

GANs, Optimal Transport, and Implicit Distribution Estimation

Tengyuan Liang

Econometrics and Statistics



The University of Chicago Booth School of Business

Cowles Foundation for Research in Economics

[Home](#) [People](#) [About Us](#) [Research Programs](#) [Publications](#) [News](#) [Events](#) [Resources](#) [Contact](#)

[HOME](#) > [FROM THE ARCHIVES](#) > [PEOPLE](#) > [TJALLING C. KOOPMANS \(1910-1985\)](#)

From the Archives

[From the Archives](#)

[Events](#)

[People](#)

[Photo Gallery](#)

[Research Reports](#)

Tjalling C. Koopmans (1910-1985)

Director: July 1948-1955; 1961-1964; 1965-1967

Ph.D. University of Leiden, 1936

Tjalling C. Koopmans lectured at the Rotterdam School of Economics and served on the staff of the Netherlands Economic Institute, 1936-37. From 1938 to 1940 he was engaged in business-cycle research at the League of Nations in Geneva. In 1910-41 he was on the staff of the Local and State Government Section of the School for Public and International Affairs, Princeton University, and also taught statistics at New York University. In 1941-42, he was economist with the Penn Mutual Life Insurance Company, and in 1942-44 he was statistician to the Combined Shipping Adjustment Board at Washington. Koopmans joined the staff of the Cowles Commission in July 1944, as a research associate. In 1946 he also became an associate professor in the Department of Economics at the University of Chicago. In 1948 he was appointed director of research of the Commission and professor of economics at the University of Chicago. He was elected a Fellow of the Econometric Society in 1940, of the Institute of Mathematical Statistics in 1941, of the American Statistical Association in 1949, and a member of the International



OUTLINE

Implicit Distribution Estimation

Given i.i.d. $Y_1, \dots, Y_n \sim \nu$. Use **transformation** $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ to **represent and learn** unknown dist. $Y \sim \nu$ via simple $Z \sim \mu$ (say Uniform or Gaussian).

$$T(Z) \stackrel{\text{close in dist.}}{\approx} Y$$

OUTLINE

Implicit Distribution Estimation

Given i.i.d. $Y_1, \dots, Y_n \sim \nu$. Use **transformation** $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ to **represent and learn** unknown dist. $Y \sim \nu$ via simple $Z \sim \mu$ (say Uniform or Gaussian).

$$T(Z) \stackrel{\text{close in dist.}}{\approx} Y$$

equivalently

$$T_{\#} \mu \stackrel{?}{\approx} \nu$$

OUTLINE

Implicit Distribution Estimation

Generative Adversarial Networks

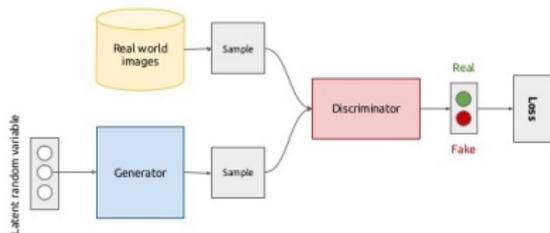
- statistical rates
- pair regularization
- optimization

Optimal Transport

- estimate the Wasserstein metric
- vs.
- estimate under the Wasserstein metric

GENERATIVE ADVERSARIAL NETWORKS

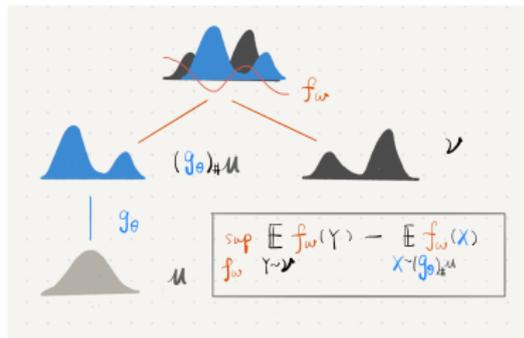
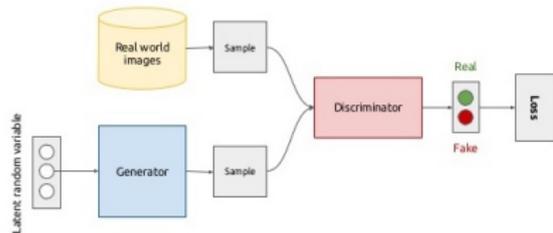
Generative adversarial networks (conceptual)



- GAN Goodfellow et al. (2014)
- WGAN Arjovsky et al. (2017); Arjovsky and Bottou (2017)
- MMD GAN Li, Swersky, and Zemel (2015); Dziugaite, Roy, and Ghahramani (2015); Arbel, Sutherland, Birkowski, and Gretton (2018)
- f -GAN Nowozin, Cseke, and Tomioka (2016)
- Sobolev GAN Mroueh et al. (2017)
- many others... Liu, Bousquet, and Chaudhuri (2017); Tolstikhin, Gelly, Bousquet, Simon-Gabriel, and Schölkopf (2017)

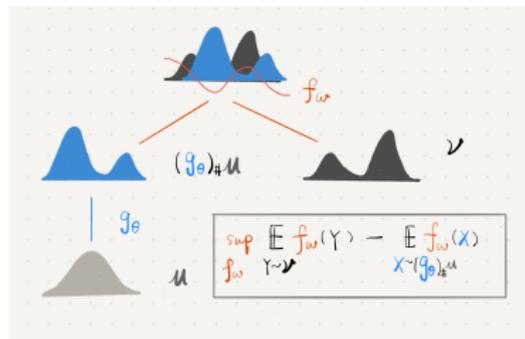
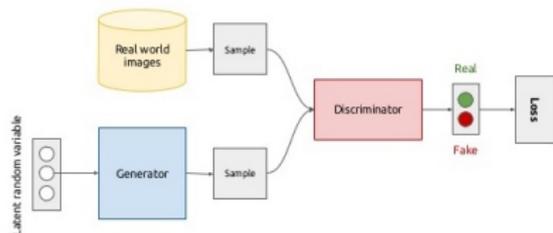
GENERATIVE ADVERSARIAL NETWORKS

Generative adversarial networks (conceptual)



GENERATIVE ADVERSARIAL NETWORKS

Generative adversarial networks (conceptual)



Generator g_θ , Discriminator f_ω

$$U(\theta, \omega) = \underbrace{\mathbb{E}_{Y \sim \nu} [f_\omega(Y)]}_{\text{target}} - \underbrace{\mathbb{E}_{Z \sim \mu} [f_\omega(g_\theta(Z))]}_{\text{input}}$$

$$\min_{\theta} \max_{\omega} U(\theta, \omega)$$

GANs are widely used in practice, **however**

MUCH NEEDS TO BE UNDERSTOOD, **IN THEORY**

- Approximation:
what dist. can be approximated by the generator $(g_{\theta})_{\#}(\mu)$?
- **Statistical:**
given n samples, what is the **statistical/generalization error rate?**
- Computational:
local convergence for practical **optimization**, how to **stabilize**?
- Landscape:
are local saddle points good globally?

