# Local continuity of log-concave projection, with applications to estimation under model misspecification

Rina Foygel Barber\* and Richard J. Samworth<sup>†</sup> December 18, 2020

#### Abstract

The log-concave projection is an operator that maps a d-dimensional distribution P to an approximating log-concave density. It is known that, with suitable metrics on the underlying spaces, this projection is continuous, but not uniformly continuous. In this work we prove a local uniform continuity result for log-concave projection—in particular, establishing that this map is locally Hölder-(1/4) continuous. A matching lower bound verifies that this exponent cannot be improved. We also examine the implications of this continuity result for the empirical setting—given a sample drawn from a distribution P, we bound the squared Hellinger distance between the log-concave projection of the empirical distribution of the sample, and the log-concave projection of P. In particular, this yields interesting statistical results for the misspecified setting, where P is not itself log-concave.

# 1 Introduction

In nonparametric statistics and inference, many problems are formulated in terms of shape constraints. Examples include isotonic regression and convex regression (for supervised learning problems, placing constraints on the shape of the regression function relating the response to the covariates), and monotone or log-concave density estimation (for unsupervised learning problems, placing constraints on a distribution that is the target we wish to estimate).

Among these examples, log-concave density estimation is especially challenging in that it cannot be formulated as an  $L_2$ -projection onto a convex constraint set. Remarkably, projection onto the space of log-concave densities can still be uniquely defined, but unlike a convex projection, this operation is not uniformly continuous (Dümbgen et al., 2011) and its mathematical and statistical properties are therefore difficult to analyze. In this work, we examine the continuity properties of log-concave projection more closely to establish locally uniform convergence, and study the statistical implications of these results.

<sup>\*</sup>Department of Statistics, University of Chicago

<sup>&</sup>lt;sup>†</sup>Statistical Laboratory, University of Cambridge

## 1.1 Background

We begin by establishing some notation used throughout the paper, and then give background on log-concave projection and its known properties.

#### 1.1.1 Notation

Throughout the paper,  $\|\cdot\|$  denotes the usual Euclidean norm. For a distribution P, we write  $\mathbb{E}_P[\cdot]$  and  $\mathbb{P}_P\{\cdot\}$  to denote expectation or probability taken with respect to a random variable or vector X drawn from distribution P, and  $\mu_P := \mathbb{E}_P[X]$  denotes its mean. We will analogously write  $\mathbb{E}_f[\cdot]$ ,  $\mathbb{P}_f\{\cdot\}$ , and  $\mu_f$  for a density f. We say a distribution, density, or random vector is isotropic if it has zero mean and identity covariance matrix. Given  $x \in \mathbb{R}^d$  and r > 0, we write  $\mathbb{B}_d(x,r) := \{y \in \mathbb{R}^d : \|y - x\| \le r\}$  for the closed Euclidean ball of radius r centered at x,  $\mathbb{B}_d(r) = \mathbb{B}_d(0,r)$  for the closed Euclidean ball of radius r centered at zero, and  $\mathbb{S}_{d-1}(r) := \{y \in \mathbb{R}^d : \|y\| = r\}$  for the sphere of radius r centered at zero. For the unit ball and unit sphere we write  $\mathbb{B}_d = \mathbb{B}_d(1)$  and  $\mathbb{S}_{d-1} = \mathbb{S}_{d-1}(1)$ . For  $x \in \mathbb{R}$ ,  $(x)_+$  denotes  $\max\{x,0\}$ , and  $(x)_-$  denotes  $\max\{-x,0\}$ . For independent observations  $X_1,\ldots,X_n \in \mathbb{R}^d$ , we will write  $\widehat{P}_n$  to denote the empirical distribution. We write  $\operatorname{Leb}_d$  for Lebesgue measure on  $\mathbb{R}^d$ .

The  $L_1$ -Wasserstein distance  $d_W$  is defined for two distributions P, Q on  $\mathbb{R}^d$  as

$$\mathrm{d}_{\mathrm{W}}(P,Q) := \inf \left\{ \mathbb{E}_{\tilde{P}} \left[ \|X - Y\| \right] : \quad \begin{array}{l} \text{Distributions } \tilde{P} \text{ on } (X,Y) \in \mathbb{R}^d \times \mathbb{R}^d \\ \text{such that marginally } X \sim P \text{ and } Y \sim Q \end{array} \right\} \in [0,+\infty].$$

For any distributions P, Q on  $\mathbb{R}^d$ , this infimum is attained for some coupling  $\tilde{P}$  (Villani, 2008, Theorem 4.1). We will also use the Hellinger distance  $d_H$ , defined for densities f, g on  $\mathbb{R}^d$  as

$$\mathrm{d}^2_\mathrm{H}(f,g) := \int_{\mathbb{R}^d} \left( \sqrt{f(x)} - \sqrt{g(x)} \right)^2 \; \mathrm{d}x.$$

The Hellinger distance is known to satisfy  $0 \le d_H^2(f,g) \le \min\{2, d_{KL}(f||g)\}$  for any densities f, g, where  $d_{KL}(f||g) := \mathbb{E}_f \left[ \log \left( f(X)/g(X) \right) \right]$  is the Kullback–Leibler divergence. Both  $d_W$  and  $d_H$  satisfy the triangle inequality, while  $d_{KL}$  does not.

#### 1.1.2 The log-concave projection

For any  $d \in \mathbb{N}$ , let  $\mathcal{P}_d$  denote the set of probability distributions P on  $\mathbb{R}^d$  satisfying  $\mathbb{E}_P[\|X\|] < \infty$  and  $\mathbb{P}_P\{X \in H\} < 1$  for every hyperplane  $H \subseteq \mathbb{R}^d$ , that is, P does not place all its mass in any hyperplane. Further, let  $\mathcal{F}_d$  denote the set of all upper semi-continuous, log-concave densities on  $\mathbb{R}^d$ . Then, by Dümbgen et al. (2011, Theorem 2.2), there exists a well-defined projection  $\psi^* : \mathcal{P}_d \to \mathcal{F}_d$ , given by

$$\psi^*(P) := \underset{f \in \mathcal{F}_d}{\operatorname{argmax}} \mathbb{E}_P \left[ \log f(X) \right].$$

When  $P \in \mathcal{P}_d$  has a (Lebesgue) density  $f_P$  satisfying  $\mathbb{E}_{f_P}[|\log f_P(X)|] < \infty$ , we can see that  $\psi^*(P)$  is the (unique) minimizer over  $f \in \mathcal{F}_d$  of the Kullback–Leibler divergence from  $f_P$  to f—since the KL divergence acts as a sort of distance, we can think of  $f = \psi^*(P)$  as the

"closest" log-concave density to  $f_P$ , which explains the use of the terminology 'projection' to describe this map. In particular, if  $f_P$  itself is log-concave, then  $\psi^*(P) = f_P$ .

To see the gain of defining  $\psi^*$  more broadly (i.e., on all distributions  $P \in \mathcal{P}_d$ , rather than only on distributions with densities), consider the empirical setting, where  $\widehat{P}_n$  is the empirical distribution of a sample. Then the result of Dümbgen et al. (2011, Theorem 2.2) tells us that, provided the convex hull of the data is d-dimensional, there exists a unique log-concave maximum likelihood estimator. We can therefore carry out log-concave density estimation via maximum likelihood in much the same way as if the class  $\mathcal{F}_d$  were a standard parametric model. To understand the estimation properties of this procedure, suppose we metrise  $\mathcal{P}_d$  with the  $L_1$ -Wasserstein distance  $d_W$ , and metrise  $\mathcal{F}_d$  with the Hellinger distance  $d_H$ . Then, by Dümbgen et al. (2011, Theorem 2.15), the map  $\psi^*$  is continuous. For the empirical distribution  $\widehat{P}_n$  obtained by drawing a sample  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} P$ , we therefore have

$$d_{\mathrm{H}}(\psi^*(\widehat{P}_n), \psi^*(P)) \stackrel{\mathrm{a.s.}}{\to} 0.$$

(This follows from the above continuity result because, by Varadarajan's theorem (Dudley, 2002, Theorem 11.4.1) and the strong law of large numbers, it holds that  $d_{\mathbf{W}}(\widehat{P}_n, P) \stackrel{\text{a.s.}}{\to} 0$ .) Thus, if  $P \in \mathcal{P}_d$  has a log-concave density, then the log-concave maximum likelihood estimator is strongly consistent—and moreover, even if the log-concavity is misspecified, then the estimator  $\psi^*(\widehat{P}_n)$  still converges to the log-concave projection  $\psi^*(P)$  of P. In this sense, then, the log-concave maximum likelihood estimator converges to the closest element of  $\mathcal{F}_d$  to P, so can be regarded as robust to misspecification.

Despite these positive results establishing continuity and consistency of  $\psi^*$ , however, the situation appears much less promising when it comes to obtaining rates of convergence (e.g., via a Lipschitz-type property of the map). Indeed, we cannot hope for Lipschitz continuity of this map, since the review article by Samworth (2018) gives the following example to show that  $\psi^*$  is not even uniformly continuous: let  $P^{(n)} = \text{Unif}[-1/n, 1/n]$  and  $Q^{(n)} = \text{Unif}[-1/n^2, 1/n^2]$ . Then  $d_{\mathbf{W}}(P^{(n)}, Q^{(n)}) \to 0$ , but since  $P^{(n)}$  and  $Q^{(n)}$  have log-concave densities  $f^{(n)} := \frac{n}{2} \mathbb{1}_{[-1/n,1/n]}$  and  $g^{(n)} := \frac{n^2}{2} \mathbb{1}_{[-1/n^2,1/n^2]}$  respectively, we deduce that

$$d_{H}(\psi^{*}(P^{(n)}), \psi^{*}(Q^{(n)})) = d_{H}(f^{(n)}, g^{(n)}) \to 0.$$
(1)

Summary of contributions While we have seen that log-concave projection does not satisfy uniform continuity, a natural question is whether it may be possible to place further restrictions on the class  $\mathcal{P}_d$  to obtain a result of this type. Moreover, from the statistical point of view, we would like to find a uniform rate of convergence for  $d_H(\psi^*(\hat{P}_n), \psi^*(P))$ , where  $\hat{P}_n$  is the empirical distribution of a sample of size n drawn from  $P \in \mathcal{P}_d$ , which again might require stronger assumptions than simply  $P \in \mathcal{P}_d$ .

The first main result of this paper (Theorem 2) reveals that the metric space map  $\psi^*: (\mathcal{P}_d, d_{\mathrm{W}}) \to (\mathcal{F}_d, d_{\mathrm{H}})$  is locally Hölder-(1/4) continuous, which establishes a precise understanding of the continuity properties of log-concave projection. Theorem 4 establishes a matching lower bound, revealing that the exponent 1/4 cannot be improved. Next, we specialise to the empirical setting, proving a bound on  $\mathbb{E}_P\left[d_{\mathrm{H}}^2(\psi^*(\widehat{P}_n), \psi^*(P))\right]$  in Theorem 5. For  $d \geq 2$ , this result is a straightforward consequence of combining our main result

in Theorem 2 with the recent work of Lei (2020), which bounds  $d_W(\widehat{P}_n, P)$  in expectation, while the case d = 1 requires a completely different approach. To the best of our knowledge, this work provides the first understanding of the range of possible rates of convergence of the log-concave maximum likelihood estimator in the misspecified setting.

## 1.2 Outline of paper

The remainder of the paper is organized as follows. In Section 2 we present our main results, establishing the local Hölder continuity of log-concave projection, and examining the empirical setting, as described above. We review prior work on log-concave projection and related problems in Section 3. The proofs of our main results are presented in Section 4, with technical details deferred to Appendix A.

## 2 Main results

As mentioned in Section 1, Dümbgen et al. (2011, Theorem 2.15) show that the log-concave projection operator  $\psi^*$  satisfies continuity with respect to appropriate metrics:

The log-concave projection 
$$\psi^*: (\mathcal{P}_d, d_W) \to (\mathcal{F}_d, d_H)$$
 is a continuous map. (2)

Our main results examine the continuity of the log-concave projection operator  $\psi^*$  more closely, and establish local uniform continuity results. To do this, we first introduce, for any distribution P on  $\mathbb{R}^d$  with  $\mathbb{E}_P[\|X\|] < \infty$ , the quantity

$$\epsilon_P := \inf_{u \in \mathbb{S}_{d-1}} \mathbb{E}_P \left[ \left| u^\top (X - \mu_P) \right| \right].$$

The quantity  $\epsilon_P$  can be thought of as a robust analogue of the minimum eigenvalue of the covariance matrix of the distribution P (note that its definition does not require P to have a finite second moment). We can also interpret  $\epsilon_P$  as measuring the extent to which P avoids placing all its mass on a single hyperplane.

First, we verify that  $\epsilon_P$  is positive for all  $P \in \mathcal{P}_d$ , and is Lipschitz with respect to the Wasserstein distance.

**Proposition 1.** We have  $\epsilon_P > 0$  for any  $P \in \mathcal{P}_d$ . Furthermore,  $|\epsilon_P - \epsilon_Q| \leq 2d_W(P,Q)$  for any distributions P, Q on  $\mathbb{R}^d$  with  $\mathbb{E}_P[||X||], \mathbb{E}_Q[||X||] < \infty$ .

We now present our first main result, which shows that  $\epsilon_P$  allows for a more detailed analysis of the continuity of the map  $\psi^*$ .

**Theorem 2.** For any  $d \ge 1$  and  $P, Q \in \mathcal{P}_d$ ,

$$d_{\mathrm{H}}(\psi^*(P), \psi^*(Q)) \le C_d \cdot \left[\frac{d_{\mathrm{W}}(P, Q)}{\max\{\epsilon_P, \epsilon_Q\}}\right]^{1/4},$$

where  $C_d > 0$  depends only on d.

This upper bound immediately implies the continuity result (2), but more importantly, to the best of our knowledge, this is the first general, quantitative statement about the local continuity of log-concave projection. Another consequence is that, when d=1, the uniform continuity counterexample in (1) is in some sense canonical: if  $(P^{(n)})$  and  $(Q^{(n)})$  are sequences in  $\mathcal{P}_1$  satisfying  $d_W(P^{(n)}, Q^{(n)}) \to 0$  and  $\lim \inf_{n \to \infty} \max\{\epsilon_{P^{(n)}}, \epsilon_{Q^{(n)}}\} > 0$ , then  $d_H(\psi^*(P^{(n)}), \psi^*(Q^{(n)})) \to 0$ .

#### 2.1 Extension to affine transformations

By Dümbgen et al. (2011, Remark 2.4), log-concave projection commutes with affine transformations; i.e., if  $\psi^*(P) = f$  then  $\psi^*(\mathbf{A} \circ P) = \mathbf{A} \circ f$  for any invertible matrix  $\mathbf{A}$ , where  $\mathbf{A} \circ P$  denotes the distribution obtained by drawing  $X \sim P$  and returning  $\mathbf{A}X$ , and similarly  $\mathbf{A} \circ f$  denotes the density of the random variable obtained by drawing X according to density f and returning  $\mathbf{A}X$ .

Turning to the terms appearing in Theorem 2, the Hellinger distance is invariant to affine transformations, but the terms on the right-hand side—namely,  $d_W(P,Q)$  and  $\max\{\epsilon_P, \epsilon_Q\}$ —are not. By considering affine transformations, we obtain the following corollary to Theorem 2, which we state without further proof:

Corollary 3. For any  $d \geq 1$  and  $P, Q \in \mathcal{P}_d$ ,

$$d_{\mathrm{H}}(\psi^*(P), \psi^*(Q)) \leq C_d \cdot \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}, \mathrm{rank}(\mathbf{A}) = d} \left[ \frac{d_{\mathrm{W}}(\mathbf{A} \circ P, \mathbf{A} \circ Q)}{\max\{\epsilon_{\mathbf{A} \circ P}, \epsilon_{\mathbf{A} \circ Q}\}} \right]^{1/4},$$

where  $C_d > 0$  depends only on d.

## 2.2 A matching lower bound

To see that our main result in Theorem 2 is optimal in terms of its dependence on the Wasserstein distance  $d_W(P,Q)$  and on the terms  $\epsilon_P, \epsilon_Q$ , we now construct an explicit example to provide a matching lower bound.

**Theorem 4.** Fix any  $d \ge 1$ ,  $\epsilon > 0$ , and  $\delta > 0$ . Then there exist distributions  $P, Q \in \mathcal{P}_d$  with  $\epsilon_P, \epsilon_Q \ge \epsilon$  and  $d_W(P, Q) \le \delta$ , such that

$$d_{\mathrm{H}}(\psi^*(P), \psi^*(Q)) \ge c_d \cdot \min\{1, (\delta/\epsilon)^{1/4}\},$$

where  $c_d > 0$  depends only on dimension d.

The theorem will be proved using the following construction: Let  $P \in \mathcal{P}_d$  be the uniform distribution on the sphere  $\mathbb{S}_{d-1}(\rho)$ , where  $\rho \propto \epsilon$ , and let  $Q \in \mathcal{P}_d$  be the mixture distribution that, with probability  $\beta \propto \delta/\epsilon$ , draws uniformly from  $\mathbb{S}_{d-1}(2\rho)$ , and with probability  $1 - \beta$  draws uniformly from  $\mathbb{S}_{d-1}(\rho)$ . Then  $d_W(P,Q) = \rho\beta \propto \delta$ , and we will see that  $d_H(\psi^*(P), \psi^*(Q)) \propto (\delta/\epsilon)^{1/4}$ , as desired.

## 2.3 Bounds for empirical processes

Now let  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} P \in \mathcal{P}_d$ , with corresponding empirical distribution function  $\widehat{P}_n$ . Under an additional moment assumption on P, we consider the problem of bounding  $d_H^2(\psi^*(\widehat{P}_n), \psi^*(P))$ . However, to be fully precise, we need to consider the possibility that  $\psi^*(\widehat{P}_n)$  may not be defined—specifically, if P places positive probability on some hyperplane  $H \subseteq \mathbb{R}^d$ , then it is possible that the empirical distribution  $\widehat{P}_n$  may place all its mass on this hyperplane, in which case we have  $\widehat{P}_n \notin \mathcal{P}_d$  and  $\psi^*(\widehat{P}_n)$  is not defined. In a slight abuse of notation, for such a case we will interpret  $d_H^2(\psi^*(\widehat{P}_n), \psi^*(P))$  as the maximum possible squared Hellinger distance (i.e., 2).

