



Sharp Inequalities between Total Variation and Hellinger Distances for Gaussian Mixtures

Joonhyuk Jung¹ Chao Gao¹

¹University of Chicago



Abstract

We study the relation between the total variation (TV) and Hellinger distances between two Gaussian location mixtures. Our first result establishes a general upper bound: for any two mixing distributions supported on a compact set, the Hellinger distance between the two mixtures is controlled by the TV distance raised to a power $1 - o(1)$, where the $o(1)$ term is of order $1/\log \log(1/\text{TV})$. We also construct two sequences of mixing distributions that demonstrate the sharpness of this bound. Taken together, our results resolve an open problem raised in Jia et al. [4] and thus lead to an entropic characterization of learning Gaussian mixtures in total variation. Our inequality also yields optimal robust estimation of Gaussian mixtures in Hellinger distance, which has a direct implication for bounding the minimax regret of empirical Bayes under Huber contamination.

Introduction

Given a probability measure π supported on \mathbb{R}^d , we define the Gaussian mixture density by

$$f_\pi(x) := \int_{\mathbb{R}^d} \phi_d(x - \theta) d\pi(\theta),$$

where ϕ_d is the d -dimensional standard Gaussian density function. We study the relation between the total variation distance $\text{TV}(p, q) := \frac{1}{2} \int |p - q|$ and the Hellinger distance $H(p, q) := \sqrt{\frac{1}{2} \int (\sqrt{p} - \sqrt{q})^2}$ of two Gaussian mixture densities.

For mixing distributions π and η supported on a bounded Euclidean ball $\{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq M\}$, it was proved by Jia et al. [4] that $H^2(f_\pi, f_\eta) \asymp \text{KL}(f_\pi \| f_\eta)$ holds up to constant factors depending on M and d . However, whether

$$\text{TV}(f_\pi, f_\eta) \asymp H(f_\pi, f_\eta)$$

holds was explicitly listed as an open question in the paper. We disprove the linear comparability by showing

$$H(f_\pi, f_\eta) \leq \text{TV}^{1-o(1)}(f_\pi, f_\eta), \quad (1)$$

where the $o(1)$ exponent is of order

$$\frac{\Theta(1)}{\log \log(1/\text{TV}(f_\pi, f_\eta))},$$

as stated in Theorem 1. In addition, Theorem 2 shows that the $o(1)$ term is indeed necessary.

Main Result

Theorem 1 (Inequality between TV distance and χ^2 -divergence)

Let π and η be probability measures supported on the d -dimensional cube $[-M, M]^d$. Let $\delta > 0$. Then, there exists $C_0 = C_0(\delta, M, d) > 0$, not depending on π or η , such that

$$\sqrt{\chi^2(f_\pi \| f_\eta)} \leq \left(C_0 \vee \text{TV}^{-\alpha(\text{TV}(f_\pi, f_\eta))} \right) \text{TV}(f_\pi, f_\eta),$$

where we define

$$\alpha(t) := \frac{2 + \delta}{\log(\log(1/t) \vee e)}, \quad t > 0.$$

- The Theorem 1 immediately implies the upper bound (1) on the Hellinger distance as $H^2(p, q) \leq \chi^2(p \| q)$ holds in general.
- A key step in establishing the above inequality is to relate the $L^1(\phi_d)$ and $L^2(\phi_d)$ norms of the ratio

$$g := \frac{f_\pi - f_\eta}{\phi_d}.$$

Indeed, the $L^1(\phi_d)$ -norm of g is exactly twice the total variation distance, while both the squared Hellinger distance and the χ^2 -divergence are closely related to the squared $L^2(\phi_d)$ -norm.

