# Nonparametric mixture MLEs under Gaussian-smoothed optimal transport distance

Fang Han[*], Zhen Miao[†], and Yandi Shen[‡]

December 7, 2021

## Abstract

The Gaussian-smoothed optimal transport (GOT) framework, pioneered in Goldfeld et al. (2020) and followed up by a series of subsequent papers, has quickly caught attention among researchers in statistics, machine learning, information theory, and related fields. One key observation made therein is that, by adapting to the GOT framework instead of its unsmoothed counterpart, the curse of dimensionality for using the empirical measure to approximate the true data generating distribution can be lifted. The current paper shows that a related observation applies to the estimation of nonparametric mixing distributions in discrete exponential family models, where under the GOT cost the estimation accuracy of the nonparametric MLE can be accelerated to a polynomial rate. This is in sharp contrast to the classical sub-polynomial rates based on unsmoothed metrics, which cannot be improved from an information-theoretical perspective. A key step in our analysis is the establishment of a new Jackson-type approximation bound of Gaussian-convoluted Lipschitz functions. This insight bridges existing techniques of analyzing the nonparametric MLEs and the new GOT framework.

**Keywords:** GOT distance, nonparametric mixture models, nonparametric maximum likelihood estimation, rate of convergence, function approximation.

## 1 Introduction

Let $f(x \,|\, \theta)$ be a known parametric density function with respect to a certain (counting or continuous) measure and $X_1, \ldots, X_n$ be $n$ i.i.d. observations drawn from the following *mixture density function*,

$$h_Q(x) := \int f(x \,|\, \theta)\mathsf{d}Q(\theta), \tag{1.1}$$

where $Q$ is unspecified and termed the *mixing distribution* of $\theta$. Our goal is to estimate the unknown $Q$ based on $X_1, \ldots, X_n$. This is the celebrated nonparametric mixing distribution estimation problem, which has been extensively studied in literature (Lindsay, 1995). The focus of this paper

---

[*]Department of Statistics, University of Washington, Seattle, WA 98195, USA; e-mail: `fanghan@uw.edu`

[†]Department of Statistics, University of Washington, Seattle, WA 98195, USA; e-mail: `zhenm@uw.edu`

[‡]Department of Statistics, University of Chicago, Chicago, IL 60637, USA. E-mail: `ydshen@uchicago.edu`

is on studying the estimation of $Q$ in the case of (identifiable) *discrete exponential family models* (Zhang, 1995), i.e., $f(x\,|\,\theta)$ taking the following form:

$$f(x\,|\,\theta) = g(\theta)w(x)\theta^x, \quad \text{with } x = 0, 1, 2, \dots, \text{ and}$$

$$w(x) > 0 \quad \text{for all } x \geq 0,$$

$$0 \leq \theta \leq \text{ (a known fixed constant) } \theta_* < \min\{\theta_r, \infty\}, \qquad (1.2)$$

where $\theta_r \in (0, \infty]$ is the radius of convergence of the power series $\theta \mapsto \sum_{x=0}^{\infty} w(x)\theta^x$ and $g(\cdot)$ is analytic in a neighborhood of 0. This model includes, among many others, Poisson and negative binomial distributions.

Estimation of $Q$ under the discrete exponential family models has been extensively investigated in literature through, e.g., the use of nonparametric maximum likelihood estimators (MLEs) (Simar, 1976), method of moments (Tucker, 1963), Fourier and kernel methods (Zhang, 1995; Loh and Zhang, 1996, 1997), and projection methods (Walter and Hamedani, 1991; Hengartner, 1997; Roueff and Rydén, 2005). Of particular interest to us is the MLE-based approach, partly due to its asymptotic efficiency under regular parametric models. In the case of nonparametric mixture models, the MLE can be written as

$$\widehat{Q} := \operatorname*{argmax}_{Q \text{ on } [0,\theta^*]} \sum_{i=1}^{n} \log h_Q(X_i), \qquad (1.3)$$

which is a convex problem with efficient solving algorithms (Simar, 1976).

Although a proof of the consistency of $\widehat{Q}$ has been standard now (cf. Chen (2017)), of central importance to statisticians and machine learning scientists in making inference based on $\widehat{Q}$ is its rate of convergence. In this regard, Zhang (1995) established the first minimax lower bound, indicating that, at the worst case, it is impossible for the MLEs to achieve a polynomial rate if measured using regular metrics such as the total variation distance and the optimal transport distance (OT; in this paper restricted to the Wasserstein-1 distance $W_1$). It is now known that, for Poisson mixtures, the minimax rate of convergence under OT distance is $\log \log n / \log n$ and could indeed be achieved by MLEs (Miao et al., 2021, Theorems 6.1 and 6.2). This slow rate demonstrates that the estimation of $Q$ suffers severely from its nonparametric structure.

Interestingly, a similar fundamental "curse" also exists in using the empirical measure $P_n$ of an independent and identically distributed (i.i.d.) sample of size $n$ to approximate the true data generating distribution $P$ in $\mathbb{R}^d$, for which the minimax rate under the cost of OT was shown to be $n^{-1/d}$ as $d > 2$ (Fournier and Guillin, 2015, Theorem 1). Partly motivated by a problem of estimating information flows in deep neural networks (Goldfeld et al., 2019), Goldfeld et al. (2020) introduced a new measure of distance $W_1^\sigma$, which is termed the *Gaussian-smoothed OT (GOT)* distance. The GOT distance, like the unsmoothed OT one, is a metric on the probability measure space with finite first moment that metrizes the weak topology. In addition, both $W_1^\sigma$ and the corresponding optimal transport plan converge weakly to the corresponding unsmoothed versions as the smoothing parameter $\sigma \to 0$ (cf. Goldfeld and Greenewald (2020, Theorems 2, 3, and 4)).

Under this new distance and with some further moment conditions on $P$, Goldfeld et al. was able to prove an upper bound of $W_1^\sigma(P_n, P)$ that is of the best possible root-$n$ order and thus overcomes the curse of dimensionality faced with the classical unsmoothed scenario. Subsequent developments

2

establish the weak convergence of $W_1^\sigma(P_n, P)$ to a functional of a Gaussian process (Sadhu et al., 2021), weaken the moment assumption (Zhang et al., 2021), and study high noise limit as $\sigma \to \infty$ (Chen and Niles-Weed, 2021).

One of the main contributions of this paper is to prove that an observation similar to what was made in Goldfeld et al. (2020) occurs to the nonparametric mixture MLEs, i.e., under some conditions on $w(\cdot)$, we have

$$\sup_{Q \text{ on } [0,\theta^*]} \mathbb{E}W_1^\sigma(\widehat{Q}, Q) \le C(\sigma, \theta_*, w)n^{-\eta(\theta_*, w)}, \tag{1.4}$$

where $C$ and $\eta$ are two positive constants only depending on $\{\sigma, \theta_*, w\}$ and $\{\theta_*, w\}$, respectively. Our result thus bridges two distinct areas, namely nonparametric mixing distribution estimation and empirical approximation to population distribution; in the earlier case, GOT is shown to boost the convergence rate to polynomial, while in the latter case GOT overcomes the curse of dimensionality.

The main technical step in our proof of (1.4) is a new Jackson-type bound on the error of degree-$k$ polynomial (for an arbitrary positive integer $k$) approximation to Gaussian-convoluted Lipschitz functions with a bounded support. Our result thus extends the classic Jackson's Theorem (Jackson (1921); see Lemma 5.6) and paves a way to leverage existing technical tools of analyzing the nonparametric MLEs, devised in an early draft written by some of the authors in this paper (Miao et al., 2021, Section 6).

**Notation.** For any positive integer $n$, let $[n] := \{1, \ldots, n\}$. For any two distributions $Q_1$ and $Q_2$ over $\mathbb{R}^d$, let $Q_1 * Q_2$ represent the convolution of $Q_1$ and $Q_2$, i.e., $Q_1 * Q_2(A) = \int \int \mathbb{1}_A(x + y)\mathsf{d}Q_1(x)\mathsf{d}Q_2(y)$, with $\mathbb{1}.(\cdot)$ standing for the indicator function. For any two measurable functions $f, g$ on $\mathbb{R}^d$, $f * g$ represents their convolution, i.e., $f * g(x) = \int f(x - y)g(y)\mathsf{d}y$. For any function $f : \mathbb{R} \to \mathbb{R}$ and $\alpha > 0$, let $f^{(\alpha)}$ represent the $\alpha$-time derivative of $f$. The OT (i.e., Wasserstein $W_1$) distance between $Q_1$ and $Q_2$ is defined as

$$W_1(Q_1, Q_2) := \sup_{\ell \in \mathrm{Lip}_1} \int \ell(\mathsf{d}Q_1 - \mathsf{d}Q_2), \quad \text{(Kantorovich-Rubinstein formula)}$$

where the supremum is over all 1-Lipschitz functions (under the Euclidean metric $\| \cdot \|$) on $\mathbb{R}^d$. Let $\mathcal{N}_\sigma$ represent the Gaussian distribution with mean 0 and covariance matrix $\sigma^2 I_d$, where $I_d$ stands for the $d$-dimensional identity matrix. Let $\phi_\sigma$ denote the corresponding density function. The GOT distance $W_1^\sigma$ is defined as

$$W_1^\sigma(Q_1, Q_2) := W_1(Q_1 * \mathcal{N}_\sigma, Q_2 * \mathcal{N}_\sigma).$$

Let $\mathcal{P}(\mathbb{R}^d)$ represent the set of all Borel probability measures on $\mathbb{R}^d$ and $\mathcal{P}_1(\mathbb{R}^d)$ be the subset of $\mathcal{P}(\mathbb{R}^d)$ with elements of finite first moment. Throughout the paper, $C, C', C'', c, c'$ are used to represent generic positive constants whole values may change in different locations.

**Paper organization.** The rest of this paper is organized as follows. Section 2 gives the preliminaries on the studied nonparametric mixture models and the MLEs. Section 3 delivers the main results, including the key technical insights to the proof. Section 4 collects the main proofs, with auxiliary proofs relegated to Section 5.

# 2 Preliminaries

## 2.1 Nonparametric mixture MLEs

Estimating the mixing distribution is known to be statistically challenging in a variety of nonparametric mixture models including the Gaussian (Wu and Yang, 2020), binomial (Tian et al., 2017; Vinayak et al., 2019), and Poisson (Miao et al., 2021) ones. Specific to the discrete exponential family models in the form of (1.2), the following $\log n$-scale information-theoretical lower bound formalized this difficulty under the standard unsmoothed $W_1$ distance.

**Theorem 2.1** (Minimax lower bounds under $W_1$ distance). *Let $\{X_i, i \in [n]\}$ be a random sample generated from the mixture density function $h_Q$ defined in (1.1) and $n \geq 2$. We then have*

*(a) For any $f(x \mid \theta)$ taking the form (1.2), we have*

$$\inf_{\widetilde{Q}} \sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widetilde{Q}) \geq \frac{c}{\log n}, \tag{2.1}$$

*where the infimum is taken over all measurable estimators of the mixing distribution $Q$ with support on $[0, \theta_*]$ and $c = c(\theta_*) > 0$ is a constant only depending on $\theta_*$.*

*(b) (Miao et al., 2021, Theorem 6.2) Suppose further $f(x \mid \theta) = e^{-\theta}\theta^x/x!$ to be the probability mass function of the Poisson with natural parameter $\theta$. We can further tighten the lower bound in (2.1) to be*

$$\inf_{\widetilde{Q}} \sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widetilde{Q}) \geq \frac{c' \log \log n}{\log n},$$

*for a constant $c' = c'(\theta_*) > 0$ depending only on $\theta_*$.*

**Remark 2.1.** In (1.2), the condition that "$w(x) > 0$ for all nonnegative integer $x$" is a (simplified) identifiability condition. Similar conditions were also posed in, e.g., Zhang (1995, Corollary 1) and Loh and Zhang (1996, Corollary 1). As a matter of fact, Stoyanov and Lin (2011, Theorem 1(a)) showed that, if there exists a constant $C > 0$ such that $f(x \mid \theta) = 0$ for all $x \geq C$ and $\theta \in [0, \theta_*]$, then $Q$ is not identifiable, i.e., there exist at least two distinct mixing distributions $Q_1, Q_2$ over $[0, \theta_*]$ such that $h_{Q_1} = h_{Q_2}$. On the other hand, it is straightforward to generalize the above result to the case of "$w(x) > 0$ for all $x \geq x_0$ for some nonnegative integer $x_0$ that is known to us".