**Theorem 5.** Fix any  $P \in \mathcal{P}_d$ , and assume that

$$\mathbb{E}_P\left[\|X\|^q\right]^{1/q} \le M_q$$

for some q > 1. Let  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} P$  for some  $n \geq 2$ , and let  $\widehat{P}_n$  denote the corresponding empirical distribution. Then

$$\mathbb{E}\left[d_{\mathrm{H}}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P))\right] \leq C_{d,q} \cdot \sqrt{\frac{M_{q}}{\epsilon_{P}}} \cdot \frac{\log^{3/2} n}{n^{\min\left\{\frac{1}{2d}, \frac{1}{2} - \frac{1}{2q}\right\}}},$$

where  $C_{d,q} > 0$  depends only on d and q.

Proof of Theorem 5. First we consider the case  $d \ge 2$ . The result will follow by combining the bound (4), obtained from Theorem 2, together with a bound on the expected Wasserstein distance between  $\widehat{P}_n$  and P (Lei, 2020). Specifically, Lei (2020, Theorem 3.1) establishes that

$$\mathbb{E}\left[\mathrm{d}_{\mathrm{W}}(\widehat{P}_{n}, P)\right] \leq \widetilde{C}_{q} M_{q} \cdot \frac{\log^{2} n}{n^{\min\left\{\frac{1}{2}, \frac{1}{d}, 1 - \frac{1}{q}\right\}}}$$
(3)

for some  $\tilde{C}_q > 0$  depending only on q. Furthermore, on the event that  $\hat{P}_n \in \mathcal{P}_d$  (i.e.,  $\hat{P}_n$  does not place all its mass in any hyperplane), then by applying Theorem 2 with  $Q = \hat{P}_n$  we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)) \leq C_{d}^{2} \cdot \frac{d_{\mathrm{W}}^{1/2}(\widehat{P}_{n}, P)}{\max\{\epsilon_{P}^{1/2}, \epsilon_{\widehat{P}_{n}}^{1/2}\}}.$$

If instead  $\widehat{P}_n$  does place all its mass in a hyperplane and so  $\psi^*(\widehat{P}_n)$  is undefined, then in this case we have  $\epsilon_{\widehat{P}_n} = 0$ , and so by Proposition 1,  $2d_W(\widehat{P}_n, P) \geq |\epsilon_{\widehat{P}_n} - \epsilon_P| = \epsilon_P$ . Recalling from above that we interpret  $d_H^2(\psi^*(\widehat{P}_n), \psi^*(P))$  as equal to 2 in the case where  $\widehat{P}_n \notin \mathcal{P}_d$ , we can see that in either case, it holds that

$$d_{H}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)) \leq \max \left\{ C_{d}^{2}, \sqrt{8} \right\} \cdot \frac{d_{W}^{1/2}(\widehat{P}_{n}, P)}{\epsilon_{P}^{1/2}}.$$
 (4)

In fact, Lei (2020, Theorem 3.1) shows that the  $\log^2 n$  term may be reduced to  $(\log n) \mathbb{1}_{\{d=1,q=2\}} + (\log n) \mathbb{1}_{\{d=2,q>2\}} + (\log^2 n) \mathbb{1}_{\{d=2,q=2\}} + (\log n) \mathbb{1}_{\{d\geq 3,q=d/(d-1)\}}$ . Since poly-logarithmic factors are not our primary concern in this work, however, we will present simpler bounds based on (3).

Now, taking the expected value and combining the bounds (3) and (4), we obtain

$$\begin{split} & \mathbb{E}\left[\mathrm{d}_{\mathrm{H}}^{2}\left(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)\right)\right] \leq \mathbb{E}\left[\max\left\{C_{d}^{2}, \sqrt{8}\right\} \cdot \frac{\mathrm{d}_{\mathrm{W}}^{1/2}(\widehat{P}_{n}, P)}{\epsilon_{P}^{1/2}}\right] \\ & \leq \max\left\{C_{d}^{2}, \sqrt{8}\right\} \cdot \left[\frac{\mathbb{E}\left[\mathrm{d}_{\mathrm{W}}(\widehat{P}_{n}, P)\right]}{\epsilon_{P}}\right]^{1/2} \leq \max\left\{C_{d}^{2}, \sqrt{8}\right\} \sqrt{\widetilde{C}_{q}} \cdot \sqrt{\frac{M_{q}}{\epsilon_{P}}} \cdot \frac{\log n}{n^{\min\left\{\frac{1}{4}, \frac{1}{2d}, \frac{1}{2} - \frac{1}{2q}\right\}}}. \end{split}$$

Choosing  $C_{d,q} = \max \{C_d^2, \sqrt{8}\} \cdot \sqrt{\tilde{C}_q}$ , this proves the desired result for the case  $d \geq 2$ .

For the case d=1, the result cannot be proved with the same argument, as the exponent on n in the bound above is at best 1/4, which does not lead to the desired scaling if q>2. We establish the desired bound for d=1 in Section 4.4, using a more technical argument.  $\square$ 

We remark that, if X is additionally assumed to be subexponential, then Lei (2020, Corollary 5.2) establishes exponential tail bounds for  $d_W(\widehat{P}_n, P)$ ; under this stronger assumption, the results of Theorem 5 could then be strengthened to give a tail bound for  $d_H^2(\psi^*(\widehat{P}_n), \psi^*(P))$ , in place of the bound on expected value.

#### 2.3.1 Lower bounds for the empirical setting

Our final main result studies the optimality of the power of n appearing in Theorem 5.

**Theorem 6.** For any  $d \ge 1$  and q > 1, there exist  $\epsilon_d^*$ ,  $c_d > 0$ , depending only on d, such that

$$\sup_{P \in \mathcal{P}_d: \mathbb{E}_P[\|X\|^q] \le 1, \, \epsilon_P \ge \epsilon_d^*} \mathbb{E}\left[ d_H^2(\psi^*(\widehat{P}_n), \psi^*(P)) \right] \ge c_d \cdot n^{-\min\left\{\frac{2}{d+1}, \frac{1}{2} - \frac{1}{2q}\right\}}.$$

Ignoring a logarithmic factor in n, the first term, namely  $n^{-\frac{2}{d+1}}$ , is the known minimax rate for any estimator under the well-specified case where P is itself log-concave, for any  $d \geq 2$  (Kim and Samworth, 2016; Kur et al., 2019). The second term is a new result and will be proved via a misspecified construction where P is not log-concave: the distribution is given by  $X = R \cdot U$ , where U is drawn uniformly from the unit sphere  $\mathbb{S}_{d-1}$ , while the radius R is drawn independently with

$$R = \begin{cases} 1/2, & \text{with probability } 1 - 1/2n, \\ n^{1/q}, & \text{with probability } 1/2n. \end{cases}$$

The intuition is that, with positive probability, the empirical distribution  $\widehat{P}_n$  (and, therefore, its log-concave projection  $\psi^*(\widehat{P}_n)$ ), is supported on the ball of radius 1/2; on the other hand, we will see in the proof that  $\psi^*(P)$  places  $\sim n^{-\frac{1}{2} + \frac{1}{2q}}$  mass outside this ball, leading to a lower bound on the Hellinger distance between these two log-concave projections.

A consequence of this last result in dimension d = 1 is that rates of convergence in log-concave density estimation can be much slower in the misspecified setting, with a minimax rate of  $n^{-1/2}$  at best, as compared to the well-specified setting when P is assumed to have a log-concave density, where the corresponding rate is  $n^{-4/5}$  (Kim and Samworth, 2016).

#### **2.3.2** A gap for dimension $d \ge 2$

Comparing the lower bound established in Theorem 6 with the upper bound given in Theorem 5, we see that for the case d=1 the two bounds match, as they both scale as  $n^{-\frac{1}{2}+\frac{1}{2q}}$  (ignoring poly-logarithmic factors). For  $d \geq 2$ , however, there is a gap—for sufficiently large q (i.e., a sufficiently strong moment condition), the upper bound scales as  $n^{-\frac{1}{2d}}$  (up to poly-logarithmic factors) while the lower bound has the faster rate  $n^{-\frac{2}{d+1}}$ . We also remark that the optimal dependence of the minimax rate on d remains unknown as well.

# 3 Relationship with prior work

Log-concave density estimation is a central problem within the field of nonparametric inference under shape constraints. Entry points to the field include the book by Groeneboom and Jongbloed (2014), as well as the 2018 special issue of the journal Statistical Science (Samworth and Sen, 2018). Other important shape-constrained problems that could benefit from the perspective taken in this work include decreasing density estimation (Grenander, 1956; Rao, 1969; Groeneboom, 1985; Birgé, 1989; Jankowski, 2014), isotonic regression (Brunk et al., 1972; Zhang, 2002; Chatterjee et al., 2015; Durot and Lopuhaä, 2018; Bellec, 2018; Yang and Barber, 2019; Han et al., 2019) and convex regression (Hildreth, 1954; Seijo and Sen, 2011; Cai and Low, 2015; Guntuboyina and Sen, 2015; Han and Wellner, 2016b; Fang and Guntuboyina, 2019), among many others. In these cases, the analysis is likely to be more straightforward, since the canonical least squares/maximum likelihood estimator can be characterised as an  $L_2$ -projection onto a convex set. By contrast, the class  $\mathcal{F}_d$  is not convex, and the Kullback–Leibler projection  $\psi^*$  is considerably more involved.

Early work on log-concave density estimation includes Walther (2002), Pal et al. (2007), Dümbgen and Rufibach (2009), Walther (2009), Cule et al. (2010), Cule and Samworth (2010), Schuhmacher et al. (2011), Samworth and Yuan (2012) and Chen and Samworth (2013). Sometimes, the class is considered as a special case of the class of s-concave densities (Koenker and Mizera, 2010; Seregin and Wellner, 2010; Han and Wellner, 2016a; Doss and Wellner, 2016; Han, 2019). For the case of correct model specification, where P has density  $f_P \in \mathcal{F}_d$  and  $\hat{f}_n := \psi^*(\hat{P}_n)$ , it is now known (Kim and Samworth, 2016; Kur et al., 2019) that

$$\sup_{f_P \in \mathcal{F}_d} \mathbb{E}\left[d_{\mathrm{H}}^2(\widehat{f_n}, f_P)\right] \le K_d \cdot \left\{ \begin{array}{ll} n^{-4/5} & \text{when } d = 1 \\ n^{-2/(d+1)} \log n & \text{when } d \ge 2, \end{array} \right.$$

where  $K_d > 0$  depends only on d, and that this risk bound is minimax optimal (up to the logarithmic factor when  $d \geq 2$ ). See also Carpenter et al. (2018) for an earlier result in the case  $d \geq 4$ , and Xu and Samworth (2020) for an alternative approach to high-dimensional log-concave density estimation that seeks to evade the curse of dimensionality in the additional presence of symmetry constraints. It is further known that when  $d \leq 3$ , the log-concave maximum likelihood estimator can adapt to certain subclasses of log-concave densities, including log-concave densities whose logarithms are piecewise affine (Kim et al., 2018; Feng et al., 2020). Although these recent works provide a relatively complete picture of the behaviour of the log-concave maximum likelihood estimator when the true distribution has a log-concave density, there is almost no prior work on risk bounds under model misspecification. The

only exception of which we are aware is Kim et al. (2018, Theorem 1), which considers a univariate case where the true distribution has a density that is very close to log-affine on its support.

One feature that distinguishes our contributions from earlier work on rates of convergence in log-concave density estimation in the correctly specified setting is that our arguments avoid entirely notions of bracketing entropy, as well as empirical process arguments that control the behaviour of M-estimators in terms of the entropy of a relevant function class (e.g. van der Vaart and Wellner, 1996; van de Geer, 2000). It turns out that, for non-convex classes of densities, these ideas are not well suited to the misspecified setting. Instead, our main tool is a detailed and delicate analysis of the Lipschitz approximations to concave functions introduced in Dümbgen et al. (2011). In their original usage, these were employed in conjunction with asymptotic results such as Skorokhod's representation theorem to derive the consistency and robustness results described above. By contrast, our analysis facilitates the direct inequality established in Theorem 2.

Another role of this work is to advocate for the benefits of regarding an estimator as a function of the empirical distribution, as opposed to the more conventional view where it is seen as a function on the sample space. The empirical distribution  $\widehat{P}_n$  of a sample  $X_1, \ldots, X_n$  encodes all of the information in the data when we regard it as a multi-set  $\{X_1, \ldots, X_n\}$ , i.e. when we discard information in the ordering of the indices. It follows that any statistic  $\widehat{\theta}_n = \widehat{\theta}_n(X_1, \ldots, X_n)$  that is invariant to permutation of its arguments can be thought of as a functional  $\theta(\widehat{P}_n)$  of the empirical distribution. Frequently, the definition of  $\theta$  can be extended to a more general class of distributions  $\mathcal{P}$ , and we may regard  $\theta$  as a projection from  $\mathcal{P}$  onto a model, or parameter space,  $\Theta$ . This perspective, which was pioneered by Richard von Mises in the 1940s (von Mises, 1947) and described in Serfling (1980, Chapter 6), offers many advantages to the statistician. In particular, once the analytical properties (e.g. continuity, differentiability) of  $\theta$  are understood, key statistical properties of the estimator (consistency, robustness to misspecification, rates of convergence), can often be deduced as simple corollaries of basic facts about the convergence of empirical distributions.

# 4 Proofs of upper bounds

In this section we prove Theorem 2 (for arbitrary dimension d), and complete the proof of Theorem 5 (for the remaining case of dimension d = 1). In Section 4.1 we review some known properties of log-concave projection, and in Section 4.2 we establish a key lemma that will be used in both proofs. In Section 4.3 we complete the proof of Theorem 2, and in Section 4.4 we complete the proof of Theorem 5 for the remaining case d = 1.

<sup>&</sup>lt;sup>2</sup>See Patilea (2001, Proposition 4.1) for applications of entropy methods to studying rates of convergence of maximum likelihood estimators for convex classes of densities. However, the class of densities f that are log-concave is not a convex class; if we instead consider the class of concave log-densities (i.e., log f, where f is a log-concave density), then this class is also not convex, because of the need for the exponentials of these log-densities to integrate to 1.

## 4.1 Background on log-concave projection

We begin by reviewing some known properties of log-concave projection, and computing some new bounds.

#### 4.1.1 Moment inequalities

The log-concave projection  $\psi^*$  is known to satisfy a useful convex ordering property (Dümbgen et al., 2011, Eqn. (3)): for any  $P \in \mathcal{P}_d$  and for  $f = \psi^*(P)$ ,

$$\mathbb{E}_f[h(X)] \le \mathbb{E}_P[h(X)] \text{ for any convex function } h: \mathbb{R}^d \to (-\infty, \infty].$$
 (5)

In particular, this implies that

$$\mathbb{E}_f\left[|v^\top(X-\mu_P)|\right] \le \mathbb{E}_P\left[|v^\top(X-\mu_P)|\right] \text{ for all } v \in \mathbb{R}^d.$$

The following lemma establishes that, up to a constant, this inequality is tight for all vectors  $v \in \mathbb{R}^d$ .

**Lemma 7.** Fix any  $P \in \mathcal{P}_d$ , and let  $f = \psi^*(P)$ . Then

$$\mathbb{E}_f \left[ |v^\top (X - \mu_P)| \right] \ge c_d \cdot \mathbb{E}_P \left[ |v^\top (X - \mu_P)| \right] \text{ for all } v \in \mathbb{R}^d,$$

where  $c_d \in (0,1]$  depends only on d.

By Dümbgen et al. (2011, Eqn. (4)), log-concave projection preserves the mean, i.e.,

$$\mu_P = \mathbb{E}_P \left[ X \right] = \mathbb{E}_f \left[ X \right].$$

We can also define the covariance matrix  $\Sigma = \text{Cov}_f(X)$ , which is finite (since all moments of a log-concave distribution are finite) and strictly positive definite. Lemma 7 immediately implies bounds on the eigenvalues of  $\Sigma$ :

**Corollary 8.** Fix any  $P \in \mathcal{P}_d$ , let  $f = \psi^*(P)$ , and let  $\Sigma = \operatorname{Cov}_f(X)$  be the covariance matrix of the distribution with density f. Then for all  $v \in \mathbb{R}^d$ ,

$$c_d^2 \{ \mathbb{E}_P [|v^\top (X - \mu_P)|] \}^2 \le v^\top \Sigma v \le 16 \{ \mathbb{E}_P [|v^\top (X - \mu_P)|] \}^2,$$

where  $c_d \in (0,1]$  is taken from Lemma 7. In particular, this implies that

$$\lambda_{\min}(\Sigma) \ge (c_d \epsilon_P)^2$$

where  $\lambda_{\min}(\Sigma)$  denotes the smallest eigenvalue of  $\Sigma$ .

Proof of Corollary 8. First, for the lower bound, by Lemma 7 and Cauchy–Schwarz,

$$c_d^2 \{ \mathbb{E}_P \left[ |v^{\top}(X - \mu_P)| \right] \}^2 \le \{ \mathbb{E}_f \left[ |v^{\top}(X - \mu_P)| \right] \}^2 \le \mathbb{E}_f \left[ |v^{\top}(X - \mu_P)|^2 \right] = v^{\top} \Sigma v.$$

Next, for the upper bound,

$$v^{\top} \Sigma v = \mathbb{E}_f \left[ |v^{\top} (X - \mu_P)|^2 \right] \le 16 \left\{ \mathbb{E}_f \left[ |v^{\top} (X - \mu_P)| \right] \right\}^2 \le 16 \left\{ \mathbb{E}_P \left[ |v^{\top} (X - \mu_P)| \right] \right\}^2$$

where the first inequality is due to Lovász and Vempala (2007, Theorem 5.22) while the second is by (5) (Dümbgen et al., 2011, Eqn. (3)).

#### 4.1.2 A lower bound on a ball

Next we show that for any P, its log-concave projection  $f = \psi^*(P)$  is lower bounded on a ball of radius of order  $\epsilon_P$ .

