\ &\leq \frac{1}{2} \left| \log \frac{1-Q_0}{1-P_0} \right| + \frac{1}{2} \left| \log \frac{q(x)}{p(x)} \right| + \frac{I}{2}. \end{aligned}$$

Therefore,

$$\begin{aligned} (1-Y_0) \int y(x) \left( \log \frac{1-Y_0}{1-P_0} \frac{y(x)}{p(x)} \right)^2 dx &\leq (1-Y_0) \int y(x) \left\{ \frac{1}{2} \left| \log \frac{1-Q_0}{1-P_0} \right| + \frac{1}{2} \left| \log \frac{q(x)}{p(x)} \right| + \frac{I}{2} \right\}^2 dx \\ &\leq (1-Y_0) \int y(x) \left\{ C \left| \log \frac{1-Q_0}{1-P_0} \right|^2 + C \left| \log \frac{q(x)}{p(x)} \right|^2 + CI^2 \right\} dx, \end{aligned} \tag{H.6}$$

where we have used the fact that  $(x+y+z)^2 \leq 9x^2 + 9y^2 + 9z^2$  in the last inequality. Define

$$\begin{aligned} \mathcal{A} &:= (1-Y_0) \int y(x) \left| \log \frac{1-Q_0}{1-P_0} \right|^2 dx, \\ \mathcal{B} &:= (1-Y_0) \int y(x) \left| \log \frac{q(x)}{p(x)} \right|^2 dx, \\ \mathcal{C} &:= (1-Y_0) \int y(x) I^2 dx. \end{aligned}$$

We bound each term separately, beginning with  $\mathcal{A}$ . Note that

$$\left| \log \frac{1-Q_0}{1-P_0} \right| = \left| \log \left( 1 - \frac{Q_0-P_0}{1-P_0} \right) \right| \stackrel{(a)}{\leq} \left| \frac{Q_0-P_0}{1-P_0} \right| + C \left( \frac{Q_0-P_0}{1-P_0} \right)^2 \stackrel{(b)}{\leq} C \left| \frac{Q_0-P_0}{1-P_0} \right|,$$

where (a) follows from Lemma I.2 and the fact that  $\frac{1-Q_0}{1-P_0}$  is bounded, and (b) follows from the fact that  $\left| \frac{Q_0-P_0}{1-P_0} \right| = \left| 1 - \frac{1-Q_0}{1-P_0} \right| \leq 1 + \left| \frac{1-Q_0}{1-P_0} \right| \leq C$ . Therefore,

$$\begin{aligned} \mathcal{A} &\leq C(1-Y_0) \int y(x) \left( \frac{Q_0-P_0}{1-P_0} \right)^2 dx \\ &= C \left( \frac{Q_0-P_0}{1-P_0} \right)^2 \int \frac{\sqrt{(1-P_0)p(x)(1-Q_0)q(x)}}{Z} dx \\ &\stackrel{(a)}{\leq} C(1+o(1)) \left( \frac{Q_0-P_0}{1-P_0} \right)^2 \int \sqrt{\frac{(1-Q_0)q(x)}{(1-P_0)p(x)}} (1-P_0)p(x) dx \\ &\stackrel{(b)}{\leq} C(1+o(1)) \left( \frac{Q_0-P_0}{1-P_0} \right)^2 (1-P_0) \\ &\leq C(1+o(1)) \frac{(Q_0-P_0)^2}{1-P_0}, \end{aligned}$$

where in (a), we use the fact that  $\frac{1}{Z} = (1+o(1))$  since  $Z \rightarrow 1$ , and in (b), we use the fact that  $\frac{1-Q_0}{1-P_0}$  and  $\frac{q(x)}{p(x)}$  are bounded. Note that

$$I \geq (1+o(1))(\sqrt{1-P_0} - \sqrt{1-Q_0})^2 = (1+o(1)) \frac{(P_0-Q_0)^2}{(\sqrt{1-P_0} + \sqrt{1-Q_0})^2} = C(1+o(1)) \frac{(P_0-Q_0)^2}{1-P_0},$$

implying that  $\mathcal{A} \leq C(1+o(1))I \leq CI$ .