In the past several decades, methods that provably (nearly) achieve the above minimax lower bound have been proposed; cf. Zhang (1995), Hengartner (1997), and Roueff and Rydén (2005) among many others. However, none of the above methods is likelihood-based, partly due to the theoretical challenges faced with analyzing the nonparametric MLEs. A major breakthrough towards understanding the rate of convergence of nonparametric mixture MLEs was made in Vinayak et al. (2019) for the binomial case and Miao et al. (2021) for the Poisson case.

The following theorem provides an extension of Miao et al. (2021, Theorem 6.1) to cover those mixture models of the general form (1.2).

**Theorem 2.2** (Minimax upper bounds of MLEs under $W_1$ distance). *Let $\{X_i, i \in [n]\}$ be a random sample generated from the mixture density function $h_Q$ defined in (1.1), and $\widehat{Q}$ be the MLE introduced in (1.3). The following are then true.*

(a) *If there exists a universal constant $C \geq 1$ such that $1/w(x) \leq C^x$ for all $x \geq 1$, then there exist some $n' = n'(\theta_*, C)$ and $C' = C'(\theta_*, C)$ such that*

$$\sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widehat{Q}) \leq \frac{C'}{\log n} \quad \text{for all } n \geq n'.$$

(b) *If there exists a universal constant $C \geq 1$ such that $1/w(x) \leq (Cx)^{Cx}$ for all $x \geq 1$, then there exists $n' = n'(\theta_*, C)$ and $C' = C'(\theta_*, C)$ such that*

$$\sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widehat{Q}) \leq \frac{C' \log\log n}{\log n} \quad \text{for all } n \geq n'.$$

**Remark 2.2.** The conditions enforced for $w(x)$ in Theorem 2.2 are classic and related to the identifiability of $Q$ discussed in Remark 2.1. Similar conditions were posed in Zhang (1995, Theorem 4), Loh and Zhang (1996, Corollary 1), Loh and Zhang (1997, Theorem 1), and Roueff and Rydén (2005, Corollary 1).

**Remark 2.3.** It is straightforward to verify that, after some standard operations including location shift, point mass inflation, and reparametrization, Theorem 2.2(a) applies to, e.g., the (zero-inflated or $C$-truncated) negative binomial, the logarithmic (Noack, 1950), the lost games (Gupta, 1984), as well as the generalized Poisson, negative binomial, and logarithmic (Janardan, 1982) distributions; Theorem 2.2(b) applies to, e.g., the (zero-inflated or $C$-truncated) Poisson as well as the Poisson polynomial (Cameron and Trivedi, 2013) distributions.

**Remark 2.4.** It may be of some theoretical interest to note that Theorem 2.2 can be further generalized to cover the following two cases.

(i) If the following bound holds,

$$1/w(x) \leq \underbrace{\exp \circ \exp \circ \cdots \exp}_{L}(Cx) \text{ for all } x,$$

then there exists a constant $C' = C'(C, \theta_*)$ such that

$$\sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widehat{Q}) \leq C' / \underbrace{\log \circ \log \circ \cdots \circ \log}_{L}(n) =: C' / \log_L(n) \quad \text{for all } n \geq n',$$

where $n' = n'(C, \theta_*)$ is a sufficiently large integer.

(ii) If the following bound holds,

$$1/w(x) \leq (Cx)^{\cdots^{(Cx)}} \quad (L \text{ times power}) \quad \text{for all } x,$$

then there exists a constant $C' = C'(C, \theta_*)$ such that

$$\sup_{Q \text{ on } [0,\theta_*]} \mathbb{E}W_1(Q, \widehat{Q}) \leq C' \log_{L-1}(n) / \log_L(n) \quad \text{for all } n \geq n',$$

where $n' = n'(C, \theta_*)$ is a sufficiently large integer.

5

## 2.2 The GOT distance

Theorem 2.1 suggests that, under the $W_1$ cost, the sub-polynomial rate in estimating the mixing distribution of a nonparametric mixture model is information-theoretically optimal. As a matter of fact, the conclusion of Theorem 2.1 goes beyond the discrete exponential family models studied in this paper; cf. Wu and Yang (2020, Proposition 8) for a similar phenomenon in the nonparametric Gaussian mixture models.

Revising the Wasserstein distance through convolution/smoothing has a long and rich history. In probability theory, this is interestingly related to heat semigroup operators on Riemannian manifold (von Renesse and Sturm, 2005), which reveals its connection to the Ricci curvature. More recently, stemming from the interest in estimating the mutual information of deep networks, Goldfeld et al. (2019) initiated the study of GOT distances, introduced as a smoothed alternative to the classic OT metric.

Indeed, the GOT distance is now known to be able to effectively alleviate some undesired issues associated with the OT distances. Let us start with the following simple fact, that the $W_1$-distance is non-increasing under convolution.

**Lemma 2.1.** *Consider $\mu_1, \mu_2, \nu \in \mathcal{P}_1(\mathbb{R}^d)$ be arbitrary three Borel probability measures on $\mathbb{R}^d$ with finite first moment. We then have*

$$W_1(\mu_1 * \nu, \mu_2 * \nu) \le W_1(\mu_1, \mu_2).$$

Lemma 2.1 confirms that the GOT distance is no greater than the original OT distance, but it does not quantify the difference. For that purpose, the existing literature has provided us with an interesting example, i.e., in approximating the population measure using the empirical one. In detail, suppose $P_n$ is the empirical measure of $P$, and both are supported on $\mathbb{R}^d$ with some integer $d \ge 3$. Fournier and Guillin (2015, Theorem 1) showed that

$$\sup_{P:\mathbb{E}_P X^2 < \infty} \mathbb{E} W_1(P, P_n) \asymp n^{-1/d},$$

which is faced with severe curse of dimensionality as the dimension $d$ becomes larger. In a recent paper of Goldfeld et al. (2020), the authors showed that, via appealing to the GOT one, this curse can be effectively handled. More specifically, they proved that, as long as $P$ is sub-gaussian with a fixed subgaussian constant, we have

$$\mathbb{E} W_1^\sigma(P, P_n) \lesssim n^{-1/2},$$

which is the parametric rate of convergence. See also Sadhu et al. (2021) for the limiting distribution of $\sqrt{n} W_1^\sigma(P, P_n)$ as well as Zhang et al. (2021) for the relaxation of the moment conditions on $P$.

The purpose of this paper is to present the second and also a statistically interesting example, for which adopting the GOT distance can significantly accelerate the convergence rate of a statistical procedure.

## 3 Main results

The following theorem is the main result of this paper.

**Theorem 3.1.** *Let $\{X_i, i \in [n]\}$ be a random sample generated from the mixture density function $h_Q$ defined in (1.1), and $\widehat{Q}$ be the MLE introduced in (1.3). Suppose that there exist some universal positive constants $c_1, c_2, c_3, C_1, C_2, C_3$ such that one of the following two conditions holds,*

*(i) $1/w(x) \leq C_1 C_2^x$ for all $x = 1, 2, \ldots$;*

*(ii) $c_1 c_2^x x^{c_3 x} \leq 1/w(x) \leq C_1 C_2^x x^{C_3 x}$ for all $x = 1, 2, \ldots$.*

*Then we have*

$$\sup_{Q \ on \ [0,\theta_*]} \mathbb{E}W_1^\sigma(Q, \widehat{Q}) \leq C \cdot n^{-c}.$$

*Here $C = C(\sigma, \theta_*, c_1, c_2, c_3, C_1, C_2, C_3)$ and $c = c(\theta_*, c_3, C_2, C_3)$ are two positive constants.*

**Remark 3.1.** Let us point out some results in the nonparametric mixture model literature that are relevant to ours. Lambert and Tierney (1984) studied the convergence of $h_{\widehat{Q}}$ to $h_Q$ in a specific nonparametric Poisson mixture model based on the regular unsmoothed distance. They observed that the convergence rate can be nearly parametric; cf. Proposition 3.1 therein. This observation is particularly relevant to ours as the $h$. operation is intrinsically also "smoothing" the probability measure. A similar observation was made in Wu and Yang (2020), who studied approximating the mixing distributions in Gaussian mixture models via moment matching. In particular, their Lemma 8 considers bounding the Gaussian-smoothed chi-squared distance between two subgaussian distributions whose first $k$ moments are matched. Their bound suggests a similar exponential-order improvement as ours. However, it is clear from the context that the proof techniques in Lambert and Tierney (1984) and Wu and Yang (2020) are distinct from the current paper, where, as we detail next, the conclusion is arrived via a new Jackson-type bound.

**Remark 3.2.** In Theorem 3.1 the explicit value of $c$ was not exposed. For readers of interest, considering $\epsilon \in (0, 1)$ to be an arbitrarily small positive constant, the largest possible $c$ we can obtain for the Poisson mixture is

$$1/10 - \epsilon$$

and for negative binomial mixture is

$$\left[2\Big\{1 + 2 \cdot \frac{\log(e/\theta_*)}{\log(1/\theta_*)}\Big\}\right]^{-1} - \epsilon, \text{ recalling that } \theta_* \in (0, 1) \text{ in this case.}$$

Although it is certainly not the main interest of this paper to devise the sharpest possible value of $c$, it is our conjecture that for any fixed $\sigma$, (at the worst case) $c$ would have to be strictly smaller than $1/2$. In other words, the parametric root-$n$ rate as was observed in Goldfeld et al. (2020) cannot be recovered in the setting of nonparametric mixture MLEs considered in this paper. A detailed investigation of the lower bound of $\mathbb{E}W_1^\sigma(Q, \widehat{Q})$ is beyond the scope of this paper, but will be the topic of a subsequent future work.

