

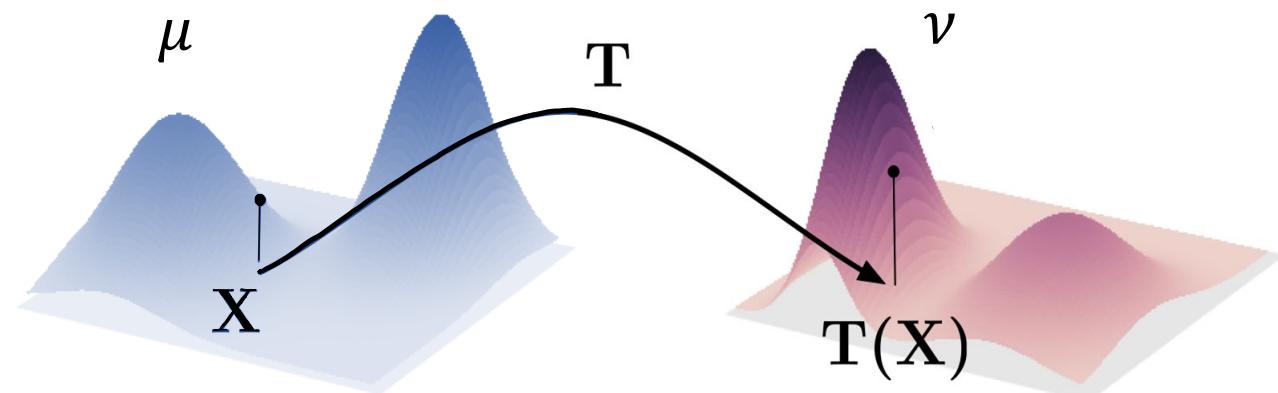
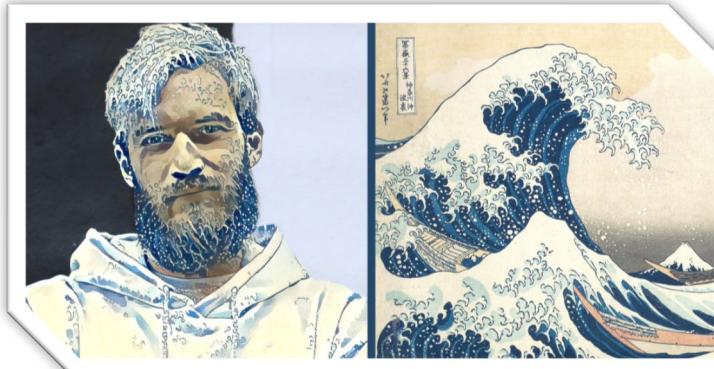
Statistical and Computational Aspects of Sliced Optimal Transport

Ziv Goldfeld
Cornell University

Machine Learning Lunch Seminar, Vanderbilt University
October 17th, 2022

Statistical Divergences for Learning & Inference

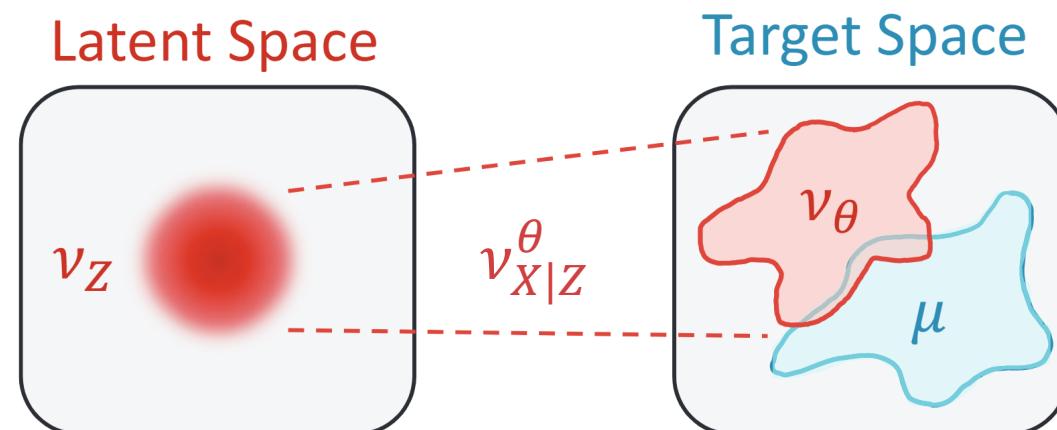
Modeling & Analysis: ML tasks boil down to comparing/transforming probability distributions



Statistical Divergences for Generative Modeling

Example: Implicit generative modeling aims to learn $\nu_\theta \approx \mu \in \mathcal{P}(\mathbb{R}^d)$

- Map $Z \sim \nu_Z \in \mathcal{P}(\mathbb{R}^p)$, $p \ll d$, to \mathbb{R}^d via parametrized transformation $\nu_{X|Z}^\theta$
- ⇒ **Generative model:** $\nu_\theta(\cdot) := \int_{\mathbb{R}^p} \nu_{X|Z}^\theta(\cdot | z) d\nu_Z(z)$



