

Nonparametric Density Estimation under IPM Losses with Statistical Convergence Rates for Generative Adversarial Networks (GANs)

Ananya Uppal, Shashank Singh, Barnabás Póczos



↑Link to paper↑

Carnegie
Mellon
University

Introduction

- Nonparametric distribution estimation: Given n IID samples $X_{1:n} = X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$ from an unknown distribution P , we want to estimate P .
 - Important sub-routine of many statistical methods
 - Usually analyzed in terms of \mathcal{L}^2 loss
 - * Severe curse of dimensionality
- We provide unified minimax-optimal estimation rates under large family of losses called Integral Probability Metrics (IPMs), for many function classes (Sobolev, Besov, RKHS).
 - Includes most common metrics on probability distributions
 - Implicitly used in Generative Adversarial Networks (GANs)
 - * Allows us to derive statistical guarantees for GANs
 - Reduced curse of dimensionality

Integral Probability Metrics (IPMs)

Definition 1 (IPM). Let \mathcal{P} be a class of probability distributions on a sample space \mathcal{X} , and \mathcal{F} a class of (bounded) functions on \mathcal{X} . Then, the metric $\rho_{\mathcal{F}} : \mathcal{P} \times \mathcal{P} \rightarrow [0, \infty]$ on \mathcal{P} is defined by

$$\rho_{\mathcal{F}}(P, Q) := \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X \sim P} [f(X)] - \mathbb{E}_{X \sim Q} [f(X)] \right|.$$

Definition 2 (Besov Ball). Let $\beta_{j,k}$ denote coefficients of a function f in a wavelet basis indexed by $j \in \mathbb{N}$, $k \in [2^j]$. For parameters $\sigma \geq 0$, $p, q \in [1, \infty]$, $f \in \mathcal{L}^2$ lies in the Besov ball $B_{p,q}^{\sigma}$ iff

$$\|f\|_{B_{p,q}^{\sigma}} := \left\| \left\{ 2^{j(\sigma + D(1/2 - 1/p))} \|\{\beta_{\lambda}\}_{\lambda \in \Lambda_j}\|_p \right\}_{j \in \mathbb{N}} \right\|_q \leq 1.$$

The parameter q affects convergence rates only by logarithmic factors, so we omit it in sequel.

Examples of IPMs

Distance	\mathcal{F}
\mathcal{L}^p (including Total Variation/ \mathcal{L}^1)	$B_{p,1}^0$, with $p' = \frac{p}{p-1}$
Wasserstein (“earth-mover”)	$B_{\infty,1}^1$ (1-Lipschitz class)
Kolmogorov-Smirnov	B_1^1 (total variation ≤ 1)
Max. mean discrepancy (MMD)	RKHS ball
GAN discriminator	parameterized by neural network

