

Quantifying the Empirical Wasserstein Distance to a Set of Measures: Beating the Curse of Dimensionality

Nian Si

Joint work with Jose Blanchet, Soumyadip Ghosh, and Mark Squillante

NeurIPS 2020



October 22, 2020

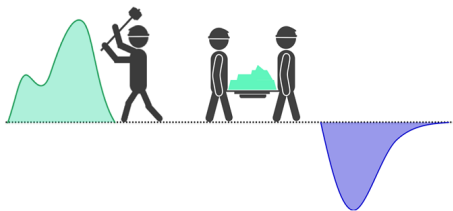
Road map

- 1 Wasserstein Distances and the Curse of Dimensionality
- 2 Robust Wasserstein Profile Function
- 3 Duality Results
 - Connections with the the Integral Probability Metric (IPM)
 - Examples
- 4 Statistical convergence

Wasserstein Distances and the Curse of Dimensionality

Definition of the Wasserstein Distance (earth mover's distance, optimal cost): for any measure P, Q ,

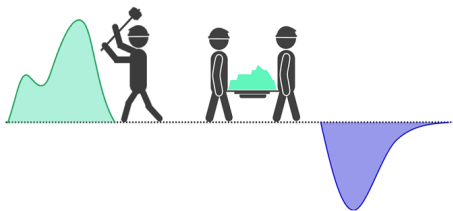
$$\mathcal{D}_c(P, Q) = \min_{\pi \in \mathcal{P}(\Omega \times \Omega)} \left\{ \left(\int c(x, w) \pi(dx, dw) \right) \right. \\ \left. : \int_{w \in \mathbb{R}^d} \pi(dx, dw) = P(dx), \int_{x \in \mathbb{R}^d} \pi(dx, dw) = Q(dw) \right\}.$$



Wasserstein Distances and the Curse of Dimensionality

Definition of the Wasserstein Distance (earth mover's distance, optimal cost): for any measure P, Q ,

$$\mathcal{D}_c(P, Q) = \min_{\pi \in \mathcal{P}(\Omega \times \Omega)} \left\{ \left(\int c(x, w) \pi(dx, dw) \right) \right. \\ \left. : \int_{w \in \mathbb{R}^d} \pi(dx, dw) = P(dx), \int_{x \in \mathbb{R}^d} \pi(dx, dw) = Q(dw) \right\}.$$

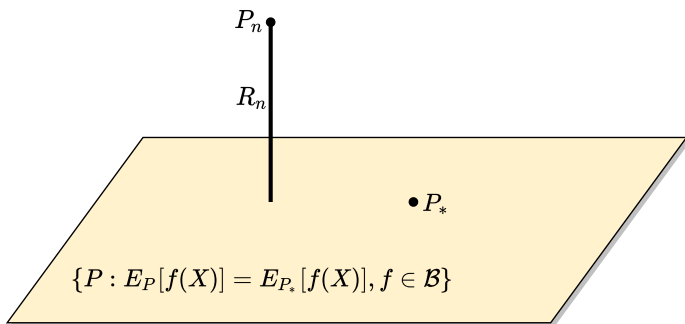


Curse of the dimensionality: $\mathcal{D}_c(P_*, P_n) = O_p(n^{-1/(d \vee 2)})$.

- How to explain the good empirical performance, e.g., Wasserstein GAN?

Robust Wasserstein Profile Function

$$R_n = \inf_{P \in \mathcal{P}(\Omega)} \{ \mathcal{D}_c(P, P_n) : \mathbb{E}_P[f(X)] = \mathbb{E}_{P_*}[f(X)], \text{ for all } f \in \mathcal{B}(\Omega) \}.$$



Duality Results

Theorem (Strong Duality)

