

Wasserstein Distance

Tong Zhou

November 30, 2020

Contents

| | | |
|----------|---|----------|
| 1 | Definitions | 2 |
| 2 | Properties of Wasserstein distance | 2 |
| 3 | Relation to Other Metrics | 2 |

1 Definitions

The p -Wasserstein distance between probability measures μ and ν on \mathbb{R}^d is defined as

$$W_p(\mu, \nu) = \inf_{\substack{X \sim \mu \\ Y \sim \nu}} (\mathbb{E} [\|X - Y\|^p])^{1/p}, \quad p \geq 1.$$

Usually, $W_p(\mu, \nu)$ is also denoted by $W_p(X, Y)$.

Remark.

- If a general metric space (\mathcal{X}, ρ) is complete and separable, the above norm in the definition is replaced by $\rho(X, Y)$.
- W_p is a proper metric: it is nonnegative, symmetric and satisfies triangle inequality.

2 Properties of Wasserstein distance

Let X and Y be random variables taking values in $\mathcal{X} = \mathbb{R}^d$; the notation $(\mathcal{X}, \|\cdot\|)$ is maintained.

- For any real number a , $W_p(aX, aY) = |a| W_p(X, Y)$.
- For any fixed vector $x \in \mathcal{X}$, $W_p(X + x, Y + x) = W_p(X, Y)$.
- For any fixed $x \in \mathcal{X}$, we have $W_2^2(X + x, Y) = \|x + \mathbb{E}[X] - \mathbb{E}[Y]\|^2 + W_2^2(X, Y)$.
- For product measures and when $p = 2$, we have $W_2^2(\otimes_{i=1}^n \mu_i, \otimes_{i=1}^n \nu_i) = \sum_{i=1}^n W_2^2(\mu_i, \nu_i)$ in the analytic notation.

3 Relation to Other Metrics

For random variables X and Y on \mathcal{X} , let Ω be the union of their ranges and set

$$D = \sup_{x, y \in \Omega} \|x - y\|, \quad d_{\min} = \inf_{x \neq y \in \Omega} \|x - y\|$$

It possesses the following properties:

- If $p \leq q$, then $W_p \leq W_q$, by Jensen's inequality.
- On the other hand, $W_q^q \leq W_p^p D^{q-p}$.
- Duality arguments yield the particularly useful *Kantorovich-Rubinstein* representation for W_1 as

$$W_1(X, Y) = \sup_{\|f\|_{\text{Lip}} \leq 1} |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]|, \quad \|f\|_{\text{Lip}} = \sup_{x \neq y} \frac{|f(x) - f(y)|}{\|x - y\|},$$

- This shows that W_1 is larger than the *Bounded Lipschitz* (BL) metric

$$W_1(X, Y) \geq \text{BL}(X, Y) = \sup_{\|f\|_\infty + \|f\|_{\text{Lip}} \leq 1} |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]|$$

that metrises convergence in distribution.

- Let P denote the Prokhorov distance. Then $P^2(X, Y) \leq W_1(X, Y) \leq (D + 1)P(X, Y)$.
- For the class of random variables supported on a fixed bounded subset $K \subset \mathcal{X}$, BL and W_1 are equivalent up to constant, and all metrics W_p are topologically equivalent.
- The Wasserstein distances W_p can be bounded by a version of total variation TV. A weaker but more explicit bound for $p = 1$ is $W_1(X, Y) \leq D \times \text{TV}(X, Y)$.
- For discrete random variables, there is an opposite bound $\text{TV} \leq \frac{W_1}{d_{\min}}$.
- The total variation between convolutions with a sufficiently smooth measure is bounded above by W_1 .
- The *Toscani* (or *Toscani-Fourier*) distance is also bounded above by W_1 .

$$\beta_\lambda(\mathcal{T}, \mathbf{W}) = (\mathbf{W} \mathbf{X} \mathbf{X}^\top \mathbf{W}^\top + \lambda \mathbf{I})^{-1} \mathbf{W} \mathbf{X} \mathbf{y},$$

We may now calculate the bias and variance of the model described above via the following formulations:

$$\begin{aligned} \text{Bias}_\lambda(\theta)^2 &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{T}, \mathbf{W}} [f_\lambda(\mathbf{x}; \mathcal{T}, \mathbf{W})] - f_0(\mathbf{x}) \right]^2, \\ \text{Variance}_\lambda(\theta) &= \mathbb{E}_{\mathbf{x}} \left[\text{Var}_{\mathcal{T}, \mathbf{W}} (f_\lambda(\mathbf{x}; \mathcal{T}, \mathbf{W})) \right]. \end{aligned}$$