**Lemma 9.** Fix any  $P \in \mathcal{P}_d$ , and let  $f = \psi^*(P)$ . Then there exist  $b_d, r_d \in (0, 1]$ , depending only on d, such that

$$f(x) \ge b_d \cdot \sup_{x' \in \mathbb{R}^d} f(x')$$
 for all  $x \in \mathbb{B}_d(\mu_P, r_d \epsilon_P)$ .

Proof of Lemma 9. Let  $\Sigma = \text{Cov}_f(X)$ , and define the isotropic, log-concave density  $g(x) = f(\Sigma^{1/2}x + \mu_P) \det^{1/2}(\Sigma)$ . By Lovász and Vempala (2007, Theorem 5.14(a) and (b)),

$$\inf_{x:\|x\| \le 1/9} g(x) \ge b_d \sup_{x \in \mathbb{R}^d} g(x),$$

where  $b_d \in (0,1]$  depends only on d. This immediately implies that

$$f(x) \ge b_d \sup_{x' \in \mathbb{R}^d} f(x')$$
 for all  $x \in \mathbb{R}^d$  with  $\|\Sigma^{-1/2}(x - \mu_P)\| \le 1/9$ .

But  $\|\Sigma^{-1/2}(x-\mu_P)\| \leq \lambda_{\min}^{-1/2}(\Sigma)\|x-\mu_P\| \leq \|x-\mu_P\|/(c_d\epsilon_P)$  by Corollary 8, so the result holds with  $r_d = c_d/9$ .

## 4.2 Key lemma: the Lipschitz majorization

Let

$$\Phi_d := \bigg\{ \phi : \mathbb{R}^d \to [-\infty, \infty) \ : \quad \begin{array}{c} \phi \text{ is a proper concave, upper semi-continuous function,} \\ \text{and } \phi(x) \to -\infty \text{ as } \|x\| \to \infty \end{array} \bigg\},$$

and define the function  $\phi^*: \mathcal{P}_d \to \Phi_d$  that maps a distribution P to the log-density  $\phi = \phi^*(P)$  given by  $\phi(x) = \log \left[ \psi^*(P) \right](x)$ . Dümbgen et al. (2011, Theorem 2.2) establishes that the log-density  $\phi = \phi^*(P)$  maximizes  $\ell(\phi, P) := \mathbb{E}_P \left[ \phi(X) \right] - \int_{\mathbb{R}^d} e^{\phi(x)} \, dx + 1$  over  $\Phi_d$ . We now show that this maximum can be nearly attained by a Lipschitz function. In particular, for any  $\phi \in \Phi_d$  and any L > 0, define its L-Lipschitz majorization  $\phi^L : \mathbb{R}^d \to \mathbb{R}$  by

$$\phi^{L}(x) := \sup_{y \in \mathbb{R}^{d}} \{ \phi(y) - L \|x - y\| \}.$$
 (6)

It can easily be verified that this function is concave, L-Lipschitz, and satisfies  $\phi^L(x) \geq \phi(x)$  for all  $x \in \mathbb{R}^d$ . Furthermore, it holds that  $\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x < \infty$  (this follows from the fact that there exist constants  $a \in \mathbb{R}$ , b > 0 such that  $\phi(y) \leq a - b \|y\|$  for all  $y \in \mathbb{R}^d$  (Dümbgen et al., 2011)), and moreover  $\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x > 0$ .

Next we normalize to produce a log-density. For any  $\phi \in \Phi_d$ , we define

$$\tilde{\phi}^L(x) := \phi^L(x) - \log\left(\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x\right). \tag{7}$$

The following result proves that, if  $\phi = \phi^*(P)$ , then for L sufficiently large,  $\tilde{\phi}^L \in \Phi_d$  is nearly optimal for P (in the sense of maximizing  $\ell(\cdot, P)$ ).

**Lemma 10.** Fix any  $P \in \mathcal{P}_d$ , let  $\phi = \phi^*(P)$ , and let  $\phi^L$  and  $\tilde{\phi}^L$  be defined as in (6) and (7). Then for any  $L \geq \frac{2d}{r_d \epsilon_P}$ ,

$$\ell(\tilde{\phi}^L, P) \ge \ell(\phi^L, P) \ge \ell(\phi, P) - \frac{4d}{Lb_d r_d \epsilon_P}$$

where  $r_d, b_d \in (0,1]$  are taken from Lemma 9. In particular, this implies that

$$\mathbb{E}_P\left[\tilde{\phi}^L(X)\right] \ge \mathbb{E}_P\left[\phi(X)\right] - \frac{4d}{Lb_d r_d \epsilon_P}.$$

#### 4.2.1 Bounding the Hellinger distance

Now we apply Lemma 10 to the problem of bounding Hellinger distance.

Corollary 11. Fix any  $P, Q \in \mathcal{P}_d$ , and define  $\epsilon = \min\{\epsilon_P, \epsilon_Q\} > 0$ . Let  $\phi_P = \phi^*(P)$  and  $\phi_Q = \phi^*(Q)$ , and let  $f_P = \psi^*(P)$  and  $f_Q = \psi^*(Q)$  be the corresponding density functions. Let  $\phi_P^L$  and  $\phi_Q^L$  be the L-Lipschitz majorizations of  $\phi_P$  and  $\phi_Q$ , respectively, as defined in (6), for some  $L \geq \frac{2d}{r_d \epsilon}$ , where  $r_d \in (0,1]$  is taken from Lemma 9. Then

$$d_{\mathrm{H}}^{2}(f_{P}, f_{Q}) \leq \frac{16d}{Lb_{d}r_{d}\epsilon} + \left(\mathbb{E}_{P}\left[\phi_{P}^{L}(X)\right] - \mathbb{E}_{Q}\left[\phi_{P}^{L}(X)\right]\right) + \left(\mathbb{E}_{Q}\left[\phi_{Q}^{L}(X)\right] - \mathbb{E}_{P}\left[\phi_{Q}^{L}(X)\right]\right),$$

where  $b_d \in (0,1]$  is taken from Lemma 9.

Proof of Corollary 11. Let  $\tilde{\phi}_P^L$ ,  $\tilde{\phi}_Q^L$  be defined as in (7), and let  $\tilde{f}_P^L$ ,  $\tilde{f}_Q^L$  be the corresponding densities, i.e.,  $\tilde{f}_P^L(x) = e^{\tilde{\phi}_P^L(x)}$  and similarly for  $\tilde{f}_Q^L$ . We first calculate

$$d_{\mathrm{KL}}(f_P||\tilde{f}_P^L) = \mathbb{E}_{f_P}\left[\phi_P(X) - \tilde{\phi}_P^L(X)\right] \le \mathbb{E}_P\left[\phi_P(X) - \tilde{\phi}_P^L(X)\right]$$

and

$$d_{\mathrm{KL}}(f_P||\tilde{f}_Q^L) = \mathbb{E}_{f_P} \left[ \phi_P(X) - \tilde{\phi}_Q^L(X) \right] \le \mathbb{E}_P \left[ \phi_P(X) - \tilde{\phi}_Q^L(X) \right],$$

where the inequalities hold by Dümbgen et al. (2011, Remark 2.3). The same bounds hold with the roles of P and Q reversed. Furthermore, by the triangle inequality,

$$\begin{split} \mathbf{d}_{\mathrm{H}}^{2}(f_{P}, f_{Q}) &= \frac{1}{2} \mathbf{d}_{\mathrm{H}}^{2}(f_{P}, f_{Q}) + \frac{1}{2} \mathbf{d}_{\mathrm{H}}^{2}(f_{P}, f_{Q}) \\ &\leq \frac{1}{2} \left\{ \mathbf{d}_{\mathrm{H}}(f_{P}, \tilde{f}_{P}^{L}) + \mathbf{d}_{\mathrm{H}}(f_{Q}, \tilde{f}_{P}^{L}) \right\}^{2} + \frac{1}{2} \left\{ \mathbf{d}_{\mathrm{H}}(f_{P}, \tilde{f}_{Q}^{L}) + \mathbf{d}_{\mathrm{H}}(f_{Q}, \tilde{f}_{Q}^{L}) \right\}^{2} \\ &\leq \mathbf{d}_{\mathrm{H}}^{2}(f_{P}, \tilde{f}_{P}^{L}) + \mathbf{d}_{\mathrm{H}}^{2}(f_{Q}, \tilde{f}_{P}^{L}) + \mathbf{d}_{\mathrm{H}}^{2}(f_{P}, \tilde{f}_{Q}^{L}) + \mathbf{d}_{\mathrm{H}}^{2}(f_{Q}, \tilde{f}_{Q}^{L}) \\ &\leq \mathbf{d}_{\mathrm{KL}}(f_{P}||\tilde{f}_{P}^{L}) + \mathbf{d}_{\mathrm{KL}}(f_{Q}||\tilde{f}_{P}^{L}) + \mathbf{d}_{\mathrm{KL}}(f_{P}||\tilde{f}_{Q}^{L}) + \mathbf{d}_{\mathrm{KL}}(f_{Q}||\tilde{f}_{Q}^{L}), \end{split}$$

where the last step holds by the standard inequality relating KL divergence with Hellinger distance (i.e.,  $d_H^2 \leq d_{KL}$ ). Combining all these calculations, and then rearranging terms, we

see that<sup>3</sup>

$$\begin{aligned} \mathrm{d}_{\mathrm{H}}^{2}(f_{P},f_{Q}) &\leq \mathbb{E}_{P} \left[ \phi_{P}(X) - \tilde{\phi}_{P}^{L}(X) \right] + \mathbb{E}_{Q} \left[ \phi_{Q}(X) - \tilde{\phi}_{P}^{L}(X) \right] \\ &+ \mathbb{E}_{P} \left[ \phi_{P}(X) - \tilde{\phi}_{Q}^{L}(X) \right] + \mathbb{E}_{Q} \left[ \phi_{Q}(X) - \tilde{\phi}_{Q}^{L}(X) \right] \\ &= 2 \left( \mathbb{E}_{P} \left[ \phi_{P}(X) - \tilde{\phi}_{P}^{L}(X) \right] + \mathbb{E}_{Q} \left[ \phi_{Q}(X) - \tilde{\phi}_{Q}^{L}(X) \right] \right) \\ &+ \left( \mathbb{E}_{P} \left[ \tilde{\phi}_{P}^{L}(X) \right] - \mathbb{E}_{Q} \left[ \tilde{\phi}_{P}^{L}(X) \right] \right) + \left( \mathbb{E}_{Q} \left[ \tilde{\phi}_{Q}^{L}(X) \right] - \mathbb{E}_{P} \left[ \tilde{\phi}_{Q}^{L}(X) \right] \right) \\ &= 2 \left( \mathbb{E}_{P} \left[ \phi_{P}(X) - \tilde{\phi}_{P}^{L}(X) \right] + \mathbb{E}_{Q} \left[ \phi_{Q}(X) - \tilde{\phi}_{Q}^{L}(X) \right] \right) \\ &+ \left( \mathbb{E}_{P} \left[ \phi_{P}^{L}(X) \right] - \mathbb{E}_{Q} \left[ \phi_{P}^{L}(X) \right] \right) + \left( \mathbb{E}_{Q} \left[ \phi_{Q}^{L}(X) \right] - \mathbb{E}_{P} \left[ \phi_{Q}^{L}(X) \right] \right), \end{aligned}$$

where the last step holds since  $\tilde{\phi}_P^L$ ,  $\tilde{\phi}_Q^L$  are simply shifts of the functions  $\phi_P^L$ ,  $\phi_Q^L$ , respectively. Finally, applying Lemma 10 concludes the proof.

## 4.3 Completing the proof of Theorem 2

We will now apply Corollary 11 to prove Theorem 2, bounding  $d_H^2(f_P, f_Q)$  in terms of the Wasserstein distance. Define

$$L = \sqrt{\frac{8d}{r_d b_d \min\{\epsilon_P, \epsilon_Q\} d_W(P, Q)}},$$

where  $r_d, b_d \in (0, 1]$  are taken from Lemma 9. Take a coupling (X, Y) of d-dimensional random vectors with marginal distributions  $X \sim P$  and  $Y \sim Q$ , such that  $\mathbb{E}[\|X - Y\|] = d_W(P, Q)$ , which is guaranteed to exist by Villani (2008, Theorem 4.1). Then, since  $\phi_P^L$  is L-Lipschitz, we have

$$\mathbb{E}\left[\phi_P^L(X)\right] - \mathbb{E}\left[\phi_P^L(Y)\right] \le \mathbb{E}\left[L \|X - Y\|\right] = Ld_{\mathbf{W}}(P, Q),$$

and similarly

$$\mathbb{E}\left[\phi_Q^L(Y)\right] - \mathbb{E}\left[\phi_Q^L(X)\right] \le L\mathrm{d}_{\mathrm{W}}(P,Q).$$

If  $L \geq \frac{2d}{r_d \min\{\epsilon_P, \epsilon_Q\}}$ , then applying Corollary 11, we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq \frac{16d}{Lb_{d}r_{d}\min\{\epsilon_{P}, \epsilon_{Q}\}} + 2Ld_{\mathrm{W}}(P, Q) = \sqrt{\frac{128dd_{\mathrm{W}}(P, Q)}{r_{d}b_{d}\min\{\epsilon_{P}, \epsilon_{Q}\}}}.$$

If instead  $L < \frac{2d}{r_d \min\{\epsilon_P, \epsilon_Q\}}$ , then  $\frac{db_d d_W(P,Q)}{2r_d \min\{\epsilon_P, \epsilon_Q\}} > 1$ . Since Hellinger distance is always bounded by  $\sqrt{2}$ , we then have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq 2 \leq \sqrt{\frac{2db_{d}d_{\mathrm{W}}(P, Q)}{r_{d}\min\{\epsilon_{P}, \epsilon_{Q}\}}} \leq \sqrt{\frac{2dd_{\mathrm{W}}(P, Q)}{r_{d}b_{d}\min\{\epsilon_{P}, \epsilon_{Q}\}}},$$

<sup>&</sup>lt;sup>3</sup>All expectations in this display are finite, because, e.g.,  $\sup_{x \in \mathbb{R}^d} \phi_P(x) = \sup_{x \in \mathbb{R}^d} \phi_P^L(x) < \infty$ ; moreover,  $\mathbb{E}_P\left[\phi_P^L(X)\right] \geq \mathbb{E}_P\left[\phi_P(X)\right] > -\infty$  because  $P \in \mathcal{P}_d$ , and  $\mathbb{E}_P\left[\phi_Q^L(X)\right] > -\infty$  because  $\phi_Q^L$  is Lipschitz and P has a finite first moment.

where the last step holds trivially since  $b_d \leq 1$ . Thus, in either case, we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq \sqrt{\frac{128d}{r_{d}b_{d}}} \cdot \sqrt{\frac{d_{\mathrm{W}}(P, Q)}{\min\{\epsilon_{P}, \epsilon_{Q}\}}}.$$

We now split into cases. If  $d_W(P,Q) \leq \max\{\epsilon_P, \epsilon_Q\}/4$ , then

$$\frac{\mathrm{d_W}(P,Q)}{\min\{\epsilon_P,\epsilon_Q\}} = \frac{\mathrm{d_W}(P,Q)}{\max\{\epsilon_P,\epsilon_Q\} - |\epsilon_P - \epsilon_Q|} \leq \frac{\mathrm{d_W}(P,Q)}{\max\{\epsilon_P,\epsilon_Q\} - 2\mathrm{d_W}(P,Q)} \leq \frac{2\mathrm{d_W}(P,Q)}{\max\{\epsilon_P,\epsilon_Q\}},$$

where the second step applies Proposition 1. If instead  $d_W(P,Q) > \max\{\epsilon_P, \epsilon_Q\}/4$  then we will instead use the trivial bound

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq 2 \leq 4\sqrt{\frac{d_{\mathrm{W}}(P, Q)}{\max\{\epsilon_{P}, \epsilon_{Q}\}}} \leq 4\sqrt{\frac{d}{r_{d}b_{d}}} \cdot \sqrt{\frac{d_{\mathrm{W}}(P, Q)}{\max\{\epsilon_{P}, \epsilon_{Q}\}}}$$

where the last step is trivial since  $d \ge 1$  and  $r_d, b_d \in (0, 1]$ . Thus, in both cases, we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq 16\sqrt{\frac{d}{r_{d}b_{d}}} \cdot \sqrt{\frac{d_{\mathrm{W}}(P, Q)}{\max\{\epsilon_{P}, \epsilon_{Q}\}}}.$$

This proves the theorem, when we choose  $C_d = 4\left(\frac{d}{r_d b_d}\right)^{1/4}$ .

## 4.4 Completing the proof of Theorem 5: the case d=1

Before proving the theorem, we first state several supporting lemmas. First we state a deterministic result:

**Lemma 12.** Let  $P, Q \in \mathcal{P}_1$  satisfy  $\max\{\mathbb{E}_P[|X|^q]^{1/q}, \mathbb{E}_Q[|X|^q]^{1/q}\} \leq M_q$  for some q > 1. Define

$$\Delta_{\mathrm{CDF}}(P,Q)$$

$$:= \max \left\{ \sup_{t \in \mathbb{R}} \left| \sqrt{\mathbb{P}_P\{X > t\}} - \sqrt{\mathbb{P}_Q\{X > t\}} \right|, \sup_{t \in \mathbb{R}} \left| \sqrt{\mathbb{P}_P\{X < t\}} - \sqrt{\mathbb{P}_Q\{X < t\}} \right| \right\}.$$

Then

$$\mathrm{d}_{\mathrm{H}}^{2}\big(\psi^{*}(P),\psi^{*}(Q)\big) \leq C_{*}\sqrt{\frac{M_{q}}{\max\{\epsilon_{P},\epsilon_{Q}\}}} \cdot \Big\{\Delta_{\mathrm{CDF}}(P,Q) \cdot \log\big(e/\Delta_{\mathrm{CDF}}(P,Q)\big)\Big\}^{1-1/q},$$

for a universal constant  $C_* > 0$ .