Next we give a proof sketch of Theorem 2.2. Invoking the same trick that was used in the proof of Theorem 6.1(a) in Miao et al. (2021), for any given 1-Lipschitz function $\ell(\cdot)$ such that $\ell(0) = 0$,

we introduce the following function to approximate it,

$$\widehat{\ell}_k(\theta) := \sum_{x=0}^{k} b_{x,\ell} f(x|\theta), \quad \text{for } b_{x,\ell} \in \mathbb{R} \text{ and } \theta \in [0, \theta_*].$$

Some straightforward manipulations (see Section 4.2 for details) then yield

$$W_1(Q, \widehat{Q}) \leq \sup_{\ell \in \text{Lip}_1, \ell(0)=0} \left\{ 2 \sup_{\theta \in [0,\theta_*]} \left| \ell(\theta) - \widehat{\ell}_k(\theta) \right| + \sum_{x=0}^{k} b_{x,\ell} \left( h_Q(x) - h_Q^{\text{obs}}(x) \right) \right.$$
$$\left. + \sum_{x=0}^{k} b_{x,\ell} \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right\}, \quad (3.1)$$

where $h_Q^{\text{obs}}(x) := n^{-1} \sum_{i=1}^{n} \mathbb{1}(x = X_i)$. The sub-polynomial bound of $\mathbb{E} W_1(Q, \widehat{Q})$ could then be explained by the following fact (detailed proofs to be presented in Section 4.2): For any function $\ell$ considered above and any $k = 1, 2, \ldots$, there exists an $\widehat{\ell}_k$ and a constant $C = C(\theta_*) > 0$ only depending on $\theta_*$ such that the following two inequalities hold. First, a Jackson-type bound (see Lemma 5.4):

$$\sup_{\ell \in \text{Lip}_1, \ell(0)=0} \sup_{\theta \in [0,\theta_*]} \left| \ell(\theta) - \widehat{\ell}_k(\theta) \right| \leq C/k; \quad (3.2)$$

second,

$$\sup_{\ell \in \text{Lip}_1, \ell(0)=0} \max_{0 \leq x \leq k} \left| b_{x,\ell} \right| \leq C^k \max_{1 \leq x \leq k} \left\{ \frac{1}{w(x)} \right\}.$$

Plugging different bounds of $\max\{1/w(x)\}$ into the above two inequalities then gives us the desired results in Theorem 2.2.

With these concepts in mind, let us then move on to examine the case when the GOT distance is used. Similar to the derivation of (3.1) and further noting that $\int \ell \mathsf{d}(Q * \mathcal{N}_\sigma) = \int \ell_\sigma \mathsf{d}Q$ with

$$\ell_\sigma := \ell * \phi_\sigma,$$

one can show that

$$W_1^\sigma(Q, \widehat{Q}) \leq \sup_{\ell \in \text{Lip}_1, \ell(0)=0} \left\{ 2 \sup_{\theta \in [0,\theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - \widehat{\ell}_k(\theta) \right| + \sum_{x=0}^{k} b_{x,\ell} \left( h_Q(x) - h_Q^{\text{obs}}(x) \right) \right.$$
$$\left. + \sum_{x=0}^{k} b_{x,\ell} \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right\}. \quad (3.3)$$

The last two terms in (3.3) can be similarly handled as in (3.1), and it remains to control the first term. To this end, we introduce the following lemma, which turns out to be an extension of the Jackson-type one.

**Lemma 3.1** (Polynomial approximation of Gaussian-convoluted Lipschitz functions). *Let $0 \in [a, b] \subset \mathbb{R}$ be a bounded interval and let $\ell(\cdot)$ be a 1-Lipschitz function over $[a, b]$. For any $\sigma > 0$ and integer $k > 1$, there exist a constant $C = C(a, b) > 0$ only depending on $a, b$ and a polynomial $p_k(\cdot)$*

*of degree at most $k$ such that*

$$\sup_{\theta \in [a,b]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta) \right| \leq C_1 e\sigma \cdot \left[ \frac{2\sqrt{e}\sigma\sqrt{k}}{b-a} \right]^{-k} k^{-1/4},$$

*where we recall that $\ell_\sigma := \ell * \phi_\sigma$ with $\phi_\sigma$ standing for the density function of $\mathcal{N}_\sigma$.*

In striking contrast to the linear convergence in Jackson-type bounds such as (3.2), Lemma 3.1 states that approximation to Gaussian-convoluted Lipschitz functions by degree-$k$ polynomials is super-exponentially fast, hinting a substantial gain of convergence speed whence GOT distances are used to quantify the distance. We refer to Section 4.3 for the complete proof of Theorem 3.1.

## 4   Proofs

In the subsequent proofs, we sometimes drop the track of dependence on $C, C'$ for simplicity.

### 4.1   Proof of Theorem 2.1

*Proof.* The proof is based on Le Cam's two-point method (cf. Tsybakov (2009, Chapter 2.3)) and uses the following proposition.

**Proposition 4.1** (Lemma 3 in Tian et al. (2017), Proposition 4.3 in Vinayak et al. (2019).). *For any positive integer $k$ and any $M > 0$, there exist two distributions $P_1, P_2$ with support in $[0, M]$ such that $P_1, P_2$ have first $k$ moments identical and $W_1(P_1, P_2) \geq M/(2k)$.*

We first upper bound $g^{(x)}(0)\theta_*^x/x!$. To this end, define $\widetilde{g}(\theta) := g(\theta_* \theta)$. We then have

$$\frac{1}{\widetilde{g}(\theta)} = \sum_{x=0}^{\infty} w(x)\theta_*^x \theta^x \quad \text{for all } \theta \in [0,1]$$

by the definition of the mixing density function in (1.2). Furthermore, the radius of convergence of $\sum_{x=0}^{\infty} w(x)\theta_*^x \theta^x$ is $\theta_r/\theta_* > 1$. Accordingly, by Krantz and Parks (2002, Corollary 1.1.10), there exists some universal constant $C > 0$ such that

$$w(x)\theta_*^x \leq C \quad \text{for all } x = 0, 1, 2, \ldots.$$

The proof of Krantz and Parks (2002, Corollary 1.1.12) then yields that $\widetilde{g}(\theta)$, of the form $\widetilde{g}(\theta) = \sum_{x=0}^{\infty} \widetilde{g}^{(x)}(0)\theta^x/x!$, has a radius of convergence at least $1/(C+1)$. Invoking Krantz and Parks (2002, Corollary 1.1.10) again shows that there exists another universal constant $C' > 0$ such that

$$|g^{(x)}(0)\theta_*^x/x!| = |\widetilde{g}^{(x)}(0)/x!| \leq C'(C+1)^x, \quad \text{for all } x = 0, 1, 2, \ldots. \tag{4.1}$$

We then combine (4.1) with Proposition 4.1 to finish the proof. On one hand, for any $k = 1, 2, \ldots$, Proposition 4.1 guarantees the existence of two distributions $Q_1, Q_2$ over $[0, \theta_*/(C+3)]$ such that

$$\int \theta^x dQ_1(\theta) = \int \theta^x dQ_2(\theta), \quad \text{for all } x \in [k] \quad \text{and} \quad W_1(Q_1, Q_2) \geq \theta_*/(2(C+3)k).$$

On the other hand, the total variance distance between $h_{Q_1}$ and $h_{Q_2}$ satisfies

$$\mathrm{TV}(h_{Q_1}, h_{Q_2}) \le \frac{1}{2} \sum_{x=0}^{\infty} \Big| \int_0^{\theta_*} g(\theta) w(x) \theta^x \mathsf{d}Q_1(\theta) - \int_0^{\theta_*} g(\theta) w(x) \theta^x \mathsf{d}Q_2(\theta) \Big|$$

$$\le \sum_{x=0}^{\infty} w(x) \sum_{m: m+x \ge k+1} \frac{|g^{(m)}(0)|}{m!} \Big( \frac{\theta_*}{C+3} \Big)^{m+x}$$

$$\le 2(C+3)^2 C' \Big( \frac{C+2}{C+3} \Big)^k.$$

Picking $k = k(n)$ so that

$$2(C+3)^2 C' \Big( \frac{C+2}{C+3} \Big)^k = 1/(2n),$$

it follows from Le Cam's lower bound for two hypotheses that, denoting $Q^{\otimes n}$ to be the $n$-time product measure of $Q$,

$$\inf_{\widetilde{Q}} \sup_Q \mathbb{E} W_1(Q, \widetilde{Q}) \ge \frac{1}{2} W_1(Q_1, Q_2) \Big\{ 1 - \mathrm{TV}(h_{Q_1}^{\otimes n}, h_{Q_2}^{\otimes n}) \Big\}$$

$$\ge \frac{1}{2} W_1(Q_1, Q_2) \{ 1 - n/(2n) \} = \frac{1}{4} W_1(Q_1, Q_2),$$

with $W_1(Q_1, Q_2) \ge \theta_*/(2(C+3)k)$ by the construction. This completes the proof. $\quad\square$

## 4.2   Proof of Theorem 2.2

*Proof of Theorem 2.2.* By definition of $W_1$, we have

$$W_1(Q_1, Q_2) = \sup_{\ell \in \mathrm{Lip}_1} \int \ell(\mathsf{d}Q_1 - \mathsf{d}Q_2) = \sup_{\ell \in \mathrm{Lip}_1, \ell(0)=0} \int \ell(\mathsf{d}Q_1 - \mathsf{d}Q_2).$$

To control each $\int \ell(\mathsf{d}Q_1 - \mathsf{d}Q_2)$, define the following approximation function of $\ell(\theta)$:

$$\theta \mapsto \widehat{\ell}(\theta) := \sum_{x=0}^{\infty} b_x f(x|\theta), \text{ where } b_x \in \mathbb{R} \text{ and } \theta \in [0, \theta_*],$$

Recall that $h_Q(x) = \int f(x|\theta) \mathsf{d}Q(\theta)$. Then direct calculation yields that

$$\int_0^{\theta_*} \ell(\theta) \mathsf{d}\big(Q(\theta) - \widehat{Q}(\theta)\big) = \int_0^{\theta_*} \big(\ell(\theta) - \widehat{\ell}(\theta)\big) \mathsf{d}\big(Q(\theta) - \widehat{Q}(\theta)\big) + \sum_{x=0}^{\infty} b_x \big(h_Q(x) - h_{\widehat{Q}}(x)\big)$$

$$\le 2\|\ell - \widehat{\ell}\|_\infty + \Big| \sum_{x=0}^{\infty} b_x \big(h_Q(x) - h_Q^{\mathrm{obs}}(x)\big) \Big| + \Big| \sum_{x=0}^{\infty} b_x \big(h_Q^{\mathrm{obs}}(x) - h_{\widehat{Q}}(x)\big) \Big|,$$

where $\|\ell - \widehat{\ell}\|_\infty \equiv \sup_{\theta \in [0,\theta_*]} |\ell(\theta) - \widehat{\ell}(\theta)|$ and $h_Q^{\mathrm{obs}}(x) \equiv n^{-1} \sum_{i=1}^n \mathbf{1}_{X_i = x}$. This implies

$$W_1(Q, \widehat{Q}) \le \sup_{\ell \in \mathrm{Lip}(1)} \Big\{ 2\|\ell - \widehat{\ell}\|_\infty + \Big| \sum_{x=0}^{\infty} b_x \big(h_Q(x) - h_Q^{\mathrm{obs}}(x)\big) \Big| + \Big| \sum_{x=0}^{\infty} b_x \big(h_Q^{\mathrm{obs}}(x) - h_{\widehat{Q}}(x)\big) \Big| \Big\}. \quad (4.2)$$

By Lemmas 5.2 and 5.3, for an arbitrary $\delta \in (0, 1/2)$ and an arbitrary $\epsilon \in (0, 1)$, there exist constants $n_1 = n_1(\epsilon)$ and $C_1 = C_1(\epsilon, \theta_*)$ such that the sum of the last two terms in (4.2) is upper bounded by

$$C_1 \max_{x \ge 0} |b_x| / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}}$$

10

for all $n \geq n_1$ with probability at least $1 - 2\delta$. The bound on $\max_{x \geq 0} |b_x|$ depends on the tail of $1/w(x)$.

