

# Nonparametric Density Estimation & Convergence of GANs under Besov IPM Losses

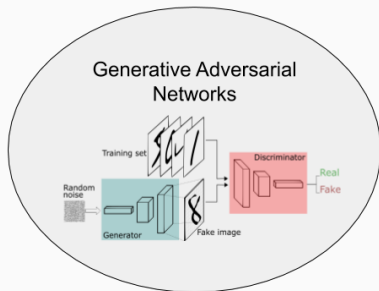
---

Ananya Uppal, **Shashank Singh\***, & Barnabás Póczos  
Carnegie Mellon University

NeurIPS 2019 Oral Presentation  
December 12, Vancouver

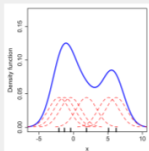
\* Now at Google

## Theoretical Guarantees for GANs

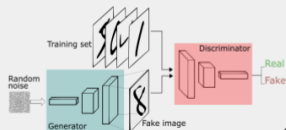


## Theoretical Guarantees for GANs through the lens of nonparametric density estimation

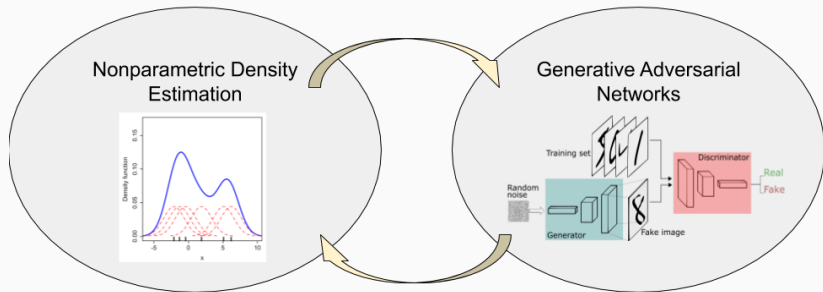
### Nonparametric Density Estimation



### Generative Adversarial Networks



## Theoretical Guarantees for GANs through the lens of nonparametric density estimation



# Contributions

1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators

# Contributions

1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality

# Contributions

1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality
3. Many classical estimators are provably sub-optimal
  - e.g., kernel density estimator
  - gap increases with dimension

# Contributions

1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality
3. Many classical estimators are provably sub-optimal
  - e.g., kernel density estimator
  - gap increases with dimension
4. Certain GANs are minimax optimal



# Contributions

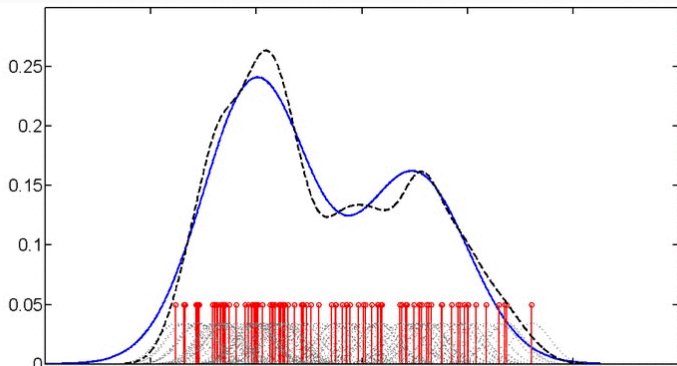
1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality
3. Many classical estimators are provably sub-optimal
  - e.g., kernel density estimator
  - gap increases with dimension
4. Certain GANs are minimax optimal
5. Besov IPMs also have important theoretical roles in math-stats.
  - Unify several previous works in nonparametric density est.

# Contributions

1. New generalization of density estimation
  - **Besov IPMs** – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality
3. Many classical estimators are provably sub-optimal
  - e.g., kernel density estimator
  - gap increases with dimension
4. **Certain GANs are minimax optimal**
5. Besov IPMs also have important theoretical roles in math-stats.
  - Unify several previous works in nonparametric density est.

# Density Estimation

- Observe  $n$  independent samples  $X_1, \dots, X_n \sim P$ .
- Assume  $P \in \mathcal{P}$ .
- Want to estimate  $P$ .