Minimum distance estimation: Solve

$$\theta^* \in \operatorname{argmin}_\theta \delta(\mu, \nu_\theta)$$

Optimal Transport Distances

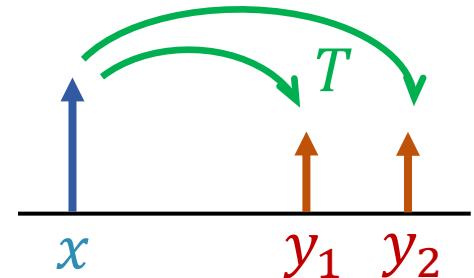
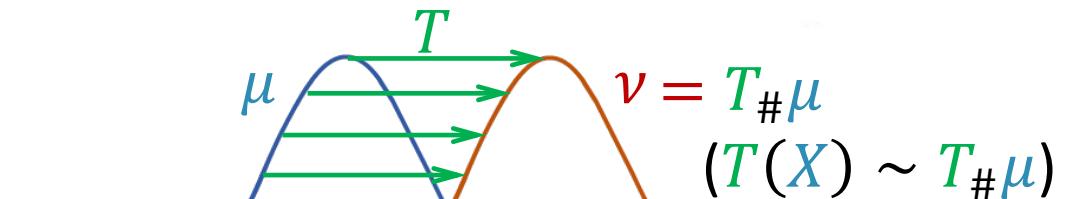
Distributions: $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$

Cost: $c: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, \infty)$

Transport map: $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ s.t. $T_\# \mu = \nu$

OT problem (Monge 1781): $M_c(\mu, \nu) := \inf_{T: T_\# \mu = \nu} \int_{\mathbb{R}^d} c(x, T(x)) d\mu(x)$

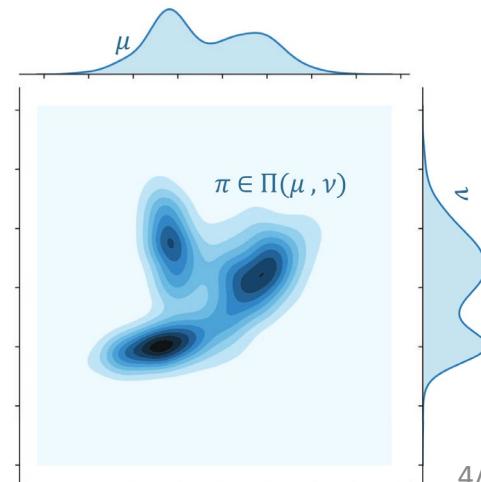
🚫 $\{T: T_\# \mu = \nu\}$ may be empty, not closed, non-linear problem, ...



Coupling: $\Pi(\mu, \nu) = \{\pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d): \pi(\cdot \times \mathbb{R}^d) = \mu, \pi(\mathbb{R}^d \times \cdot) = \nu\}$

OT problem (Kantorovich '42)

$$K_c(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} c(x, y) d\pi(x, y)$$



The Wasserstein Metric

p -Wasserstein metric

For $p \in [1, \infty)$ and $\mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$: $W_p(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^p d\pi(x, y) \right)^{1/p}$

Comments: Kantorovich OT with distance cost (or power thereof) $W_p(\mu, \nu) := \left(K_{\|\cdot\|_p}(\mu, \nu) \right)^{\frac{1}{p}}$

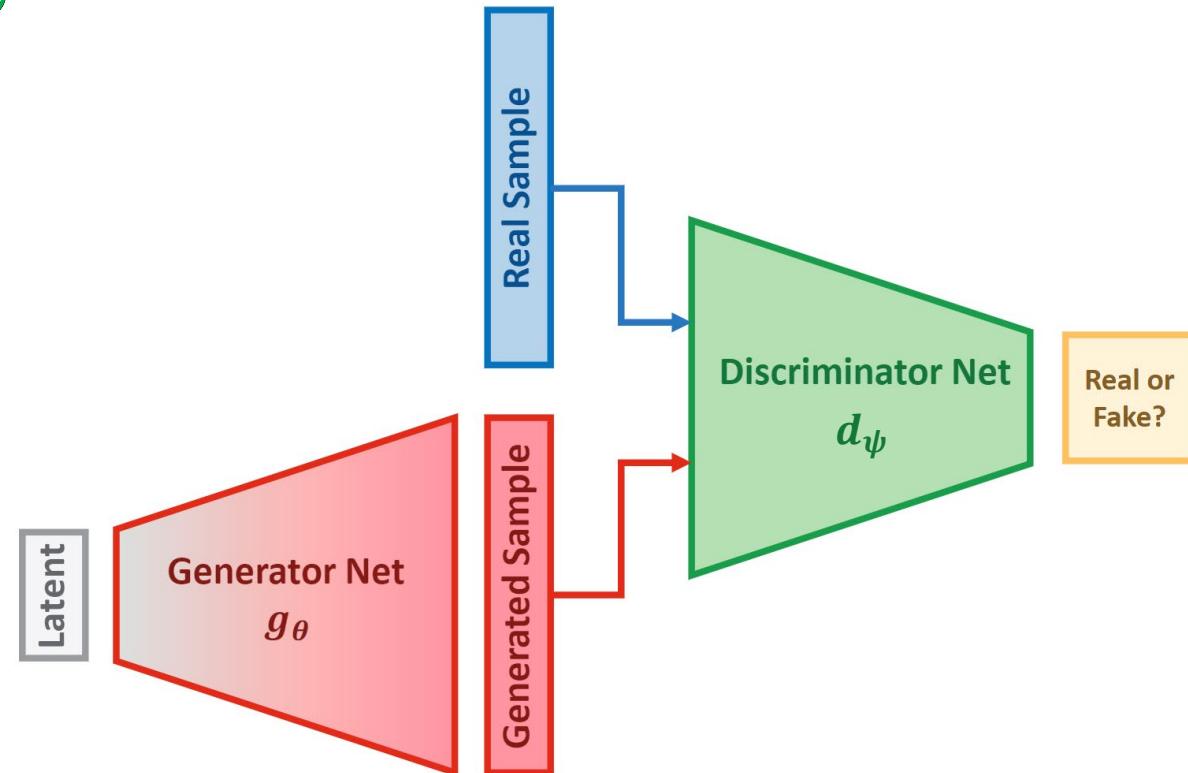
- **Wasserstein space:** $\mathfrak{W}_p = (\mathcal{P}_p(\mathbb{R}^d), W_p)$ is a metric space & rich geometric structure
- **Robust to support mismatch:** $W_p(\mu, \nu) < \infty, \forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$
- **Duality:** $W_1(\mu, \nu) = \sup_{\varphi \in \text{Lip}_1(\mathbb{R}^d)} \int_{\mathbb{R}^d} \varphi d(\mu - \nu)$