(i) If $1/w(x) \leq C_2^x$ for some universal constant $C_2 > 1$ and all $x \geq 1$, it follows from Lemma 5.4 that any 1-Lipschitz function $\ell(\theta)$ on $[0, \theta_*]$ can be approximated by $\widehat{\ell}(\theta) = \sum_{x=0}^{k} b_x f(x|\theta)$, such that $\max_{\theta \in [0, \theta_*]} |\ell(\theta) - \widehat{\ell}(\theta)| \leq C_3/k$ and

$$\max_{x \geq 0} |b_x| = \max_{x \in [0,k]} |b_x| \leq C_3^k/w(k) \leq (C_2 C_3)^k$$

for $k \geq 1$, where $C_3 = C_3(\theta_*) > 1$ is a constant. Hence it follows from (4.2) that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k + C_1 (C_2 C_3)^k / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}},$$

for any $n \geq 1$ with probability at least $1 - 2\delta$. Taking $k = k(n)$ such that $(C_2 C_3)^k = n^c$ for some small positive constant $c$ specified later, it follows that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k(n) + C_1 n^c / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}} = 2C_3/k(n) + C_1 n^{c+\epsilon/2-1/2}/\sqrt{\delta^{1+\epsilon}}.$$

Note that $(C_2 C_3)^{k(n)} = n^c$ implies $k(n) = c \log n / \log(C_2 C_3)$. Letting $\epsilon = 1/4$ and $c = 1/8$, it follows that

$$W_1(Q, \widehat{Q}) \leq 2C_3 \log(C_2 C_3)/(c \log n) + C_1 n^{c+\epsilon/2-1/2}/\sqrt{\delta^{1+\epsilon}}$$
$$\leq 16 C_3 \log(C_2 C_3)/\log n + C_1 n^{-1/4}/\delta^{5/8}.$$

Therefore, for sufficiently large $n$ (depending on $\theta_*$), there exists a positive constant $C_4 = C_4(\theta_*)$ such that $\mathbb{E} W_1(Q, \widehat{Q}) \lesssim \log n$ by integrating the tail estimate.

(ii) If $1/w(x) \leq (C_5 x)^{C_5 x}$ for some universal constant $C_5$ and all $x \geq 1$, it follows from Lemma 5.4 that any 1-Lipschitz function $\ell(\theta)$ on $[0, \theta_*]$ can be approximated by $\widehat{\ell}(\theta) = \sum_{x=0}^{k} b_x f(x|\theta)$ such that $\max_{\theta \in [0, \theta_*]} |\ell(\theta) - \widehat{\ell}(\theta)| \leq C_3/k$, and

$$\max_x |b_x| \leq C_3^k/w(k) \leq (C_5(C_3)^{1/C_5} k)^{C_5 k} \leq (C_6 k)^{C_6 k}$$

for $k \geq 1$, where $C_6 = C_6(\theta_*)$ is a constant. Hence it follows that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k + (C_6 k)^{C_6 k} C_1 / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}},$$

for any $n \geq 1$ with probability at least $1 - 2\delta$. Taking $k = k(n)$ satisfying $(C_6 k)^{C_6 k} = n^c$ for a small positive constant $c$ specified later, it follows that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k(n) + C_1 n^{c+\epsilon/2-1/2}/\sqrt{\delta^{1+\epsilon}}.$$

Since $(C_6 k)^{C_6 k} = n^c$ is equivalent to $\log(C_6 k) \exp(\log(C_6 k)) = c \log n$, it follows that $\log(C_6 k(n)) = W(c \log n)$ and hence $k(n) = \exp(W(c \log n))/C_6$, where $W(\cdot)$ is the Lambert W function. Using the expansion

$$W(x) = \log x - \log \log x + o(1), \text{ as } x \to \infty,$$

and hence there exists a universal constant $C_7 > 0$ such that

$$\exp(W(x)) \geq x/(2 \log x) \text{ for } x \geq C_7.$$

Therefore, for sufficiently large $n$, we have

$$k(n) \geq \frac{c \log n}{2C_6 \log(c \log n)}. \tag{4.3}$$

As a result,

$$W_1(Q, \widehat{Q}) \leq \{4C_3 C_6 \log(c \log n)\}/(c \log n) + C_1 n^{c+\epsilon/2-1/2}/\sqrt{\delta^{1+\epsilon}},$$

with probability at least $1 - 2\delta$. Letting $c = 1/8, \epsilon = 1/4$, we have

$$W_1(Q, \widehat{Q}) \lesssim \log \log n / \log n + n^{-1/4} \delta^{-5/8}.$$

Therefore, for sufficiently large $n$ (depending on $\theta_*$), it follows that $\mathbb{E}W_1(Q, \widehat{Q}) \lesssim \log \log n / \log n$ by integrating the tail estimate. $\qquad\qquad\square$

## 4.3 Proof of Theorem 3.1

*Proof of Theorem 3.1.* The proof is separated to three steps.

**Step 1.** In the first step, we prove that for any $\sigma > 0$, integer $k > 1$, and any $\ell \in \text{Lip}(1)$ on $[-\theta_*, \theta_*]$ with $\ell(0) = 0$, there exist a positive constant $C_4 = C_4(\theta_*, \sigma)$ and a set of coefficients

$$\left\{ b_x \in \mathbb{R}, x = 0, \ldots, 2k \right\}$$

such that

$$\sup_{\theta \in [0, \theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - \sum_{0 \leq x \leq 2k} b_x f(x|\theta) \right| \leq C_3 \left\{ \left[ \theta_* \sigma \sqrt{ek} \right]^{-k} + \sum_{x \geq k+1} w(x) \theta_*^x \right\},$$

where we recall that $\ell_\sigma(\theta) := [\ell * \phi_\sigma](\theta)$ and $\phi_\sigma$ is the probability density function of $\mathcal{N}_\sigma$.

For any $k = 1, 2, \ldots$, let $q_k(\theta) := \sum_{x=0}^{k} w(x) \theta^x$ be an approximation of the function $\theta \mapsto 1/g(\theta) = \sum_{x=0}^{\infty} w(x) \theta^x$ on $[0, \theta_*]$. Then one can readily verify that

$$R_k(\theta) := g(\theta) \cdot \left\{ \frac{1}{g(\theta)} - q_k(\theta) \right\} = g(\theta) \cdot \sum_{x \geq k+1} w(x) \theta^x \leq g(0) \cdot \sum_{x \geq k+1} w(x) \theta_*^x$$

whenever $\theta \in [0, \theta_*]$.

Let $p_k(\theta)$ be the degree-$k$ polynomial achieving the approximation bound in Lemma 3.1. We then have

$$\sup_{\theta \in [-\theta_*, \theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta) \right| \leq C_5 e \sigma \cdot \left[ 2\theta_*^{-1} \sigma \sqrt{ek} \right]^{-k}, \tag{4.4}$$

where $C_5 > 0$ is a universal constant.

Let's construct $\{b_x \in \mathbb{R}, x \in [2k]\}$ to be coefficients such that

$$p_k(\theta) q_k(\theta) = \sum_{x=0}^{2k} b_x w(x) \theta^x.$$

Then have $g(\theta)p_k(\theta)q_k(\theta) = \sum_{x=0}^{2k} b_x f(x|\theta)$, and the proof in this step is complete by noting that

$$\sup_{\theta\in[0,\theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta)q_k(\theta)g(\theta) \right|$$

$$= \sup_{\theta\in[0,\theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta)\big[1 - R_k(\theta)\big] \right|$$

$$\leq 2 \sup_{\theta\in[0,\theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta) \right| + \sup_{\theta\in[0,\theta_*]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) \right| \cdot \sup_{\theta\in[0,\theta_*]} \left| R_k(\theta) \right|$$

$$\overset{(*)}{\leq} C_5 2e\sigma \cdot \left[ 2\theta_*^{-1}\sigma\sqrt{ek} \right]^{-k} + 2(\theta_* + \sigma)g(0) \cdot \sum_{x\geq k+1} w(x)\theta_*^x \tag{4.5}$$

and

$$\sup_{\theta\in[0,\theta_*]} g(\theta)p_k(0)q_k(0) \leq g(0)p_k(0)q_k(0).$$

Here in $(*)$ we use the fact that, as $\ell(0) = 0$ and $\ell \in \mathrm{Lip}(1)$,

$$\sup_{\theta\in[-\theta_*,\theta_*]} \left| \ell_\sigma(\theta) \right| = \sup_{\theta\in[-\theta_*,\theta_*]} \left| \int \ell(\theta - \theta_1)\phi_\sigma(\theta_1)\mathrm{d}\theta_1 \right| \leq \int (\theta_* + |\theta_1|)\phi_\sigma(\theta_1)\mathrm{d}\theta_1 \leq \theta_* + \sigma.$$

**Step 2.** In this step, we upper bound $\max_{x\in[2k]} |b_x|$. Let

$$\widetilde{r}(\theta) := p_k(\theta_*\theta)q_k(\theta_*\theta) := \sum_{x=1}^{2k} \widetilde{b}_x w(x)\theta^x$$

be a rescaled version of $p_k(\theta)q_k(\theta)$, so that $\widetilde{b}_x = \theta_*^x b_x$. Then by Lemma 5.5, it holds that for each $1 \leq x \leq 2k$,

$$\left| \widetilde{b}_x \right| w(x) \leq \frac{(2k)^x}{x!} \sup_{|\theta|\leq 1} \left| \widetilde{r}(\theta) \right| \leq \frac{(2k)^x}{x!} \sup_{|\theta|\leq\theta_*} p_k(\theta) \cdot \sup_{|\theta|\leq\theta_*} q_k(\theta).$$

Since

$$\sup_{|\theta|\leq\theta_*} q_k(\theta) \leq 1/g(\theta_*)$$

and by (4.4),

$$\sup_{|\theta|\leq\theta_*} p_k(\theta) \leq C$$

for some positive constant $C$ only depending on $\theta_*$ and $\sigma$, it follows that

$$\max_{1\leq x\leq 2k} \left| b_x \right| \leq C_6 \max_{1\leq x\leq 2k} \frac{(2k)^x}{w(x)\theta_*^x x!} \leq C_6 \max_{1\leq x\leq 2k} \frac{1}{w(x)} \cdot \max_{1\leq x\leq 2k} \frac{1}{\theta_*^x} \cdot \max_{1\leq x\leq 2k} \frac{(2k)^x}{x!}$$

where $C_6 = C_6(\theta_*, \sigma) > 0$. Combining the above inequality with

$$\max_{1\leq x\leq 2k} 1/\theta_*^x \leq \left( \max\{1, 1/\theta_*\} \right)^{2k}$$

and

$$\max_{1\leq x\leq 2k} (2k)^x/x! \leq e^{2k},$$

13

it follows that

$$\max_{1 \leq x \leq 2k} \left| b_x \right| \leq C_6 \cdot \left( e \cdot \max\{1, 1/\theta_*\} \right)^{2k} \cdot \max_{1 \leq x \leq 2k} \frac{1}{w(x)}.$$

**Step 3.** In this step we prove the claim of the theorem. Recall that

$$W_1^\sigma(\widehat{Q}, Q) = \sup_\ell \int \ell \mathsf{d}[\widehat{Q} * \mathcal{N}_\sigma] - \ell \mathsf{d}[Q * \mathcal{N}_\sigma],$$

where $\ell \in \mathrm{Lip}(1)$ with $\ell(0) = 0$. It further holds that

$$
\begin{aligned}
W_1^\sigma(\widehat{Q}, Q) &= \sup_{\ell \in \mathrm{Lip}(1):\ell(0)=0} \int \ell \mathsf{d}[\widehat{Q} * \mathcal{N}_\sigma] - \ell \mathsf{d}[Q * \mathcal{N}_\sigma] \\
&= \sup_{\ell \in \mathrm{Lip}(1):\ell(0)=0} \int (\ell_\sigma(\theta) - \ell_\sigma(0))[\mathsf{d}\widehat{Q} - \mathsf{d}Q] \\
&= \sup_{\ell \in \mathrm{Lip}(1):\ell(0)=0} \int \left\{ \ell_\sigma(\theta) - \ell_\sigma(0) - \sum_{0 \leq x \leq 2k} b_x f(x|\theta) \right\} [\mathsf{d}\widehat{Q} - \mathsf{d}Q] + \\
&\quad \sup_{\ell \in \mathrm{Lip}(1):\ell(0)=0} \int \sum_{0 \leq x \leq 2k} b_x f(x|\theta)[\mathsf{d}\widehat{Q} - \mathsf{d}Q] \\
&:= (I) + (II).
\end{aligned}
$$