# GANs

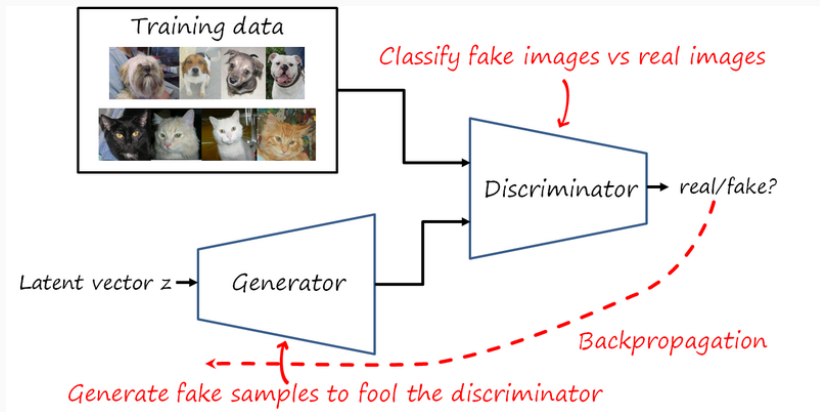
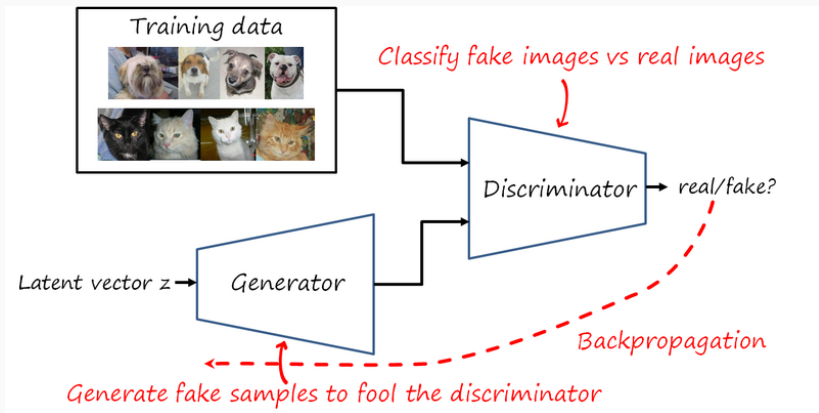


Figure from <http://www.lherranz.org/2018/08/07/imagetranslation/>.

# GANs



$$\hat{P}_{\text{GAN}} := \underbrace{\operatorname{argmin}}_{Q \in \mathcal{P}} \underbrace{\sup}_{f \in \mathcal{F}} \mathbb{E}_{X \sim Q} [f(X)] - \mathbb{E}_{X \sim P_n} [f(X)],$$

Generator      Discriminator

Figure from <http://www.lherranz.org/2018/08/07/imagetranslation/>.

# GANs as Regularized ERM Density Estimators

$$\hat{P}_{\text{GAN}} := \operatorname{argmin}_{Q \in \mathcal{P}} \sup_{f \in \mathcal{F}} \mathbb{E}_{X \sim Q} [f(X)] - \mathbb{E}_{X \sim P_n} [f(X)]$$

# GANs as Regularized ERM Density Estimators

$$\begin{aligned}\hat{P}_{\text{GAN}} &:= \operatorname{argmin}_{Q \in \mathcal{P}} \sup_{f \in \mathcal{F}} \underbrace{\mathbb{E}_{X \sim Q} [f(X)] - \mathbb{E}_{X \sim P_n} [f(X)]}_{d_{\mathcal{F}}(Q, P_n)} \\ &= \operatorname{argmin}_{Q \in \mathcal{P}} d_{\mathcal{F}}(Q, P_n)\end{aligned}$$

Empirical Risk Minimization (ERM)

- Hypothesis class  $\mathcal{P}$
- Loss  $d_{\mathcal{F}}$
- regularize data before feeding to GAN
  - instance noise (Sønderby et al. (2017), ICLR)

# Integral Probability Metrics (IPMs)

- $\mathcal{F}$  – class of *discriminator functions*

The metric  $d_{\mathcal{F}} : \mathcal{P} \times \mathcal{P} \rightarrow [0, \infty]$  is defined by

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X \sim P} [f(X)] - \mathbb{E}_{X \sim Q} [f(X)] \right|, \quad \text{for all } P, Q \in \mathcal{P}.$$

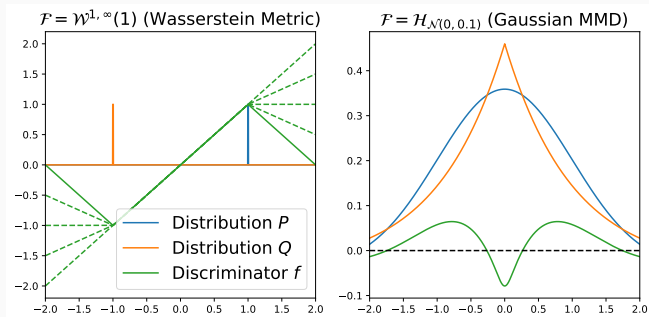


# Integral Probability Metrics (IPMs)

- $\mathcal{F}$  – class of *discriminator functions*

The metric  $d_{\mathcal{F}} : \mathcal{P} \times \mathcal{P} \rightarrow [0, \infty]$  is defined by

$$d_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X \sim P} [f(X)] - \mathbb{E}_{X \sim Q} [f(X)] \right|, \quad \text{for all } P, Q \in \mathcal{P}.$$