Sharpness

Theorem 2 (Sharpness)

There exist two sequences of probability measures $\{\pi_n\}$ and $\{\eta_n\}$ supported on $[-M, M]$ such that, if we define

$$\text{TV}_n := \text{TV}(f_{\pi_n}, f_{\eta_n}), \quad H_n := H(f_{\pi_n}, f_{\eta_n}),$$

then $\text{TV}_n \downarrow 0$ as $n \rightarrow \infty$, and moreover it holds for all n that

$$H_n \geq \text{TV}_n^{1-\alpha^*(\text{TV}_n)}, \quad (2)$$

where we define

$$\alpha^*(t) := \frac{0.33}{\log \log(1/t)}, \quad t > 0.$$

- We can also construct d -dimensional probability measures π_n and η_n satisfying (2) because we have $\text{TV}(f_\pi, f_\eta) = \text{TV}(f_{\pi^*}, f_{\eta^*})$ and $H(f_\pi, f_\eta) = H(f_{\pi^*}, f_{\eta^*})$ for

$$\pi = \pi^* \otimes \delta_0 \otimes \cdots \otimes \delta_0, \quad \eta = \eta^* \otimes \delta_0 \otimes \cdots \otimes \delta_0,$$

where δ_0 denotes the point mass at zero and \otimes the product measure.

Proof Sketch for the Main Result

Here we give a sketch of the proof for the one-dimensional setting with $d = 1$. Consider the Hermite decomposition of $g := \frac{f_\pi - f_\eta}{\phi_1}$ in $L^2(\phi_1)$. Concretely, write $g = q + r$, where

$$q = \sum_{k=0}^n \frac{\Delta_k}{\sqrt{k!}} h_k, \quad r = \sum_{k=n+1}^{\infty} \frac{\Delta_k}{\sqrt{k!}} h_k, \quad \Delta_k = \int_{\mathbb{R}} \theta^k d(\pi - \eta)(\theta), \quad h_k(x) = \frac{(-1)^k}{\sqrt{k! \phi_1(x)}} \frac{d^k}{dx^k} \phi_1(x),$$

and n is an integer to be determined later. To control the $L^1(\phi_1)$ -norm of q , we define

$$c_n := \inf \left\{ \|P\|_{L^1(\phi_1)} : P \in \Pi_n, \|P\|_{L^2(\phi_1)} = 1 \right\}. \quad (3)$$

(Let Π_n be the set of real polynomials of degree $\leq n$ so that $q \in \Pi_n$.) The Nikolskii-type inequality [6] and the restricted-range inequality (Theorem 6.2(b) of Lubinsky [5]) show that $c_n \geq cn^{-1/4}e^{-n}$ holds for some universal constant $c > 0$. In addition to c_n , another technical ingredient is to control the tail norm $\|r\|_{L^2(\phi_1)}$. Compact support implies $|\Delta_k| \leq 2(2M)^k$ and

$$\|r\|_{L^2(\phi_1)} \leq \left(\sum_{k=n+1}^{\infty} \frac{4(4M^2)^k}{k!} \right)^{1/2} \leq \left(\frac{C}{n+1} \right)^{(n+1)/2},$$

where $C > 0$ is a constant depending solely on M . Now we bound $\|g\|_{L^1(\phi_1)}$ from below:

$$\begin{aligned} \|g\|_{L^1(\phi_1)} &\geq \|q\|_{L^1(\phi_1)} - \|r\|_{L^1(\phi_1)} \\ &\geq c_n \|q\|_{L^2(\phi_1)} - \|r\|_{L^2(\phi_1)} \\ &\geq c_n \|g\|_{L^2(\phi_1)} - 2 \|r\|_{L^2(\phi_1)}, \end{aligned} \quad (\text{by (3)})$$

where the last inequality uses $c_n \leq 1$ and the decomposition $g = q + r$. Together with the lower bound on c_n and the upper bound on $\|r\|_{L^2(\phi_1)}$, we obtain

$$\|g\|_{L^1(\phi_1)} \geq \sup_{n \geq 1} \left\{ cn^{-1/4}e^{-n} \|g\|_{L^2(\phi_1)} - 2 \left(\frac{C}{n+1} \right)^{(n+1)/2} \right\}.$$

Finally, we choose

$$n \approx \frac{2 \log(1/\text{TV}(f_\pi, f_\eta))}{\log \log(1/\text{TV}(f_\pi, f_\eta))}$$

to conclude the proof. In our paper, we present multidimensional extensions of the Nikolskii-type and restricted-range inequalities. Building on these results, we provide the full proof of the theorem.

Application 1: Learning Gaussian Mixtures in TV

In this section, we consider the problem of estimating a Gaussian mixture $P \in \mathcal{P}_{M,d}$ based on i.i.d. samples drawn from P , where $\mathcal{P}_{M,d}$ denotes the collection of d -dimensional Gaussian mixtures where the mixing distributions are supported on the d -dimensional cube $[-M, M]^d$.