FORMULATION

\mathcal{T}_G class of **generator** transformations, \mathcal{F}_D class of **discriminator** functions
 ν target dist.

$$\text{population} \quad g^* \in \arg \min_{g \in \mathcal{T}_G} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#} \mu} [f(X)] - \mathbb{E}_{Y \sim \nu} [f(Y)] \right\}$$

FORMULATION

\mathcal{T}_G class of **generator** transformations, \mathcal{F}_D class of **discriminator** functions
 ν target dist.

population
$$g^* \in \arg \min_{g \in \mathcal{T}_G} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#} \mu} [f(X)] - \mathbb{E}_{Y \sim \nu} [f(Y)] \right\}$$

$\widehat{\nu}^n$ empirical dist.

empirical
$$\widehat{g} \in \arg \min_{g \in \mathcal{T}_G} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#} \mu} [f(X)] - \mathbb{E}_{Y \sim \widehat{\nu}^n} [f(Y)] \right\}$$

$\widehat{g}_{\#} \mu$ as estimate for ν

FORMULATION

\mathcal{T}_G class of **generator** transformations, \mathcal{F}_D class of **discriminator** functions
 ν target dist.

population
$$g^* \in \arg \min_{g \in \mathcal{T}_G} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#} \mu} [f(X)] - \mathbb{E}_{Y \sim \nu} [f(Y)] \right\}$$

$\widehat{\nu}^n$ empirical dist.

empirical
$$\widehat{g} \in \arg \min_{g \in \mathcal{T}_G} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#} \mu} [f(X)] - \mathbb{E}_{Y \sim \widehat{\nu}^n} [f(Y)] \right\}$$

$\widehat{g}_{\#} \mu$ as estimate for ν

- Density learning/estimation: long history nonparametric statistics
 model target density $\rho_{\nu} \in \mathbf{W}^{\alpha}$ - Sobolev space with smoothness $\alpha \geq 0$
 Stone (1982); Nemirovski (2000); Tsybakov (2009); Wassermann (2006)
- GAN statistical theory is needed
 Arora and Zhang (2017); Arora et al. (2017a,b); Liu et al. (2017)

DISCRIMINATOR METRIC

Define the critic metric (IPM)

$$d_{\mathcal{F}}(\mu, \nu) := \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X \sim \mu} f(X) - \mathbb{E}_{Y \sim \nu} f(Y) \right| .$$

DISCRIMINATOR METRIC

Define the critic metric (IPM)

$$d_{\mathcal{F}}(\mu, \nu) := \sup_{f \in \mathcal{F}} | \mathbb{E}_{X \sim \mu} f(X) - \mathbb{E}_{Y \sim \nu} f(Y) | .$$

- \mathcal{F} Lip-1: Wasserstein metric d_W
- \mathcal{F} bounded by 1: total variation/Radon metric d_{TV}
- RKHS \mathcal{H} , $\mathcal{F} = \{f \in \mathcal{H}, \|f\|_{\mathcal{H}} \leq 1\}$: MMD GAN
- \mathcal{F} Sobolev smoothness β : Sobolev GAN

Statistical question: statistical error rate with n -i.i.d samples, $\mathbb{E} d_{\mathcal{F}}(\nu, \widehat{\mu}_n)$?
for a range of \mathcal{F} and ν with certain regularity.

SUMMARY OF FIRST HALF OF TALK

Goal	Evaluation Metric	Results		Generator Class \mathcal{G}	Discriminator Class \mathcal{F}	Property
Adversarial Framework (nonparametric)	$d_{\mathcal{F}}$	Sobolev GAN	minimax optimal	Sobolev W^{α}	Sobolev W^{β}	
		MMD GAN	upper bound	smooth subspace in RKHS	RKHS \mathcal{H}	
			oracle results	any	Sobolev W^{β}	\mathcal{G}^+
Generative Adversarial Networks (parametric)	d_{TV}	leaky-ReLU GANs	upper bound	leaky-ReLU	leaky-ReLU	$\mathcal{F}^{\dagger}, m^*$
	d_{TV}, d_{KL}, d_H	any GANs	oracle results	neural networks	neural networks	$\mathcal{G}^+, \mathcal{F}^{\dagger}, m^*$
	d_W	Lipschitz GANs	oracle results	Lipschitz neural networks	Lipschitz neural networks	$\mathcal{G}^+, \mathcal{F}^{\dagger}, m^*$

SUMMARY OF FIRST HALF OF TALK

Goal	Evaluation Metric	Results		Generator Class \mathcal{G}	Discriminator Class \mathcal{F}	Property
Adversarial Framework (nonparametric)	$d_{\mathcal{F}}$	Sobolev GAN	minimax optimal	Sobolev W^{α}	Sobolev W^{β}	
		MMD GAN	upper bound	smooth subspace in RKHS	RKHS \mathcal{H}	
			oracle results	any	Sobolev W^{β}	\mathcal{G}^{\dagger}
Generative Adversarial Networks (parametric)	d_{TV}	leaky-ReLU GANs	upper bound	leaky-ReLU	leaky-ReLU	$\mathcal{F}^{\dagger}, m^*$
	d_{TV}, d_{KL}, d_H	any GANs	oracle results	neural networks	neural networks	$\mathcal{G}^{\dagger}, \mathcal{F}^{\dagger}, m^*$
	d_W	Lipschitz GANs	oracle results	Lipschitz neural networks	Lipschitz neural networks	$\mathcal{G}^{\dagger}, \mathcal{F}^{\dagger}, m^*$

The symbols: (\mathcal{G}^{\dagger}) and (\mathcal{F}^{\dagger}) to denote the mis-specification for the generator class and the discriminator class respectively, and (m^*) to indicate the dependence on the number of generator samples.

Implicit Distribution Estimator: GANs, Optimal Transport

vs.

Explicit Density Estimator: KDE, Projection/Series Estimator, ...

Adversarial Framework (nonparametric)

MINIMAX OPTIMAL RATES: SOBOLEV GAN

Consider the target $\mathcal{G} := \{\nu : \rho_\nu \in W^\alpha\}$ Sobolev space with smoothness α , and the evaluation metric $\mathcal{F} = W^\beta$ with smoothness β .

MINIMAX OPTIMAL RATES: SOBOLEV GAN

Consider the target $\mathcal{G} := \{\nu : \rho_\nu \in W^\alpha\}$ Sobolev space with smoothness α , and the evaluation metric $\mathcal{F} = W^\beta$ with smoothness β .

Theorem (L. '17 & L. '18, Sobolev).

The **minimax optimal rate** is

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} d_{\mathcal{F}}(\nu, \tilde{\nu}_n) \asymp n^{-\frac{\alpha+\beta}{2\alpha+d}} \vee n^{-\frac{1}{2}}.$$

Here $\tilde{\nu}_n$ any estimator based on n samples. d -dim.