# Examples of IPMs

- 1-Wasserstein<sup>1</sup>
- Max. Mean Discrepancy (MMD)<sup>2</sup>
- $L^r$  distances<sup>3</sup>
- Kolmogorov-Smirnov
- Hilbert-Sobolev distances
- Besov distances
- Neural net distance (GANs)

---

<sup>1</sup>a.k.a. optimal transport or earthmover's distance

<sup>2</sup>Including energy distances

<sup>3</sup>Including total variation distance

- 2-parameter family of function spaces  $\mathcal{B}_p^s$ 
  - $s \in (0, \infty)$ ,  $p \in [1, \infty]$

# Besov Spaces

- 2-parameter family of function spaces  $\mathcal{B}_p^s$ 
  - $s \in (0, \infty)$ ,  $p \in [1, \infty]$

For integer  $s$

$$\mathcal{B}_p^s \approx \{f \in \mathcal{L}^p : \|f^{(s)}\|_p \leq C\}$$

where  $f^{(s)} = s^{\text{th}}$  derivative of  $f$ .

# Besov Spaces

- 2-parameter family of function spaces  $\mathcal{B}_p^s$ 
  - $s \in (0, \infty)$ ,  $p \in [1, \infty]$

For integer  $s$

$$\mathcal{B}_p^s \approx \{f \in \mathcal{L}^p : \|f^{(s)}\|_p \leq C\}$$

where  $f^{(s)} = s^{\text{th}}$  derivative of  $f$ .

Examples:

Ex. 1: Lipschitz/Hölder spaces:  $\mathcal{B}_\infty^s \approx \mathcal{C}^s$

Ex. 2: Sobolev spaces:  $\mathcal{B}_2^s \approx \mathcal{H}^s$

Ideal ERM:

$$\operatorname{argmin}_{Q \in \mathcal{P}} d_{\mathcal{F}}(Q, P_n).$$

Ideal ERM:

$$\operatorname{argmin}_{Q \in \mathcal{P}} d_{\mathcal{F}}(Q, P_n).$$

$\mathcal{P}$  and  $\mathcal{F}$  are  $\infty$ -dimensional... How to approximate?

# Neural Network GANs

Ideal ERM:

$$\operatorname{argmin}_{Q \in \mathcal{P}} d_{\mathcal{F}}(Q, P_n).$$

$\mathcal{P}$  and  $\mathcal{F}$  are  $\infty$ -dimensional... How to approximate?

ReLU Neural Networks (Suzuki (2019), ICLR):

$$\mathcal{B}_p^S \approx \Phi(L, W, S, B)$$

$\Phi(L, W, S, B)$  = class of fully-connected ReLU networks of size:

- $L$  = # of layers (depth)
- $W$  = # neurons/layer (width)
- $S$  = # nonzero weights/layer (sparsity)
- $B$  = largest weight value



# Neural Network GANs

Ideal ERM:

$$\operatorname{argmin}_{Q \in \mathcal{P}} d_{\mathcal{F}}(Q, P_n).$$

$\mathcal{P}$  and  $\mathcal{F}$  are  $\infty$ -dimensional... How to approximate?

ReLU Neural Networks (Suzuki (2019), ICLR):

$$\mathcal{B}_p^S \approx \Phi(L, W, S, B)$$

$\Phi(L, W, S, B)$  = class of fully-connected ReLU networks of size  $L \in O(\log n)$ ,  $W, S, B \in O(\text{poly}(n))$ .

$$\hat{P}_{\text{GAN}} = \operatorname{argmin}_{Q \in \Phi(L_g, W_g, S_g, B_g)} d_{\Phi(L_d, W_d, S_d, B_d)}(Q, P_n).$$

# GANs are optimal\*\*\*

$$\hat{P}_{\text{GAN}} = \underset{Q \in \Phi(L_g, W_g, S_g, B_g)}{\text{argmin}} d_{\Phi(L_d, W_d, S_d, B_d)}(Q, P_n).$$

is optimal for estimating Besov distributions under Besov IPMs.

# GANs are optimal\*\*\*

$$\hat{P}_{\text{GAN}} = \underset{Q \in \Phi(L_g, W_g, S_g, B_g)}{\text{argmin}} d_{\Phi(L_d, W_d, S_d, B_d)}(Q, P_n).$$

is optimal for estimating Besov distributions under Besov IPMs.