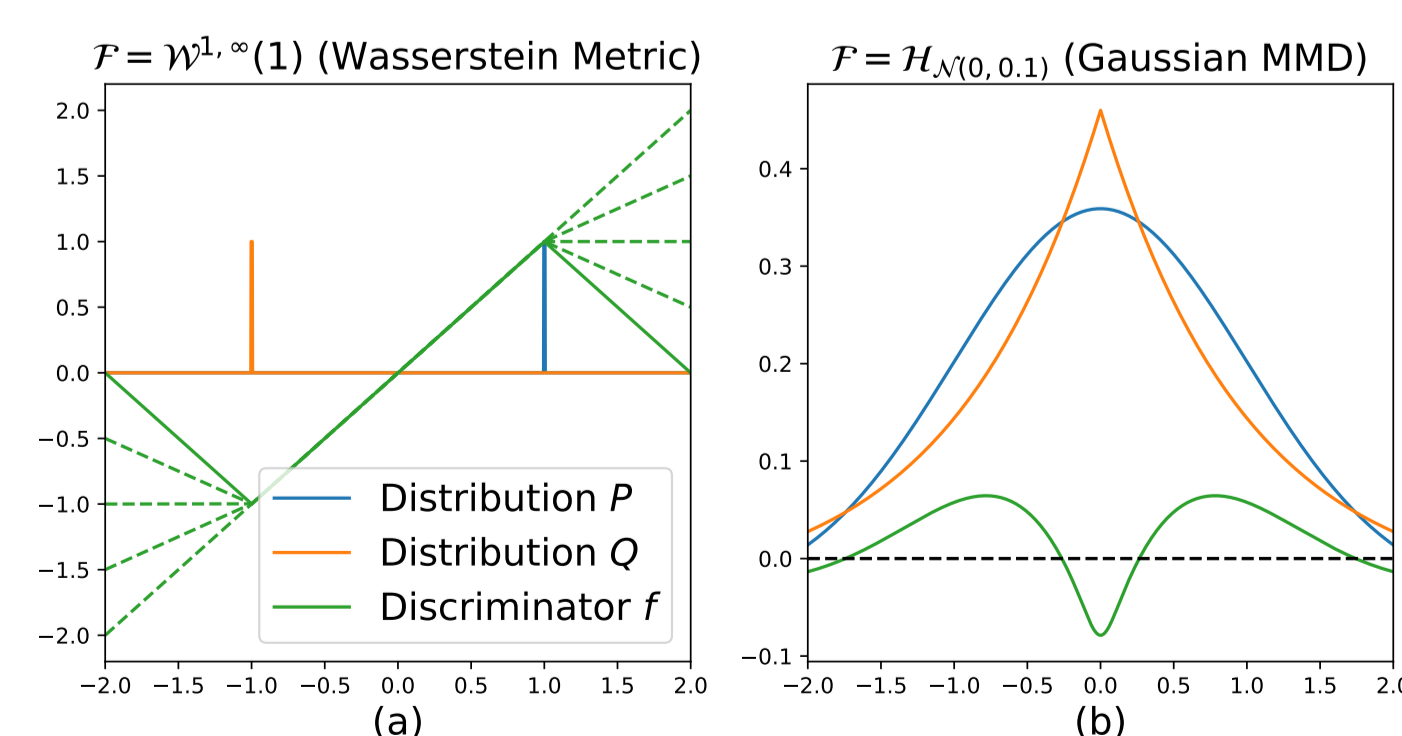


Figure 1: Examples of probability distributions P and Q and corresponding discriminator functions f^* . In (a), P and Q are Dirac masses at $+1$ and -1 , resp., and \mathcal{F} is the 1-Lipschitz class, so that $\rho_{\mathcal{F}}$ is the Wasserstein metric. In (b), P and Q are standard Gaussian and Laplace distributions, resp., and \mathcal{F} is a ball in an RKHS with a Gaussian kernel.

Minimax Rates for General Estimators

Theorem 1. Suppose $\sigma_g \geq D/p_g$, $p'_d > p_g$. Then, up to polylog factors in n ,

$$M(B_{p_g}^{\sigma_g}, B_{p'_d}^{\sigma_d}) := \inf_{\hat{P}} \sup_{p \in B_{p_g}^{\sigma_g}} \mathbb{E} \left[\rho_{B_{p'_d}^{\sigma_d}}(p, \hat{P}(X_{1:n})) \right] \asymp n^{-\frac{\sigma_g + \sigma_d}{2\sigma_g + D}} + n^{-\frac{\sigma_g + \sigma_d + D - D/p_g - D/p'_d}{2\sigma_g + D - 2D/p_g}} + n^{-\frac{1}{2}}.$$

Moreover, this rate is achieved by the wavelet thresholding estimator of Donoho et al. [1].

Minimax Rates for Linear Estimators

Definition 3 (Linear Estimator). A distribution estimate \hat{P} is said to be linear if there exist measures $T_i(X_i, \cdot)$ such that for all measurable A ,

$$\hat{P}(A) = \sum_{i=1}^n T_i(X_i, A).$$

Examples: empirical distribution, kernel density estimate, or orthogonal series estimate.