Liang (2017); Singh et al. (2018); Weed and Berthet (2019)

MINIMAX OPTIMAL RATES: MMD GAN

Consider a reproducing kernel Hilbert space (RKHS) \mathcal{H}

- integral operator \mathcal{T} with eigenvalue decay $t_i \asymp i^{-\kappa}$, $0 < \kappa < \infty$
- evaluation metric $\mathcal{F} = \{f \in \mathcal{H} \mid \|f\|_{\mathcal{H}} \leq 1\}$
- target density ρ_{ν} in $\mathcal{G} = \{\nu \mid \|\mathcal{T}^{-\frac{\alpha-1}{2}} \rho_{\nu}\|_{\mathcal{H}} \leq 1\}$ with smoothness α

MINIMAX OPTIMAL RATES: MMD GAN

Consider a reproducing kernel Hilbert space (RKHS) \mathcal{H}

- integral operator \mathcal{T} with eigenvalue decay $t_i \asymp i^{-\kappa}$, $0 < \kappa < \infty$
- evaluation metric $\mathcal{F} = \{f \in \mathcal{H} \mid \|f\|_{\mathcal{H}} \leq 1\}$
- target density ρ_{ν} in $\mathcal{G} = \{\nu \mid \|\mathcal{T}^{-\frac{\alpha-1}{2}} \rho_{\nu}\|_{\mathcal{H}} \leq 1\}$ with smoothness α

Theorem (L. '18, RKHS).

The **minimax optimal rate** is

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} d_{\mathcal{F}}(\nu, \tilde{\nu}_n) \lesssim n^{-\frac{(\alpha+1)\kappa}{2\alpha\kappa+2}} \vee n^{-\frac{1}{2}} .$$

MINIMAX OPTIMAL RATES: MMD GAN

Consider a reproducing kernel Hilbert space (RKHS) \mathcal{H}

- integral operator \mathcal{T} with eigenvalue decay $t_i \asymp i^{-\kappa}$, $0 < \kappa < \infty$
- evaluation metric $\mathcal{F} = \{f \in \mathcal{H} \mid \|f\|_{\mathcal{H}} \leq 1\}$
- target density ρ_{ν} in $\mathcal{G} = \{\nu \mid \|\mathcal{T}^{-\frac{\alpha-1}{2}} \rho_{\nu}\|_{\mathcal{H}} \leq 1\}$ with smoothness α

Theorem (L. '18, RKHS).

The **minimax optimal rate** is

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} d_{\mathcal{F}}(\nu, \tilde{\nu}_n) \lesssim n^{-\frac{(\alpha+1)\kappa}{2\alpha\kappa+2}} \vee n^{-\frac{1}{2}}.$$

$\kappa > 1$: intrinsic dim. $\sum_{i \geq 1} t_i = \sum_{i \geq 1} i^{-\kappa} \leq C$, parametric rate $n^{-\frac{(\alpha+1)\kappa}{2\alpha\kappa+2}} \vee n^{-\frac{1}{2}} = n^{-1/2}$.

$\kappa < 1$: sample complexity scales $n = e^{2 + \frac{2}{\alpha+1} \left(\frac{1}{\kappa} - 1\right)}$, effective dim. $\frac{1}{\kappa}$.

ORACLE INEQUALITY FOR GANS

Generator class \mathcal{G} may not contain the **target** ν : oracle approach.

ORACLE INEQUALITY FOR GANS

Generator class \mathcal{G} may not contain the **target** ν : oracle approach.

Let $\mathcal{T}_{\mathcal{G}}$ be any **generator transformation**. The **discriminator metric** $\mathcal{F}_D = W^\beta$, target density $\rho_\nu \in W^\alpha$.

$\widehat{g}_{\#}\mu, \widetilde{g}_{\#}\mu$ are **Implicit Density Estimators!**

Corollary (L. '17).

With empirical $\widehat{\nu}^n$ as plug-in, GAN

$$\widehat{g} \in \arg \min_{g \in \mathcal{T}_{\mathcal{G}}} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#}\mu} [f(X)] - \mathbb{E}_{Y \sim \widehat{\nu}^n} [f(Y)] \right\},$$

attains a **sub-optimal rate**

$$\mathbb{E} d_{\mathcal{F}_D}(\widehat{g}_{\#}\mu, \nu) \leq \min_{g \in \mathcal{T}_{\mathcal{G}}} d_{\mathcal{F}_D}(g_{\#}\mu, \nu) + \boxed{n^{-\frac{\beta}{d}} \vee \frac{\log n}{\sqrt{n}}}.$$

Canas and Rosasco (2012): $\beta = 1$

ORACLE INEQUALITY FOR GANS

Generator class \mathcal{G} may not contain the target ν : oracle approach.

Let $\mathcal{T}_{\mathcal{G}}$ be any generator transformation. The discriminator metric $\mathcal{F}_D = W^\beta$, target density $\rho_\nu \in W^\alpha$.

$\widehat{g}_{\#}\mu, \widetilde{g}_{\#}\mu$ are Implicit Densitators!

Corollary (L. '17).

With empirical $\widehat{\nu}^n$ as plug-in

$$\widehat{g} \in \arg \min_{g \in \mathcal{T}_{\mathcal{G}}} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#}\mu} [f(X)] \right\}$$

attains a sub-optimal rate

$$\mathbb{E} d_{\mathcal{F}_D}(\widehat{g}_{\#}\mu, \nu) \leq \min_{g \in \mathcal{T}_{\mathcal{G}}} d_{\mathcal{F}_D}(g_{\#}\mu, \nu) + n^{-\frac{\beta}{d}} \vee \frac{\log n}{\sqrt{n}}.$$

Corollary (L. '17).

In contrast, a regularized empirical $\widehat{\nu}^n$ as plug-in

$$\widetilde{g} \in \arg \min_{g \in \mathcal{T}_{\mathcal{G}}} \max_{f \in \mathcal{F}_D} \left\{ \mathbb{E}_{X \sim g_{\#}\mu} [f(X)] - \mathbb{E}_{Y \sim \widehat{\nu}^n} [f(Y)] \right\},$$

a faster rate is attainable

$$\mathbb{E} d_{\mathcal{F}_D}(\widetilde{g}_{\#}\mu, \nu) \leq \min_{g \in \mathcal{T}_{\mathcal{G}}} d_{\mathcal{F}_D}(g_{\#}\mu, \nu) + n^{-\frac{\alpha+\beta}{2\alpha+d}} \vee \frac{1}{\sqrt{n}}.$$

SUB-OPTIMALITY AND REGULARIZATION

Regularization helps achieve faster rate

Use $\tilde{\nu}^n$ “smoothed” empirical estimate, that serves as **regularization**

For example, kernel smoothing: $\tilde{\nu}^n(x) = \frac{1}{nh_n} K\left(\frac{x-x_i}{h_n}\right)$, SGD works

Turns out, this is used in practice, called “instance noise” or “data augmentation”

Sønderby et al. (2016); Liang et al. (2017); Arjovsky and Bottou (2017); Mescheder et al. (2018)

Generative Adversarial Networks and Pair Regularization (parametric)

Consider the parametrized GAN estimator

$$\widehat{\theta}_{m,n} \in \arg \min_{\theta: g_{\theta} \in \mathcal{G}} \max_{\omega: f_{\omega} \in \mathcal{F}} \{ \widehat{\mathbb{E}}_m f_{\omega}(g_{\theta}(Z)) - \widehat{\mathbb{E}}_n f_{\omega}(Y) \},$$

with m generator samples and n target samples.