Next, in order to prove Theorem 5, we will want to apply this result with  $Q = \widehat{P}_n$ , i.e., we want to bound  $\Delta_{\text{CDF}}(\widehat{P}_n, P)$ . Let F denote the distribution function of P, and, for  $t \in (0,1)$ , let  $F^{-1}(t) := \inf\{x : F(x) \ge t\}$ . Then, with  $U \sim \text{Unif}[0,1]$ , we know that

 $F^{-1}(U) \sim P$ . We may therefore assume that  $X_1, \ldots, X_n$  are generated as  $X_i = F^{-1}(U_i)$ , where  $U_1, \ldots, U_n \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$ . Since  $F^{-1}$  is monotonic, we have

$$\Delta_{\text{CDF}}(\widehat{P}_n, P) \le \Delta_{\text{CDF}}(\widehat{U}_n, \text{Unif}[0, 1]),$$
(8)

where  $\widehat{U}_n$  is the empirical distribution of  $U_1, \ldots, U_n$ . Therefore, it suffices to consider the case that P is the uniform distribution. We now apply results from Shorack and Wellner (2009) to prove a tail bound on  $\Delta_{\text{CDF}}(\widehat{U}_n, \text{Unif}[0, 1])$ .

**Lemma 13.** Fix any  $n \geq 2$ , and let  $\widehat{U}_n$  be the empirical distribution of  $U_1, \ldots, U_n \stackrel{\text{iid}}{\sim} \text{Unif}[0,1]$ . Then, for any c > 0,

$$\mathbb{P}\left\{\Delta_{\mathrm{CDF}}(\widehat{U}_n, \mathrm{Unif}[0, 1]) \le c' \sqrt{\frac{\log n}{n}}\right\} \ge 1 - n^{-c},$$

where c' > 0 depends only on c.

With these lemmas in place, we are now in a position to prove Theorem 5. Let  $M_{q,n} = \left(\frac{1}{n}\sum_{i=1}^{n}|X_i|^q\right)^{1/q}$  and  $\Delta = \Delta_{\text{CDF}}(\widehat{P}_n, P)$ . If  $\widehat{P}_n \in \mathcal{P}_1$  (that is,  $\widehat{P}_n$  does not place all its mass on a single point), then we have

$$d_{\mathrm{H}}^{2}\left(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)\right) \leq \min\left\{2, C_{*}\sqrt{\frac{\max\{M_{q}, M_{q, n}\}}{\max\{\epsilon_{P}, \epsilon_{\widehat{P}_{n}}\}}} \cdot \left(\Delta \log(e/\Delta)\right)^{1-1/q}\right\}$$
(9)

by applying Lemma 12 with  $Q = \widehat{P}_n$ . On the other hand, if  $\widehat{P}_n$  does place all its mass on one point, then recall that  $\psi^*(\widehat{P}_n)$  is not defined but we take  $d^2_H(\psi^*(\widehat{P}_n), \psi^*(P)) = 2$  by convention. For this case, we can trivially calculate

$$\Delta \ge \min \{ \sqrt{\mathbb{P}_P \{ X > \mu_P \}}, \sqrt{\mathbb{P}_P \{ X < \mu_P \}} \}.$$

We will now need an additional lemma:

**Lemma 14.** Fix any  $P \in \mathcal{P}_1$  and any q > 1. Suppose  $M_q = \mathbb{E}_P[|X|^q]^{1/q} < \infty$ . Then

$$\min \left\{ \mathbb{P}_P \{ X > \mu_P \}, \mathbb{P}_P \{ X < \mu_P \} \right\} \ge \left( \frac{\epsilon_P}{4M_q} \right)^{\frac{q}{q-1}}.$$

This implies

$$\Delta \ge \left(\frac{\epsilon_P}{4M_q}\right)^{\frac{q}{2(q-1)}}$$

for the case where  $\widehat{P}_n \notin \mathcal{P}_1$  (i.e.,  $\widehat{P}_n$  is supported on a single point). Since also  $\Delta \leq 1$  by definition, this means that

$$\sqrt{\frac{\max\{M_q, M_{q,n}\}}{\epsilon_P}} \cdot \left(\Delta \log(e/\Delta)\right)^{1-1/q} \ge \frac{1}{2} = \frac{\mathrm{d}_{\mathrm{H}}^2\left(\psi^*(\widehat{P}_n), \psi^*(P)\right)}{4}.$$

Combining this with (9) for the case  $\widehat{P}_n \in \mathcal{P}_1$ , we see that

$$d_{\mathrm{H}}^{2}\left(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)\right) \leq \min\left\{2, \max\{C_{*}, 4\}\sqrt{\frac{\max\{M_{q}, M_{q, n}\}}{\epsilon_{P}}} \cdot \left(\Delta \log(e/\Delta)\right)^{1-1/q}\right\}$$

holds for both cases.

Next, we will combine this calculation with Lemma 13, applied with c=1/2. Let c' be the constant from Lemma 13. First, if  $c'\sqrt{\frac{\log n}{n}} > 1$ , then

$$\mathbb{E}\left[d_{H}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P))\right] \leq 2 \leq 2\left(c'\sqrt{\frac{\log n}{n}}\right)^{1-1/q} \leq \frac{2c'^{1-1/q}}{(\log 2)^{1-1/q}} \frac{\log^{\frac{3}{2}(1-1/q)} n}{n^{\frac{1}{2}-\frac{1}{2q}}} \\ \leq \frac{2c'^{1-1/q}}{(\log 2)^{1-1/q}} \cdot \sqrt{\frac{2M_{q}}{\epsilon}} \cdot \frac{\log^{\frac{3}{2}(1-1/q)} n}{n^{\frac{1}{2}-\frac{1}{2q}}},$$

where the last step holds since

$$\epsilon_P = \mathbb{E}_P[|X - \mu_P|] \le \mathbb{E}_P[|X|] + |\mu_P| \le 2\mathbb{E}_P[|X|] \le 2\{\mathbb{E}_P[|X|^q]\}^{1/q} \le 2M_q.$$
 (10)

If instead  $c'\sqrt{\frac{\log n}{n}} \le 1$ , then we have

$$\begin{split} &\mathbb{E}\left[\mathrm{d}_{\mathrm{H}}^{2}\left(\psi^{*}(\widehat{P}_{n}),\psi^{*}(P)\right)\right] \\ &\leq \mathbb{E}\left[\min\left\{2,\max\{C_{*},4\}\sqrt{\frac{\max\{M_{q},M_{q,n}\}}{\max\{\epsilon_{P},\epsilon_{\widehat{P}_{n}}\}}}\cdot\left(\Delta\log(e/\Delta)\right)^{1-1/q}\right\}\right] \\ &\leq 2\mathbb{P}\left\{\Delta > c'\sqrt{\frac{\log n}{n}}\right\} + \mathbb{E}\left[\max\{C_{*},4\}\sqrt{\frac{M_{q}+M_{q,n}}{\epsilon_{P}}}\cdot\left\{c'\sqrt{\frac{\log n}{n}}\log\left(\frac{e}{c'\sqrt{\frac{\log n}{n}}}\right)\right\}^{1-1/q}\right] \\ &\leq 2n^{-1/2} + \max\{C_{*},4\}\sqrt{\frac{M_{q}+\mathbb{E}\left[M_{q,n}\right]}{\epsilon_{P}}}\cdot\left\{c'\sqrt{\frac{\log n}{n}}\log\left(\frac{e}{c'\sqrt{\frac{\log n}{n}}}\right)\right\}^{1-1/q} \\ &\leq 2n^{-1/2} + \max\{C_{*},4\}\sqrt{\frac{2M_{q}}{\epsilon_{P}}}\cdot\left\{c'\sqrt{\frac{\log n}{n}}\log\left(\frac{e}{c'\sqrt{\frac{\log n}{n}}}\right)\right\}^{1-1/q} \\ &\leq \sqrt{\frac{2M_{q}}{\epsilon_{P}}}\cdot\left[2n^{-1/2} + \max\{C_{*},4\}\left\{c'\sqrt{\frac{\log n}{n}}\log\left(\frac{e}{c'\sqrt{\frac{\log n}{n}}}\right)\right\}^{1-1/q}\right], \end{split}$$

where the third-to-last step applies Jensen's inequality, the second-to-last step holds because  $\mathbb{E}[M_{q,n}] \leq M_q$ , and the last step holds by (10). After simplifying, we obtain

$$\mathbb{E}\left[d_{H}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P))\right] \leq C_{1,q} \sqrt{\frac{M_{q}}{\epsilon_{P}}} \cdot \frac{\log^{\frac{3}{2}(1-1/q)} n}{n^{\frac{1}{2}-\frac{1}{2q}}}$$

for all  $n \geq 2$  when  $C_{1,q}$  is chosen appropriately. This completes the proof of Theorem 5 for the case d = 1.

# A Additional proofs

## A.1 Proof of Proposition 1

First fix any distribution P on  $\mathbb{R}^d$  with  $\mathbb{E}_P[\|X\|] < \infty$ . Observe that  $u \mapsto \mathbb{E}_P[\|u^\top (X - \mu_P)\|]$  is a continuous function on  $\mathbb{S}_{d-1}$ , since for any  $u, v \in \mathbb{S}_{d-1}$ , we have

$$\begin{aligned} \left| \mathbb{E}_{P} \left[ |u^{\top} (X - \mu_{P})| \right] - \mathbb{E}_{P} \left[ |v^{\top} (X - \mu_{P})| \right] \right| &\leq \mathbb{E}_{P} \left[ \left| (u - v)^{\top} (X - \mu_{P}) \right| \right] \\ &\leq \|u - v\| \cdot \mathbb{E}_{P} \left[ \|X - \mu_{P}\| \right] \leq \|u - v\| \cdot 2\mathbb{E}_{P} \left[ \|X\| \right], \end{aligned}$$

and  $\mathbb{E}_P[||X||] < \infty$  by assumption. Therefore,  $u \mapsto \mathbb{E}_P[|u^\top(X - \mu_P)|]$  must attain its infimum, that is,

$$\epsilon_P = \inf_{u \in \mathbb{S}_{d-1}} \mathbb{E}_P \left[ \left| u^\top (X - \mu_P) \right| \right] = \mathbb{E}_P \left[ \left| u_0^\top (X - \mu_P) \right| \right]$$

for some  $u_0 \in \mathbb{S}_{d-1}$ .

Next suppose  $P \in \mathcal{P}_d$ . We will show that  $\epsilon_P > 0$ . As above, we have  $\epsilon_P = \mathbb{E}_P \left[ \left| u_0^\top (X - \mu_P) \right| \right]$  for some  $u_0 \in \mathbb{S}_{d-1}$ . If  $\epsilon_P = 0$ , then this implies that  $u_0^\top (X - \mu_P) = 0$  with probability 1, meaning that P places all its mass on a single hyperplane  $H = \{x \in \mathbb{R}^d : u_0^\top x = u_0^\top \mu_P\}$ . This contradicts the assumption  $P \in \mathcal{P}_d$ , thus proving the first claim.

Finally, consider distributions P, Q on  $\mathbb{R}^d$  with  $\mathbb{E}_P[\|X\|]$ ,  $\mathbb{E}_Q[\|X\|] < \infty$ . By Villani (2008, Theorem 4.1), we can find a pair of d-dimensional random vectors X and Y such that marginally  $X \sim P$ ,  $Y \sim Q$  and  $\mathbb{E}[\|X - Y\|] = \mathrm{d}_W(P, Q)$ . Let  $u_0$  be defined as above, so that  $\epsilon_P = \mathbb{E}[|u_0^\top (X - \mu_P)|]$ . Then

$$\begin{split} \epsilon_{Q} - \epsilon_{P} &= \inf_{u \in \mathbb{S}_{d-1}} \mathbb{E} \left[ \left| u^{\top}(Y - \mu_{Q}) \right| \right] - \mathbb{E} \left[ \left| u_{0}^{\top}(X - \mu_{P}) \right| \right] \\ &\leq \mathbb{E} \left[ \left| u_{0}^{\top}(Y - \mu_{Q}) \right| \right] - \mathbb{E} \left[ \left| u_{0}^{\top}(X - \mu_{P}) \right| \right] \\ &\leq \mathbb{E} \left[ \left| u_{0}^{\top}(X - Y) \right| \right] + \left| u_{0}^{\top}(\mu_{P} - \mu_{Q}) \right| \\ &\leq \mathbb{E} \left[ \|X - Y\| \right] + \|\mu_{P} - \mu_{Q}\| \\ &\leq 2\mathbb{E} \left[ \|X - Y\| \right] \\ &= 2d_{W}(P, Q). \end{split}$$

An identical argument proves the reverse bound, and we deduce that  $|\epsilon_P - \epsilon_Q| \leq 2d_W(P, Q)$ , as desired.

#### A.2 Proof of Lemma 7

Let  $\Sigma = \text{Cov}_f(X)$  and define an isotropic log-concave density g on  $\mathbb{R}^d$  by  $g(x) = f(\Sigma^{1/2}x + \mu_P) \det^{1/2}(\Sigma)$ . Note that, if  $X \sim f$ , then  $\Sigma^{-1/2}(X - \mu_P) \sim g$ . Hence

$$\begin{split} \mathbb{E}_f \left[ |v^{\top} (X - \mu_P)| \right] &= \mathbb{E}_f \left[ |(\Sigma^{1/2} v)^{\top} (\Sigma^{-1/2} (X - \mu_P))| \right] = \mathbb{E}_g \left[ |(\Sigma^{1/2} v)^{\top} X| \right] \\ &\geq \frac{1}{4} \left( \mathbb{E}_g \left[ ((\Sigma^{1/2} v)^{\top} X)^2 \right] \right)^{1/2} = \frac{1}{4} \left\| \Sigma^{1/2} v \right\|, \end{split}$$

where the inequality applies Lovász and Vempala (2007, Theorem 5.22), and the last step holds because q is isotropic.

Next, define a distribution Q obtained by drawing  $X \sim P$  and then taking the affine transformation  $\Sigma^{-1/2}(X - \mu_P)$ . By definition of Q, we have

$$\mathbb{E}_{P}\left[|v^{\top}(X - \mu_{P})|\right] = \mathbb{E}_{P}\left[|(\Sigma^{1/2}v)^{\top}(\Sigma^{-1/2}(X - \mu_{P}))|\right] = \mathbb{E}_{Q}\left[|(\Sigma^{1/2}v)^{\top}X|\right] \leq \|\Sigma^{1/2}v\| \cdot \mathbb{E}_{Q}\left[\|X\|\right].$$

Since log-concave projection commutes with affine transformations, we have

$$\psi^*(Q) = g,$$

which is an isotropic log-concave density. Lemma 15 below establishes that  $\mathbb{E}_Q[||X||] \leq a_d$ , where  $a_d > 0$  depends only on d. Therefore, we have proved that, for any  $v \in \mathbb{R}^d$ ,

$$\mathbb{E}_P\left[|v^{\top}(X-\mu_P)|\right] \le \left\|\Sigma^{1/2}v\right\| \cdot a_d$$

while

$$\mathbb{E}_f\left[|v^{\top}(X-\mu_P)|\right] \ge \frac{1}{4} \left\|\Sigma^{1/2}v\right\|.$$

Setting  $c_d = \frac{1}{4a_d}$  establishes the desired result.

#### A.2.1 Supporting lemma for Lemma 7

**Lemma 15.** There exists  $a_d > 0$ , depending only on d, such that, for any isotropic log-concave density f on  $\mathbb{R}^d$  and any  $P \in \mathcal{P}_d$  with  $\psi^*(P) = f$ ,

$$\mathbb{E}_P[||X||] \le a_d.$$

*Proof of Lemma 15.* By Fresen (2013, Lemma 13), since f is an isotropic log-concave density, it holds that

$$f(x) \le e^{\beta_d - \alpha_d ||x||}$$
 for all  $x \in \mathbb{R}^d$ ,

where  $\alpha_d > 0$  and  $\beta_d \in \mathbb{R}$  depend only on d. We can therefore calculate

$$\mathbb{E}_{P}\left[\log f(X)\right] \leq \mathbb{E}_{P}\left[\beta_{d} - \alpha_{d} \|X\|\right] = \beta_{d} - \alpha_{d} \mathbb{E}_{P}\left[\|X\|\right].$$

On the other hand, consider the log-concave density

$$g(x) = \left(\frac{d^d}{\mathbb{E}_P[\|X\|]^d (d-1)! S_{d-1}}\right) \cdot \exp\left\{-\frac{d \|x\|}{\mathbb{E}_P[\|X\|]}\right\},\,$$

where  $S_{d-1}$  denotes the surface area of the unit sphere  $\mathbb{S}_{d-1}$  in  $\mathbb{R}^d$  (with  $S_0 = 2$ ). We have

$$\mathbb{E}_{P}\left[\log g(X)\right] = \log\left(\frac{d^{d}}{\mathbb{E}_{P}\left[\left\|X\right\|\right]^{d} (d-1)! S_{d-1}}\right) - d.$$

But, since  $f = \psi^*(P)$ , it must hold that

$$\mathbb{E}_P \left[ \log f(X) \right] \ge \mathbb{E}_P \left[ \log g(X) \right],$$

and so

$$\beta_d - \alpha_d \mathbb{E}_P[\|X\|] \ge \log\left(\frac{(d/e)^d}{(d-1)!S_{d-1}}\right) - d\log \mathbb{E}_P[\|X\|].$$

The result follows.

## A.3 Proof of Lemma 10

We will prove below that, when  $L \ge \frac{2d}{r_d \epsilon_P}$ , the function  $\phi^L(x) = \sup_{y \in \mathbb{R}^d} \{\phi(x) - L \|x - y\|\}$  satisfies

$$\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x \le 1 + \frac{4d}{Lb_d r_d \epsilon_P}. \tag{11}$$

Assuming this holds, we then have

$$\begin{split} \ell(\phi^L,P) &= \mathbb{E}_P\left[\phi^L(X)\right] - \int_{\mathbb{R}^d} e^{\phi^L(x)} \; \mathrm{d}x + 1 \geq \mathbb{E}_P\left[\phi^L(X)\right] - \frac{4d}{Lb_d r_d \epsilon_P} \\ &\geq \mathbb{E}_P\left[\phi(X)\right] - \frac{4d}{Lb_d r_d \epsilon_P} = \ell(\phi,P) - \frac{4d}{Lb_d r_d \epsilon_P}, \end{split}$$

where the last inequality holds since  $\phi^L \ge \phi$  pointwise. Finally, normalizing to  $\tilde{\phi}^L$  can only improve the objective function, since

$$\ell(\tilde{\phi}^L, P) = \mathbb{E}_P\left[\tilde{\phi}^L(X)\right] = \mathbb{E}_P\left[\phi^L(X)\right] - \log\left(\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x\right) \ge \ell(\phi^L, P),$$

because  $\log t \le t - 1$  for all t > 0.