By Step 2, we have

$$(I) \leq 2C_4 \left\{ \left[ 2\theta_*^{-1} \sigma \sqrt{ek} \right]^{-k} + \sum_{x \geq k+1} w(x)\theta_*^x \right\}. \tag{4.6}$$

Next we bound $(II)$. Recall that

$$h^{\mathrm{obs}}(x) := \sum_{i=1}^n \mathbb{1}(X_i = x)/n.$$

We have

$$\int \sum_{0 \leq x \leq 2k} b_x f(x|\theta)[\mathsf{d}\widehat{Q}(\theta) - \mathsf{d}Q(\theta)] \leq \left| \sum_{0 \leq x \leq 2k} b_x[h_{\widehat{Q}}(x) - h^{\mathrm{obs}}(x)] \right| + \left| \sum_{0 \leq x \leq 2k} b_x[h^{\mathrm{obs}}(x) - h_Q(x)] \right|.$$

It follows from Lemma 5.2 that for any $\delta > 0$, it holds with probability $1 - \delta$ that

$$\left| \sum_{0 \leq x \leq 2k} b_x \left[ h^{\mathrm{obs}}(x) - h_Q(x) \right] \right| \leq \max_{0 \leq x \leq 2k} |b_x| \sqrt{\frac{\log(2/\delta)}{2n}}.$$

Moreover, it follows from Lemma 5.3 that for an arbitrary $\delta \in (0,1)$ and an arbitrary $\epsilon \in (0,1)$, there exists a constant $C_7 = C_7(\epsilon, \theta_*) > 0$ such that

$$\left| \sum_{x=0}^{2k} b_x \left\{ h_Q^{\mathrm{obs}}(x) - h_{\widehat{Q}}(x) \right\} \right| \leq C_7 \max_{0 \leq x \leq 2k} |b_x| \sqrt{\frac{1}{n^{1-\epsilon}\delta^{1+\epsilon}}}$$

holds with probability at least $1 - \delta$.

Consequently, we have

$$(II) \leq C_8 \max_{0 \leq x \leq 2k} |b_x| \Big/ \sqrt{n^{1-\epsilon}\delta^{1+\epsilon}}$$

14

with probability at least $1 - \delta$ for some constant $C_8 = C_8(\epsilon, \theta_*) > 0$. Note that $\max_{0 \leq x \leq 2k} |b_x|$ have been upper bounded in Step 2.

Putting together the estimates for $(I)$ and $(II)$, we have that with probability at least $1 - \delta$, $W_1^\sigma(Q, \widehat{Q})$ is upper bounded by

$$\left[2\theta_*^{-1}\sigma\sqrt{ek}\right]^{-k} + \sum_{x \geq k+1} w(x)\theta_*^x + \left(e \cdot \max\{1, 1/\theta_*\}\right)^{2k} \cdot \max_{1 \leq x \leq 2k} \frac{1}{w(x)} \Big/ \sqrt{n^{1-\epsilon}\delta^{1+\epsilon}}. \tag{4.7}$$

up to a constant depending on $\sigma, \theta_*$ and $\epsilon$.

(i) If $c_1 c_2^x \leq 1/w(x) \leq C_1 C_2^x$, (4.7) becomes

$$[2\theta_*^{-1}\sigma\sqrt{ek}]^{-k} + \sum_{x \geq k+1} w(x)\theta_*^x + C_9^{2k}/\sqrt{n^{1-\epsilon}\delta^{1+\epsilon}},$$

where $C_9 := e \cdot \max\{1, 1/\theta_*\} \cdot \max\{1, C_2\}$ is a positive constant. For the second term, it follows from Krantz and Parks (2002, Corollary 1.1.10) that for any $R \in (\theta_*, \theta_r)$ there exists some constant $C_{10} = C_{10}(R) > 0$ such that $w(x) \leq C_{10}/R^x$ for all $x = 0, 1, 2\ldots$, and hence

$$\sum_{x \geq k+1} w(x)\theta_*^x \leq C_{10} \sum_{x \geq k+1} (\theta_*/R)^x \leq C_{10} \cdot (\theta_*/[R - \theta_*]) \cdot [\theta_*/R]^k \text{ for any } k = 1, 2, \ldots.$$

Therefore, the second term dominates the first term in (4.7), and (4.7) becomes

$$[\theta_*/R]^k + C_9^{2k}/\sqrt{n^{1-\epsilon}\delta^{1+\epsilon}}.$$

The proof is then complete by letting $C_9^{2k} = n^\alpha$ for some $\alpha \in (0, 1/2 - \epsilon/2)$. The final bound is then $n^{-\frac{(1-\epsilon)\log(R/\theta_*)}{2\log(R/\theta_*)+4\log C_9}}$ for any $\epsilon \in (0, 1)$.

(ii) If $c_1 c_2^x x^{c_3 x} \leq 1/w(x) \leq C_1 C_2^x x^{C_3 x}$, (4.7) becomes

$$[2\theta_*^{-1}\sigma\sqrt{ek}]^{-k} + (C_{11}k)^{-c_3 k} + (C_{12}k)^{2C_3 k}/\sqrt{n^{1-\epsilon}\delta^{1+\epsilon}}$$

for some positive constants $C_{11}$ and $C_{12}$ and the proof is then complete by letting $(C_{12}k)^{2C_3 k} = n^\alpha$ for some $\alpha \in (0, 1/2 - \epsilon/2)$. The final bound is then $n^{-\frac{(1-\epsilon)/2}{1+\max\{4C_3, 2C_3/c_3\}}}$. $\qquad\square$

# 5 Auxiliary results

## 5.1 Auxiliary lemmas

**Lemma 5.1** (Theorem 6.2 in Chapter 7, DeVore (1976)). *For any integer $r \geq 1$, let*

$$W_\infty^r([-1, 1]) := \Big\{\psi : [-1, 1] \to \mathbb{R} : \psi^{(r-1)} \text{ is absolutely continuous and}$$

$$\text{the supremum of } \psi^{(r)} \text{ on } [-1, 1] \text{ is finite}\Big\}$$

*be the Sobolev space on $[-1, 1]$. For functions $f \in W_\infty^r([-1, 1])$ and any integer $k > r$, there exists a polynomial $p_k$ of degree at most $k$ such that*

$$\sup_{\theta \in [-1,1]} \left|f(\theta) - p_k(\theta)\right| \leq Ck^{-r}\omega\big(f^{(r)}, k^{-1}\big),$$

*where $C > 0$ is a universal constant and*

$$\omega(f^{(r)}, k^{-1}) := \sup_{\theta_1, \theta_2 : |\theta_1 - \theta_2| \leq k^{-1}} \left| f^{(r)}(\theta_1) - f^{(r)}(\theta_2) \right|.$$

Recall that for a sample $\{X_i, i \in [n]\}$ and $x \in \mathbb{N}$, $h_Q^{\text{obs}}(x) = n^{-1} \sum_{i=1}^n \mathbf{1}_{X_i = x}$. The following lemmas provides the concentration of $h_Q^{\text{obs}}$ around $h_Q$.

**Lemma 5.2** (Lemma A.1 in Miao et al. (2021)). *Let $\{X_i, i \in [n]\}$ be an i.i.d. sample generated from the probability mass function $h_Q$ in (1.1). Then for any $\delta \in (0,1)$ the following inequality holds with probability at least $1 - \delta$,*

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_Q(x) \right) \right| \leq \max_{x \geq 0} |b_x| \sqrt{\frac{\log(2/\delta)}{2n}},$$

*where $b_x \in \mathbb{R}$ for all $x \in \mathbb{N}$.*

**Lemma 5.3** (A generalized version of Lemma A.2 in Miao et al. (2021)). *Let $\{X_i, i \in [n]\}$ be an i.i.d. sample generated from the mixture distribution $h_Q$ in (1.1). Then for an arbitrary $\delta \in (0,1)$ and an arbitrary $\epsilon \in (0,1)$, there exists a constant $C = C(\epsilon, \theta_*) > 0$ such that for any $n \geq 1$,*

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right| \leq C_1 \max_{x \geq 0} |b_x| \sqrt{\frac{1}{n^{1-\epsilon} \delta^{1+\epsilon}}}$$

*holds with probability at least $1 - \delta$. Here $b_x \in \mathbb{R}$ for all $x \in \mathbb{N}$.*

**Lemma 5.4** (A generalized version of Proposition A.2. in Miao et al. (2021)). *For any 1-Lipschitz function $\theta \mapsto \ell(\theta)$ on $[0, \theta_*]$ with $\ell(0) = 0$, there exists some $\widehat{\ell}(\theta) = \sum_{x=0}^{k} b_x f(x|\theta)$ such that $\max_{\theta \in [0, \theta_*]} |\ell(\theta) - \widehat{\ell}| \leq C/k$, and*

$$\max_{x \in [0,k]} |b_x| \leq C^k \cdot \max_{1 \leq x \leq k} 1/w(x) \text{ for } k \geq 1,$$

*where $C = C(\theta_*)$ is a positive constant. It can be further proved that there exists some universal constant $C' > 0$ such that*

$$C^x / w(x) \geq e^x / C'$$

*for all nonnegative integer $x$.*

**Lemma 5.5** (Chapter 2.6 Equation 9 in Timan (2014)). *Suppose $k$ is a non-negative integer and $\theta \mapsto p_k(\theta) \equiv \sum_{x=0}^{k} c_x \theta^x$. Then it follows that coefficients $\{c_x\}_{x=0}^{k}$ satisfy*

$$|c_x| \leq \frac{k^x}{x!} \max_{|\theta| \leq 1} |p_k(\theta)|.$$

**Lemma 5.6** (Jackson's theorem, Lemma 10 of Han and Shiragur (2021) or see DeVore (1976)). *Let $k > 0$ be any integer, and $[a, b] \subseteq \mathbb{R}$ be any bounded interval. For any 1-Lipschitz function $\ell(\cdot)$ on $[a, b]$, there exists a universal constant $C$ independent of $k, \ell$ such that there exists a polynomial $p_k(\cdot)$ of degree at most $k$ such that*

$$|\ell(\theta) - p_k(\theta)| \leq C \sqrt{(b-a)(\theta - a)}/k, \ \forall \theta \in [a, b]. \tag{5.1}$$

*In particular, the following norm bound holds:*

$$\sup_{\theta \in [a,b]} |\ell(\theta) - p_k(\theta)| \leq C(b-a)/k. \tag{5.2}$$

## 5.2 Proofs of Remarks

*Proof of Remark 2.4:* (1) If $1/w(x) \leq \exp_L(C_9 x)$ for some universal constant $C_9$ and all $x \geq 1$, it follows from Lemma 5.4 that any 1-Lipschitz function $\ell(\theta)$ on $[0, \theta_*]$ can be approximated by $\widehat{\ell}(\theta) = g(\theta) \sum_{x=0}^{k} b_x w(x) \theta^x$ with an uniform approximation error of $C_3/k$ with

$$\max_x |b_x| \leq C_3^k / w(k) \leq C_3^k \exp_L(C_9 k) \leq \exp_L(C_{10} k)$$

for $k \geq 1$, where $C_{10} = C_{10}(\theta_*)$ is a constant. Hence it follows from the first steps in the proof of Theorem 2.2 that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k + C_1 \exp_L(C_{10} k) / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}},$$

for $n \geq n_1$ with probability at least $1 - 2\delta$. Analogously, by letting $k = k(n)$ such that $\exp_L(C_{10} k) = n^c$ for a small $c$, we then have $E\{W_1(Q, \widehat{Q})\} \leq C_{11}/\log_L(n)$, where $C_{11} = C_{11}(\theta_*)$ is a constant.