How well GANs learn the distribution, under **objective evaluation metric**, say

$$d_{TV} \left((g_{\widehat{\theta}_{m,n}})_{\#} \mu, \nu \right) ?$$

GENERALIZED ORACLE INEQUALITY

$$\text{approx. err.} \quad \boxed{A_1(\mathcal{F}, \mathcal{G}, \nu) := \sup_{\theta} \inf_{\omega} \left\| \log \frac{\rho_{\nu}}{\rho_{\mu_{\theta}}} - f_{\omega} \right\|}, \quad A_2(\mathcal{G}, \nu) := \inf_{\theta} \left\| \log \frac{\rho_{\mu_{\theta}}}{\rho_{\nu}} \right\|^{1/2},$$

$$\text{sto. err.} \quad S_{n,m}(\mathcal{F}, \mathcal{G}) := \sqrt{\text{Pdim}(\mathcal{F}) \frac{\log(m \wedge n)}{m \wedge n}} \vee \sqrt{\text{Pdim}(\mathcal{F} \circ \mathcal{G}) \frac{\log(m)}{m}},$$

$\text{Pdim}(\cdot)$ the pseudo-dimension of the neural network function.

GENERALIZED ORACLE INEQUALITY

$$\text{approx. err.} \quad \boxed{A_1(\mathcal{F}, \mathcal{G}, \nu) := \sup_{\theta} \inf_{\omega} \left\| \log \frac{\rho_{\nu}}{\rho_{\mu_{\theta}}} - f_{\omega} \right\|}, \quad A_2(\mathcal{G}, \nu) := \inf_{\theta} \left\| \log \frac{\rho_{\mu_{\theta}}}{\rho_{\nu}} \right\|^{1/2},$$

$$\text{sto. err.} \quad S_{n,m}(\mathcal{F}, \mathcal{G}) := \sqrt{\text{Pdim}(\mathcal{F}) \frac{\log(m \wedge n)}{m \wedge n}} \vee \sqrt{\text{Pdim}(\mathcal{F} \circ \mathcal{G}) \frac{\log(m)}{m}},$$

$\text{Pdim}(\cdot)$ the pseudo-dimension of the neural network function.

Theorem (L. '18, generalized oracle inequality).

$$\begin{aligned} & \mathbb{E} d_{TV}^2 \left(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right), \mathbb{E} d_W^2 \left(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right), \\ & \mathbb{E} d_{KL} \left(\nu \| (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right) + \mathbb{E} d_{KL} \left((\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \| \nu \right) \\ & \leq A_1(\mathcal{F}, \mathcal{G}, \nu) + A_2(\mathcal{G}, \nu) + S_{n,m}(\mathcal{F}, \mathcal{G}). \end{aligned}$$

GENERALIZED ORACLE INEQUALITY

$$\text{approx. err.} \quad A_1(\mathcal{F}, \mathcal{G}, \nu) := \sup_{\theta} \inf_{\omega} \left\| \log \frac{\rho_{\nu}}{\rho_{\mu_{\theta}}} - f_{\omega} \right\|, \quad A_2(\mathcal{G}, \nu) := \inf_{\theta} \left\| \log \frac{\rho_{\mu_{\theta}}}{\rho_{\nu}} \right\|^{1/2},$$

$$\text{sto. err.} \quad S_{n,m}(\mathcal{F}, \mathcal{G}) := \sqrt{\text{Pdim}(\mathcal{F}) \frac{\log(m \wedge n)}{m \wedge n}} \vee \sqrt{\text{Pdim}(\mathcal{F} \circ \mathcal{G}) \frac{\log(m)}{m}},$$

$\text{Pdim}(\cdot)$ the pseudo-dimension of the neural network function.

Theorem (L. '18, generalized oracle inequality).

$$\begin{aligned} & \mathbb{E} d_{TV}^2 \left(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right), \mathbb{E} d_W^2 \left(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right), \\ & \mathbb{E} d_{KL} \left(\nu \| (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \right) + \mathbb{E} d_{KL} \left((\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu \| \nu \right) \\ & \leq A_1(\mathcal{F}, \mathcal{G}, \nu) + A_2(\mathcal{G}, \nu) + S_{n,m}(\mathcal{F}, \mathcal{G}) . \end{aligned}$$

We emphasize on the interplay between $(\mathcal{G}, \mathcal{F})$ as a **pair** of tuning parameters for **regularization**.

$$\text{approx. err. } \boxed{A_1(\mathcal{F}, \mathcal{G}, \nu) := \sup_{\theta} \inf_{\omega} \left\| \frac{\sqrt{\rho_{\nu}} - \sqrt{\rho_{\mu_{\theta}}}}{\sqrt{\rho_{\nu}} + \sqrt{\rho_{\mu_{\theta}}}} - f_{\omega} \right\|},$$

$$A_2(\mathcal{G}, \nu) := \inf_{\theta} \left\| \frac{\sqrt{\rho_{\nu}} - \sqrt{\rho_{\mu_{\theta}}}}{\sqrt{\rho_{\nu}} + \sqrt{\rho_{\mu_{\theta}}}} \right\|,$$

Theorem (L. '18, generalized oracle inequality).

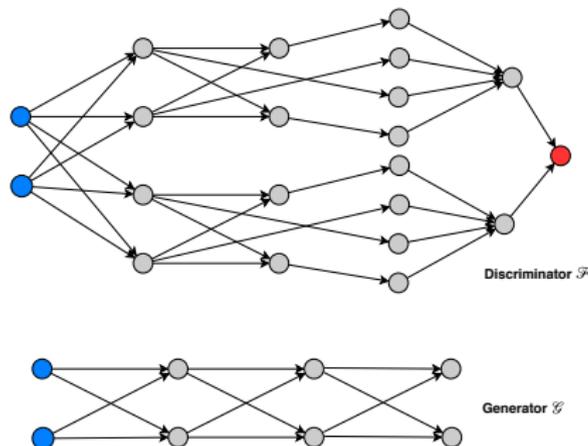
$$\begin{aligned} & \mathbb{E} d_{TV}^2(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu), \mathbb{E} d_H^2(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}})_{\#} \mu), \\ & \leq A_1(\mathcal{F}, \mathcal{G}, \nu) + A_2(\mathcal{G}, \nu) + S_{n,m}(\mathcal{F}, \mathcal{G}). \end{aligned}$$

similar result for Hellinger d_H , for non-absolutely continuous $(\mathcal{G}_{\theta})_{\#} \mu$ and ν .