From this point on, we only need to prove (11) in order to complete the proof of the lemma. For any  $x \in \mathbb{R}^d$ , we will write  $y_x$  to denote a point attaining the supremum, i.e.,  $\phi^L(x) = \phi(y_x) - L ||x - y_x||$  (Lemma 16 below verifies the existence and measurability of such a map  $x \mapsto y_x$ ).

We now derive the desired bound (11). We have

$$\begin{split} \int_{\mathbb{R}^d} e^{\phi^L(x)} \; \mathrm{d}x &= \int_{\mathbb{R}^d} e^{\phi(y_x)} \cdot e^{-L\|x-y_x\|} \; \mathrm{d}x \\ &= \int_{\mathbb{R}^d} \left( \int_{-\infty}^{\phi(y_x)} e^t \; \mathrm{d}t \right) \cdot \left( \int_{L\|x-y_x\|}^{\infty} e^{-s} \; \mathrm{d}s \right) \, \mathrm{d}x \\ &= \int_{-\infty}^{M_\phi} \int_{0}^{\infty} e^{t-s} \left( \int_{\mathbb{R}^d} \mathbbm{1} \left\{ \phi(y_x) \geq t, \|x-y_x\| \leq s/L \right\} \; \mathrm{d}x \right) \, \mathrm{d}s \; \mathrm{d}t, \end{split}$$

where the last step follows by Fubini's theorem, and where  $M_{\phi} = \sup_{x \in \mathbb{R}^d} \phi(x)$  (note that we must have  $M_{\phi} < \infty$  by definition of  $\Phi_d$ ). We now examine this indicator function. For  $t \in \mathbb{R}$  define the super-level set  $D_t = \{x : \phi(x) \ge t\}$ . Note that  $D_t$  is convex for any t by concavity of  $\phi$ , and furthermore is bounded since  $\phi$  is a log-density. Moreover, we can observe that  $D_t$  has non-empty interior for any  $t < M_{\phi}$ , since  $\phi$  is concave and is a log-density.

Now, for any compact, convex set  $C \subseteq \mathbb{R}^d$  and any  $\delta > 0$ , define the  $\delta$ -neighborhood of C by

$$Nbd(C, \delta) := \{x \in \mathbb{R}^d : dist(x, C) \le \delta\},\$$

where  $\operatorname{dist}(x,C) := \min_{y \in C} \|x - y\|$ . (If C is the empty set, then this neighborhood is also defined to be the empty set.) If  $x \in \mathbb{R}^d$  is such that  $\phi(y_x) \geq t$ , then  $y_x \in D_t$ , and if, furthermore,  $\|x - y_x\| \leq s/L$ , then

$$x \in \mathrm{Nbd}(D_t, s/L).$$

Hence,

$$\int_{\mathbb{R}^d} e^{\phi^L(x)} \, \mathrm{d}x \le \int_{-\infty}^{M_\phi} \int_0^\infty e^{t-s} \cdot \mathrm{Leb}_d \Big( \mathrm{Nbd}(D_t, s/L) \Big) \, \mathrm{d}s \, \mathrm{d}t.$$

On the other hand, we have

$$\int_{-\infty}^{M_{\phi}} \int_{0}^{\infty} e^{t-s} \cdot \operatorname{Leb}_{d}(D_{t}) \, ds \, dt = \int_{-\infty}^{M_{\phi}} e^{t} \cdot \operatorname{Leb}_{d}(D_{t}) \, dt = \int_{-\infty}^{M_{\phi}} e^{t} \left( \int_{\mathbb{R}^{d}} \mathbb{1} \left\{ \phi(x) \ge t \right\} \, dx \right) dt$$
$$= \int_{\mathbb{R}^{d}} \int_{-\infty}^{\phi(x)} e^{t} \, dt \, dx = \int_{\mathbb{R}^{d}} e^{\phi(x)} \, dx = 1, \tag{12}$$

by again applying Fubini's theorem. Therefore, to prove (11), we only need to show that

$$\int_{-\infty}^{M_{\phi}} \int_{0}^{\infty} e^{t-s} \cdot \text{Leb}_{d} \Big( \text{Nbd}(D_{t}, s/L) \backslash D_{t} \Big) \, ds \, dt \le \frac{4d}{Lb_{d}r_{d}\epsilon_{P}}. \tag{13}$$

Next we will use a basic result about neighborhoods of convex sets—Lemma 17 verifies that

$$\delta \mapsto \frac{\mathrm{Leb}_d(\mathrm{Nbd}(C,\delta) \setminus C)}{\delta}$$

is a non-decreasing function for any compact, convex set  $C \subseteq \mathbb{R}^d$  with non-empty interior. Therefore, for any  $t < M_{\phi}$ , it holds that

$$\operatorname{Leb}_{d}\left(\operatorname{Nbd}(D_{t}, s/L)\backslash D_{t}\right) \leq \frac{2d}{Lr_{d}\epsilon_{P}} \cdot \operatorname{Leb}_{d}\left(\operatorname{Nbd}\left(D_{t}, \frac{sr_{d}\epsilon_{P}}{2d}\right)\backslash D_{t}\right)$$

since we have assumed  $L \ge \frac{2d}{r_d \epsilon_P}$ . We also have  $D_t \subseteq D_{t+\log b_d}$ , where  $b_d \in (0,1]$  is the constant appearing in Lemma 9, and so

$$\operatorname{Leb}_d\left(\operatorname{Nbd}\left(D_t, \frac{sr_d\epsilon_P}{2d}\right)\setminus D_t\right) \leq \operatorname{Leb}_d\left(\operatorname{Nbd}\left(D_t, \frac{sr_d\epsilon_P}{2d}\right)\right) \leq \operatorname{Leb}_d\left(\operatorname{Nbd}\left(D_{t+\log b_d}, \frac{sr_d\epsilon_P}{2d}\right)\right).$$

Recall from Lemma 9 that  $D_{M_{\phi} + \log b_d}$  contains  $\mathbb{B}_d(\mu_P, r_d \epsilon_P)$ . Therefore, for any  $t < M_{\phi}$ ,  $D_{t + \log b_d} \supseteq D_{M_{\phi} + \log b_d}$  also contains this ball, and so

$$\operatorname{Nbd}\left(D_{t+\log b_d}, \frac{sr_d\epsilon_P}{2d}\right) = D_{t+\log b_d} + \frac{s}{2d} \cdot \mathbb{B}_d(\mu_P, r_d\epsilon_P)$$

$$\subseteq D_{t+\log b_d} + \frac{s}{2d} \cdot D_{t+\log b_d}$$

$$= \left(1 + \frac{s}{2d}\right) \cdot D_{t+\log b_d},$$

where for two sets  $A, B \subseteq \mathbb{R}^d$ , we write  $A + B := \{x + y : x \in A, y \in B\}$  to denote their Minkowski sum. Therefore,

$$\operatorname{Leb}_d\left(\operatorname{Nbd}\left(D_{t+\log b_d}, \frac{sr_d\epsilon_P}{2d}\right)\right) \le \operatorname{Leb}_d(D_{t+\log b_d}) \cdot \left(1 + \frac{s}{2d}\right)^d \le \operatorname{Leb}_d(D_{t+\log b_d}) \cdot e^{s/2}$$

for any  $t < M_{\phi}$ . Combining this with our work above, we obtain

$$\operatorname{Leb}_d\left(\operatorname{Nbd}(D_t, s/L)\backslash D_t\right) \le \frac{2d}{Lr_d\epsilon_P} \cdot \operatorname{Leb}_d(D_{t+\log b_d}) \cdot e^{s/2}$$
 (14)

for any  $t < M_{\phi}$ . Therefore,

$$\begin{split} \int_{-\infty}^{M_{\phi}} \int_{0}^{\infty} e^{t-s} \cdot \operatorname{Leb}_{d} \Big( \operatorname{Nbd}(D_{t}, s/L) \backslash D_{t} \Big) \, \mathrm{d}s \, \mathrm{d}t \\ & \leq \int_{-\infty}^{M_{\phi}} \int_{0}^{\infty} e^{t-s} \cdot \frac{2d}{Lr_{d}\epsilon_{P}} \cdot \operatorname{Leb}_{d}(D_{t+\log b_{d}}) \cdot e^{s/2} \, \mathrm{d}s \, \mathrm{d}t \\ & = \frac{2d}{Lr_{d}\epsilon_{P}} \cdot \left( \int_{-\infty}^{M_{\phi}} e^{t} \cdot \operatorname{Leb}_{d}(D_{t+\log b_{d}}) \, \mathrm{d}t \right) \cdot \left( \int_{0}^{\infty} e^{-s} \cdot e^{s/2} \, \mathrm{d}s \right) \\ & = \frac{4d}{Lr_{d}\epsilon_{P}} \cdot \int_{-\infty}^{M_{\phi}} e^{t} \cdot \operatorname{Leb}_{d}(D_{t+\log b_{d}}) \, \mathrm{d}t \\ & = \frac{4d}{Lb_{d}r_{d}\epsilon_{P}} \cdot \int_{-\infty}^{M_{\phi}} e^{t+\log b_{d}} \cdot \operatorname{Leb}_{d}(D_{t+\log b_{d}}) \, \mathrm{d}t \\ & = \frac{4d}{Lb_{d}r_{d}\epsilon_{P}} \cdot \int_{-\infty}^{M_{\phi}+\log b_{d}} e^{t} \cdot \operatorname{Leb}_{d}(D_{t}) \, \mathrm{d}t \\ & \leq \frac{4d}{Lb_{d}r_{d}\epsilon_{P}} \cdot \int_{-\infty}^{M_{\phi}} e^{t} \cdot \operatorname{Leb}_{d}(D_{t}) \, \mathrm{d}t \\ & = \frac{4d}{Lb_{d}r_{d}\epsilon_{P}}, \end{split}$$

where for the last step we again apply (12). This completes the proof of Lemma 10.

## A.3.1 Supporting lemmas for Lemma 10

**Lemma 16.** For any  $x \in \mathbb{R}^d$  and any  $\phi \in \Phi_d$ , there exists a Borel measurable map  $x \mapsto y_x$  such that  $y_x$  attains  $\sup_{y \in \mathbb{R}^d} {\{\phi(y) - L || x - y||\}}$ .

Proof of Lemma 16. Let  $M_{\phi} := \sup_{x \in \mathbb{R}^d} \phi(x)$ , and let  $x_{\phi} \in \operatorname{argmax}_{x \in \mathbb{R}^d} \phi(x)$  (note that, by definition of  $\Phi_d \ni \phi$ ,  $M_{\phi}$  must be finite, and  $x_{\phi}$  must exist). Define

$$\mathcal{Y} = \left\{ y \in \mathbb{R}^d : \phi(y) \ge \phi(y') - L \|y - y'\| \text{ for all } y' \in \mathbb{R}^d \right\}.$$

Note that  $\mathcal{Y}$  is non-empty, since trivially  $x_{\phi} \in \mathcal{Y}$ .

Next define  $h: \mathbb{R}^d \times \mathcal{Y} \to \mathbb{R}$  as  $h(x,y) = \phi(y) - L ||x-y||$ . For each  $x \in \mathbb{R}^d$ , define

$$S(x) = \mathcal{Y} \cap \mathbb{B}_d(x, ||x - x_{\phi}||).$$

Note that, for any x, we have  $x_{\phi} \in S(x)$  by definition.

Now we will apply Aliprantis and Border (2006, Theorem 18.19), which guarantees the existence of a Borel measurable function  $x \mapsto y_x \in S(x)$  such that, for each x,

$$y_x \in \operatorname*{argmax}_{y \in S(x)} h(x, y),$$

as long as we verify the following conditions:

ullet R<sup>d</sup> is a measurable space, and  ${\mathcal Y}$  is a separable metrizable space. This holds trivially.

- h is a Carathéodory function (i.e.,  $x \mapsto h(x,y)$  is measurable for any  $y \in \mathcal{Y}$ , and  $y \mapsto h(x,y)$  is continuous for almost every  $x \in \mathbb{R}^d$ ). It holds trivially that  $x \mapsto h(x,y)$  is measurable. To check that  $y \mapsto h(x,y)$  is continuous for any fixed x, it is sufficient to verify that  $\phi$  is continuous on  $\mathcal{Y}$ . In fact, examining the definition of  $\mathcal{Y}$ , we can see that  $\phi$  is L-Lipschitz on  $\mathcal{Y}$  by definition, thus ensuring continuity.
- S(x) is non-empty and compact for any  $x \in \mathbb{R}^d$ . We have already seen that  $x_{\phi} \in S(x)$  for all x. To check compactness, it is sufficient to verify that  $\mathcal{Y}$  is closed, which follows immediately from the definition of  $\mathcal{Y}$  along with the fact that  $\phi$  is upper semi-continuous (by definition of  $\phi \in \Phi_d$ ).
- In the terminology of Aliprantis and Border (2006), the correspondence  $\mathcal{X} \to \mathcal{Y}$ , mapping  $x \mapsto S(x) \subseteq \mathcal{Y}$ , is weakly measurable, meaning that the set  $X_A := \{x \in \mathbb{R}^d : S(x) \cap A \neq \emptyset\}$  is measurable for any open subset  $A \subseteq \mathcal{Y}$ . Aliprantis and Border (2006, Lemma 18.2) establishes that, since  $\mathcal{Y}$  is metrizable, this is implied by the stronger condition that  $X_A$  is measurable for every *closed* subset  $A \subseteq \mathcal{Y}$ , so we will check this stronger condition.

Let  $A \subseteq \mathcal{Y}$  be a closed subset. Consider any  $x, x_1, x_2, \ldots \in \mathbb{R}^d$  such that  $x_i \in X_A$  for all  $i \geq 1$  and such that  $\lim_{i \to \infty} x_i = x$ . Let  $R = \sup_i \|x_i - x_\phi\|$ , which is finite since the sequence converges. This means that  $S(x_i) \subseteq \mathbb{B}_d(x_\phi, 2R)$  for all i. For each  $i, x_i \in X_A$  implies that  $S(x_i) \cap A \neq \emptyset$ , and so we can find some  $y_i \in S(x_i) \cap A \subseteq \mathbb{B}_d(x_\phi, 2R)$ . Therefore, we can find some convergent subsequence, i.e.,  $i_1, i_2, \ldots$  such that  $\lim_{j \to \infty} y_{i_j} = y$  for some  $y \in \mathbb{R}^d$ . By assumption, A is a closed subset of  $\mathcal{Y}$ , and we have already shown that  $\mathcal{Y}$  is a closed subset of  $\mathbb{R}^d$ . Therefore,  $A \subseteq \mathbb{R}^d$  is closed, and so we must have  $y \in A$ . Now we check that  $y \in S(x)$ . We know that  $y \in A \subseteq \mathcal{Y}$ , and so we only need to check that  $y \in \mathbb{B}_d(x, \|x - x_\phi\|)$ . This holds because, for each  $j \geq 1, y_{i_j} \in S(x_{i_j}) \subseteq \mathbb{B}_d(x_{i_j}, \|x_{i_j} - x_\phi\|)$ , and so

$$||y - x|| = \lim_{i \to \infty} ||y_{i_j} - x_{i_j}|| \le \lim_{i \to \infty} ||x_{i_j} - x_{\phi}|| = ||x - x_{\phi}||.$$

We have now seen that  $y \in S(x) \cap A$ , proving that  $S(x) \cap A \neq \emptyset$  and so  $x \in X_A$ . Therefore, we have established that  $X_A$  is closed, and is therefore measurable.

Finally we check that, for any x,

$$\sup_{y \in \mathbb{R}^d} \{ \phi(y) - L \|x - y\| \} = \sup_{y \in S(x)} \{ \phi(y) - L \|x - y\| \}.$$

First, for any  $y \notin \mathbb{B}_d(x, ||x - x_{\phi}||)$ , we have  $||x - y|| > ||x - x_{\phi}||$ , and so since  $\phi(y) \leq \phi(x_{\phi})$  by definition of  $x_{\phi}$ , it holds that

$$\phi(y) - L \|x - y\| < \phi(x_{\phi}) - L \|x - x_{\phi}\|.$$

Therefore,

$$\sup_{y \in \mathbb{R}^d} \{ \phi(y) - L \|x - y\| \} = \sup_{y \in \mathbb{B}_d(x, \|x - x_\phi\|)} \{ \phi(y) - L \|x - y\| \}.$$

Next, since  $\phi$  is upper semi-continuous, the supremum on the right-hand side is attained, i.e., there exists some  $y_1 \in \mathbb{B}_d(x, ||x - x_{\phi}||)$  such that

$$\phi(y_1) - L \|x - y_1\| = \sup_{y \in \mathbb{B}_d(x, \|x - x_\phi\|)} \{\phi(y) - L \|x - y\|\} = \sup_{y \in \mathbb{R}^d} \{\phi(y) - L \|x - y\|\}.$$

Now we verify that  $y_1 \in \mathcal{Y}$ . To see this, fix any  $y' \in \mathbb{R}^d$ . Then

$$\phi(y') - L \|x - y'\| \le \sup_{y \in \mathbb{R}^d} \{\phi(y) - L \|x - y\|\} = \phi(y_1) - L \|x - y_1\|$$

and so

$$\phi(y_1) \ge \phi(y') - L \|x - y'\| + L \|x - y_1\| \ge \phi(y') - L \|y_1 - y'\|.$$

Since this holds for all  $y' \in \mathbb{R}^d$ , we have established that  $y_1 \in \mathcal{Y}$ . Therefore,  $y_1 \in S(x)$ , which verifies  $\sup_{y \in \mathbb{R}^d} \{\phi(y) - L \|x - y\|\} = \sup_{y \in S(x)} \{\phi(y) - L \|x - y\|\}$ .

**Lemma 17.** Let  $C \subseteq \mathbb{R}^d$  be any compact, convex set with non-empty interior. Then

$$\delta \mapsto \frac{\mathrm{Leb}_d(\mathrm{Nbd}(C,\delta)\backslash C)}{\delta}$$

is a non-decreasing function of  $\delta > 0$ .