From Duality to Generative Modeling

Dual representation: $W_1(\mu, \nu) = \sup_{\varphi \in \text{Lip}_1(\mathbb{R}^d)} \mathbb{E}_{\mu}[\varphi(X)] - \mathbb{E}_{\nu}[\varphi(Y)]$

W-GAN [Arjovsky et al '17]:

- μ (X data sample)
- $\nu = \nu_{\theta}$ ($Y = g_{\theta}(Z)$ generated sample)
- $\varphi = d_{\psi}$ ($\text{Lip}_1(\mathbb{R}^d)$ constraint)



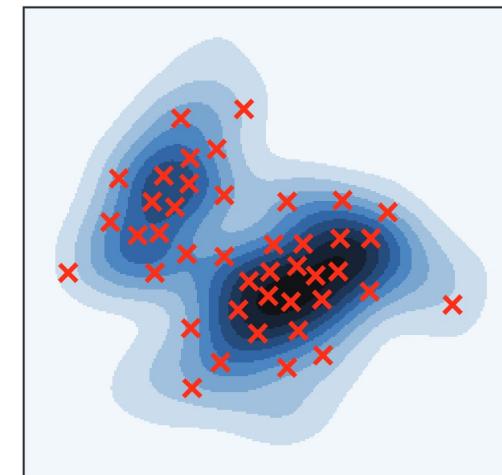
$$\inf_{\theta} W_1(\mu, \nu_{\theta}) \cong \inf_{\theta} \sup_{\psi: d_{\psi} \in \text{Lip}_1(\mathbb{R}^d)} \mathbb{E}[d_{\psi}(X)] - \mathbb{E}[d_{\psi}(g_{\theta}(Z))]$$

Empirical Estimation in High Dimensions

Sampling: μ, ν are unknown & we get $X_1, \dots, X_n \sim \mu$ and $Y_1, \dots, Y_n \sim \nu$

- **Empirical measures:** $\hat{\mu}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and $\hat{\nu}_n := \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$

→ Can we approximate $W_p(\mu, \nu) \approx W_p(\hat{\mu}_n, \hat{\nu}_n)$?



Theorem (Dudley '69, Boissard-Le Gouic'14, Niles Weed-Bach'19,...)

For $\mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$ and $d > 2p$: $\mathbb{E}[|W_p(\mu, \nu) - W_p(\hat{\mu}_n, \hat{\nu}_n)|] \asymp n^{-\frac{1}{d}}$

Core issue: For $\mu \ll \text{Leb}(\mathbb{R}^d)$: $W_p(\hat{\mu}_n, \mu) \gtrsim n^{-\frac{1}{d}}$ a.s.

🚫 **Statistical implication:** Too slow for $d \gg 1$!

🚫 **Computation:** $W_p(\hat{\mu}_n, \hat{\nu}_n)$ is LP w/ complexity $O(n^3)$ → infeasible for large n

Regularized OT Distances

Approach: Regularize the OT problem to improve statistical/computational properties

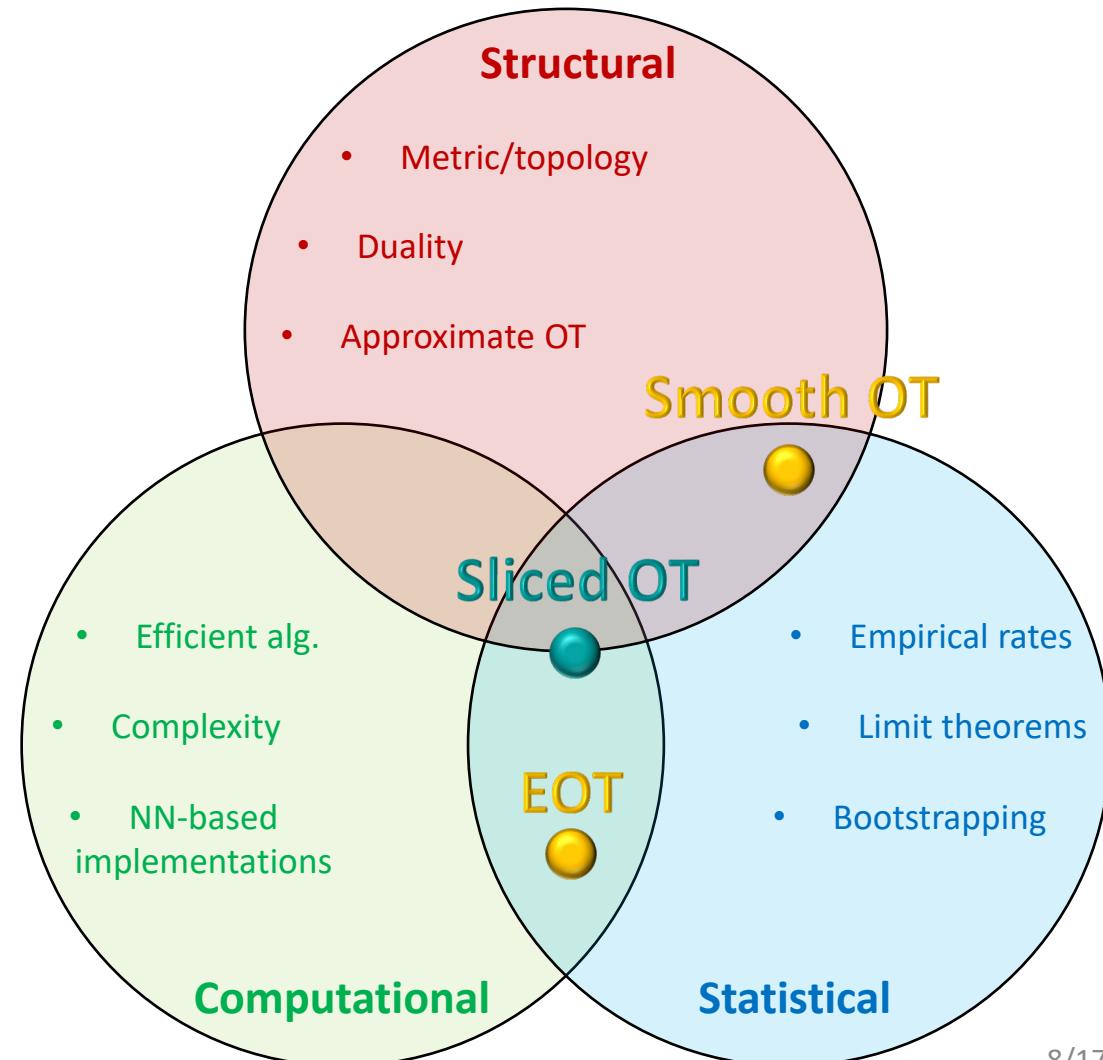
1. **Entropic OT:** Let $D_{\text{KL}}(\cdot \parallel \cdot)$ be the KL divergence:

$$E_c^\epsilon(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int c \, d\pi + \epsilon D_{\text{KL}}(\pi \parallel \mu \otimes \nu)$$

2. **Smooth OT:** For a smoothing kernel $\kappa_\sigma \in \mathcal{C}^\infty$

$$S_c^\sigma(\mu, \nu) := K_c(\mu * \kappa_\sigma, \nu * \kappa_\sigma)$$

3. **Sliced OT:** Focus of this talk

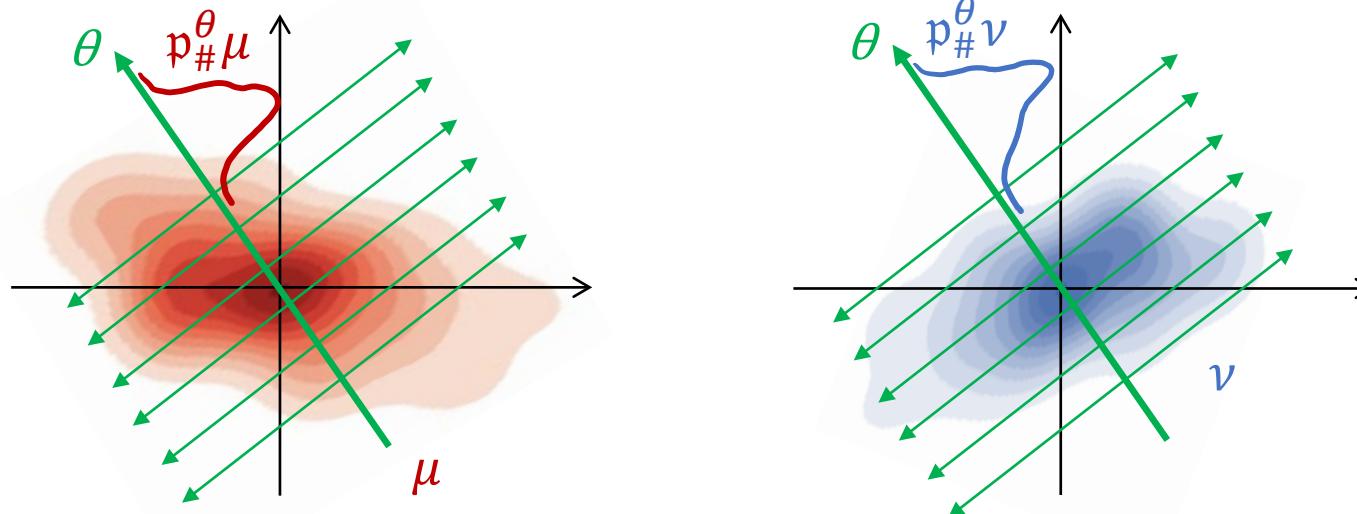


Sliced Wasserstein Distances

Definition (Rabin-Peyré-Delon-Bernot '11)

$$\underline{W}_p(\mu, \nu) := \left[\int_{\mathbb{S}^{d-1}} W_p^p(p_\#^\theta \mu, p_\#^\theta \nu) d\sigma(\theta) \right]^{1/p} \quad \text{and} \quad \overline{W}_p(\mu, \nu) := \sup_{\theta \in \mathbb{S}^{d-1}} W_p(p_\#^\theta \mu, p_\#^\theta \nu)$$

where $p^\theta(x) = \theta^T x$ and σ is the uniform measure on the unit sphere \mathbb{S}^{d-1}



Sliced Wasserstein Distances

Definition (Rabin-Peyré-Delon-Bernot '11)

$$\underline{W}_p(\mu, \nu) := \left[\int_{\mathbb{S}^{d-1}} W_p^p(p_\#^\theta \mu, p_\#^\theta \nu) d\sigma(\theta) \right]^{1/p} \quad \text{and} \quad \overline{W}_p(\mu, \nu) := \sup_{\theta \in \mathbb{S}^{d-1}} W_p(p_\#^\theta \mu, p_\#^\theta \nu)$$

where $p^\theta(x) = \theta^T x$ and σ is the uniform measure on the unit sphere \mathbb{S}^{d-1}