Applications of pair regularization

APPLICATION I: PARAMETRIC RATES FOR LEAKY RELU NETWORKS

When the generator \mathcal{G} and discriminator \mathcal{F} are both **leaky ReLU** networks with depth L (width properly chosen depends on dimension).



When the target density is realizable by the generator.

$$\log p_{(g_\theta)_\# \mu}(x) = c_1 \sum_{l=1}^{L-1} \sum_{i=1}^d \mathbf{1}_{m_{li}(x) \geq 0} + c_0,$$

APPLICATION I: PARAMETRIC RATES FOR LEAKY RELU NETWORKS

When the generator \mathcal{G} and discriminator \mathcal{F} are both **leaky ReLU** networks with depth \mathbf{L} (width properly chosen depends on dimension).

Theorem (L. '18, leaky ReLU).

$$\mathbb{E} d_{TV}^2 \left(\nu, (\mathcal{G}\widehat{\theta}_{m,n})_{\#} \mu \right) \lesssim \sqrt{d^2 \mathbf{L}^2 \log(d\mathbf{L}) \left(\frac{\log m}{m} \vee \frac{\log n}{n} \right)}.$$

The results hold for **very deep networks** with depth $\mathbf{L} = o(\sqrt{n/\log n})$.

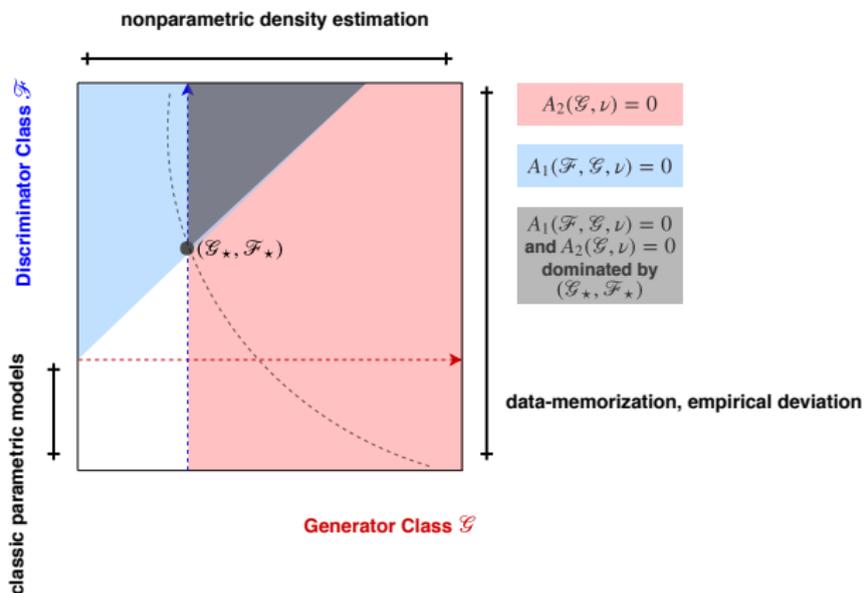
APPLICATION II: LEARNING MULTIVARIATE GAUSSIAN

Corollary (L. '18, Gaussian).

Consider $\nu \sim N(\mu, \Sigma)$. GANs enjoy near optimal sampling complexity (w.r.t. dim. d), with proper choices of the architecture and activation,

$$\mathbb{E} d_{TV}^2 \left(\nu, (\mathcal{G}_{\widehat{\theta}_{m,n}}) \# \mu \right) \lesssim \sqrt{\frac{d^2 \log d}{n \wedge m}}.$$

PAIR REGULARIZATION: WHY GANs MIGHT BE BETTER



Optimization

(local convergence)

FORMULATION

Generator g_θ , Discriminator f_ω

$$U(\theta, \omega) = \underbrace{\mathbb{E}_{Y \sim \nu} [h_1 \circ f_\omega(Y)]}_{\text{target}} - \underbrace{\mathbb{E}_{Z \sim \mu} [h_2 \circ f_\omega(g_\theta(Z))]}_{\text{input}}$$

$$\min_{\theta} \max_{\omega} U(\theta, \omega)$$

- global optimization for general $U(\theta, \omega)$ is hard [Singh et al. \(2000\)](#); [Pfau and Vinyals \(2016\)](#); [Salimans et al. \(2016\)](#)

FORMULATION

Generator g_θ , Discriminator f_ω

$$U(\theta, \omega) = \underbrace{\mathbb{E}_{Y \sim \nu} [h_1 \circ f_\omega(Y)]}_{\text{target}} - \underbrace{\mathbb{E}_{Z \sim \mu} [h_2 \circ f_\omega(g_\theta(Z))]}_{\text{input}}$$

$$\min_{\theta} \max_{\omega} U(\theta, \omega)$$

- global optimization for general $U(\theta, \omega)$ is hard [Singh et al. \(2000\)](#); [Pfau and Vinyals \(2016\)](#); [Salimans et al. \(2016\)](#)

Local saddle point (θ_*, ω_*) such that no incentive to deviate locally

$$U(\theta_*, \omega) \leq U(\theta_*, \omega_*) \leq U(\theta, \omega_*) ,$$

for (θ, ω) in an open neighborhood of (θ_*, ω_*) .

- also called local Nash Equilibrium (NE)
- modest goal: initialized properly, algorithm converges to a local NE

INTERACTION MATTERS: $\frac{\partial^2}{\partial \theta \partial \omega} \mathcal{U}(\theta, \omega)$

Geometrically fast local convergence to **stable equilibrium**

However, “**interaction term**” matters, slows down the convergence \Leftarrow curse

INTERACTION MATTERS: $\frac{\partial^2}{\partial \theta \partial \omega} U(\theta, \omega)$

Geometrically fast local convergence to **stable equilibrium**

However, “**interaction term**” matters, slows down the convergence \Leftarrow **curse**

Unstable equilibrium? turns out “**interaction term**” matters, utilize it renders geometrically fast convergence \Leftarrow **blessing**

Motivation for: *optimistic mirror descent, extra-gradients, negative-momentum ...*

*“However, no guarantees are known beyond the convex-concave setting and, more importantly for the paper, even in convex-concave games, no guarantees are known for **the last-iterate pair.**”*

— Daskalakis, Ilyas, Syrgkanis, and Zeng (2017)

GEOMETRICALLY FAST CONVERGENCE TO UNSTABLE EQUILIBRIUM

OMD proposed in [Daskalakis et al. \(2017\)](#)

$$\begin{aligned}\theta_{t+1} &= \theta_t - 2\eta \nabla_{\theta} U(\theta_t, \omega_t) + \boxed{\eta \nabla_{\theta} U(\theta_{t-1}, \omega_{t-1})} \\ \omega_{t+1} &= \omega_t + 2\eta \nabla_{\omega} U(\theta_t, \omega_t) - \boxed{\eta \nabla_{\omega} U(\theta_{t-1}, \omega_{t-1})}\end{aligned}$$