*Proof of Lemma 17.* This result follows immediately from Steiner's formula (Schneider, 2014, Chapter 4), which states that for all  $\epsilon \geq 0$ ,

$$\operatorname{Leb}_d(\operatorname{Nbd}(C,\epsilon)) = \operatorname{Leb}_d(C) + \sum_{k=1}^d V_{d-k}(C) \cdot \operatorname{Leb}_k(\mathbb{B}_k) \cdot \epsilon^k,$$

where  $V_{d-k}(C) \geq 0$  is the (d-k)-th intrinsic volume of C. Rearranging, we have

$$\frac{\operatorname{Leb}_d(\operatorname{Nbd}(C,\epsilon)\backslash C)}{\epsilon} = \sum_{k=1}^d V_{d-k}(C) \cdot \operatorname{Leb}_k(\mathbb{B}_k) \cdot \epsilon^{k-1},$$

which is a non-decreasing function of  $\epsilon$ .

#### A.4 Proof of Lemma 12

First we consider the bounded case. Suppose that P and Q are both supported on [-R, R] for some R > 0. Write  $\Delta = \Delta_{\text{CDF}}(P, Q)$  and  $\epsilon = \min\{\epsilon_P, \epsilon_Q\}$ . Let  $r_1, b_1 \in (0, 1]$  be the universal constants defined in Lemma 9 (for dimension d = 1), and fix any  $L \geq \frac{4}{r_1 \epsilon}$ . By Corollary 11, we have

$$d_{\mathrm{H}}^{2}\left(\psi^{*}(P), \psi^{*}(Q)\right) \leq \frac{16}{Lb_{1}r_{1}\epsilon} + \left(\mathbb{E}_{P}\left[\phi_{P}^{L}(X)\right] - \mathbb{E}_{Q}\left[\phi_{P}^{L}(X)\right]\right) + \left(\mathbb{E}_{Q}\left[\phi_{Q}^{L}(X)\right] - \mathbb{E}_{P}\left[\phi_{Q}^{L}(X)\right]\right).$$

Now we bound the two differences. For any  $\phi \in \Phi_d$  define  $M_{\phi} = \sup_{x \in \mathbb{R}^d} \phi(x)$  (note that  $M_{\phi}$  is finite by definition of  $\Phi_d$ ). We note that  $M_{\phi_P} = M_{\phi_P^L}$  by definition of  $\phi_P^L$ , and that

 $\phi_P^L(X) \ge M_{\phi_P} - 2LR$  with probability 1 under either P or Q, since the distributions are supported on [-R,R] and so  $\phi_P$  must attain its maximum somewhere in this range. We then have

$$\begin{split} \mathbb{E}_P\left[\phi_P^L(X)\right] - \mathbb{E}_Q\left[\phi_P^L(X)\right] &= \mathbb{E}_Q\left[M_{\phi_P} - \phi_P^L(X)\right] - \mathbb{E}_P\left[M_{\phi_P} - \phi_P^L(X)\right] \\ &= \int_0^{2LR} \left(\mathbb{P}_Q\big\{M_{\phi_P} - \phi_P^L(X) \geq t\big\} - \mathbb{P}_P\big\{M_{\phi_P} - \phi_P^L(X) \geq t\big\}\right) \; \mathrm{d}t. \end{split}$$

It is trivial to verify that

$$\left| \sqrt{\mathbb{P}_P\{X \not\in C\}} - \sqrt{\mathbb{P}_Q\{X \not\in C\}} \right| \le \Delta\sqrt{2}$$

for any convex set (i.e., an interval)  $C \subseteq \mathbb{R}$ , by definition of  $\Delta$  (this follows from the fact that  $|\sqrt{a+c}-\sqrt{b+d}|^2 \le |\sqrt{a}-\sqrt{b}|^2 + |\sqrt{c}-\sqrt{d}|^2$  for any  $a,b,c,d \ge 0$ ). Since  $\phi_P^L$  is concave, the set  $\{x: M_{\phi_P} - \phi_P^L(x) < t\}$  is convex, and so

$$\mathbb{P}_Q \left\{ M_{\phi_P} - \phi_P^L(X) \ge t \right\} \le \left( \sqrt{\mathbb{P}_P \{ M_{\phi_P} - \phi_P^L(X) \ge t \}} + \Delta \sqrt{2} \right)^2$$

and so, since it also holds that  $\phi_P^L \ge \phi_P$  pointwise, we have

$$\mathbb{P}_Q\big\{M_{\phi_P} - \phi_P^L(X) \geq t\big\} - \mathbb{P}_P\big\{M_{\phi_P} - \phi_P^L(X) \geq t\big\} \leq \Delta\sqrt{8} \cdot \sqrt{\mathbb{P}_P\{M_{\phi_P} - \phi_P(X) \geq t\}} + 2\Delta^2.$$

Lemma 18 below will establish that, for  $t \geq \frac{8R}{r_1\epsilon}$ , we have  $\mathbb{P}_P\{M_{\phi_P} - \phi_P(X) \geq t\} \leq \frac{32}{b_1r_1\epsilon} \cdot \frac{R}{t^2}$  Applying this bound, we have

$$\begin{split} &\mathbb{E}_{P}\left[\phi_{P}^{L}(X)\right] - \mathbb{E}_{Q}\left[\phi_{P}^{L}(X)\right] \\ &\leq \int_{0}^{2LR} \left(\Delta\sqrt{8} \cdot \sqrt{\mathbb{P}_{P}\{M_{\phi_{P}} - \phi_{P}(X) \geq t\}} + 2\Delta^{2}\right) \, \mathrm{d}t \\ &= \Delta\sqrt{8} \int_{0}^{2LR} \sqrt{\mathbb{P}_{P}\{M_{\phi_{P}} - \phi_{P}(X) \geq t\}} \, \, \mathrm{d}t + 4LR\Delta^{2} \\ &= \Delta\sqrt{8} \left(\int_{0}^{\frac{8R}{r_{1}\epsilon}} \sqrt{\mathbb{P}_{P}\{M_{\phi_{P}} - \phi_{P}(X) \geq t\}} \, \, \mathrm{d}t + \int_{\frac{8R}{r_{1}\epsilon}}^{2LR} \sqrt{\mathbb{P}_{P}\{M_{\phi_{P}} - \phi_{P}(X) \geq t\}} \, \, \mathrm{d}t\right) + 4LR\Delta^{2} \\ &\leq \Delta\sqrt{8} \sqrt{\frac{8R}{r_{1}\epsilon}} \cdot \left(\int_{0}^{\frac{8R}{r_{1}\epsilon}} \mathbb{P}_{P}\{M_{\phi_{P}} - \phi_{P}(X) \geq t\} \, \, \mathrm{d}t\right)^{1/2} + \Delta\sqrt{8} \int_{\frac{8R}{r_{1}\epsilon}}^{2LR} \sqrt{\frac{32}{b_{1}r_{1}\epsilon} \cdot \frac{R}{t^{2}}} \, \mathrm{d}t + 4LR\Delta^{2} \\ &\leq \Delta\sqrt{8} \sqrt{\frac{8R}{r_{1}\epsilon}} \cdot \sqrt{\mathbb{E}_{P}\left[M_{\phi_{P}} - \phi_{P}(X)\right]} + \Delta\sqrt{8} \sqrt{\frac{32R}{b_{1}r_{1}\epsilon}} \log\left(Lr_{1}\epsilon/4\right) + 4LR\Delta^{2} \\ &\leq \Delta\sqrt{8} \sqrt{\frac{8Rh_{1}}{r_{1}\epsilon}} + \Delta\sqrt{8} \sqrt{\frac{32R}{b_{1}r_{1}\epsilon}} \log\left(Lr_{1}\epsilon/4\right) + 4LR\Delta^{2}, \end{split}$$

where the last step applies Lemma 19 below, which will establish that  $\mathbb{E}_P \left[ \phi(X) \right] \geq M_\phi - h_1$  for a universal constant  $h_1$ . By symmetry the same bound holds for  $\mathbb{E}_Q \left[ \phi_Q^L(X) \right] - \mathbb{E}_P \left[ \phi_Q^L(X) \right]$ . Combining all our work so far, then,

$$d_{\mathrm{H}}^{2}\left(\psi^{*}(P), \psi^{*}(Q)\right) \leq \frac{16}{Lb_{1}r_{1}\epsilon} + 2\left\{\Delta\sqrt{8}\left(\sqrt{\frac{8Rh_{1}}{r_{1}\epsilon}} + \sqrt{\frac{32R}{b_{1}r_{1}\epsilon}}\log\left(Lr_{1}\epsilon/4\right)\right) + 4LR\Delta^{2}\right\}.$$

Next we split into cases. If  $\frac{1}{\Delta\sqrt{R\epsilon}} \geq \frac{4}{r_1\epsilon}$ , then setting  $L = \frac{1}{\Delta\sqrt{R\epsilon}}$  we apply this bound to obtain

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C' \Delta \sqrt{R/\epsilon} \max \left\{1, \log\left(\frac{1}{\Delta \sqrt{R/\epsilon}}\right)\right\},$$

for a universal constant C'. Since  $\epsilon \leq 2R$  by definition, and  $\Delta \leq 1$ , we can relax this to

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C' \Delta \sqrt{R/\epsilon} \log(e/\Delta).$$

If instead  $\frac{1}{\Delta\sqrt{R\epsilon}} < \frac{4}{r_1\epsilon}$ , then

$$d_{\mathrm{H}}^2(\psi^*(P), \psi^*(Q)) \le 2 \le \frac{8}{r_1} \Delta \sqrt{R/\epsilon}.$$

Therefore, combining both cases, we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C'' \Delta \sqrt{\frac{R}{\min\{\epsilon_{P}, \epsilon_{Q}\}}} \log(e/\Delta)$$
(15)

for a universal constant  $C'' = \max\{C', 8/r_1\}$ . Next we will need to relate  $\min\{\epsilon_P, \epsilon_Q\}$  with  $\max\{\epsilon_P, \epsilon_Q\}$ . Without loss of generality, suppose that  $\mu_P \ge \mu_Q$ . We then have

$$\begin{split} &\frac{\epsilon_Q}{2} = \frac{1}{2} \mathbb{E}_Q \left[ |X - \mu_Q| \right] = \mathbb{E}_Q \left[ (X - \mu_Q)_+ \right] \\ &\geq \mathbb{E}_Q \left[ (X - \mu_P)_+ \right] = \int_{\mu_P}^R \mathbb{P}_Q \{ X > t \} \ \mathrm{d}t \\ &\geq \int_{\mu_P}^R \mathbb{P}_P \{ X > t \} - 2\Delta \sqrt{\mathbb{P}_P \{ X > t \}} \ \mathrm{d}t \\ &\geq \int_{\mu_P}^R \mathbb{P}_P \{ X > t \} \ \mathrm{d}t - 2\Delta \sqrt{R - \mu_P} \sqrt{\int_{\mu_P}^R \mathbb{P}_P \{ X > t \}} \ \mathrm{d}t \\ &\geq \mathbb{E}_P \left[ (X - \mu_P)_+ \right] - 2\Delta \sqrt{2R} \sqrt{\mathbb{E}_P \left[ (X - \mu_P)_+ \right]} \\ &= \frac{\epsilon_P}{2} - 2\Delta \sqrt{R \cdot \epsilon_P}, \end{split}$$

where the final inequality follows because  $|\mu_P| \leq R$ . We can similarly calculate

$$\frac{\epsilon_P}{2} = \frac{1}{2} \mathbb{E}_P \left[ |X - \mu_P| \right] = \mathbb{E}_P \left[ (X - \mu_P)_- \right] \ge \frac{\epsilon_Q}{2} - 2\Delta \sqrt{R \cdot \epsilon_Q}.$$

Combining these two bounds, then,

$$\max\{\epsilon_P, \epsilon_Q\} = \min\{\epsilon_P, \epsilon_Q\} + |\epsilon_P - \epsilon_Q| \le \min\{\epsilon_P, \epsilon_Q\} + 4\Delta_{CDF}(P, Q) \cdot \sqrt{R \cdot \max\{\epsilon_P, \epsilon_Q\}}.$$
(16)

Now we work with the general case, where P, Q may not have bounded support. Fix any R > 0. For any  $x \in \mathbb{R}$  define

$$[x]_R := \begin{cases} -R, & x < -R, \\ x, & |x| \le R, \\ R, & x > R, \end{cases}$$
 (17)

the truncation of x to the range [-R, R]. Let  $[P]_R$  denote the distribution of  $[X]_R$  when  $X \sim P$ , and same for  $[Q]_R$ . Lemma 20 below calculates that  $d_W(P, [P]_R) \leq \frac{M_q^q}{R^{q-1}}$ . Applying Theorem 2 to compare the distributions P and  $[P]_R$ , then, we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}([P]_{R})) \leq C_{1}^{2} \sqrt{\frac{d_{\mathrm{W}}(P, [P]_{R})}{\max\{\epsilon_{P}, \epsilon_{[P]_{R}}\}}} \leq C_{1}^{2} \sqrt{\frac{M_{q}^{q}}{\epsilon_{[P]_{R}} R^{q-1}}},$$

and the same bound holds with Q in place of P. Therefore, by the triangle inequality,

$$d_{H}^{2}(\psi^{*}(P),\psi^{*}(Q)) \leq \left\{ d_{H}(\psi^{*}(P),\psi^{*}([P]_{R})) + d_{H}(\psi^{*}(Q),\psi^{*}([Q]_{R})) + d_{H}(\psi^{*}([P]_{R}),\psi^{*}([Q]_{R})) \right\}^{2} \\ \leq 3d_{H}^{2}(\psi^{*}(P),\psi^{*}([P]_{R})) + 3d_{H}^{2}(\psi^{*}(Q),\psi^{*}([Q]_{R})) + 3d_{H}^{2}(\psi^{*}([P]_{R}),\psi^{*}([Q]_{R})) \\ \leq 6C_{1}^{2}\sqrt{\frac{M_{q}^{q}}{\min\{\epsilon_{[P]_{R}},\epsilon_{[Q]_{R}}\}R^{q-1}}} + 3d_{H}^{2}(\psi^{*}([P]_{R}),\psi^{*}([Q]_{R})).$$

$$(18)$$

We now need to apply the bound (15) to the bounded distributions  $[P]_R$  and  $[Q]_R$ , in order to bound this last term. Combining (15) with (18), we obtain

$$d_{H}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq 6C_{1}^{2} \sqrt{\frac{M_{q}^{q}}{\min\{\epsilon_{[P]_{R}}, \epsilon_{[Q]_{R}}\}R^{q-1}}} + 3C'' \Delta_{CDF}([P]_{R}, [Q]_{R}) \sqrt{\frac{R}{\min\{\epsilon_{[P]_{R}}, \epsilon_{[Q]_{R}}\}}} \log(e/\Delta_{CDF}([P]_{R}, [Q]_{R})).$$

Now fix

$$R = M_q \Big\{ \Delta_{\text{CDF}}([P]_R, [Q]_R) \log(e/\Delta_{\text{CDF}}([P]_R, [Q]_R)) \Big\}^{-2/q}.$$

This yields

$$\begin{aligned} \mathrm{d}_{\mathrm{H}}^2 \Big( \psi^*(P), \psi^*(Q) \Big) \\ &\leq C'_* \sqrt{\frac{M_q}{\min\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}}} \cdot \left\{ \Delta_{\mathrm{CDF}}([P]_R, [Q]_R) \log \left( \frac{e}{\Delta_{\mathrm{CDF}}([P]_R, [Q]_R)} \right) \right\}^{1-1/q}, \end{aligned}$$

when the universal constant  $C'_* > 0$  is chosen appropriately. Next, it holds trivially that  $\Delta_{\text{CDF}}([P]_R, [Q]_R) \leq \Delta_{\text{CDF}}(P, Q)$ , and since  $t \mapsto t \log(e/t)$  is increasing on  $t \in (0, 1]$ , we therefore have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C'_{*} \sqrt{\frac{M_{q}}{\min\{\epsilon_{[P]_{R}}, \epsilon_{[Q]_{R}}\}}} \cdot \left\{\Delta_{\mathrm{CDF}}(P, Q) \log(e/\Delta_{\mathrm{CDF}}(P, Q))\right\}^{1-1/q}.$$

Finally, we need to lower bound  $\epsilon_{[P]_R}$  and  $\epsilon_{[Q]_R}$ . First we relate  $\min\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}$  to  $\max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}$ . Applying (16) from above, along with the fact that  $\Delta_{\text{CDF}}([P]_R, [Q]_R) \leq \Delta_{\text{CDF}}(P, Q)$ , we have

$$\max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\} \le \min\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\} + 4\Delta_{\mathrm{CDF}}(P, Q)\sqrt{R \cdot \max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}}.$$

If  $8\Delta_{\text{CDF}}(P,Q)\sqrt{R} \leq \sqrt{\max\{\epsilon_{[P]_R},\epsilon_{[Q]_R}\}}$ , then this proves that

$$\max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\} \le 2\min\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}$$

and so

$$\mathrm{d}_{\mathrm{H}}^2\big(\psi^*(P),\psi^*(Q)\big) \leq C'_*\sqrt{\frac{2M_q}{\max\{\epsilon_{[P]_R},\epsilon_{[Q]_R}\}}} \cdot \Big\{\Delta_{\mathrm{CDF}}(P,Q)\log\big(e/\Delta_{\mathrm{CDF}}(P,Q)\big)\Big\}^{1-1/q}.$$

If instead  $8\Delta_{\text{CDF}}(P,Q)\sqrt{R} > \sqrt{\max\{\epsilon_{[P]_R},\epsilon_{[Q]_R}\}}$ , then we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq 2 \leq \frac{16\Delta_{\mathrm{CDF}}(P, Q)\sqrt{R}}{\sqrt{\max\{\epsilon_{[P]_{R}}, \epsilon_{[Q]_{R}}\}}}.$$

Plugging in the definition of R and combining both cases, we obtain

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C_{*}'' \sqrt{\frac{M_{q}}{\max\{\epsilon_{[P]_{R}}, \epsilon_{[Q]_{R}}\}}} \cdot \left(\Delta_{\mathrm{CDF}}(P, Q) \log(e/\Delta_{\mathrm{CDF}}(P, Q))\right)^{1-1/q}$$

for an appropriately chosen universal constant  $C_*''$ . The last step is to relate  $\max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\}$  to  $\max\{\epsilon_P, \epsilon_Q\}$ . Applying Proposition 1 together with the bound on  $d_W(P, [P]_R)$  from Lemma 20, we have

$$\epsilon_{[P]_R} \ge \epsilon_P - 2d_W(P, [P]_R) \ge \epsilon_P - 2 \cdot \frac{M_q^q}{R^{q-1}},$$

and the same bound holds for Q in place of P. If  $\frac{2M_q^q}{R^{q-1}} \leq \frac{\max\{\epsilon_P, \epsilon_Q\}}{2}$ , then

$$\max\{\epsilon_{[P]_R}, \epsilon_{[Q]_R}\} \ge \frac{\max\{\epsilon_P, \epsilon_Q\}}{2}$$

and so we obtain

$$\mathrm{d}_{\mathrm{H}}^{2}\big(\psi^{*}(P),\psi^{*}(Q)\big) \leq C_{*}''\sqrt{\frac{2M_{q}}{\max\{\epsilon_{P},\epsilon_{Q}\}}} \cdot \Big\{\Delta_{\mathrm{CDF}}(P,Q)\log\big(e/\Delta_{\mathrm{CDF}}(P,Q)\big)\Big\}^{1-1/q}.$$

If instead  $\frac{2M_q^q}{R^{q-1}} > \frac{\max\{\epsilon_P, \epsilon_Q\}}{2}$ , then it trivially holds that

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \le 2 \le 2\sqrt{\frac{4M_{q}^{q}}{\max\{\epsilon_{P}, \epsilon_{Q}\}R^{q-1}}}.$$

Plugging in the definition of R, and combining the two cases, we obtain

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \leq C_{*}\sqrt{\frac{M_{q}}{\max\{\epsilon_{P}, \epsilon_{Q}\}}} \cdot \left\{\Delta_{\mathrm{CDF}}(P, Q) \log(e/\Delta_{\mathrm{CDF}}(P, Q))\right\}^{1-1/q}$$

for appropriately chosen universal constant  $C_*$ , which completes the proof of Lemma 12.