Motivation: OT between 1D distributions $\mu, \nu \in \mathcal{P}(\mathbb{R})$:

- **General:** $W_p(\mu, \nu) = \|F_\mu^{-1} - F_\nu^{-1}\|_{L^p([0,1])} \stackrel{p=1}{=} \|F_\mu - F_\nu\|_{L^1(\mathbb{R})}$
- **Discrete:** $W_p^p(\mu, \nu) = \frac{1}{n} \sum_{i=1}^n |X_{(i)} - Y_{(i)}|^p,$

where $X_{(1)} \leq \dots \leq X_{(n)}$ and $Y_{(1)} \leq \dots \leq Y_{(n)}$ are order statistics

→ Computing projected distance $W_p(p_\#^\theta \mu, p_\#^\theta \nu)$ is cheap!

Structural Properties

Metric/topological structure [Nadjahi *et al* '20 , Bayrakta-Guo '21]:

\underline{W}_p and \overline{W}_p are metrics on $\mathcal{P}_p(\mathbb{R}^d)$ and generate the same topology as W_p

Duality: $\overline{W}_1(\mu, \nu) = \sup_{\theta \in \mathbb{S}^{d-1}} W_1(p_\#^\theta \eta, p_\#^\theta \nu)$

Sliced W-GAN [Deshpande *et al* '18]



Approximation: Exponential gap between sliced and classic W_p

Lemma (Bonnotte '13)

There exists $c_{p,d}, C_{p,d} > 0$ s.t. for any $\mu, \nu \in \mathcal{P}(\mathbb{B}_d(0, R))$, we have

$$\underline{W}_p^p(\mu, \nu) \leq c_{p,d} W_p^p(\mu, \nu) \leq C_{p,d} R^{p - \frac{1}{d+1}} \underline{W}_p^{\frac{1}{d+1}}(\mu, \nu)$$

→ **Cannot** use sliced distances to approximate classic ones to arbitrary precision

Computational Aspects: Average-Sliced

Recall: $\mathfrak{w}_n^p(\theta) := W_p^p(p_\#^\theta \hat{\mu}_n, p_\#^\theta \hat{\nu}_n)$ computable via sorting & $\underline{W}_p^p(\hat{\mu}_n, \hat{\nu}_n) = \int_{\mathbb{S}^{d-1}} \mathfrak{w}_n^p(\theta) d\sigma(\theta)$

Average-sliced: MC integration $\underline{W}_p^p(\mu, \nu) \approx \widehat{W}_{MC}^p := \frac{1}{L} \sum_{\ell=1}^L \mathfrak{w}_n^p(\Theta_\ell)$ [Nadjahi *et al* '20]

Theorem (Nietert-Sadhu-G.-Kato '22)

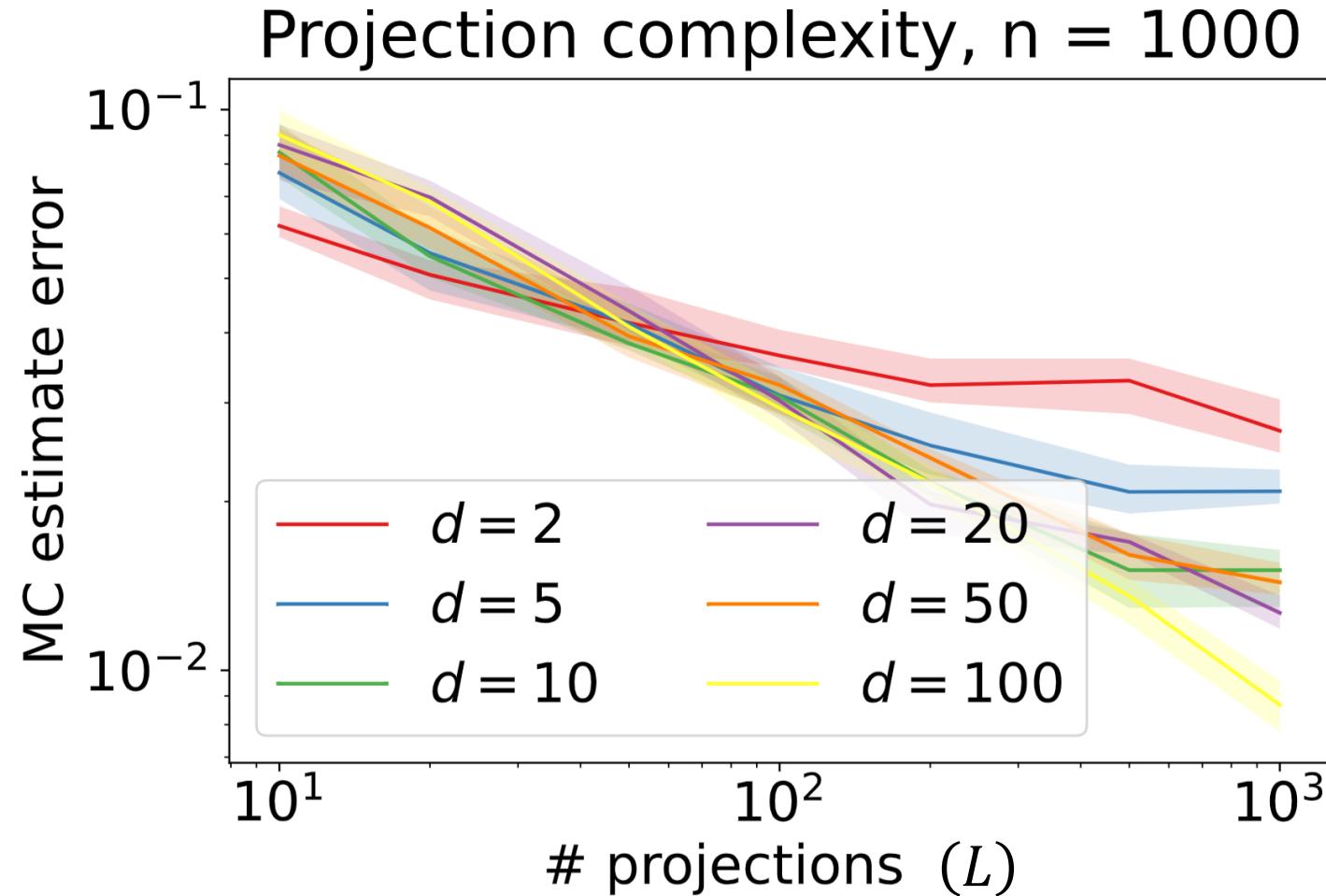
For log-concave $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ with means m_μ, m_ν and covariances Σ_μ, Σ_ν , we have