[Rakhlin and Sridharan \(2013\)](#)

For bi-linear game $U(\theta, \omega) = \theta^T C \omega$, to obtain ϵ -close solution

shown in [Daskalakis et al. \(2017\)](#) : $T \gtrsim \boxed{\epsilon^{-4} \log \frac{1}{\epsilon}} \cdot \text{Poly} \left(\frac{\lambda_{\max}(CC^T)}{\lambda_{\min}(CC^T)} \right)$

GEOMETRICALLY FAST CONVERGENCE TO UNSTABLE EQUILIBRIUM

OMD proposed in [Daskalakis et al. \(2017\)](#)

$$\begin{aligned}\theta_{t+1} &= \theta_t - 2\eta \nabla_{\theta} U(\theta_t, \omega_t) + \boxed{\eta \nabla_{\theta} U(\theta_{t-1}, \omega_{t-1})} \\ \omega_{t+1} &= \omega_t + 2\eta \nabla_{\omega} U(\theta_t, \omega_t) - \boxed{\eta \nabla_{\omega} U(\theta_{t-1}, \omega_{t-1})}\end{aligned}$$

[Rakhlin and Sridharan \(2013\)](#)

For bi-linear game $U(\theta, \omega) = \theta^T C \omega$, to obtain ϵ -close solution

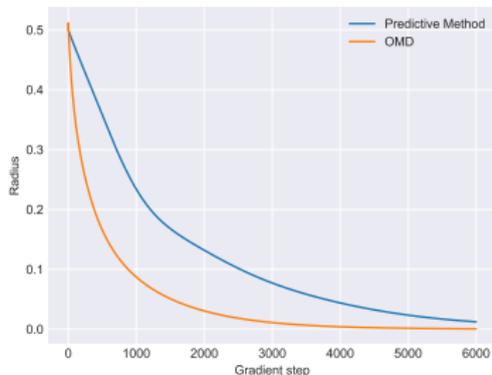
shown in [Daskalakis et al. \(2017\)](#) : $T \gtrsim \boxed{\epsilon^{-4} \log \frac{1}{\epsilon}}$ · Poly $\left(\frac{\lambda_{\max}(CC^T)}{\lambda_{\min}(CC^T)} \right)$

Theorem (L. & Stokes, '18).

we proved : $T \gtrsim \boxed{\log \frac{1}{\epsilon}}$ · $\frac{\lambda_{\max}(CC^T)}{\lambda_{\min}(CC^T)}$

further generalized beyond bi-linear game in [Mokhtari et al. \(2019\)](#).

GEOMETRICALLY FAST CONVERGENCE TO UNSTABLE EQUILIBRIUM



Theorem (L. & Stokes, '18).

$$\text{we proved : } T \gtrsim \left\lceil \log \frac{1}{\epsilon} \right\rceil \cdot \frac{\lambda_{\max}(CC^T)}{\lambda_{\min}(CC^T)}$$

further generalized beyond bi-linear game in [Mokhtari et al. \(2019\)](#).

- **Statistical** :-)
given n samples, what is the **statistical/generalization error rate**?
- **Approximation** :-(
what dist. can be approximated by the generator $g_{\theta}(Z)$?
- **Computational** :-O
local convergence for practical **optimization**, how to **stabilize**?
- **Landscape** :-(
are local saddle points good globally?

Other approach? theory of **optimal transport** \Rightarrow GANs?

OPTIMAL TRANSPORT

Wasserstein- p metric,

$$W_p(\mu, \nu) := \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} \|x - y\|^p d\pi \right)^{1/p} \quad \Pi(\mu, \nu) \text{ all couplings}$$

OPTIMAL TRANSPORT

Wasserstein- p metric,

$$W_p(\mu, \nu) := \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} \|x - y\|^p d\pi \right)^{1/p} \quad \Pi(\mu, \nu) \text{ all couplings}$$

Theorem (Brenier, '87, $p = 2$).

Let $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$. Let μ, ν absolutely continuous w.r.t. Lebesgue measure. There exists a unique **convex** $\psi_{\text{opt}} : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\begin{aligned} \frac{1}{2} W_2^2(\mu, \nu) &= \inf_{\pi \in \Pi(\mu, \nu)} \int \frac{1}{2} \|x - y\|^2 d\pi \\ &= \int \left(\frac{\|x\|^2}{2} - \psi_{\text{opt}}(x) \right) \mu(dx) + \int \left(\frac{\|y\|^2}{2} - \psi_{\text{opt}}^*(y) \right) \nu(dy) \end{aligned}$$

Here $\psi^*(y) = \sup_x \{ \langle y, x \rangle - \psi(x) \}$ is the Legendre-Fenchel conjugate of ψ .

OPTIMAL TRANSPORT

Approximation :-)

Consider $[0, 1]^d$, $Z \sim \text{Unif}([0, 1]^d)$, with a convex ψ
 $(\nabla\psi)(Z)$ can represent distribution ν !

Theorem (Brenier, '87, $p = 2$).

Let $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$. Let μ, ν absolutely continuous w.r.t. Lebesgue measure. There exists a unique **convex** $\psi_{\text{opt}} : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\begin{aligned} \frac{1}{2} W_2^2(\mu, \nu) &= \inf_{\pi \in \Pi(\mu, \nu)} \int \frac{1}{2} \|x - y\|^2 d\pi \\ &= \int \left(\frac{\|x\|^2}{2} - \psi_{\text{opt}}(x) \right) \mu(dx) + \int \left(\frac{\|y\|^2}{2} - \psi_{\text{opt}}^*(y) \right) \nu(dy) \\ &= \int \frac{1}{2} \|x - (\nabla\psi_{\text{opt}})(x)\|^2 \mu(dx), \quad \boxed{\nu = (\nabla\psi_{\text{opt}})_{\#} \mu} \end{aligned}$$

OPTIMAL TRANSPORT

Computation :-)

 linear program, or smooth convex program simple

Landscape

Recall input measure μ given, empirical target measure $\widehat{\nu}^n$

$$\frac{1}{2}W_2^2(\mu, \widehat{\nu}^n) = \sup_{\phi} \left\{ \int \phi^c(x) \mu(dx) + \int \phi(y) \widehat{\nu}^n(dy) \right\}$$

where $\phi^c(x) := \inf_y \left\{ \frac{1}{2} \|x - y\|^2 - \phi(y) \right\}$.