(2) If $1/w(x) \leq (C_{12} x)^{\cdots (C_{12} x)}$ ($L \in \mathbb{N}^+$ times power) for some universal constant $C_{12}$ and all $x \geq 1$, it follows from Lemma 5.4 that any 1-Lipschitz function $\ell(\theta)$ on $[0, \theta_*]$ can be approximated by $\widehat{\ell}(\theta) = g(\theta) \sum_{x=0}^{k} b_x w(x) \theta^x$ with an uniform approximation error of $C_3/k$ with

$$\max_x |b_x| \leq C_3^k / w(k) \leq C_3^k (C_{12} k)^{\cdots (C_{12} k)} \leq (C_{13} k)^{\cdots (C_{13} k)}$$

for $k \geq 1$, where $C_{13} = C_{13}(\theta_*)$ is a constant. Hence it follows from the first steps in the proof of Theorem 2.2 that

$$W_1(Q, \widehat{Q}) \leq 2C_3/k + C_1 (C_{13} k)^{\cdots (C_{13} k)} / \sqrt{n^{1-\epsilon} \delta^{1+\epsilon}},$$

for $N \geq N_1$ with probability at least $1 - 2\delta$. By letting $k = k(n)$ such that $(C_{13} k)^{\cdots (C_{13} k)} = n^c$ for a small $c$, we have $E\{W_1(Q, \widehat{Q})\} \leq C_{14} \log_L(n)/\log_{L-1}(n)$, where $C_{14} = C_{14}(\theta_*)$ is a constant. $\qquad \square$

## 5.3 Proofs of Lemmas

*Proof of Lemma 2.1.* Recall the duality definition of $W_1(\mu_1, \mu_2)$ as

$$W_1(\mu_1, \mu_2) := \inf \mathbb{E} \|X - Y\|,$$

with the infimum taken over all couplings of $(X, Y)$ such that $X \sim \mu_1$ and $Y \sim \mu_2$. We then consider any such $(X, Y)$ and assume $Z$ to be independent of $(X, Y)$ and follows the distribution of $\nu$. Then it is immediate that

$$W_1(\mu_1 * \nu, \mu_2 * \nu) \leq \mathbb{E} \|(X + Z) - (Y + Z)\| = \mathbb{E} \|X - Y\|,$$

and accordingly (by taking infimum over all such $(X, Y)$)

$$W_1(\mu_1 * \nu, \mu_2 * \nu) \leq W_1(\mu_1, \mu_2).$$

This completes the proof. $\qquad \square$

*Proof of Lemma 3.1.* By rescaling, we assume that $a = -1$ and $b = 1$. For any integer $r \geq 1$, let

$$W_\infty^r([a,b]) := \Big\{ \psi : [a,b] \to \mathbb{R} : \psi^{(r-1)} \text{ is absolutely continuous and}$$

$$\text{the essential supremum of } \psi^{(r)} \text{ on } [a,b] \text{ is finite} \Big\}$$

be the Sobolev space on $[a,b]$. Then it is readily verifiable that for any $\ell \in \text{Lip}(1)$ and $\sigma^2 > 0$,

$$\ell_\sigma(\theta) - \ell_\sigma(0) = (\ell * \phi_\sigma)(\theta) - (\ell * \phi_\sigma)(0),$$

when restricted on $[a,b]$, belongs to $W_\infty^r([a,b])$. Hence by Lemma 5.1, we have that for any integer $k > r$, there exists some polynomial $p_k$ of degree $k$ such that

$$\sup_{\theta \in [a,b]} \left| \ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta) \right| \leq C_1 k^{-r} \omega\big(\ell_\sigma^{(r)}, k^{-1}\big),$$

In the above inequality, $C_1 = C_1(a,b) > 0$ is a constant and

$$\omega(\psi, t) := \sup_{\theta_1, \theta_2 : |\theta_1 - \theta_2| \leq t} |\psi(\theta_1) - \psi(\theta_2)|$$

is the modulus of continuity of function $\psi$ at radius $t$. To bound the righthand side of the above display, note that, with $H_n(\cdot)$ denoting the $n$-th Hermite polynomial, we have

$$\ell_\sigma^{(r)}(\theta) = \int \ell(\theta_1) \phi_\sigma^{(r)}(\theta - \theta_1) d\theta_1 = \sigma^{-r}(-1)^r \int \ell(\theta - \theta_1) \phi_\sigma(\theta_1) H_r(\theta_1/\sigma) d\theta_1.$$

Hence for any $\theta_1, \theta_2$ such that $|\theta_1 - \theta_2| \leq k^{-1}$, we have

$$\left| \ell_\sigma^{(r)}(\theta_1) - \ell_\sigma^{(r)}(\theta_2) \right| \leq \sigma^{-r} \int |\ell(\theta_1 - \theta) - \ell(\theta_2 - \theta)| \phi_\sigma(\theta) |H_r(\theta/\sigma)| d\theta$$

$$\leq \sigma^{-r} k^{-1} \int \phi_\sigma(\theta) |H_r(\theta/\sigma)| d\theta = \sigma^{-r} k^{-1} \int \phi_1(\theta) |H_r(\theta)| d\theta$$

$$\leq \sigma^{-r} k^{-1} \Big[ \int \phi_1(\theta) H_r^2(\theta) d\theta \Big]^{1/2} = \sigma^{-r} k^{-1} \sqrt{r!}.$$

It further follows from the Sterling formula $\sqrt{r!} \leq \sqrt{e r^{r+1/2} e^{-r}}$ that

$$\left| \ell_\sigma^{(r)}(\theta_1) - \ell_\sigma^{(r)}(\theta_2) \right| \leq \sigma^{-r} k^{-1} \sqrt{e r^{r+1/2} e^{-r}}.$$

Using $r < k$, we hence obtain

$$\sup_{\theta \in [a,b]} |\ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta)| \leq C_1 \sqrt{e} (\sqrt{e} \sigma k / \sqrt{r})^{-r} r^{1/4} k^{-1} \leq C_1 \sqrt{e} (\sqrt{e} \sigma \sqrt{k})^{-r} k^{-3/4}.$$

By rescaling, we then have for any $a \leq 0, b \geq 0$ it follows that

$$\sup_{\theta \in [a,b]} |\ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta)| \leq C_1 ([b-a]/2)^{r+1} \sqrt{e} (\sqrt{e} \sigma \sqrt{k})^{-r} k^{-3/4}$$

$$\leq \frac{C_1 (b-a) \sqrt{e}}{2} \cdot \Big[ \frac{2\sqrt{e} \sigma \sqrt{k}}{b-a} \Big]^{-r} k^{-3/4}.$$

Now taking $r = k - 1$, we have

$$\sup_{\theta \in [a,b]} |\ell_\sigma(\theta) - \ell_\sigma(0) - p_k(\theta)| \leq C_1 e \sigma \cdot \Big[ \frac{2\sqrt{e} \sigma \sqrt{k}}{b-a} \Big]^{-k} k^{-1/4}$$

and accordingly complete the proof. $\qquad\square$

*Proof of Lemma 5.3.* Whenever there is no ambiguity, let $h_Q^{\text{obs}}$, $h_{\widehat{Q}}$, and $h_Q$ also represent distributions with respect to corresponding probability mass functions $x \mapsto h_Q^{\text{obs}}(x)$, $x \mapsto h_{\widehat{Q}}(x)$, and $x \mapsto h_Q(x)$.

This proof consists of two steps. In the first step, we prove that

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right|$$

can be upper bounded by $\text{KL}(h_Q^{\text{obs}}, h_Q)$, where KL is the Kullback–Leibler divergence. In the second step, we upper bound $\text{KL}(h_Q^{\text{obs}}, h_Q)$ by truncation arguments.

**Step 1.** It follows from the triangle inequality that

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right| \leq \max_{x \geq 0} |b_x| \sum_{x=0}^{\infty} \left| h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right| = \max_{x \geq 0} |b_x| \cdot \left\| h_Q^{\text{obs}} - h_{\widehat{Q}} \right\|_1,$$

where $\left\| h_Q^{\text{obs}} - h_{\widehat{Q}} \right\|_1$ represents the total variation distance between distributions $h_Q^{\text{obs}}$ and $h_{\widehat{Q}}$. It further follows from Pinsker's inequality that

$$\left\| h_Q^{\text{obs}} - h_{\widehat{Q}} \right\|_1 \leq \sqrt{\frac{1}{2} \cdot \text{KL}(h_Q^{\text{obs}}, h_{\widehat{Q}})},$$

and hence

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right| \leq \max_{x \geq 0} |b_x| \sqrt{\frac{1}{2} \cdot \text{KL}(h_Q^{\text{obs}}, h_{\widehat{Q}})} \leq \max_{x \geq 0} |b_x| \sqrt{\frac{1}{2} \cdot \text{KL}(h_Q^{\text{obs}}, h_Q)},$$

by noting that maximum likelihood estimators maximize likelihood functions.

**Step 2.** Suppose $C_1 = C_1(\theta_*)$ is the smallest positive integer larger than $\theta_* g(0)(1/g)'(\theta_*)$. Define

$$T_i := X_i \mathbb{1}(X_i \leq C_1 - 1) + C_1 \mathbb{1}(X_i \geq C_1) \text{ for all } i \in [N].$$

Let $t_Q$ be the probability mass function of $T_1$ and let $t_Q^{\text{obs}}$ be the sample version of $t_Q$, i.e.

$$x \mapsto t_Q(x) := P(T_1 = x) \text{ and } x \mapsto t_Q^{\text{obs}}(x) := \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(T_i = x), \text{ for } x \in \{0, \ldots, C_1\}.$$

Note that

$$t_Q(x) = h_Q(x) \quad \text{and} \quad t_Q^{\text{obs}}(x) = h_Q^{\text{obs}}(x) \text{ for } x = 0, \ldots, C_1 - 1$$

and

$$t_Q(C_1) = \sum_{x \geq C_1} h_Q(x), \quad t_Q^{\text{obs}}(C_1) = \sum_{x \geq C_1} h_Q^{\text{obs}}(x).$$

Hence it follows that

$$\begin{aligned}
\text{KL}(h_Q^{\text{obs}}, h_Q) &= \sum_{x=0}^{C_1-1} t_Q^{\text{obs}}(x) \log \frac{t_Q^{\text{obs}}(x)}{t_Q(x)} + \sum_{x \geq C_1} h_Q^{\text{obs}}(x) \log \frac{h_Q^{\text{obs}}(x)}{h_Q(x)} \\
&= \text{KL}(t_Q^{\text{obs}}, t_Q) - t_Q^{\text{obs}}(C_1) \log \frac{t_Q^{\text{obs}}(C_1)}{t_Q(C_1)} + \sum_{x \geq C_1} h_Q^{\text{obs}}(x) \log \frac{h_Q^{\text{obs}}(x)}{h_Q(x)},
\end{aligned}$$

where $t_Q^{\mathrm{obs}}$ and $t_Q$ are viewed as distributions with respect to corresponding probability mass functions of $x \mapsto t_Q(x)$ and $x \mapsto t_Q^{\mathrm{obs}}(x)$.