$$\mathbb{E}[|\underline{W}_p^p(\mu, \nu) - \widehat{W}_{MC}^p|] \lesssim_p \frac{\|m_\mu - m_\nu\|^p + \|\Sigma_\mu\|_{op}^{p/2} + \|\Sigma_\nu\|_{op}^{p/2}}{\sqrt{dL}} + \underbrace{\delta_n(\mu, \nu)}_{\text{empirical est. error}}$$

Analysis: 1st term corresponds to $L^{-1/2} \sqrt{\text{var}(W_p^p(p_\#^\Theta \mu, p_\#^\Theta \nu))}$, for $\Theta \sim \text{Unif}(\mathbb{S}^{d-1})$

1. Show that $W_p^p(p_\#^\theta \mu, p_\#^\theta \nu)$ is a Lipschitz function of $\theta \in \mathbb{S}^{d-1}$
2. Use concentration of Lipschitz functions over the unit sphere to bound variance

Computational Aspects: Average-Sliced (Empirical)



Computational Aspects: Max-Sliced

Max-sliced: Rewrite $\overline{W}_p^p(\hat{\mu}_n, \hat{\nu}_n) = - \min_{\theta \in \mathbb{B}_d} \max_{\sigma \in S_n} \underbrace{\left(-\frac{1}{n} \sum_{i=1}^n |\theta^T (X_i - Y_{\sigma(i)})|^p \right)}_{:= \rho(\sigma, \theta)}$

Subgradient method: Subgradient of $\tilde{w}_n^p(\theta) := \max_{\sigma \in S_n} \rho(\sigma, \theta)$ is easy to evaluate

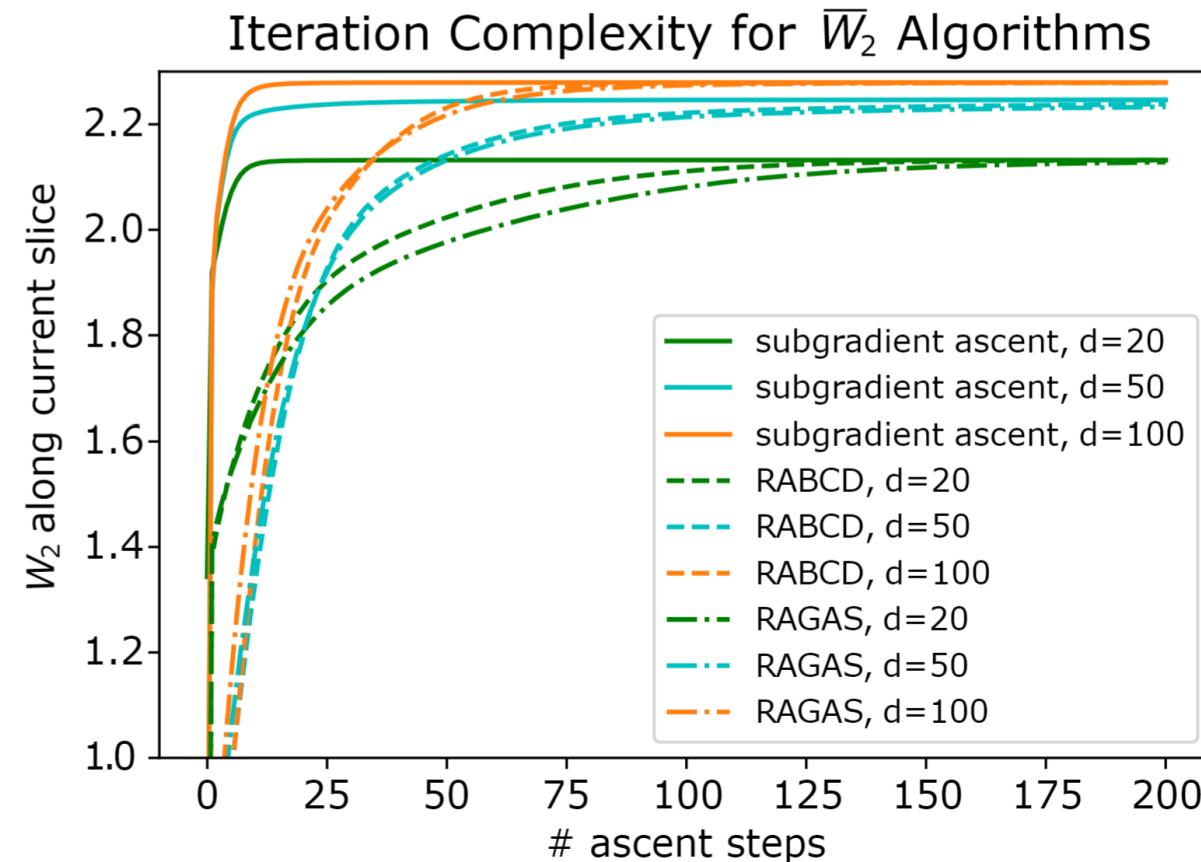
1. Sudifferential of max-function $\partial \tilde{w}_n^p(\theta) = \text{Conv}(\{\partial_\theta \rho(\sigma^*, \theta) : \sigma^* \in \text{argmax}_{\sigma \in S_n} \rho(\sigma, \theta)\})$
2. For fixed θ compute optimal permutation $\sigma^* \in S_n$ via order statistic
3. Evaluate the subgradient vector in $\partial_\theta \rho(\sigma^*, \theta)$