Genevay, Cuturi, Peyré, and Bach (2016)

OPTIMAL TRANSPORT

Computation :-)

linear program, or smooth convex program simple Landscape

Add ϵ -entropic regularization

$$\frac{1}{2} W_{2,\epsilon}^2(\mu, \widehat{\nu}^n) = \sup_{\phi} \left\{ \int \phi_{\epsilon}^c(x) \mu(dx) + \int \phi(y) \widehat{\nu}^n(dy) \right\}$$

where $\phi_{\epsilon}^c(x) := -\epsilon \log \left[\int \exp \left(-\frac{\frac{1}{2} \|x-y\|^2 - \phi(y)}{\epsilon} \right) \widehat{\nu}^n(dy) \right]$.

On data y_1, \dots, y_n
 optimization reduces to **SGD** on $[\phi(y_1), \dots, \phi(y_n)] \in \mathbb{R}^n$

Genevay, Cuturi, Peyré, and Bach (2016)

varying ϵ , solving $W_{2,\epsilon}^2(\mu, \widehat{\nu}^n)$ induced **transportation map**

$$(Id - \nabla \phi^c)(x) = \frac{\sum_{i=1}^n y_i \exp\left(-\frac{\frac{1}{2}\|x-y_i\|^2 - \phi(y_i)}{\epsilon}\right)}{\sum_{i=1}^n \exp\left(-\frac{\frac{1}{2}\|x-y_i\|^2 - \phi(y_i)}{\epsilon}\right)}$$



On data y_1, \dots, y_n
 optimization reduces to **SGD** on $[\phi(y_1), \dots, \phi(y_n)] \in \mathbb{R}^n$

OPTIMAL TRANSPORT AND PAIR REGULARIZATION

Recall input measure μ given, empirical target measure $\widehat{\nu}^n$

$$\frac{1}{2}W_2^2(\mu, \widehat{\nu}^n) = \sup_{\phi} \left\{ \int \phi^c(x) \mu(dx) + \int \phi(y) \widehat{\nu}^n(dy) \right\}$$

where $\phi^c(x) := \inf_y \left\{ \frac{1}{2} \|x - y\|^2 - \phi(y) \right\}$.

Analogy to GANs:

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ as **discriminator function**

$Id - \nabla \phi^c : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as **generator transformation**

OPTIMAL TRANSPORT AND PAIR REGULARIZATION

Recall input measure μ given, empirical target measure $\widehat{\nu}^n$

$$\frac{1}{2}W_2^2(\mu, \widehat{\nu}^n) = \sup_{\phi} \left\{ \int \phi^c(x) \mu(dx) + \int \phi(y) \widehat{\nu}^n(dy) \right\}$$

where $\phi^c(x) := \inf_y \left\{ \frac{1}{2} \|x - y\|^2 - \phi(y) \right\}$.

Analogy to GANs:

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ as **discriminator function**

$Id - \nabla \phi^c : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as **generator transformation**

However, $(Id - \nabla \phi^c)_{\#} \mu = \widehat{\nu}^n$ data memorization

$$W_2 \left((Id - \nabla \phi^c)_{\#} \mu, \nu \right) = W_2 \left(\widehat{\nu}^n, \nu \right) \asymp n^{-\frac{1}{d}}$$

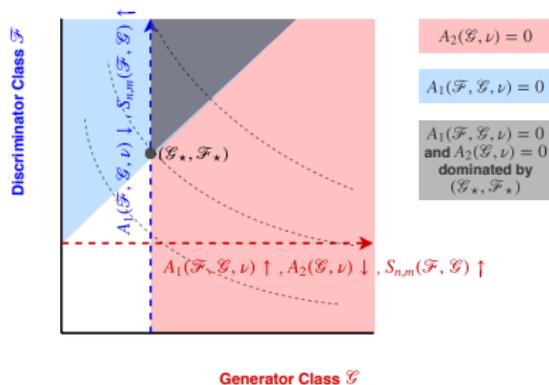
PAIR REGULARIZATION, AGAIN

Analogy to GANs:

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ as **discriminator function**

$Id - \nabla \phi^c : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as **generator transformation**

Solution: **pair regularization** $\mathcal{F}_* = \{\phi, \text{regular}\}$, $\mathcal{G}_* = \{Id - \nabla \phi^c, \text{regular}\}$ for better **statistical rate**



Estimating Transportation Cost

ANOTHER APPLICATION OF PAIR REGULARIZATION

Regularity in OT [Caffarelli \(1992, 1991\)](#): $\mu, \nu \in \mathbf{C}^\alpha$ Hölder.

Statistical question: estimate “transportation cost” $W_2^2(\mu, \nu)$ based on n -i.i.d. samples $y_1, \dots, y_n \sim \nu$. Suppose $\mu \sim \text{Unif}([0, 1]^d)$ known.

Lemma (L. & Sadhanala, '19).

$$\sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} |\tilde{W}_n - W_2^2(\mu, \nu)| \lesssim n^{-\frac{2\alpha+2}{2\alpha+d}} + n^{-\frac{1}{2}}$$

ANOTHER APPLICATION OF PAIR REGULARIZATION

Regularity in OT Caffarelli (1992, 1991): $\mu, \nu \in C^\alpha$ Hölder.

Statistical question: estimate “transportation cost” $W_2^2(\mu, \nu)$ based on n -i.i.d. samples $y_1, \dots, y_n \sim \nu$. Suppose $\mu \sim \text{Unif}([0, 1]^d)$ known.

Lemma (L. & Sadhanala, '19).

$$\sup_{\nu \in C^\alpha} \mathbb{E} |\tilde{W}_n - W_2^2(\mu, \nu)| \lesssim n^{-\frac{2\alpha+2}{2\alpha+d}} + n^{-\frac{1}{2}}$$

Elbow phenomenon: $\alpha \geq \frac{d}{2} - 2$, one gets parametric rate

ANOTHER APPLICATION OF PAIR REGULARIZATION

Regularity in OT [Caffarelli \(1992, 1991\)](#): $\mu, \nu \in C^\alpha$ Hölder.

Statistical question: estimate “transportation cost” $W_2^2(\mu, \nu)$ based on n -i.i.d. samples $y_1, \dots, y_n \sim \nu$. Suppose $\mu \sim \text{Unif}([0, 1]^d)$ known.

Lemma (L. & Sadhanala, '19).

$$\sup_{\nu \in C^\alpha} \mathbb{E} |\tilde{W}_n - W_2^2(\mu, \nu)| \lesssim n^{-\frac{2\alpha+2}{2\alpha+d}} + n^{-\frac{1}{2}}$$

Pair regularization: $\phi \in C^{\alpha+2}$, $Id - \nabla \phi^c \in C^{\alpha+1}$, by [Caffarelli \(1992, 1991\)](#)

ANOTHER APPLICATION OF PAIR REGULARIZATION

Regularity in OT [Caffarelli \(1992, 1991\)](#): $\mu, \nu \in \mathbf{C}^\alpha$ Hölder.

Statistical question: estimate “transportation cost” $W_2^2(\mu, \nu)$ based on n -i.i.d. samples $y_1, \dots, y_n \sim \nu$. Suppose $\mu \sim \text{Unif}([0, 1]^d)$ known.

