

Background

The classic Wasserstein distance W_p , defined by

$$W_p(\mu, \nu) := \left(\inf_{\pi \in \Pi(\mu, \nu)} \int |x - y|^p d\pi(x, y) \right)^{1/p}$$

↑
set of couplings between μ and ν

is a popular discrepancy measure between probability measures with many applications in statistics and ML.

Motivation

Despite its proven utility, W_p suffers from a sensitivity to outliers, with its strict marginal constraints allowing a small amount of distant mass to contribute greatly to the measured distance¹. E.g., for any $\varepsilon > 0$,

$$\lim_{|x| \rightarrow \infty} W_p(\mu, (1 - \varepsilon)\mu + \varepsilon \delta_x) = \infty.$$

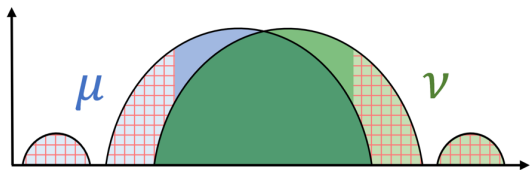
Object of Study

Outlier-robust Wasserstein distance:

robustness radius

$$W_p^\varepsilon(\mu, \nu) := \inf_{\substack{\mu', \nu' \in \mathcal{M}_+(\mathbb{R}^d) \\ \mu' \leq \mu, \nu' \leq \nu \\ \mu'(\mathbb{R}^d) = \nu'(\mathbb{R}^d) = 1 - \varepsilon}} W_p\left(\frac{\mu'}{1 - \varepsilon}, \frac{\nu'}{1 - \varepsilon}\right), \quad (1)$$

i.e. we remove an ε -fraction of mass from both μ and ν (and renormalize) to minimize their OT cost



The gridded light blue and green regions have μ and ν mass ε , respectively, and are removed to obtain optimal μ' and ν' for W_1^ε

Population-Limit Robustness Guarantees

Contamination model:

$$\|\tilde{\mu} - \mu\|_{TV}, \|\tilde{\nu} - \nu\|_{TV} \leq \varepsilon,$$

clean distributions
↓
↑ contaminated distributions

Distributional assumptions:

$$\mu, \nu \in \mathcal{D} \text{ for } \mathcal{D} \in \{\mathcal{D}_q, \mathcal{D}_2^{\text{cov}}\} \text{ where}$$

$$\mathcal{D}_q := \{\kappa \in \mathcal{P}(\mathbb{R}^d) : \mathbb{E}_\kappa[\|X - x\|^q] \leq M \text{ for some } x \in \mathbb{R}^d\},$$

$$\mathcal{D}_2^{\text{cov}} := \{\kappa \in \mathcal{P}(\mathbb{R}^d) : \Sigma_\kappa \preceq M^2 I_d\}$$

Minimax risk:

$$\widehat{W} : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R},$$

$$R_\infty(\mathcal{D}, \varepsilon) := \inf_{\widehat{W}} \sup_{\mu, \nu \in \mathcal{D}} \sup_{\substack{\|\tilde{\mu} - \mu\|_{TV} \leq \varepsilon \\ \|\tilde{\nu} - \nu\|_{TV} \leq \varepsilon}} |\widehat{W}(\tilde{\mu}, \tilde{\nu}) - W_p(\mu, \nu)|$$

Optimality of W_p^ε :

$$p < q, \quad R_\infty(\mathcal{D}_q, \varepsilon) \asymp M \varepsilon^{1/p-1/q}$$

$$p < 2, \quad R_\infty(\mathcal{D}_2^{\text{cov}}, \varepsilon) \asymp M \sqrt{d} \varepsilon^{1/p-1/2}$$

achieved by $\widehat{W} = W_p^\varepsilon$

Finite-Sample Robustness Guarantees

Contamination model:

$\mathcal{M}^{\text{AC}}(\mu, \nu, \varepsilon)$ — n i.i.d. samples from μ and ν , ε -fraction arbitrarily corrupted to obtain $\tilde{X}_1, \dots, \tilde{X}_n$ and $\tilde{Y}_1, \dots, \tilde{Y}_n$ w/ distribution P_n

Minimax risk:

Estimator \widehat{W}_n determined by corrupted samples

$$R_n(\mathcal{D}, \varepsilon) := \inf_{\widehat{W}_n} \sup_{\mu, \nu \in \mathcal{D}} \sup_{P_n \in \mathcal{M}^{\text{AC}}(\mu, \nu, \varepsilon)} \mathbb{E}_{P_n} |\widehat{W}_n - W_p(\mu, \nu)|$$

Optimality of W_p^ε :

In general, $R_n(\mathcal{D}, \varepsilon) \asymp R_\infty(\mathcal{D}, \varepsilon) + \tilde{O}(R_n(\mathcal{D}, 0))$

$$d > d_0(p, q)$$

$$p < q, \quad R_n(\mathcal{D}_q, \varepsilon) \asymp M \varepsilon^{1/p-1/q} + \tilde{O}(n^{-1/d})$$

$$p < 2, \quad R_n(\mathcal{D}_2^{\text{cov}}, \varepsilon) \asymp M \sqrt{d} \varepsilon^{1/p-1/2} + \tilde{O}(n^{-1/d})$$

achieved by $\widehat{W}_n = W_p^\varepsilon$

Duality

Kantorovich dual for classic W_p :

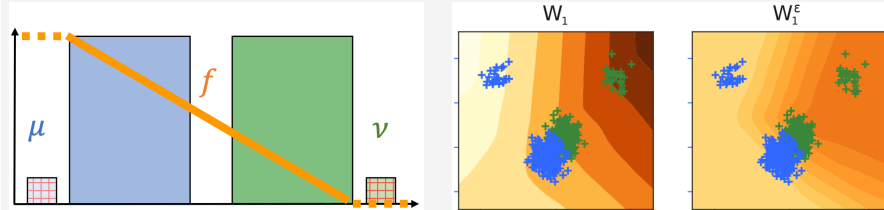
$$W_p(\mu, \nu)^p = \sup_{f \in C_b(\mathbb{R}^d)} \int f d\mu + \int f^c d\nu$$

c-transform of f , equal to $-f$ when $p = 1$
← continuous, bounded real functions

Dual for our W_p^ε :

$$(1 - \varepsilon)W_p(\mu, \nu)^p = \sup_{f \in C_b(\mathbb{R}^d)} \int f d\mu + \int f^c d\nu - 2\varepsilon \|f\|_\infty \quad (2)$$

This elegant dual form is useful for analysis and enables robustification of popular duality-based OT solvers via a simple modification.



1D densities plotted with their optimal potential for the W_1^ε dual problem

Contour plots for optimal dual potentials to W_1 and W_1^ε between 2D Gaussian mixtures

Properties

1. The infimum in (1) and the supremum in (2) are achieved
2. If f is an optimal potential for (2), then there are $\mu' = \mu - \alpha$ and $\nu' = \nu - \beta$ minimizing (1) s.t. $\text{supp}(\alpha) \subseteq \text{argmax}(f)$ and $\text{supp}(\beta) \subseteq \text{argmin}(f)$ i.e. the max and min level sets of f encode outlier locations



Samples generated by a robustified GAN (left), inspired by the dual formulation (2), alongside samples generated by standard Wasserstein GAN, after training on corrupted MNIST dataset.

1. L. Chizat, G. Peyré, B. Schmitzer, and F.-X. Vialard. Unbalanced optimal transport: dynamic and Kantorovich formulations. *Journal of Functional Analysis*, 274(11):3090–3123, 2018.