\*\*\*Caveats:

1. Well-optimized (maybe computationally challenging)

# GANs are optimal\*\*\*

$$\hat{P}_{\text{GAN}} = \underset{Q \in \Phi(L_g, W_g, S_g, B_g)}{\text{argmin}} d_{\Phi(L_d, W_d, S_d, B_d)}(Q, P_n).$$

is optimal for estimating Besov distributions under Besov IPMs.

\*\*\*Caveats:

1. Well-optimized (maybe computationally challenging)
2. Well-tuned (neural network sizes)

# GANs are optimal\*\*\*

$$\hat{P}_{\text{GAN}} = \underset{Q \in \Phi(L_g, W_g, S_g, B_g)}{\text{argmin}} d_{\Phi(L_d, W_d, S_d, B_d)}(Q, P_n).$$

is optimal for estimating Besov distributions under Besov IPMs.

\*\*\*Caveats:

1. Well-optimized (maybe computationally challenging)
2. Well-tuned (neural network sizes)
3. Assumes fully-connected ReLU networks

# Summary

1. New generalization of density estimation
  - Besov IPMs – new losses motivated partly by GAN discriminators
2. New minimax rates under these losses
  - Reduced curse of dimensionality
3. Many classical estimators are provably sub-optimal
  - e.g., kernel density estimator
4. Certain GANs are minimax optimal
5. Besov IPMs also have important theoretical roles in math. stats.
  - Unify several previous works in nonparametric density est.

Poster #243 tonight



## Further Reading

- Liang, Tengyuan. (2019) “On how well generative adversarial networks learn densities: Nonparametric and parametric results.” arXiv
- Bauer & Kohler. “On deep learning as a remedy for the curse of dimensionality in nonparametric regression.” *Annals of Statistics*.
- Johannes Schmidt-Hieber. “Nonparametric regression using deep neural networks with ReLU activation function.” *Annals of Statistics*.



# A Minimax Framework for Implicit Generative Modeling

Implicit generative model (sampler):

$$\hat{X}: \underbrace{\mathcal{X}^n}_{\text{Training Data}} \times \underbrace{\mathcal{Z}}_{\text{Randomness}} \rightarrow \underbrace{\mathcal{X}}_{\text{Novel Sample}}$$

Output distribution: conditional distribution  $P_{\hat{X}(X_1, \dots, X_n, Z) | X_1, \dots, X_n}$  of novel sample  $X_{n+1}$  given training data  $X_1, \dots, X_n$ .

Define the *implicit risk of  $\hat{X}$  at  $P$*  by

$$R_I(P, \hat{X}) := \mathbb{E}_{X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P} \left[ \ell(P, P_{\hat{X}(X_1, \dots, X_n, Z) | X_1, \dots, X_n}) \right].$$

# Implicit versus Explicit Generative Modeling

**Theorem (When do good samplers imply good density estimators?)**

Let  $\mathcal{F}_G$  be a family of probability distributions on a sample space  $\mathcal{X}$ .

Suppose

1. Loss  $\ell : \mathcal{P} \times \mathcal{P} \rightarrow [0, \infty]$  satisfies a weak triangle inequality
2.  $M_D(\mathcal{F}_G, \ell, m) \rightarrow 0$  as  $m \rightarrow \infty$ . (i.e., there exists a uniformly consistent density estimator)
3. we can draw arbitrarily many IID samples  $Z_1, Z_2, \dots$  of the latent variable  $Z$
4. Output distributions of (nearly) minimax samplers lie in  $\mathcal{F}_G$

Then,  $M_D(\mathcal{F}_G, \ell, n) \lesssim M_I(\mathcal{F}_G, \ell, n)$ .

# Implicit versus Explicit Generative Modeling

**Theorem (When do good samplers imply good density estimators?)**

Let  $\mathcal{F}_G$  be a family of probability distributions on a sample space  $\mathcal{X}$ .

Suppose

1. Loss  $\ell : \mathcal{P} \times \mathcal{P} \rightarrow [0, \infty]$  satisfies a weak triangle inequality
2.  $M_D(\mathcal{F}_G, \ell, m) \rightarrow 0$  as  $m \rightarrow \infty$ . (i.e., there exists a uniformly consistent density estimator)
3. we can draw arbitrarily many IID samples  $Z_1, Z_2, \dots$  of the latent variable  $Z$
4. Output distributions of (nearly) minimax samplers lie in  $\mathcal{F}_G$

Then,  $M_D(\mathcal{F}_G, \ell, n) \lesssim M_I(\mathcal{F}_G, \ell, n)$ .

**Proof:** Train a new density estimator  $\hat{P}$  with  $m$  IID samples drawn from the sampler  $\hat{X}$ . Then,  $R(\hat{P}) \leq R(\hat{X}) + \varepsilon_m$ .