Suppose the underlying space Ω is compact and the cost function $c(\cdot, \cdot)$ is a non-negative continuous function with $c(x, x) = 0$, for $x \in \Omega$. Then, we have the strong duality

$$\begin{aligned} R_n &:= \inf_{P \in \mathcal{P}(\Omega)} \{ \mathcal{D}_c(P, P_n) : \mathbb{E}_P[f(X)] = \mathbb{E}_{P_n}[f(X)], \text{ for all } f \in \mathcal{B}(\Omega) \}. \\ &= \sup_{f \in \mathcal{LB}(\Omega)} \{ \mathbb{E}_{P_*}[f(X)] - \mathbb{E}_{P_n}[f^c(X)] \}, \end{aligned}$$

where $f^c(x) = \sup_{z \in \Omega} \{ f(z) - c(z, x) \}$ and $\mathcal{LB}(\Omega)$ denotes the linear span generated by $\mathcal{B}(\Omega)$, namely

$$\mathcal{LB}(\Omega) = \left\{ f(\cdot) = \sum_{i=1}^m \lambda_i f_i(\cdot) : \{f_i(\cdot)\}_{i=1}^m \subset \mathcal{B}(\Omega), \lambda \in \mathbb{R}^m, \text{ and } m \in \mathbb{Z}_+ \right\}.$$

Duality Results

Theorem (Strong Duality)

Suppose the underlying space Ω is compact and the cost function $c(\cdot, \cdot)$ is a non-negative continuous function with $c(x, x) = 0$, for $x \in \Omega$. Then, we have the strong duality

$$\begin{aligned} R_n &:= \inf_{P \in \mathcal{P}(\Omega)} \{ \mathcal{D}_c(P, P_n) : \mathbb{E}_P[f(X)] = \mathbb{E}_{P_n}[f(X)], \text{ for all } f \in \mathcal{B}(\Omega) \}. \\ &= \sup_{f \in \mathcal{LB}(\Omega)} \{ \mathbb{E}_{P_n^*}[f(X)] - \mathbb{E}_{P_n}[f^c(X)] \}, \end{aligned}$$

where $f^c(x) = \sup_{z \in \Omega} \{ f(z) - c(z, x) \}$ and $\mathcal{LB}(\Omega)$ denotes the linear span generated by $\mathcal{B}(\Omega)$, namely

$$\mathcal{LB}(\Omega) = \left\{ f(\cdot) = \sum_{i=1}^m \lambda_i f_i(\cdot) : \{f_i(\cdot)\}_{i=1}^m \subset \mathcal{B}(\Omega), \lambda \in \mathbb{R}^m, \text{ and } m \in \mathbb{Z}_+ \right\}.$$

Connections with the the Integral Probability Metric (IPM)

$$\text{IPM}_{\mathcal{F}}(P, P_n) = \sup_{f \in \mathcal{F}} \left| \int f dP - \int f dP_n \right|.$$

Connections with the the Integral Probability Metric (IPM)

$$\text{IPM}_{\mathcal{F}}(P, P_n) = \sup_{f \in \mathcal{F}} \left| \int f dP - \int f dP_n \right|.$$

- R_n is not a metric in general.

Connections with the the Integral Probability Metric (IPM)

$$\text{IPM}_{\mathcal{F}}(P, P_n) = \sup_{f \in \mathcal{F}} \left| \int f dP - \int f dP_n \right|.$$

- R_n is not a metric in general.
- We add a new modeling feature, which is the hypothesis class.

Connections with the the Integral Probability Metric (IPM)

$$\text{IPM}_{\mathcal{F}}(P, P_n) = \sup_{f \in \mathcal{F}} \left| \int f dP - \int f dP_n \right|.$$

- R_n is not a metric in general.
- We add a new modeling feature, which is the hypothesis class.
- Our expression for the strong duality uses the combination of both the function f and its c -conjugate f^c in contrast with IPM.