Theorem 2. Suppose $r > \sigma_g \geq D/p_g$. Then, up to polylog factors in n ,

$$M_{\text{lin}}(B_{p_g}^{\sigma_g}, B_{p'_d}^{\sigma_d}) := \inf_{\hat{P}} \sup_{p \in B_{p_g}^{\sigma_g}} \mathbb{E} \left[\rho_{B_{p'_d}^{\sigma_d}}(p, \hat{P}(X_{1:n})) \right] \asymp n^{-\frac{\sigma_g + \sigma_d}{2\sigma_g + D}} + n^{-\frac{\sigma_g + \sigma_d - D/p_g + D/p'_d}{2\sigma_g + D - 2D/p_g + 2D/p'_d}} + n^{-\frac{1}{2}},$$

where the inf is over all linear estimates of $p \in \mathcal{F}_g$, and μ_p is the distribution with density p .

Error Bounds for GANs

A natural statistical model for a perfectly optimized GAN as a distribution estimator is

$$\hat{P} := \operatorname{argmin}_{Q \in \mathcal{F}_g} \sup_{f \in \mathcal{F}_d} \mathbb{E} [f(X)] - \mathbb{E}_{X \sim \hat{P}_n} [f(X)], \quad (1)$$

where \mathcal{F}_d and \mathcal{F}_g are function classes parameterized by the discriminator and generator, resp [2].

Theorem 3 (Convergence Rate of a Regularized GAN). Fix a Besov density class $B_{p_g}^{\sigma_g}$ with $\sigma_g > D/p_g$ and discriminator class $B_{p'_d}^{\sigma_d}$. Then, for some constant $C > 0$ depending only on $B_{p'_d}^{\sigma_d}$ and $B_{p_g}^{\sigma_g}$, for any desired approximation error $\epsilon > 0$, one can construct a GAN \hat{P} of the form (1) (with \tilde{P}_n denoting the wavelet-thresholded distribution) whose discriminator network N_d and generator network N_g are fully-connected ReLU networks, such that

$$\sup_{P \in B_{p_g}^{\sigma_g}} \mathbb{E} \left[d_{B_{p'_d}^{\sigma_d}}(\hat{P}, P) \right] \lesssim \epsilon + n^{-\eta(D, \sigma_d, p_d, \sigma_g, p_g)},$$

where $\eta(D, \sigma_d, p_d, \sigma_g, p_g)$ is the optimal exponent in Theorem 1.

- N_d and N_g have (rate-optimal) depth $\text{polylog}(1/\epsilon)$ and width, max weight, and sparsity $\text{poly}(1/\epsilon)$.
- Proof uses recent fully-connected ReLU network for approximating Besov functions [3].