Theorem 3 (Learning Gaussian mixtures in TV distance)

Suppose \mathcal{P} is a compact subset of $\mathcal{P}_{M,d}$. Then, we have

$$\epsilon_n^{2\left(1 + \frac{\Theta(1)}{\log(\log(1/\epsilon_n) \vee e)}\right)} \lesssim \inf_{\hat{P}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\text{TV}^2(P, \hat{P}) \right] \lesssim \epsilon_n^2,$$

where

$$\inf_{\hat{P}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[H^2(P, \hat{P}) \right] \asymp \epsilon_n^2 \asymp \inf_{\epsilon > 0} \left(\epsilon^2 + \frac{1}{n} \log N_{H,loc}(\mathcal{P}, \epsilon) \right),$$

and $N_{H,loc}(\mathcal{P}, \epsilon)$ is the local Hellinger covering number of \mathcal{P} .

Application 2: Robust Density Estimation in Hellinger

In this section, we consider the problem of estimating a Gaussian mixture with contaminated data,

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P := (1 - \epsilon)P_\pi + \epsilon Q, \quad (4)$$

where the distribution $P_\pi \in \mathcal{P}_{M,d}$ has density function f_π and Q is an arbitrary distribution of contamination. The data generating process in (4) is recognized as Huber's contamination model [3].

Theorem 4 (Robust density estimation in Hellinger distance)

Consider the data generating process in (4). Then, we have

$$\inf_{\hat{f}} \sup_{\pi, Q} \mathbb{E} \left[H^2(f_\pi, \hat{f}) \right] \asymp \epsilon^2 \left(1 - \frac{\Theta(1)}{\log(\log(1/\epsilon) \vee e)} \right), \quad (5)$$

provided that $n \geq \text{poly}(1/\epsilon)$, where the expectation is under (4) and the supremum is taken over all Q and π such that $\text{supp}(\pi) \subseteq [-M, M]^d$.

Theorem 5 (Robust regret bound)

Consider the data generating process in (4). Suppose $\hat{\theta}^*(\cdot)$ is the oracle Bayes estimator given by Tweedie's formula [2]. Then,

$$\inf_{\hat{\theta}} \sup_{\pi, Q} \mathbb{E} \left[\mathbb{E}_{X \sim f_\pi} \left\| \hat{\theta}(X) - \hat{\theta}^*(X) \right\|^2 \right] \lesssim \epsilon^2 \left(1 - \frac{\Theta(1)}{\log(\log(1/\epsilon) \vee e)} \right) + \frac{1}{n^{1-o(1)}},$$

where the outer expectation is under (4) and the supremum is taken over all Q and π such that $\text{supp}(\pi) \subseteq [-M, M]^d$.

Remarks

- Applying the same argument as in Chen et al. [1], the Theorem 2 (Sharpness) immediately implies the lower bound in (5).
- We also prove that Yatracos' estimator [7] attains the upper bound in (5).
- Dependency of the minimax rate (5) on n is a long-standing open question in the clean data setting ($\epsilon = 0$), and it remains open in the "small ϵ " regime.

References

- [1] Mengjie Chen, Chao Gao, and Zhao Ren. Robust covariance and scatter matrix estimation under Huber's contamination model. *The Annals of Statistics*, 46(5):1932–1960, 2018.
- [2] Bradley Efron. Tweedie's formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- [3] Peter J Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1):73–101, 1964.
- [4] Zeyu Jia, Yuri Polyanskiy, and Yihong Wu. Entropic characterization of optimal rates for learning Gaussian mixtures. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 4296–4335. PMLR, 2023.
- [5] Doron S Lubinsky. A survey of weighted approximation for exponential weights. *arXiv preprint math/0701099*, 2007.
- [6] Paul Nevai and Vilmos Totik. Sharp Nikolskii inequalities with exponential weights. *Analysis Mathematica*, 13(4):261–267, 1987.
- [7] Yannis G Yatracos. Rates of convergence of minimum distance estimators and Kolmogorov's entropy. *The Annals of Statistics*, 13(2): 768–774, 1985.