Examples

1. When $\mathcal{B}(\Omega)$ is the space of all 1-Lipschitz functions, $f^c(x) = f(x)$ and R_n reduces to 1-Wasserstein distance. Then, we recover the Kantorovich-Rubinstein duality result:

$$R_n = \sup_{f \in \text{Lip}_1(\Omega)} \{ \mathbb{E}_{P_*} [f(X)] - \mathbb{E}_{P_n} [f(X)] \} = \mathcal{D}_1(P_*, P_n).$$

Examples

1. When $\mathcal{B}(\Omega)$ is the space of all 1-Lipschitz functions, $f^c(x) = f(x)$ and R_n reduces to 1-Wasserstein distance. Then, we recover the Kantorovich-Rubinstein duality result:

$$R_n = \sup_{f \in \text{Lip}_1(\Omega)} \{ \mathbb{E}_{P_*} [f(X)] - \mathbb{E}_{P_n} [f(X)] \} = \mathcal{D}_1(P_*, P_n).$$

2. Suppose that $\mathcal{B}(\Omega)$ is finite dimensional, such as $\mathcal{B}(\Omega) = \{f_i(x)\}_{i=1}^K$. Then, we have

$$R_n = \sup_{\lambda \in \mathbb{R}^K} \left\{ \mathbb{E}_{P_*} \left[\sum_{i=1}^K \lambda_i f_i(X) \right] - \mathbb{E}_{P_n} \left[\sup_{z \in \Omega} \left\{ \sum_{i=1}^K \lambda_i f_i(z) - c(z, X) \right\} \right] \right\},$$

which recovers the duality result obtained in Blanchet et al. (2019).

Examples

3. Infinite dimensional case: fix linearly independent unit vectors

$\theta_1, \dots, \theta_K \in \mathbb{R}^d$, and consider function class

$\mathcal{B}(\Omega) = \cup_{i=1}^K \{f(\theta_i^\top \cdot)|_\Omega : f \in \mathcal{F}_B\}$, where \mathcal{F}_B collects some 1-dimensional continuous functions, in which case

$$\mathcal{LB}(\Omega) = \left\{ f(\cdot) = \sum_{i=1}^K \lambda_i f_i(\theta_i^\top \cdot)|_\Omega : \{f_i(\cdot)\}_{i=1}^K \subset \mathcal{F}_B, \lambda \in \mathbb{R}^K \right\}.$$

Examples

3. Infinite dimensional case: fix linearly independent unit vectors

$\theta_1, \dots, \theta_K \in \mathbb{R}^d$, and consider function class

$\mathcal{B}(\Omega) = \cup_{i=1}^K \{f(\theta_i^\top \cdot)|_\Omega : f \in \mathcal{F}_B\}$, where \mathcal{F}_B collects some 1-dimensional continuous functions, in which case

$$\mathcal{LB}(\Omega) = \left\{ f(\cdot) = \sum_{i=1}^K \lambda_i f_i(\theta_i^\top \cdot)|_\Omega : \{f_i(\cdot)\}_{i=1}^K \subset \mathcal{F}_B, \lambda \in \mathbb{R}^K \right\}.$$

Theorem

Following the setting in Example 3 and for $\Omega = \mathbb{R}^d$, we have the strong duality:

$$R_n = \sup_{f \in \mathcal{LB}(\mathbb{R}^d)} \{ \mathbb{E}_{P_*} [f(X)] - \mathbb{E}_{P_n} [f^c(X)] \}.$$

Statistical convergence

Theorem

Consider function class $\mathcal{B}(\Omega) = \cup_{i=1}^K \{f(\theta_i^\top \cdot)|_\Omega : f \in \mathcal{F}_B\}$ in Example 3. We assume the space Ω is compact and some technical conditions on the function class \mathcal{F}_B , we have

$$nR_n \Rightarrow \sup_{f \in \mathcal{LB}(\Omega)} \left\{ -2H^f - \mathbb{E}_{P_*} \left[\|\nabla_X f(X)\|_2^2 \right] \right\},$$

where $\nabla_x f(x)$ is the gradient of $f(\cdot)$ evaluated at x and H^f is a Gaussian process indexed by f with

$$H^f \sim \mathcal{N}(0, \text{var}(f(X))) \text{ and } \text{cov}(H^{f_1}, H^{f_2}) = \text{cov}(f_1(X), f_2(X)).$$

Reference

Si, Nian, Jose Blanchet, Soumyadip Ghosh, and Mark Squillante. "Quantifying the Empirical Wasserstein Distance to a Set of Measures: Beating the Curse of Dimensionality" NeurIPS 2020

Thanks!