⇒ EM-like alg. + projected subgradient method to compute $\overline{W}_p^p(\hat{\mu}_n, \hat{\nu}_n)$ [Nadjahi et al '20]

Theorem (Nietert-Sadhu-G.-Kato '22)

$\tilde{w}_n^2(\theta)$ is weakly-convex in $\theta \in \mathbb{B}_d$.

Computational Aspects: Max-Sliced (Empirical)



Statistical Analysis: Empirical Convergence Rates

- [Nadjahi *et al* '20]: Rates for average-sliced distances follow 1D rates of base distance
- [Niles Weed-Rigollet '19, Lin *et al* '21]: $O_d(n^{-1/(2p)})$, under $T_q(\sigma^2)$, Poincaré, bdd moments

Theorem (Nietert-Sadhu-G.-Kato '22)

For log-concave $\mu \in \mathcal{P}(\mathbb{R}^d)$ with covariance Σ s.t. $\text{rank}(\Sigma) = k$, we have

$$\mathbb{E}[\underline{W}_p(\hat{\mu}_n, \mu)] \lesssim \frac{(\|\Sigma\|_{\text{op}}^{1/2} (\log n)^{\mathbb{I}\{p=2\}})}{n^{1/(2p)}}$$

$$\mathbb{E}[\overline{W}_p(\hat{\mu}_n, \mu)] \lesssim \frac{\|\Sigma\|_{\text{op}}^{1/2} k \log n}{n^{1/p}} + \frac{\|\Sigma\|_{\text{op}}^{1/2} \sqrt{k \log n}}{n^{1/(2p)}} + \frac{(\|\Sigma\|_{\text{op}}^{1/2} (\log n)^{\mathbb{I}\{p=2\}})}{n^{1/(2p)}}$$

Analysis (for \underline{W}_p): If $\mu \in \mathcal{P}_p(\mathbb{R}^d)$ is log-concave then so is $\mathfrak{p}_\#^\theta \mu$ & has bounded Cheeger constant

$\implies \sup_\theta \mathbb{E}[W_p(\mathfrak{p}_\#^\theta \hat{\mu}_n, \mathfrak{p}_\#^\theta \mu)] \lesssim \text{RHS via 1D } W_p \text{ rates for bdd Cheeger const [Bobkov-Ledoux '16]}$

Statistical Analysis: Limit Distribution (Method)

Limit theorems: Key for valid inference but remained an open question

Approach: Relies on extended functional delta method in normed spaces [Römisch '04]

Theorem (G.-Kato-Rioux-Sadhu '22)

Setup: \mathcal{F} class of Borel functions with finite envelope F

$\ell^\infty(\mathcal{F})$ space of bdd functionals from \mathcal{F} to \mathbb{R} equipped w/ $\|\mu\|_{\infty,\mathcal{F}} := \sup_{f \in \mathcal{F}} |\mu(f)|$
 Φ is a functional from (a subset of) $\ell^\infty(\mathcal{F})$ to \mathbb{R} $\widetilde{\mathbb{E}_\mu[f]}$

Assume:

1. The empirical process weakly converges $\sqrt{n}(\hat{\mu}_n - \mu) \xrightarrow{d} G_\mu$ in $\ell^\infty(\mathcal{F})$
2. Φ is a Lipschitz functional wrt $\|\cdot\|_{\infty,\mathcal{F}}$, i.e., $|\Phi(v) - \Phi(v')| \leq C \|v - v'\|_{\infty,\mathcal{F}}$
3. Φ is Gâteaux differentiable at μ w/ derivative $\Phi'_\mu(\rho) = \lim_{t \rightarrow 0^+} \frac{1}{t} (\Phi(\mu + t\rho) - \Phi(\mu))$

Then: $\sqrt{n}(\Phi(\hat{\mu}_n) - \Phi(\mu)) \xrightarrow{d} \Phi'_\mu(G_\mu)$

Statistical Analysis: Limit Distribution (Derivation)

Goal: Limit thm for $\sqrt{n} \left(\underline{W}_p(\hat{\mu}_n, \nu) - \underline{W}_p(\mu, \nu) \right)$ for fixed $\nu \in \mathcal{P}(\mathbb{R}^d)$ with $\mu \neq \nu$ and $p > 1$

Derivation: Assume compact supports and take $\Phi(\eta) := \underline{W}_p^p(\eta, \nu)$