If $t_Q^{\mathrm{obs}}(C_1) = 0$, then

$$t_Q^{\mathrm{obs}}(C_1) \log \frac{t_Q^{\mathrm{obs}}(C_1)}{t_Q(C_1)} = 0.$$

Otherwise it follows from the inequality

$$\log(1 + x) \leq x \text{ for } x > 0$$

that

$$-t_Q^{\mathrm{obs}}(C_1) \log \frac{t_Q^{\mathrm{obs}}(C_1)}{t_Q(C_1)} \leq \sum_{x \geq C_1} \left\{ h_Q(x) - h_Q^{\mathrm{obs}}(x) \right\}.$$

Analogously, we have

$$\sum_{x \geq C_1} h_Q^{\mathrm{obs}}(x) \log \frac{h_Q^{\mathrm{obs}}(x)}{h_Q(x)} \leq \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} + \sum_{x \geq C_1} \left\{ h_Q^{\mathrm{obs}}(x) - h_Q(x) \right\}$$

and hence

$$-t_Q^{\mathrm{obs}}(C_1) \log \frac{t_Q^{\mathrm{obs}}(C_1)}{t_Q(C_1)} + \sum_{x \geq C_1} h_Q^{\mathrm{obs}}(x) \log \frac{h_Q^{\mathrm{obs}}(x)}{h_Q(x)} \leq \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)}.$$

**Step 2(a).** We first upper bound $\sum_{x \geq C_1} (h_Q^{\mathrm{obs}}(x) - h_Q(x))^2 / h_Q(x)$. Fix an arbitrary $\epsilon \in (0, 1)$ and choose a $\gamma > 0$ in $(1 - \epsilon, 1)$. Define $A := \alpha^{(1-\gamma)/3}$, where $\alpha := (\theta_* + \theta_r)/(2\theta_*) > 1$. Note that $\alpha \theta_* < \theta_r$ and we have $1/g(\theta) = \sum_{x=0}^{\infty} w(x) \theta^x < \infty$ for all $\theta \in [0, \alpha \theta_*]$. It then follows from Hölder's inequality that

$$n^{1-\epsilon} \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} = n^{1-\epsilon} \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{-x} A^x$$

$$\leq n^{1-\epsilon} \left( \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{-x/\gamma} \right)^\gamma \left( \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{x/(1-\gamma)} \right)^{1-\gamma}.$$

It further follows from $A > 1$ that

$$n \cdot \mathbb{E}\left\{ \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{-x/\gamma} \right\} = \sum_{x \geq C_1} (1 - h_Q(x)) A^{-x/\gamma} \leq \sum_{x \geq C_1} A^{-x/\gamma} = \frac{A^{-C_1/\gamma}}{1 - A^{-1/\gamma}} < \infty$$

and hence for an arbitrary $\delta \in (0, 1)$, we have

$$n \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{-x/\gamma} \leq \frac{A^{-C_1/\gamma}}{1 - A^{-1/\gamma}} \frac{1}{\delta}$$

with probability at least $1 - \delta$. Therefore, with probability at least $1 - \delta$, we have

$$\left( \sum_{x \geq C_1} \frac{(h_Q^{\mathrm{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{-x/\gamma} \right)^\gamma \leq \left( \frac{A^{-C_1/\gamma}}{1 - A^{-1/\gamma}} \frac{1}{N\delta} \right)^\gamma \leq \frac{1}{\left( a^{\frac{1-\gamma}{3\gamma}} - 1 \right)^\gamma} \frac{1}{(N\delta)^\gamma},$$

where the last inequality follows from $C_1 \geq 1$. On the other hand,

$$\sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{x/(1-\gamma)} \leq \sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x))^2}{h_Q(x)} \alpha^{x/3} + \sum_{x \geq C_1} h_Q(x) \alpha^{x/3}.$$

We first show that the second term on the right hand side is bounded, which is true if

$$\sum_{x \geq C_1} h_Q(x) \alpha^x \leq g(\theta_*) \Big/ g(\alpha \theta_*).$$

Since $\sum_{x=0}^{\infty} g(\theta) w(x) \theta^x = 1$ and $1/g(\theta) = \sum_{x=0}^{\infty} w(x) \theta^x$, it follows from

$$(1/g)'(\theta) = \sum_{x=1}^{\infty} x w(x) \theta^{x-1} > 0 \quad \text{and} \quad (1/g)''(\theta) = \sum_{x=2}^{\infty} x(x-1) w(x) \theta^{x-2} > 0$$

that $g(\cdot)$ is monotonically decreasing on $[0, \theta_*]$ and $(1/g)'(\cdot)$ is monotonically increasing on $[0, \theta_*]$. Therefore, it follows from

$$\log f(x|\theta) = \log(1 - \pi) - \log(1/g(\theta)) + x \log \theta + \log w(x)$$

that

$$\frac{\mathrm{d}(\log f(x|\theta))}{\mathrm{d}\theta} = \frac{1}{\theta} \left( x - \theta g(\theta)(1/g)'(\theta) \right) \geq \frac{1}{\theta} \left( x - \theta_* g(0)(1/g)'(\theta_*) \right) \geq \frac{1}{\theta} (x - C_1) \geq 0$$

for all $x \geq C_1$. Therefore we have

$$h_Q(x) = \int_0^{\theta_*} f(x|\theta) dQ \leq \sup_{\theta \in [0, \theta_*]} f(x|\theta) = f(x|\theta_*)$$

and

$$\sum_{x \geq C_1} h_Q(x) \alpha^x \leq \sum_{x \geq C_1} f(x|\theta_*) \alpha^x \leq \sum_{x \geq 0} f(x|\theta_*) \alpha^x \leq \frac{g(\theta_*)}{g(\alpha \theta_*)} < \infty.$$

For any fixed $k > 0$, define $A_n$ to be the event

$$A_n := \left\{ h_Q^{\text{obs}}(x) > k h_Q(x) \alpha^{x/3} \text{ for some } x \geq C_1 \right\}.$$

Then, it follows from Markov's inequality that

$$P(A_n) \leq \sum_{x \geq C_1} P(h_Q^{\text{obs}}(x) > k h_Q(x) \alpha^{x/3}) \leq \frac{1}{k} \sum_{x \geq C_1} \mathbb{E}\{h_Q^{\text{obs}}(x)\} \frac{1}{h_Q(x) \alpha^{x/3}} \leq \frac{1}{k} \frac{1}{\alpha^{1/3} - 1}.$$

Thus, $P(A_n)$ can be made arbitrarily small by choosing $k$ large enough and on the complement of $A_n$ we have

$$\sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x))^2}{h_Q(x)} \alpha^{x/3} \leq k^2 \sum_{x \geq C_1} h_Q(x) \alpha^x \leq k^2 \frac{g(\theta_*)}{g(\alpha \theta_*)}.$$

Therefore, for an arbitrary $\delta \in (0, 1)$, we have

$$\sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x))^2}{h_Q(x)} \alpha^{x/3} \leq \frac{g(\theta_*)}{g(\alpha \theta_*)} \left( \frac{1}{\delta} \frac{1}{\alpha^{1/3} - 1} \right)^2$$

with probability at least $1 - \delta$. Thus, for an arbitrary $\delta \in (0, 1)$, with probability at least $1 - \delta$, it

follows that

$$\left\{ \sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x) - h_Q(x))^2}{h_Q(x)} A^{\frac{x}{1-\gamma}} \right\}^{1-\gamma} \leq \left\{ \frac{g(\theta_*)}{g(\alpha\theta_*)} \left( \frac{1}{\delta} \frac{1}{\alpha^{1/3} - 1} \right)^2 + \frac{g(\theta_*)}{g(\alpha\theta_*)} \right\}^{1-\gamma} \leq \frac{C_2}{\delta^{2-2\gamma}},$$

where $C_2 = C_2(\theta_*) = g(\theta_*)[1/(\alpha^{1/3} - 1)^2 + 1]/g(\alpha\theta_*)$ is a constant. For an arbitrary $\delta \in (0, 1/2)$, with probability at least $1 - 2\delta$, it follows that

$$n^{1-\epsilon} \sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x) - h_Q(x))^2}{h_Q(x)} \leq n^{1-\epsilon} \frac{1}{\left(\alpha^{\frac{1-\gamma}{3\gamma}} - 1\right)^\gamma} \frac{1}{(n\delta)^\gamma} \frac{C_2}{\delta^{2-2\gamma}} = n^{1-\epsilon-\gamma} \frac{C_2}{\left(\alpha^{\frac{1-\gamma}{3\gamma}} - 1\right)^\gamma} \frac{1}{\delta^{2-\gamma}}.$$

Thus, by letting $\gamma$ go to $1 - \epsilon$, we have

$$\sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x) - h_Q(x))^2}{h_Q(x)} \leq \frac{C_2}{\left(\alpha^{\frac{\epsilon}{3(1-\epsilon)}} - 1\right)^{1-\epsilon}} \frac{1}{n^{1-\epsilon}} \frac{1}{\delta^{1+\epsilon}}.$$

As a result, for arbitrary $\delta \in (0, 1/2)$ and $\epsilon \in (0, 1)$, with probability at least $1 - 2\delta$, we have

$$\text{KL}(h^{\text{obs}}, h_Q) \leq \text{KL}(t_Q^{\text{obs}}, t_Q) + \sum_{x \geq C_1} \frac{(h_Q^{\text{obs}}(x) - h_Q(x))^2}{h_Q(x)} \leq \text{KL}(t_Q^{\text{obs}}, t_Q) + C_3 \frac{1}{n^{1-\epsilon}} \frac{1}{\delta^{1+\epsilon}},$$

where $C_3 = C_3(\epsilon, \theta_*) = C_2/(\alpha^{\frac{\epsilon}{3(1-\epsilon)}} - 1)^{1-\epsilon}$.

**Step 2(b).** We then upper bound $\text{KL}(t_Q^{\text{obs}}, t_Q)$. It follows from Mardia et al. (2020) that with probability at least $1 - \delta$,

$$\text{KL}(t^{\text{obs}}, t_Q) \leq \frac{C_1 + 1}{2n} \log \frac{4n}{C_1 + 1} + \frac{1}{n} \log \frac{3e}{\delta},$$

and hence for any $\epsilon \in (0, 1)$ and $\delta \in (0, 1/3)$, with probability at least $1 - 3\delta$,

$$\text{KL}(h_Q^{\text{obs}}, h_Q) \leq \frac{1}{N\delta^{1+\epsilon}} \left( 3C_1 \log(2n) + C_3 n^\epsilon \right).$$

Therefore, it follows that there exists a constant $C_4 = C_4(\epsilon, \theta_*)$ such that for any $n \geq 1$

$$\text{KL}(h_Q^{\text{obs}}, h_Q) \leq \frac{C_4}{n^{1-\epsilon}\delta^{1+\epsilon}}$$

holds with probability at least $1 - 3\delta$ for any $\epsilon \in (0, 1)$ and $\delta \in (0, 1/3)$. Therefore,

$$\left| \sum_{x=0}^{\infty} b_x \left( h_Q^{\text{obs}}(x) - h_{\widehat{Q}}(x) \right) \right| \leq \max_{x \geq 0} |b_x| \sqrt{\frac{C_4}{2n^{1-\epsilon}\delta^{1+\epsilon}}}$$

holds for all $n \geq n_1$ with probability at least $1 - 3\delta$ for any $\epsilon \in (0, 1)$ and $\delta \in (0, 1/3)$. $\square$

*Proof of Lemma 5.4.* This proof consists of two steps. In the first step, we prove the existence of $\widehat{\ell}$ and upper bound the difference between $\widehat{\ell}$ and $\ell$. In the second step, we upper bound coefficients of $\widehat{\ell}$, i.e., $\max_{x \geq 0} |b_x|$.

