

Obtaining Fairness using Optimal Transport Theory

ICML 2019, Long Beach

E. del Barrio^a, F. Gamboa^b, **P. Gordaliza**^{a,b} and JM. Loubes^b

a : IMUVa (Universidad de Valladolid)

b : IMT (Université de Toulouse)



Framework for achieving Fairness

Target class

$$Y = \begin{cases} 0 & \text{failure} \\ 1 & \text{success} \end{cases}$$

Visible attributes

$$X \in \mathbb{R}^d, d \geq 1,$$

Protected attribute

$$S = \begin{cases} 0 & \text{unfavored} \\ 1 & \text{favored} \end{cases}$$

Goal: Replace X by \tilde{X} such that for all $g \in \mathcal{G}$
 $\mathcal{L}(g(\tilde{X}) | S = 0) = \mathcal{L}(g(\tilde{X}) | S = 1)$

Methodology: Find $\tilde{X} = T_S(X)$ such that
 $\mathcal{L}(T_0(X) | S = 0) = \mathcal{L}(T_1(X) | S = 1)$

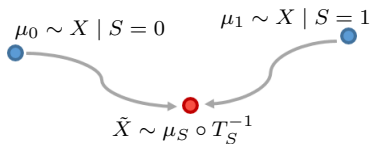
Questions:

- Best choice for the distribution $\tilde{X} \sim \nu$?
- Optimal way of transporting μ_1 and μ_0 to ν ?

Reasonable and feasible solutions: T_S optimal transport map carrying μ_S towards their Wasserstein barycenter μ_B with weights $\pi_0 = P(S = 0)$ and $\pi_1 = P(S = 1)$:

$$\mu_{S\#} T_S = \mu_B$$

$$\mu_B \in \operatorname{argmin}_{\nu \in \mathcal{P}_2} \{ \pi_0 W_2^2(\mu_0, \nu) + \pi_1 W_2^2(\mu_1, \nu) \}$$



Our proposal

- Justification for repair with Wasserstein Barycenter:** under some regularity conditions, $\mathcal{E}(\tilde{X}) := R_B(\tilde{X}) - R_B(X, S)$

$$\mathcal{E}(\tilde{X}) \leq 2\sqrt{2}K \left(\sum_{s=0,1} \pi_s \mathcal{W}_2^2(\mu_s, \mu_{s\neq} T_s) \right)^{\frac{1}{2}}, K > 0$$

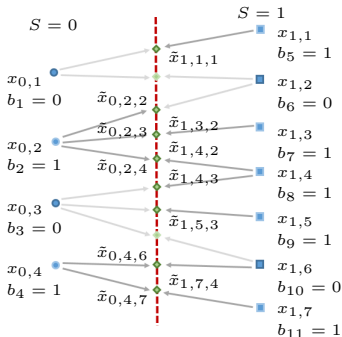
- Multidimensional extension for Total Repair**
- New Partial Repair method: Random Repair**
(Improvement of Geometric Repair)

For $s \in \{0, 1\}$,

$$\tilde{\mu}_{s,\lambda} = \mathcal{L}(BT_s(X) + (1 - B)X \mid S = s),$$

$B \sim \mathcal{B}(\lambda)$ with $\lambda \in (0, 1)$ level of repair

$\lambda = 0$ $\lambda = 1$
 Accuracy ← **Trade-off** → Blurring
 of $g(\tilde{X})$ of S



Related work



E. DEL BARRIO AND J.-M. LOUBES. (2019) Central limit theorems for empirical transportation cost in general dimension. *The Annals of Probability*, **47**, 926–951.



E. DEL BARRIO, P. GORDALIZA AND J.-M. LOUBES. (2019) A central limit theorem on the real line with application to fairness assessment in machine learning. *Information and Inference: A Journal of the IMA*.

Thanks for the attention!