1. \underline{W}_p^p is Lipschitz in \overline{W}_1 & $\overline{W}_1(\eta, \nu) = \|\eta - \nu\|_{\infty, \mathcal{F}}$ w/ $\mathcal{F} = \{\varphi \circ p^\theta : \varphi \in \text{Lip}_1(\mathbb{R}), \theta \in \mathbb{S}^{d-1}\}$
 - **Recall:** $\overline{W}_1(\eta, \nu) = \sup_{\substack{\varphi \in \text{Lip}_1(\mathbb{R}) \\ \theta \in \mathbb{S}^{d-1}}} (\eta - \nu)(\varphi \circ p^\theta) = \|\eta - \nu\|_{\infty, \mathcal{F}}$
 2. \mathcal{F} is a μ -Donsker class (via Donskerness of $\text{Lip}_1(\mathbb{R})$) $\iff \sqrt{n}(\hat{\mu}_n - \mu) \xrightarrow{d} G_\mu$ in $\ell^\infty(\mathcal{F})$
 3. Evaluate the Gâteaux derivative $\frac{d}{dt^+} \underline{W}_p^p(\mu + t\rho, \nu) = \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}^d} \varphi_\theta(\theta^T x) d\rho(x) d\sigma(\theta)$
- \implies Apply general theorem

Statistical Analysis: Limit Distribution (Result)

Theorem (G.-Kato-Rioux-Sadhu '22)

Fix $p > 1$. Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ be compactly supported s.t. μ is a.c. and $\text{spt}(\mu)$ is convex.

For $\theta \in \mathbb{S}^{d-1}$, let φ^θ be an OT potential from $\mathfrak{p}_\#^\theta \mu$ to $\mathfrak{p}_\#^\theta \nu$ (unique up to additive const.). Then:

$$\sqrt{n} \left(\underline{W}_p^p(\hat{\mu}_n, \nu) - \underline{W}_p^p(\mu, \nu) \right) \xrightarrow{d} \mathcal{N}(0, \sigma_p^2)$$

where $\sigma_p^2 = \int \int \text{Cov}(\varphi^\theta \circ \mathfrak{p}^\theta, \varphi^\psi \circ \mathfrak{p}^\psi) d\sigma(\theta) d\sigma(\psi)$.

Comments: All results extend to \overline{W}_p ; separate (more general) derivations for $p = 1$

1. Linearity of the Gâteaux derivative directly implies consistency of the bootstrap
2. p th power can be removed via standard delta method for the map $x \mapsto x^{1/p}$
3. Similar approach yields distributional limits for entropic/smooth OT

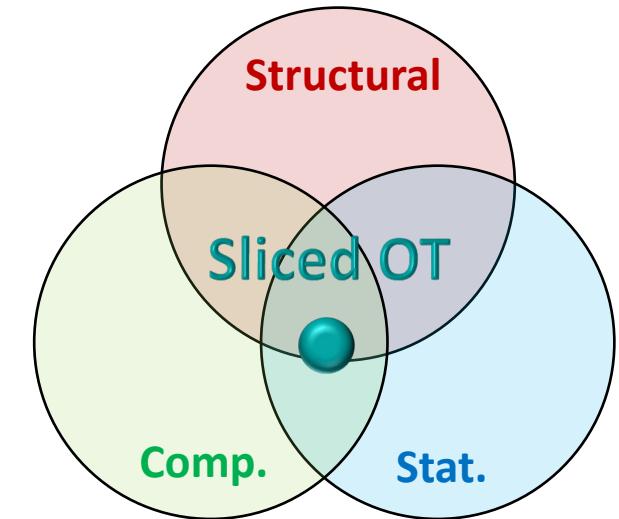
Summary

Classic W_p : Metric on $\mathcal{P}_p(\mathbb{R}^d)$ with rich & meaningful structure

- Many applications in ML and beyond
- Hard to compute & statistical curse of dimensionality $n^{-1/d}$

Sliced W_p : Project distributions to low dimension & average/maximize

- Inherits structure of W_p (metric, topology, duality, but **not** proxy)
- Easy to compute & formal guarantees [A]
- Fast empirical convergence [A] & rich limit distribution theory [B]



[A] Nietert, Sadhu, Goldfeld, Kato, “Statistical, robustness, and computational guarantees for sliced W_p ”, ArXiv:2210.09160

[B] Goldfeld, Kato, Rioux, Sadhu, “Statistical inference with regularized optimal transport”, ArXiv:2205.04283

Thank you!

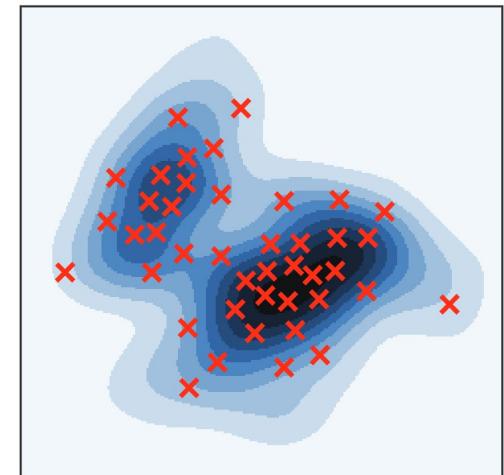
Generative Modeling: Generalization

Goal: Solve $\text{OPT} = \inf_{\theta} W_1(\mu, \nu_{\theta})$ exactly (find θ^*)

Estimation: We don't have μ and only get samples X_1, \dots, X_n

- **Empirical distribution:** $\hat{\mu}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$

→ Inherently we work with $W_1(\hat{\mu}_n, \nu_{\theta})$



Optimization: Can solve $\inf_{\theta} W_1(\hat{\mu}_n, \nu_{\theta})$ approximately

$$\text{Find } \hat{\theta}_n \text{ s.t. } W_1(\hat{\mu}_n, \nu_{\hat{\theta}_n}) \leq \inf_{\theta} W_1(\hat{\mu}_n, \nu_{\theta}) + \epsilon$$

Generalization: $W_1(\mu, \nu_{\hat{\theta}_n}) - \text{OPT} \leq 2W_1(\hat{\mu}_n, \mu) + \epsilon$

→ Boils down to empirical approximation question under W_1