#### A.4.1 Supporting lemmas for Lemma 12

**Lemma 18.** Let  $P \in \mathcal{P}_d$  and let  $\phi = \phi^*(P)$ . Let  $M_{\phi} := \sup_{x \in \mathbb{R}^d} \phi(x)$  and let  $x_{\phi} \in \underset{\text{argmax}}{\operatorname{argmax}} \phi(x)$  (which is guaranteed to exist by definition of  $\Phi_d \ni \phi$ ). Fix any R > 0 and  $t \geq \frac{8dR}{r_d \epsilon_P}$ , where  $r_d \in (0,1]$  is taken from Lemma 9. Then

$$\mathbb{P}_P\{\phi(X) \le M_\phi - t \text{ and } ||X - x_\phi|| \le 2R\} \le \frac{32d}{b_d r_d \epsilon_P} \cdot \frac{R}{t^2},$$

where  $b_d \in (0,1]$  is taken from Lemma 9.

Proof of Lemma 18. First, for any x with  $||x - x_{\phi}|| \le 2R$ ,

$$\phi^{t/4R}(x) = \sup_{y \in \mathbb{R}^d} \left\{ \phi(y) - \frac{t}{4R} \|y - x\| \right\} \ge \phi(x_\phi) - \frac{t}{4R} \|x - x_\phi\| \ge M_\phi - \frac{t}{2}.$$

Hence, if  $\phi(x) \leq M_{\phi} - t$  and  $||x - x_{\phi}|| \leq 2R$ , then

$$\phi^{t/4R}(x) - \phi(x) \ge \frac{t}{2}.$$

Moreover, by definition of  $\phi = \phi^*(P)$ , since  $\phi^{t/4R} \in \Phi_d$ , it holds that

$$\mathbb{E}_{P}\left[\phi(X)\right] = \ell(\phi, P) \ge \ell(\phi^{t/4R}, P) = \mathbb{E}_{P}\left[\phi^{t/4R}(X)\right] - \int_{\mathbb{R}^{d}} e^{\phi^{t/4R}(x)} \, \mathrm{d}x + 1$$

$$\ge \mathbb{E}_{P}\left[\phi^{t/4R}(X)\right] - \frac{4d}{\frac{t}{4R}b_{d}r_{d}\epsilon_{P}},$$

where the last step holds by (11) as calculated in the proof of Lemma 10, noting that  $\frac{t}{4R} \ge \frac{2d}{r_d \epsilon_P}$ . We deduce that

$$\mathbb{P}_{P}\{\phi(X) \leq M_{\phi} - t \text{ and } \|X - x_{\phi}\| \leq 2R\} \leq \mathbb{P}_{P}\left\{\phi^{t/4R}(X) - \phi(X) \geq \frac{t}{2}\right\}$$

$$\leq \frac{\mathbb{E}_{P}\left[\phi^{t/4R}(X) - \phi(X)\right]}{t/2} \leq \frac{\frac{4d}{\frac{t}{4R}b_{d}r_{d}\epsilon_{P}}}{t/2} = \frac{32d}{b_{d}r_{d}\epsilon_{P}} \cdot \frac{R}{t^{2}},$$

as required.

**Lemma 19.** Fix any  $P \in \mathcal{P}_d$  and let  $\phi = \phi^*(P)$ . Then

$$\mathbb{E}_P\left[\phi(X)\right] \ge M_\phi - h_d,$$

where  $M_{\phi} = \sup_{x \in \mathbb{R}^d} \phi(x)$  and where  $h_d \geq 0$  depends only on d.

Proof of Lemma 19. Write  $\mathbb{E}_{\phi}[\cdot]$  to denote the expectation with respect to the distribution with log-density  $\phi$ . Let  $\mu_{\phi} := \mathbb{E}_{\phi}[X]$  be the mean and  $\Sigma := \mathbb{E}_{\phi}[(X - \mu_{\phi})(X - \mu_{\phi})^{\top}]$  the covariance of this distribution. Let  $\bar{\phi}$  denote the log-density of the isotropic, log-concave random vector  $\Sigma^{-1/2}(X - \mu_{\phi})$ , where X has log-density  $\phi$ . Let  $M_{\bar{\phi}} := \sup_{x \in \mathbb{R}^d} \bar{\phi}(x)$ .

Since  $x \mapsto \phi(x) + \frac{1}{2} \{ M_{\phi} - \phi(x) \}$  is concave and coercive, it holds by Dümbgen et al. (2011, Remark 2.3) that

$$\mathbb{E}_P \left[ M_{\phi} - \phi(X) \right] \leq \mathbb{E}_{\phi} \left[ M_{\phi} - \phi(X) \right].$$

Next, we can trivially verify that

$$\mathbb{E}_{\phi} \left[ M_{\phi} - \phi(X) \right] = \mathbb{E}_{\bar{\phi}} \left[ M_{\bar{\phi}} - \bar{\phi}(X) \right]$$

since the log-densities  $\phi$  and  $\bar{\phi}$  are related via the linear transformation on random variables above. Furthermore,

$$\mathbb{E}_{\bar{\phi}}\left[M_{\bar{\phi}} - \bar{\phi}(X)\right] = M_{\bar{\phi}} - \int_{\mathbb{R}^d} e^{\bar{\phi}(y)} \cdot \bar{\phi}(y) \; \mathrm{d}y \leq M_{\bar{\phi}} + \frac{d}{2}\log(2\pi e),$$

where the last step holds since  $\bar{\phi}$  is the log-density of an isotropic distribution on  $\mathbb{R}^d$ , and so its entropy is bounded by that of the standard d-dimensional Gaussian (e.g. Cover and Thomas, 1991, Theorem 9.6.5). Finally, by Lovász and Vempala (2007, Theorem 5.14(e)),  $M_{\bar{\phi}} \leq m_d$  where  $m_d \in \mathbb{R}$  depends only on the dimension d. Therefore, combining everything,

$$\mathbb{E}_P\left[M_\phi - \phi(X)\right] \le m_d + \frac{d}{2}\log(2\pi e),$$

which proves the desired bound.

**Lemma 20.** Let  $P \in \mathcal{P}_1$  satisfy  $\mathbb{E}_P[|X|^q]^{1/q} \leq M_q$ , for some q > 1. Let  $[P]_R$  be the distribution of  $[X]_R$  when  $X \sim P$  (where the truncation  $[X]_R$  is defined as in (17)). Then

$$d_{\mathcal{W}}(P,[P]_R) \le \frac{M_q^q}{R^{q-1}}.$$

Proof of Lemma 20. Drawing  $X \sim P$ , note that  $(X, [X]_R)$  is a coupling of the distributions P and  $[P]_R$ . Hence

$$d_{W}(P, [P]_{R}) \le \mathbb{E}_{P}[|X - [X]_{R}|] = \mathbb{E}_{P}[(|X| - R)_{+}] \le \mathbb{E}_{P}[\frac{|X|^{q}}{R^{q-1}}] \le \frac{M_{q}^{q}}{R^{q-1}},$$

as required.  $\Box$ 

## A.5 Proof of Lemma 13

Write  $\widehat{U}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{U_i \leq t\}$ . First we calculate

$$\Delta_{\mathrm{CDF}}\big(\widehat{U}_n, \mathrm{Unif}[0,1]\big) = \max\left\{\underbrace{\sup_{t \in [0,1]} \left| \sqrt{1 - \widehat{U}_n(t)} - \sqrt{1 - t} \right|}_{=\Delta_0} \right., \underbrace{\sup_{t \in [0,1]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right|}_{=:\Delta_1}\right\},$$

by obseving that

$$\sup_{t \in [0,1]} \left| \sqrt{\frac{1}{n} \sum_{i=1}^{n} \mathbb{1} \{U_i < t\}} - \sqrt{t} \right| = \sup_{t \in [0,1]} \left| \sqrt{\frac{1}{n} \sum_{i=1}^{n} \mathbb{1} \{U_i \le t\}} - \sqrt{t} \right|$$

(i.e., the supremum is unchanged by replacing < with  $\le$ ). We can further write

$$\Delta_1 = \max \left\{ \underbrace{\sup_{t \in [0, \frac{\log n}{n}]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right|}_{=:\Delta_{1,0}} \right., \quad \underbrace{\sup_{t \in \left[\frac{\log n}{n}, 1 - \frac{\log n}{n}\right]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right|}_{=:\Delta_{1,1}} \right., \quad \underbrace{\sup_{t \in \left[1 - \frac{\log n}{n}, 1\right]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right|}_{=:\Delta_{1,2}} \right\}.$$

We have

$$\Delta_{1,0} = \sup_{t \in [0, \frac{\log n}{n}]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right| \le \sqrt{\frac{\log n}{n}} + \sqrt{\widehat{U}_n\left(\frac{\log n}{n}\right)} \le 2\sqrt{\frac{\log n}{n}} + \Delta_{1,1},$$

and

$$\Delta_{1,2} = \sup_{t \in [1 - \frac{\log n}{n}, 1]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right| \le \sqrt{\frac{\log n}{n}} + \left(1 - \sqrt{\widehat{U}_n \left(1 - \frac{\log n}{n}\right)}\right) \le 2\sqrt{\frac{\log n}{n}} + \Delta_{1,1}.$$

Furthermore,

$$\Delta_{1,1} = \sup_{t \in [\frac{\log n}{n}, 1 - \frac{\log n}{n}]} \left| \sqrt{\widehat{U}_n(t)} - \sqrt{t} \right| = \sup_{t \in [\frac{\log n}{n}, 1 - \frac{\log n}{n}]} \frac{\left| \widehat{U}_n(t) - t \right|}{\sqrt{\widehat{U}_n(t)} + \sqrt{t}} \le \sup_{t \in [\frac{\log n}{n}, 1 - \frac{\log n}{n}]} \frac{\left| \widehat{U}_n(t) - t \right|}{\sqrt{t}}.$$

Combining these calculations, we have

$$\Delta_1 \le 2\sqrt{\frac{\log n}{n}} + \sup_{t \in \lceil \frac{\log n}{n}, 1 - \frac{\log n}{n} \rceil} \frac{|\widehat{U}_n(t) - t|}{\sqrt{t}}.$$

Similarly we can calculate

$$\Delta_0 \le 2\sqrt{\frac{\log n}{n}} + \sup_{t \in \left[\frac{\log n}{2}, 1 - \frac{\log n}{2}\right]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{1 - t}},$$

and so we have

$$\begin{split} \Delta_{\text{CDF}}\big(\widehat{U}_n, \text{Unif}[0, 1]\big) &\leq 2\sqrt{\frac{\log n}{n}} + \sup_{t \in [\frac{\log n}{n}, 1 - \frac{\log n}{n}]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{\min\{t, 1 - t\}}} \\ &= 2\sqrt{\frac{\log n}{n}} + \max\left\{\sup_{t \in [\frac{\log n}{n}, \frac{1}{2}]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{t}}, \sup_{t \in [\frac{1}{2}, 1 - \frac{\log n}{n}]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{1 - t}}\right\}. \end{split}$$

Next, Shorack and Wellner (2009, Proposition 11.1.1 (part (10)) + Inequality 11.2.1) (applied with  $q(t) = \sqrt{t}$ , with  $a = \frac{\log n}{n}$ , and with  $b = \delta = \frac{1}{2}$ ) establishes that, for any  $\lambda > 0$ ,

$$\mathbb{P}\left\{\sup_{t\in [\frac{\log n}{n},\frac{1}{2}]}\frac{|\widehat{U}_n(t)-t|}{\sqrt{t}}\geq \frac{\lambda}{\sqrt{n}}\right\}\leq 12\int_{\frac{\log n}{n}}^{1/2}\frac{1}{t}\cdot \exp\left\{-\frac{\lambda^2}{8\left(1+\frac{\lambda}{3\sqrt{\log n}}\right)}\right\}\;\mathrm{d}t,$$

as long as n satisfies  $\frac{\log n}{n} \leq \frac{1}{4}$  (which holds for n > 8; for  $n \leq 8$ , by taking  $c' \geq 2$  we can ensure that the lemma's claim is trivial, since  $\Delta_{\text{CDF}}(\widehat{U}_n, \text{Unif}[0, 1]) \leq 1$  deterministically). Furthermore, clearly we see that  $\sup_{t \in [\frac{\log n}{n}, \frac{1}{2}]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{t}}$  and  $\sup_{t \in [\frac{1}{2}, 1 - \frac{\log n}{n}]} \frac{|\widehat{U}_n(t) - t|}{\sqrt{1 - t}}$  are equal in distribution. Therefore, we have

$$\mathbb{P}\left\{\Delta_{\mathrm{CDF}}\left(\widehat{U}_n, \mathrm{Unif}[0, 1]\right) \ge 2\sqrt{\frac{\log n}{n}} + \frac{\lambda}{\sqrt{n}}\right\} \le 24\log\left(\frac{n}{2\log n}\right) \cdot \exp\left\{-\frac{\lambda^2}{8\left(1 + \frac{\lambda}{3\sqrt{\log n}}\right)}\right\}$$

for any  $\lambda > 0$ . Taking  $\lambda = 5(c+2)\sqrt{\log n}$ , we can calculate  $\exp\left\{-\frac{\lambda^2}{8\left(1+\frac{\lambda}{3\sqrt{\log n}}\right)}\right\} \le \exp\left\{-(c+2)\log n\right\} = n^{-(c+2)}$ , and so we have

$$\mathbb{P}\left\{\Delta_{\mathrm{CDF}}\left(\widehat{U}_n, \mathrm{Unif}[0, 1]\right) \ge 2\sqrt{\frac{\log n}{n}} + 5(c+2)\sqrt{\frac{\log n}{n}}\right\} \le 24\log\left(\frac{n}{2\log n}\right) \cdot n^{-(c+2)} \le n^{-c}$$

where the last step holds since we have assumed that n > 8. This proves the lemma with c' = 5c + 12.

#### A.6 Proof of Lemma 14

We have

$$\begin{split} \epsilon_{P} &= \mathbb{E}_{P} \left[ |X - \mu_{P}| \right] \\ &= 2\mathbb{E}_{P} \left[ (X - \mu_{P})_{+} \right] \\ &\leq 2\mathbb{E}_{P} \left[ |X - \mu_{P}| \cdot \mathbb{1} \left\{ X > \mu_{P} \right\} \right] \\ &\leq 2\mathbb{E}_{P} \left[ |X - \mu_{P}|^{q} \right]^{1/q} \mathbb{E}_{P} \left[ \mathbb{1} \left\{ X > \mu_{P} \right\}^{\frac{q}{q-1}} \right]^{\frac{q-1}{q}} \\ &\leq 2(\mathbb{E}_{P} \left[ |X|^{q} \right]^{1/q} + (|\mu_{P}|^{q})^{1/q}) \cdot \mathbb{P}_{P} \left\{ X > \mu_{P} \right\}^{\frac{q-1}{q}} \\ &\leq 4M_{q} \cdot \mathbb{P}_{P} \left\{ X > \mu_{P} \right\}^{\frac{q-1}{q}}. \end{split}$$

Therefore,

$$\mathbb{P}_P\{X > \mu_P\} \ge \left(\frac{\epsilon_P}{4M_a}\right)^{\frac{q}{q-1}}.$$

Similarly, the same bound holds for  $\mathbb{P}_P\{X < \mu_P\}$ .