**Step 1.** It follows from $\sum_{x=0}^{\infty} f(x|\theta) = 1$ that $\sum_{x=0}^{\infty} g(\theta)w(x)\theta^x = 1$ and hence $g(\theta) > 0$ for $\theta \in [0, \theta_*]$. As a consequence, $1/g(\theta) = \sum_{x=0}^{\infty} w(x)\theta^x$ on $[0, \theta_*]$.

Since $\theta \mapsto \sum_{x=0}^{\infty} w(x)\theta^x$ is a continuous function on $[-\theta_*, \theta_*]$ with $w(0) > 0$, there exists a universal constant $\theta_0 \in (0, \theta_*]$ such that $\theta \mapsto \sum_{x=0}^{\infty} w(x)\theta^x$ is strictly positive on $[-\theta_0, \theta_*]$. For $\theta \in [-\theta_0, 0)$, define $1/g(\theta) := \sum_{x=0}^{\infty} w(x)\theta^x$ and $\ell(\theta) := -\ell(-\theta)$. Then $\theta \mapsto \ell(\theta)$ is a 1-Lipschitz

22

function on $[-\theta_0, \theta_*]$ and for any $\theta_1, \theta_2 \in [-\theta_0, \theta_*]$ we have

$$|\ell(\theta_1)/g(\theta_1) - \ell(\theta_2)/g(\theta_2)| \leq |\ell(\theta_1)/g(\theta_1) - \ell(\theta_2)/g(\theta_1)| + |\ell(\theta_2)/g(\theta_1) - \ell(\theta_2)/g(\theta_2)|$$
$$\leq |\theta_1 - \theta_2|\{1/g(\theta_*) + \theta_*(1/g)'(\theta_*)\}.$$

Therefore, it follows from Jackson's theorem (see Lemma 5.6) that there exists a polynomial $\sum_{x=0}^{k} v_x \theta^x$ of degree $k \geq 1$ such that

$$\sup_{\theta \in [-\theta_0, \theta_*]} |\ell(\theta)/g(\theta) - \sum_{x=0}^{k} v(x)\theta^x| \leq C_1/k,$$

where $C_1 = C_1(\theta_*)$ is a positive constant independent of $k$ and $\ell$ and $v_x \in \mathbb{R}$ for all $x = 0, \ldots, k$. Let $b_x = v_x/\{w(x)(1-\pi)\}$ for $x = 1, \ldots, k$ and $b_x = 0$ for $x = 0$. Then it follows from

$$|v_0| \leq C_1/k + |\ell(0)/g(0)| = C_1/k$$

that

$$\sup_{\theta \in [-\theta_0, \theta_*]} \left| \frac{\ell(\theta)}{g(\theta)} - b_0 \frac{\pi + (1-\pi)g(\theta)w(0)}{g(\theta)} - (1-\pi) \sum_{x=1}^{k} b_x w(x)\theta^x \right| \leq 2C_1/k,$$

and hence

$$\sup_{\theta \in [-\theta_0, \theta_*]} \left| \ell(\theta) - \widehat{\ell}(\theta) \right| \leq C_2/k,$$

where

$$\widehat{\ell}(\theta) := b_0\pi + (1-\pi)g(\theta) \sum_{x=0}^{k} b_x w(x)\theta^x = \sum_{x=0}^{k} b_x f(x|\theta)$$

and $C_2 = C_2(\theta_*)$ is a positive constant independent of $k$ and $\ell$.

**Step 2**. To bound the coefficients $b_x$'s, we first define a polynomial

$$\theta \mapsto r(\theta) := \sum_{x=0}^{k} v_x(\theta_0\theta)^x \text{ on } [-1, 1]$$

and note that

$$\sup_{\theta \in [-1,1]} |r(\theta)| \leq C_1/k + \sup_{\theta \in [-\theta_0, \theta_0]} |\ell(\theta)/g(\theta)| \leq C_1/k + \theta_0/g(\theta_0).$$

We then apply Lemma 5.5 on the polynomial $r(\theta)$, and it follows that

$$|v(x)|\theta_0^x \leq \max_{|\theta| \leq 1} |r(\theta)| \cdot k^x/x! \leq C_3 k^x/x!,$$

where $C_3 = C_3(\theta_*)$ is a positive constant. Hence

$$|b_x| = |v_x|/\{w(x)(1-\pi)\} \leq C_3 \cdot (k/\theta_0)^x/\{x!w(x)(1-\pi)\}$$

and

$$\max_{x \in [0,k]} |b_x| \leq \frac{C_3}{1-\pi} \cdot \max_{x \in [0,k]} \frac{(k/\theta_0)^x}{x!w(x)} \leq C_4 \cdot \frac{(k/\theta_0)^k}{k!} \cdot \max_{1 \leq x \leq k} 1/w(x) \leq C_4 \cdot (e/\theta_0)^k \cdot \max_{1 \leq x \leq k} 1/w(x),$$

where $C_4 = C_4(\theta_*)$ is a positive constant. It follows from Krantz and Parks (2002, Corollary 1.1.10)

that $w(x) \leq C_5/\theta_*^x$ for all $x \in \mathbb{N}$ and some universal constant $C_5 \geq 1$ and hence

$$(e/\theta_0)^k/w(k) \geq (\theta_* e/\theta_0)^k/C_5 \geq e^k/C_5 > 1$$

for all sufficiently large $k$. $\qquad\square$

# References

Cameron, A. C. and Trivedi, P. K. (2013). *Regression Analysis of Count Data*. Cambridge University Press.

Chen, H.-B. and Niles-Weed, J. (2021+). Asymptotics of smoothed Wasserstein distances. *Potential Analysis*, (in press).

Chen, J. (2017). Consistency of the MLE under mixture models. *Statistical Science*, 32(1):47–63.

DeVore, R. A. (1976). Degree of approximation. *Approximation theory II*, 241(242):117–161.

Fournier, N. and Guillin, A. (2015). On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3):707–738.

Goldfeld, Z. and Greenewald, K. (2020). Gaussian-smoothed optimal transport: Metric structure and statistical efficiency. In *International Conference on Artificial Intelligence and Statistics*, pages 3327–3337. PMLR.

Goldfeld, Z., Greenewald, K., Niles-Weed, J., and Polyanskiy, Y. (2020). Convergence of smoothed empirical measures with applications to entropy estimation. *IEEE Transactions on Information Theory*, 66(7):4368–4391.

Goldfeld, Z., van den Berg, E., Greenewald, K. H., Melnyk, I., Nguyen, N., Kingsbury, B., and Polyanskiy, Y. (2019). Estimating information flow in deep neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2299–2308. PMLR.

Gupta, R. C. (1984). Estimating the probability of winning (losing) in a gambler's ruin problem with applications. *Journal of Statistical Planning and Inference*, 9(1):55–62.

Han, Y. and Shiragur, K. (2021). On the competitive analysis and high accuracy optimality of profile maximum likelihood. In *Proceedings of 2021 ACM-SIAM Symposium on Discrete Algorithms*, pages 1317–1336. SIAM.

Hengartner, N. W. (1997). Adaptive demixing in Poisson mixture models. *The Annals of Statistics*, 25(3):917–928.

Jackson, D. (1921). The general theory of approximation by polynomials and trigonometric sums. *Bulletin of the American Mathematical Society*, 27(9-10):415–431.

Janardan, K. (1982). A new discrete exponential family of distributions: Properties and application to power series distributions. *American Journal of Mathematical and Management Sciences*, 2(2):145–158.

Krantz, S. G. and Parks, H. R. (2002). *A Primer of Real Analytic Functions*. Springer Science and Business Media.

Lambert, D. and Tierney, L. (1984). Asymptotic properties of maximum likelihood estimates in the mixed Poisson model. *The Annals of Statistics*, pages 1388–1399.

Lindsay, B. G. (1995). Mixture models: theory, geometry and applications. In *NSF-CBMS Regional Conference Series in Probability and Statistics*, pages i–163. JSTOR.

Loh, W.-L. and Zhang, C.-H. (1996). Global properties of kernel estimators for mixing densities in discrete exponential family models. *Statistica Sinica*, pages 561–578.

Loh, W.-L. and Zhang, C.-H. (1997). Estimating mixing densities in exponential family models for discrete variables. *Scandinavian Journal of Statistics*, 24(1):15–32.

Mardia, J., Jiao, J., Tánczos, E., Nowak, R. D., and Weissman, T. (2020). Concentration inequalities for the empirical distribution of discrete distributions: beyond the method of types. *Information and Inference: A Journal of the IMA*, 9(4):813–850.

Miao, Z., Kong, W., Vinayak, R. K., Sun, W., and Han, F. (2021). Fisher-Pitman permutation tests based on nonparametric Poisson mixtures with application to single cell genomics. *arXiv preprint arXiv:2106.03022*.

Noack, A. (1950). A class of random variables with discrete distributions. *The Annals of Mathematical Statistics*, 21(1):127–132.

Roueff, F. and Rydén, T. (2005). Nonparametric estimation of mixing densities for discrete distributions. *The Annals of Statistics*, 33(5):2066–2108.

Sadhu, R., Goldfeld, Z., and Kato, K. (2021). Limit distribution theory for the smooth 1-Wasserstein distance with applications. *arXiv preprint arXiv:2107.13494*.

Simar, L. (1976). Maximum likelihood estimation of a compound Poisson process. *The Annals of Statistics*, 4(6):1200–1209.

Stoyanov, J. and Lin, G. D. (2011). Mixtures of power series distributions: identifiability via uniqueness in problems of moments. *Annals of the Institute of Statistical Mathematics*, 63(2):291–303.

Tian, K., Kong, W., and Valiant, G. (2017). Learning populations of parameters. In *Advances in Neural Information Processing Systems*, volume 30, pages 5778–5787.

Timan, A. F. (2014). *Theory of Approximation of Functions of a Real Variable*. Elsevier.

Tsybakov, A. B. (2009). *Introduction to Nonparametric Estimation.* Springer.

Tucker, H. G. (1963). An estimate of the compounding distribution of a compound Poisson distribution. *Theory of Probability and Its Applications*, 8(2):195–200.

Vinayak, R. K., Kong, W., Valiant, G., and Kakade, S. (2019). Maximum likelihood estimation for learning populations of parameters. In *Proceedings of the 36th International Conference on Machine Learning*, pages 6448–6457. PMLR.

von Renesse, M.-K. and Sturm, K.-T. (2005). Transport inequalities, gradient estimates, entropy and Ricci curvature. *Communications on Pure and Applied Mathematics*, 58(7):923–940.

Walter, G. and Hamedani, G. (1991). Bayes empirical bayes estimation for natural exponential families with quadratic variance functions. *The Annals of Statistics*, 19(3):1191–1224.

Wu, Y. and Yang, P. (2020). Optimal estimation of Gaussian mixtures via denoised method of moments. *The Annals of Statistics*, 48(4):1981–2007.

Zhang, C.-H. (1995). On estimating mixing densities in discrete exponential family models. *The Annals of Statistics*, 23(3):929–945.

Zhang, Y., Cheng, X., and Reeves, G. (2021). Convergence of Gaussian-smoothed optimal transport distance with sub-gamma distributions and dependent samples. In *International Conference on Artificial Intelligence and Statistics*, pages 2422–2430. PMLR.