Example Phase Diagrams

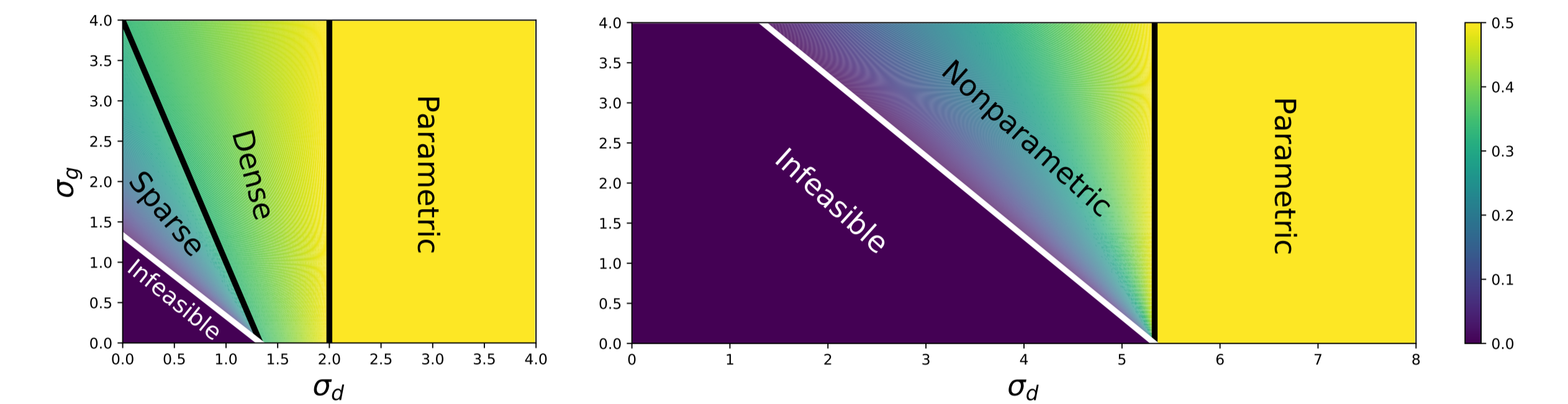


Figure 2: Minimax convergence rates as functions of discriminator smoothness σ_d and distribution function smoothness σ_g , in the case $D = 4$, $p_d = 1.2$, $p_g = 2$. Color shows exponent of minimax convergence rate (i.e., $\alpha(\sigma_d, \sigma_g)$) such that $M(B_{p'_d}^{\sigma_d}(\mathbb{R}^D), B_{p_g}^{\sigma_g}(\mathbb{R}^D)) \asymp n^{-\alpha(\sigma_d, \sigma_g)}$, ignoring polylogarithmic factors.

Applications/Examples

Example 1. (Total variation/Wasserstein-type losses) If, for some $\sigma_d > 0$, \mathcal{F} is a ball in $B_{\infty}^{\sigma_d}$, we obtain generalizations of total variation ($\sigma_d = 0$) and Wasserstein ($\sigma_d = 1$) losses. For these losses, we always have “Dense” rate

$$M(B_{p_g}^{\sigma_g}, B_{p'_d}^{\sigma_d}) \asymp n^{-\frac{\sigma_g + \sigma_d}{2\sigma_g + D}} + n^{-1/2}.$$

Example 2. (Kolmogorov-Smirnov-type losses) If, for $\sigma_d > 0$, \mathcal{F} is a ball in $B_1^{\sigma_d}$, we obtain generalizations of Kolmogorov-Smirnov loss ($\sigma_d = 0$). For these losses, we have “Sparse” rate

$$M(B_{p_g}^{\sigma_g}, B_1^{\sigma_d}) \asymp n^{-\frac{\sigma_g + \sigma_d - D/p_g}{2\sigma_g + D - 2D/p_g}} + n^{-1/2}.$$

Example 3. (Maximum Mean Discrepancy) If \mathcal{F} is a ball of radius L in a reproducing kernel Hilbert space with translation invariant kernel $K(x, y) = \kappa(x - y)$ for some $\kappa \in \mathcal{L}^2(\mathcal{X})$, then,

$$\sup_{P \text{ Borel}} \mathbb{E} \left[\rho_{\mathcal{F}}(P, \hat{P}) \right] \leq \frac{L \|\kappa\|_{\mathcal{L}^2(\mathcal{X})}}{\sqrt{n}}.$$

Example 4. (Sobolev IPMs) For $\sigma \in \mathbb{N}$, B_2^{σ} is the σ -order Hilbert-Sobolev ball $B_2^{\sigma} = \left\{ f \in \mathcal{L}^2(\mathcal{X}) : \int_{\mathcal{X}} (f^{(\sigma)}(x))^2 dx \leq \infty \right\}$, where $f^{(\sigma)}$ is the σ^{th} derivative of f . For these losses, we always have the rate

$$M(B_2^{\sigma_g}, B_2^{\sigma_d}) \asymp n^{-\frac{\sigma_g + \sigma_d}{2\sigma_g + D}} + n^{-1/2}.$$

(note that $n^{-1/2}$ dominates $\Leftrightarrow t \geq 2$).

References

- [1] David L Donoho, Iain M Johnstone, Gérard Kerkycharian, and Dominique Picard. Density estimation by wavelet thresholding. *The Annals of Statistics*, pages 508–539, 1996.
- [2] Shashank Singh, Ananya Uppal, Boyue Li, Chun-Liang Li, Manzil Zaheer, and Barnabas Poczos. Nonparametric density estimation under adversarial losses. In *NeurIPS*, 2018.
- [3] Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov spaces: optimal rate and curse of dimensionality. *arXiv preprint arXiv:1810.08033*, 2018.