## A.7 Proofs of lower bounds (Theorems 4 and 6)

We begin with some preliminary calculations that we will use for the constructions for both theorems. Fix any  $0 < \rho_0 < \rho_1$  and any  $\beta \in (0, \rho_0/\rho_1]$ . Let P be the mixture distribution drawing

$$X \sim \begin{cases} \operatorname{Unif}(\mathbb{S}_{d-1}(\rho_0)), & \text{with probability } 1 - \beta, \\ \operatorname{Unif}(\mathbb{S}_{d-1}(\rho_1)), & \text{with probability } \beta. \end{cases}$$
(19)

Defining

$$s_d = \mathbb{E}\left[|V_1|\right] \text{ for } V = (V_1, \dots, V_d) \sim \text{Unif}(\mathbb{S}_{d-1}),$$
 (20)

we can calculate

$$\epsilon_P = (1 - \beta)\rho_0 \cdot s_d + \beta\rho_1 \cdot s_d \ge s_d\rho_0.$$

We will apply Lemma 18 to this distribution P and the log-density  $\phi = \phi^*(P)$  of its log-concave projection. Observe that  $\phi$  is spherically symmetric around 0, and is constant over  $||x|| \leq \rho_0$ —in particular, this means that  $\phi(x) = M_{\phi}$  for all  $||x|| \leq \rho_0$ , where  $M_{\phi} = \sup_{x \in \mathbb{R}^d} \phi(x)$  as before. Next, let  $t_* \geq 0$  be the value of  $M_{\phi} - \phi(x)$  for points x with  $||x|| = \rho_1$  (since  $\phi$  is spherically symmetric, this is well defined). We now split into cases. If  $t_* \geq \frac{8d\rho_1}{r_d s_d \rho_0}$ , then applying Lemma 18 with  $R = \rho_1/2$ ,  $x_{\phi} = 0$ , and  $t = t_*$ , we obtain

$$\beta \le \mathbb{P}_P\{\phi(X) \le M_\phi - t_* \text{ and } ||X|| \le \rho_1\} \le \frac{16d}{b_d r_d s_d \rho_0} \cdot \frac{\rho_1}{t_*^2},$$

which proves that

$$t_* \le \sqrt{\frac{16d}{b_d r_d s_d} \cdot \frac{\rho_1}{\rho_0 \beta}}.$$

If this case does not hold, then we instead have  $t_* < \frac{8d\rho_1}{r_d s_d \rho_0}$ , so combining the two cases,

$$t_* \le \max\left\{\sqrt{\frac{16d}{b_d r_d s_d} \cdot \frac{\rho_1}{\rho_0 \beta}}, \frac{8d}{r_d s_d} \cdot \frac{\rho_1}{\rho_0}\right\} \le \max\left\{\sqrt{\frac{16d}{b_d r_d s_d}}, \frac{8d}{r_d s_d}\right\} \cdot \sqrt{\frac{\rho_1}{\rho_0 \beta}},$$

where the last step comes from our assumption on  $\beta$ . Therefore

$$\phi(x) \ge \phi(0) - \max\left\{\sqrt{\frac{16d}{b_d r_d s_d}}, \frac{8d}{r_d s_d}\right\} \cdot \sqrt{\frac{\rho_1}{\rho_0 \beta}}$$

for  $||x|| = \rho_1$  while

$$\phi(x) = \phi(0)$$

for  $||x|| \leq \rho_0$ . By concavity of  $\phi$ , then,

$$\phi(x) \ge \phi(0) - \max\left\{\sqrt{\frac{16d}{b_d r_d s_d}}, \frac{8d}{r_d s_d}\right\}$$

for all x with  $||x|| \leq \rho_0 + (\rho_1 - \rho_0) \cdot \sqrt{\frac{\rho_0 \beta}{\rho_1}}$ . Therefore, for any density f supported on  $\mathbb{B}_d(\rho_0)$ , it holds that

$$d_{H}^{2}(f, \psi^{*}(P)) \geq \int_{\mathbb{R}^{d}} e^{\phi(0) - \max\left\{\sqrt{\frac{16d}{b_{d}r_{d}s_{d}}}, \frac{8d}{r_{d}s_{d}}\right\}} \cdot \mathbb{I}\left\{\rho_{0} < \|x\| < \rho_{0} + (\rho_{1} - \rho_{0}) \cdot \sqrt{\frac{\rho_{0}\beta}{\rho_{1}}}\right\} dx$$

$$= e^{\phi(0) - \max\left\{\sqrt{\frac{16d}{b_{d}r_{d}s_{d}}}, \frac{8d}{r_{d}s_{d}}\right\}} \cdot \operatorname{Leb}_{d}\left(\mathbb{B}_{d}(\rho_{0} + (\rho_{1} - \rho_{0}) \cdot \sqrt{\rho_{0}\beta/\rho_{1}}) \backslash \mathbb{B}_{d}(\rho_{0})\right)$$

$$\geq e^{\phi(0) - \max\left\{\sqrt{\frac{16d}{b_{d}r_{d}s_{d}}}, \frac{8d}{r_{d}s_{d}}\right\}} \cdot \rho_{0}^{d-1} \cdot (\rho_{1} - \rho_{0}) \cdot \sqrt{\frac{\rho_{0}\beta}{\rho_{1}}} \cdot S_{d-1},$$

where as before  $S_{d-1}$  denotes the surface area of  $\mathbb{S}_{d-1}$ . Finally, we need to place a lower bound on  $\phi(0)$ . By Corollary 8, we know that the covariance matrix  $\Sigma$  of the distribution with log-density  $\phi$  has operator norm bounded as

$$\|\Sigma\|_{\mathrm{op}} \le 16\big((1-\beta)\rho_0 + \beta\rho_1\big)^2.$$

Furthermore,  $\tilde{\phi}(x) = \frac{1}{2} \log \det(\Sigma) + \phi(\Sigma^{1/2}x)$  is an isotropic concave log-density, and so  $\tilde{\phi}(0) \geq c'_d$  where  $c'_d > 0$  depends only on d, by Lovász and Vempala (2007, Theorem 5.14(d)). Therefore,

$$\phi(0) \ge c'_d - \frac{d}{2}\log(16) - d\log((1-\beta)\rho_0 + \beta\rho_1).$$

We conclude that

$$d_{H}^{2}(f, \psi^{*}(P)) \geq c_{d}'' \cdot \rho_{0}^{d-1} \cdot (\rho_{1} - \rho_{0}) \cdot \sqrt{\frac{\rho_{0}\beta}{\rho_{1}}} \cdot ((1 - \beta)\rho_{0} + \beta\rho_{1})^{-d}, \tag{21}$$

where  $c''_d$  depends only on d.

### A.7.1 Completing the proof of Theorem 4

To prove Theorem 4, let P be the distribution constructed in (19) with

$$\rho_0 = \epsilon/s_d, \ \rho_1 = 2\epsilon/s_d, \ \beta = \min\left\{\frac{s_d\delta}{\epsilon}, \frac{1}{2}\right\},$$

where  $s_d$  is defined as in (20). Let

$$Q = \mathrm{Unif}(\mathbb{S}_{d-1}(\rho_0)).$$

Clearly  $\epsilon_P \geq \epsilon_Q = s_d \rho_0 = \epsilon$ , and  $d_W(P,Q) = \beta(\rho_1 - \rho_0) \leq \delta$ , thus satisfying the conditions of the theorem. Since Q is supported on  $\mathbb{B}_d(\rho_0)$ ,  $\psi^*(Q)$  is also supported on this ball. Then applying our calculation (21), and plugging in our choices of  $\rho_0, \rho_1, \beta$ , after simplifying we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(P), \psi^{*}(Q)) \geq c_{d}' \cdot \frac{2^{d}}{3^{d}} \cdot \sqrt{\min\left\{\frac{s_{d}\delta}{2\epsilon}, \frac{1}{4}\right\}}.$$

This completes the proof of the theorem, when  $c_d$  is chosen appropriately.

#### A.7.2 Completing the proof of Theorem 6

The first term in the lower bound, i.e.,  $\sup_{P\in\mathcal{P}_d:\mathbb{E}_P[\|X\|^q]\leq 1, \epsilon_P\geq \epsilon_d^*}\mathbb{E}\left[\mathrm{d}^2_\mathrm{H}\left(\psi^*(\widehat{P}_n),\psi^*(P)\right)\right]\geq c_d n^{-\frac{2}{d+1}}$ , holds by Kim and Samworth (2016, Theorem 1), which establishes this as the minimax rate (for  $d\geq 2$ ) over distributions P that are log-concave (we can verify that the distribution P constructed in their proof satisfies the conditions  $\mathbb{E}_P\left[\|X\|^q\right]\leq 1, \ \epsilon_P\geq \epsilon_d^*$ , for appropriately chosen  $\epsilon_d^*$ ). If instead d=1, then the first term cannot be the minimum.

Next, to prove the second term in the lower bound, we consider a mixture model. Let P be the distribution constructed in (19) with

$$\rho_0 = \frac{1}{2}, \ \rho_1 = n^{1/q}, \ \beta = \frac{1}{2n}.$$

Then clearly  $\mathbb{E}_P[\|X\|^q] \leq 1$ , and  $\epsilon_P \geq \frac{1}{2}s_d$ , so  $\epsilon_P \geq \epsilon_d^*$  for an appropriately chosen  $\epsilon_d^*$ . Now, with probability at least 1/2, the observations  $X_1, \ldots, X_n$  are all drawn from the first component of the mixture model, i.e.,  $\psi^*(\widehat{P}_n)$  is supported on  $\mathbb{B}_d(1/2)$ . On this event, applying (21) and plugging in our choices of  $\rho_0, \rho_1, \beta$ , after simplifying we have

$$d_{\mathrm{H}}^{2}(\psi^{*}(\widehat{P}_{n}), \psi^{*}(P)) \ge c_{d}^{\prime\prime\prime} \cdot n^{-\frac{1}{2} + \frac{1}{2q}},$$

where  $c_d'''$  depends only on d. This establishes the second term in the lower bound claimed in Theorem 6, and thus completes the proof of the theorem.

# Acknowledgements

The authors thank the anonymous reviewers and Oliver Feng for helpful comments. R.F.B. was supported by the National Science Foundation via grant DMS-1654076 and by an Alfred P. Sloan fellowship. R.J.S. was supported by EPSRC grants EP/P031447/1 and EP/N031938/1.

## References

CD Aliprantis and KC Border. *Infinite Dimensional Analysis: A Hitchhiker's Guide*. Springer, 2006.

Pierre C Bellec. Sharp oracle inequalities for least squares estimators in shape restricted regression. *The Annals of Statistics*, 46(2):745–780, 2018.

Lucien Birgé. The Grenander estimator: A nonasymptotic approach. The Annals of Statistics, 17(4):1532–1549, 1989.

HD Brunk, Richard E Barlow, Daniel J Bartholomew, and James M Bremner. Statistical Inference under Order Restrictions: The Theory and Application of Isotonic Regression. John Wiley & Sons, 1972.

T Tony Cai and Mark G Low. A framework for estimation of convex functions. *Statistica Sinica*, pages 423–456, 2015.

- Timothy Carpenter, Ilias Diakonikolas, Anastasios Sidiropoulos, and Alistair Stewart. Near-optimal sample complexity bounds for maximum likelihood estimation of multivariate log-concave densities. In *COLT 2018*, 2018.
- Sabyasachi Chatterjee, Adityanand Guntuboyina, and Bodhisattva Sen. On risk bounds in isotonic and other shape restricted regression problems. *The Annals of Statistics*, 43(4): 1774–1800, 2015.
- Yining Chen and Richard J Samworth. Smoothed log-concave maximum likelihood estimation with applications. *Statistica Sinica*, pages 1373–1398, 2013.
- Thomas M Cover and Joy A Thomas. *Elements of Information Theory*. Wiley, New York, 1991.
- Madeleine Cule and Richard Samworth. Theoretical properties of the log-concave maximum likelihood estimator of a multidimensional density. *Electronic Journal of Statistics*, 4: 254–270, 2010.
- Madeleine Cule, Richard Samworth, and Michael Stewart. Maximum likelihood estimation of a multi-dimensional log-concave density. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(5):545–607, 2010.
- Charles R Doss and Jon A Wellner. Global rates of convergence of the mles of log-concave and s-concave densities. The Annals of Statistics, 44(3):954, 2016.
- Richard M Dudley. Real Analysis and Probability. Cambridge University Press, Cambridge, 2002.
- Lutz Dümbgen and Kaspar Rufibach. Maximum likelihood estimation of a log-concave density and its distribution function: Basic properties and uniform consistency. *Bernoulli*, 15(1):40–68, 2009.
- Lutz Dümbgen, Richard Samworth, and Dominic Schuhmacher. Approximation by log-concave distributions, with applications to regression. *The Annals of Statistics*, 39(2): 702–730, 2011.
- Cécile Durot and Hendrik P Lopuhaä. Limit theory in monotone function estimation. Statistical Science, 33(4):547–567, 2018.
- Billy Fang and Adityanand Guntuboyina. On the risk of convex-constrained least squares estimators under misspecification. *Bernoulli*, 25(3):2206–2244, 2019.
- Oliver Y Feng, Adityanand Guntuboyina, Arlene KH Kim, and Richard J Samworth. Adaptation in multivariate log-concave density estimation. *The Annals of Statistics, to appear*, 2020.
- Daniel Fresen. A multivariate Gnedenko law of large numbers. *The Annals of Probability*, 41(5):3051–3080, 2013.

- Ulf Grenander. On the theory of mortality measurement: part ii. Scandinavian Actuarial Journal, 1956(2):125–153, 1956.
- Piet Groeneboom. Estimating a monotone density. In *Proceedings of the Berkeley Conference* in *Honor of Jerzy Neyman and Jack Kiefer*. Wadsworth, Monterey, California, 1985.
- Piet Groeneboom and Geurt Jongbloed. *Nonparametric Estimation under Shape Constraints*, volume 38. Cambridge University Press, Cambridge., 2014.
- Adityanand Guntuboyina and Bodhisattva Sen. Global risk bounds and adaptation in univariate convex regression. *Probability Theory and Related Fields*, 163(1-2):379–411, 2015.
- Qiyang Han. Global empirical risk minimizers with "shape constraints" are rate optimal in general dimensions. arXiv preprint arXiv:1905.12823, 2019.
- Qiyang Han and Jon A Wellner. Approximation and estimation of s-concave densities via Rényi divergences. Annals of Statistics, 44(3):1332, 2016a.
- Qiyang Han and Jon A Wellner. Multivariate convex regression: global risk bounds and adaptation. arXiv preprint arXiv:1601.06844, 2016b.
- Qiyang Han, Tengyao Wang, Sabyasachi Chatterjee, and Richard J Samworth. Isotonic regression in general dimensions. *The Annals of Statistics*, 47(5):2440–2471, 2019.
- Clifford Hildreth. Point estimates of ordinates of concave functions. *Journal of the American Statistical Association*, 49(267):598–619, 1954.
- Hanna Jankowski. Convergence of linear functionals of the Grenander estimator under misspecification. *The Annals of Statistics*, 42(2):625–653, 2014.
- Arlene KH Kim and Richard J Samworth. Global rates of convergence in log-concave density estimation. *The Annals of Statistics*, 44(6):2756–2779, 2016.
- Arlene KH Kim, Adityanand Guntuboyina, and Richard J Samworth. Adaptation in log-concave density estimation. *The Annals of Statistics*, 46(5):2279–2306, 2018.
- Roger Koenker and Ivan Mizera. Quasi-concave density estimation. *The Annals of Statistics*, 38(5):2998–3027, 2010.
- Gil Kur, Yuval Dagan, and Alexander Rakhlin. The log-concave maximum likelihood estimator is optimal in high dimensions. arXiv preprint arXiv:1903.05315v3, 2019.
- Jing Lei. Convergence and concentration of empirical measures under Wasserstein distance in unbounded functional spaces. *Bernoulli*, 26(1):767–798, 2020.
- László Lovász and Santosh Vempala. The geometry of logconcave functions and sampling algorithms. Random Structures & Algorithms, 30(3):307–358, 2007.
- Jayanta Kumar Pal, Michael Woodroofe, and Mary Meyer. Estimating a Polya frequency function. *Lecture Notes-Monograph Series*, pages 239–249, 2007.

- Valentin Patilea. Convex models, mle and misspecification. *The Annals of Statistics*, 29: 94–123, 2001.
- BLS Prakasa Rao. Estimation of a unimodal density. Sankhyā: The Indian Journal of Statistics, Series A, pages 23–36, 1969.
- Richard J Samworth. Recent progress in log-concave density estimation. *Statistical Science*, 33(4):493–509, 2018.
- Richard J Samworth and Bodhisattva Sen. Editorial: Special Issue on "Nonparametric Inference under Shape Constraints". *Statistical Science*, 2018.
- Richard J Samworth and Ming Yuan. Independent component analysis via nonparametric maximum likelihood estimation. *The Annals of Statistics*, 40(6):2973–3002, 2012.
- Rolf Schneider. Convex Bodies: the Brunn-Minkowski theory. Cambridge University Press, 2014.
- Dominic Schuhmacher, André Hüsler, and Lutz Dümbgen. Multivariate log-concave distributions as a nearly parametric model. Statistics & Risk Modeling with Applications in Finance and Insurance, 28(3):277–295, 2011.
- Emilio Seijo and Bodhisattva Sen. Nonparametric least squares estimation of a multivariate convex regression function. *The Annals of Statistics*, 39(3):1633–1657, 2011.
- Arseni Seregin and Jon A Wellner. Nonparametric estimation of multivariate convex-transformed densities. *The Annals of Statistics*, 38(6):3751–3781, 2010.
- Robert J Serfling. Approximation Theorems of Mathematical Statistics. John Wiley & Sons, New York, 1980.
- Galen R Shorack and Jon A Wellner. *Empirical processes with applications to statistics*. SIAM, 2009.
- Sara A van de Geer. Empirical Processes in M-Estimation. Cambridge University Press, Cambridge, 2000.
- Aad W van der Vaart and Jon A Wellner. Weak Convergence and Empirical Processes. Springer, 1996.
- Cédric Villani. Optimal Transport: Old and New, volume 338. Springer Science & Business Media, 2008.
- Richard von Mises. On the asymptotic distribution of differentiable statistical functions. *The Annals of Mathematical Statistics*, 18(3):309–348, 1947.
- Guenther Walther. Detecting the presence of mixing with multiscale maximum likelihood. Journal of the American Statistical Association, 97(458):508–513, 2002.

- Guenther Walther. Inference and modeling with log-concave distributions. *Statistical Science*, 24(3):319–327, 2009.
- Min Xu and Richard J Samworth. High-dimensional nonparametric density estimation via symmetry and shape constraints. *The Annals of Statistics, to appear*, 2020.
- Fan Yang and Rina Foygel Barber. Contraction and uniform convergence of isotonic regression. *Electronic Journal of Statistics*, 13(1):646–677, 2019.
- Cun-Hui Zhang. Risk bounds in isotonic regression. The Annals of Statistics, 30(2):528–555, 2002.