Lemma (L. & Sadhanala, '19).

$$\sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} |\tilde{W}_n - W_2^2(\mu, \nu)| \lesssim n^{-\frac{2\alpha+2}{2\alpha+d}} + n^{-\frac{1}{2}}$$

typically an easier problem than estimating measure under W_2 , or estimating transportation map T under metric $\mathbb{E}_{X \sim \mu} \|\hat{T}(X) - T(X)\|^2$

Hütter and Rigollet (2019)

BACK TO THE ADVERSARIAL FRAMEWORK

Two related problems

Estimate under the metric/loss

Theorem (L., '17).

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} d_{\mathcal{F}}^2(\nu, \tilde{\nu}_n) \asymp n^{-\frac{2\alpha+2\beta}{2\alpha+d}} \vee n^{-1}$$

$$\mathcal{G} = W^\alpha, \mathcal{F} = W^\beta$$

No elbow phenomenon on α .

Liang (2017); Singh et al. (2018); Weed and Berthet (2019)

BACK TO THE ADVERSARIAL FRAMEWORK

Two related problems

Estimate under the metric/loss

Theorem (L., '17).

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} d_{\mathcal{F}}^2(\nu, \tilde{\nu}_n) \asymp n^{-\frac{2\alpha+2\beta}{2\alpha+d}} \vee n^{-1}$$

$$\mathcal{G} = W^\alpha, \mathcal{F} = W^\beta$$

No elbow phenomenon on α .

Estimating the metric/loss itself

Theorem (L. & Sadhanala, '19).

$$\inf_{\tilde{W}_n} \sup_{\nu \in \mathcal{G}} \mathbb{E} |\tilde{W}_n - d_{\mathcal{F}}^2(\mu, \nu)|^2 \asymp n^{-\frac{8\alpha+8\beta}{4\alpha+d}} \vee n^{-1}$$

$$\mathcal{G} = W^\alpha, \mathcal{F} = W^\beta$$

Elbow phenomenon on $\alpha = d/4 - 2\beta$.

Liang (2017); Singh et al. (2018); Weed and Berthet (2019)

typically an easier problem.

HOWEVER, FOR WASSERSTEIN METRIC

Theorem (L., '19).

Consider $d \geq 2$ and the domain $\Omega = [0, 1]^d$. Given n i.i.d. samples y_1, \dots, y_n from ν ,

$$\inf_{\tilde{W}_n} \sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} |\tilde{W}_n - W_1(\mu, \nu)| \lesssim n^{-\frac{\alpha+1}{2\alpha+d}},$$

as we know

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} W(\tilde{\nu}_n, \nu) \asymp n^{-\frac{\alpha+1}{2\alpha+d}}.$$

HOWEVER, FOR WASSERSTEIN METRIC

Theorem (L., '19).

Consider $d \geq 2$ and the domain $\Omega = [0, 1]^d$. Given n i.i.d. samples y_1, \dots, y_n from ν ,

$$\frac{\log \log(n)}{\log(n)} \cdot n^{-\frac{\alpha+1}{2\alpha+d}} \lesssim \inf_{\tilde{W}_n} \sup_{\nu \in \mathcal{C}^\alpha} \mathbb{E} |\tilde{W}_n - W_1(\mu, \nu)| \lesssim n^{-\frac{\alpha+1}{2\alpha+d}},$$

as we know

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathcal{C}^\alpha} \mathbb{E} W(\tilde{\nu}_n, \nu) \asymp n^{-\frac{\alpha+1}{2\alpha+d}}.$$

estimating the Wasserstein-1 metric itself
 is **almost as hard** as
 estimating under the Wasserstein-1 metric

HOWEVER, FOR WASSERSTEIN METRIC

Theorem (L., '19).

Consider $d \geq 2$ and the domain $\Omega = [0, 1]^d$. Given n i.i.d. samples y_1, \dots, y_n from ν ,

$$\frac{\log \log(n)}{\log(n)} \cdot n^{-\frac{\alpha+1}{2\alpha+d}} \lesssim \inf_{\tilde{W}_n} \sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} |\tilde{W}_n - W_1(\mu, \nu)| \lesssim n^{-\frac{\alpha+1}{2\alpha+d}},$$

as we know

$$\inf_{\tilde{\nu}_n} \sup_{\nu \in \mathbf{C}^\alpha} \mathbb{E} W(\tilde{\nu}_n, \nu) \asymp n^{-\frac{\alpha+1}{2\alpha+d}}.$$

- the main technicality is in deriving the lower bound: wavelets
- construct two composite/fuzzy hypotheses using delicate priors with matching $\log n$ moments
- and the Wasserstein metric differs sufficiently
- calculate total variation metric directly on the posterior of data (sum-product form), via a telescoping trick

SUMMARY

- In this talk, we study **statistical rates** for $d(\widehat{T}_{\#}\mu, \nu)$ and $\overline{d}(\mu, \nu)$, with $\nu = T_{\#}^*\mu$.

Implicit Distribution Estimation motivated from GANs, OT.

- Conceptually, to learn the distribution via transformation/transportation, vs., to estimate the transformation/transportation difficulty.
- Closely related problems in the lens of Optimal Transport.

$$\text{harder } d(\widehat{T}_{\#}\mu, \nu) \begin{array}{c} \xrightarrow{\text{induces plug-in estimate}} \\ \xleftarrow{\text{sometimes induces a transportation map}} \end{array} \overline{d}(\mu, \nu) \text{ easier}$$

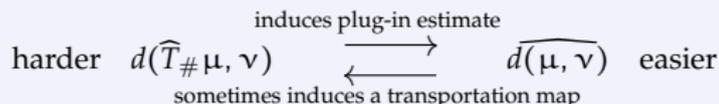
- Idea of **pair regularization**
what GANs have over classical nonparametrics.

SUMMARY

- In this talk, we study **statistical rates** for $d(\widehat{T}_{\#}\mu, \nu)$ and $\overline{d(\mu, \nu)}$, with $\nu = T_{\#}^*\mu$.

Implicit Distribution Estimation motivated from GANs, OT.

- Conceptually, to learn the distribution via transformation/transportation, vs., to estimate the transformation/transportation difficulty.
- Closely related problems in the lens of Optimal Transport.



- Idea of **pair regularization**
what GANs have over classical nonparametrics.

Many interesting open problems both **statistically** and **computationally**, with new insights on **regularization** and **adaptivity**.

Thank you!

Liang, T. (2018). — On How Well Generative Adversarial Networks Learn Densities: Nonparametric and Parametric Results. *arXiv:1811.03179* *under review*

Liang, T. & Stokes, J. (2018). — Interaction Matters: A Note on Non-asymptotic Local Convergence of Generative Adversarial Networks. *arXiv:1802.06132* *AISTATS 2019*

Liang, T. (2019). — On the Minimax Optimality of Estimating the Wasserstein Metric. *arXiv:1908.10324*

Liang, T. & Sadhanala, V. (2019). — Working Paper.