Moving onto  $\mathcal{B}$ , first suppose  $H = \Theta(1)$  and  $\max \left\{ \int p(x) \left| \log \frac{q(x)}{p(x)} \right|^2 dx, \int q(x) \left| \log \frac{q(x)}{p(x)} \right|^2 dx \right\} < \infty$ . We have

$$\begin{aligned} \mathcal{B} &\leq C \int \frac{\sqrt{(1-P_0)p(x)(1-Q_0)q(x)}}{Z} \left| \log \frac{q(x)}{p(x)} \right|^2 dx \\ &\stackrel{(a)}{\leq} C \sqrt{(1-P_0)(1-Q_0)} \int \sqrt{p(x)q(x)} \left| \log \frac{q(x)}{p(x)} \right|^2 dx \\ &\leq C \sqrt{(1-P_0)(1-Q_0)} \int (p(x) + q(x)) \left| \log \frac{q(x)}{p(x)} \right|^2 dx \\ &\stackrel{(b)}{\leq} C \sqrt{(1-P_0)(1-Q_0)} H \leq CI, \end{aligned}$$

where (a) follows because  $Z \rightarrow 1$  and (b) follows because  $H = \Theta(1)$ .

Next, we make no assumption on  $H$  but assume  $\left| \log \frac{q(x)}{p(x)} \right|$  is bounded by a constant. We have

$$\begin{aligned} \left| \log \frac{q(x)}{p(x)} \right| &= \left| \log \left( 1 - \frac{p(x) - q(x)}{p(x)} \right) \right| \\ &\stackrel{(a)}{\leq} \left| \frac{p(x) - q(x)}{p(x)} \right| + \left( \frac{p(x) - q(x)}{p(x)} \right)^2 C \\ &\stackrel{(b)}{\leq} C \left| \frac{p(x) - q(x)}{p(x)} \right|, \end{aligned}$$

where (a) follows from Lemma I.2 and the fact that  $\frac{q(x)}{p(x)}$  is bounded, and (b) follows from the fact that  $\left| \frac{p(x)-q(x)}{p(x)} \right| = \left| 1 - \frac{q(x)}{p(x)} \right| \leq 1 + \left| \frac{q(x)}{p(x)} \right| \leq C$ . Then

$$\begin{aligned} \mathcal{B} &\leq \frac{C}{Z} \int \sqrt{\frac{1-Q_0}{1-P_0} \frac{q(x)}{p(x)}} (1-P_0)p(x) \left( \frac{p(x) - q(x)}{p(x)} \right)^2 dx \\ &\stackrel{(a)}{\leq} \frac{C}{Z} (1-P_0) \int p(x) \left( \frac{p(x) - q(x)}{p(x)} \right)^2 dx, \end{aligned}$$

where in (a), we use the facts that  $\frac{1}{Z} = (1 + o(1))$  and  $\frac{1-Q_0}{1-P_0}$ , and  $\frac{p(x)}{q(x)}$  are both bounded by assumption. Now, note that  $H = \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx = \int \frac{(p(x)-q(x))^2}{(\sqrt{p(x)}+\sqrt{q(x)})^2} dx = C \int \frac{(p(x)-q(x))^2}{p(x)} dx$ . Therefore,

$$\mathcal{B} \leq C(1-P_0)H \leq C \sqrt{(1-P_0)(1-Q_0)} H \leq CI.$$

Finally, note that  $\mathcal{C} = (1-Y_0)CI^2 \leq CI$ . Substituting back into inequality (H.6), we therefore obtain

$$(1-Y_0) \int y(x) \left( \log \frac{1-Y_0}{1-P_0} \frac{y(x)}{p(x)} \right)^2 dx \leq CI,$$

so substituting back into inequality (H.5), we obtain the desired bound.  $\square$

## I Additional useful lemmas

**Lemma I.1.** *Let*

$$I = -2 \log \left( \sqrt{P_0 Q_0} + \int \sqrt{(1-P_0)(1-Q_0)p(x)q(x)} dx \right),$$

$$I^h = (\sqrt{P_0} - \sqrt{Q_0})^2 + \int \left( \sqrt{(1-P_0)p(x)} - \sqrt{(1-Q_0)q(x)} \right)^2 dx.$$

If  $I^h < 2 - 2\epsilon$ , then  $I = I^h(1 + \eta)$ , where  $|\eta| \leq \frac{I^h}{2\epsilon}$ . Thus,  $I \rightarrow 0$  if and only if  $I^h \rightarrow 0$ , in which case  $I = I^h(1 + o(1))$ .

*Proof.* We have

$$\begin{aligned} I &= -2 \log \left( \sqrt{P_0 Q_0} + \int \sqrt{(1-P_0)(1-Q_0)p(x)q(x)} dx \right) \\ &= -2 \log \left( 1 - \frac{1}{2} \left( (\sqrt{P_0} - \sqrt{Q_0})^2 + \int (\sqrt{(1-P_0)p(x)} - \sqrt{(1-Q_0)q(x)})^2 dx \right) \right) \\ &= -2 \log \left( 1 - \frac{1}{2} I^h \right) \\ &= 2 \cdot \frac{1}{2} I^h (1 + \eta), \end{aligned}$$

where  $|\eta| \leq \frac{I^h}{2\epsilon}$ . The last equality follows from Lemma I.2.  $\square$

**Lemma I.2.** Suppose  $0 < \epsilon \leq 1$ . For all  $0 \leq x < 1 - \epsilon$ , we have  $\log(1 - x) = -(1 + \eta)x$ , where  $|\eta| \leq \frac{x}{2\epsilon}$ .

*Proof.* This follows by taking the Taylor expansion of  $\log(1 - x)$  around  $x = 0$ .  $\square$

**Lemma I.3.** Let  $f(z) = \frac{1 - \frac{z}{2} - \sqrt{1-z}}{z}$ , for  $z \leq 1$  and  $z \neq 0$ , and define  $f(0) = 0$ . Then  $|f(z)| \leq |z|$ , for all  $z \leq 1$ .

*Proof.* Note that  $f$  is continuous, with derivative

$$f'(z) = -\frac{1}{z^2} - \frac{z-2}{2z^2\sqrt{1-z}}.$$

It is straightforward to check that  $f'(z) \geq 0$  for all  $z < 1$ , and we may define  $f'(0) = \frac{1}{4}$  such that  $f'(z)$  is continuous. Therefore,  $f(z)$  is monotonic and maximized at  $z = 1$ , yielding  $f(1) = \frac{1}{2}$ , and minimized at  $\lim_{z \rightarrow -\infty} f(z) = -\frac{1}{2}$ .

We now split into cases. If  $z < -\frac{1}{2}$ , then  $|f(z)| \leq \frac{1}{2} < |z|$ . If  $-1/2 \leq z \leq 1/2$ , a Taylor expansion gives

$$\sqrt{1-z} = 1 - \frac{1}{2}z - \frac{1}{8}z^2 - \frac{1}{16}z^3 - \dots - \frac{(n+1)!!}{2^n n!} z^n - \dots$$

Hence,

$$\begin{aligned} \left| \sqrt{1-z} - \left(1 - \frac{z}{2}\right) \right| &\leq \frac{1}{8}(|z|^2 + |z|^3 + \dots) \\ &\leq \frac{1}{8}|z|^2(1 + |z| + |z|^2 + \dots) \\ &\leq \frac{1}{8}|z|^2 \frac{1}{1-|z|} \leq \frac{1}{4}|z|^2, \end{aligned}$$

implying that  $|f(z)| \leq \frac{1}{4}|z|$ . Finally, if  $z > 1/2$ , we have  $|f(z)| \leq \frac{1}{2} < z$ .  $\square$