

---

## Fair Regression with Wasserstein Barycenters

*for Sorbonne Université*

Eyal Cohen<sup>1</sup>   Eline Pot<sup>1</sup>

<sup>1</sup>LPSM

*Supervisor(s):*

Christophe Denis<sup>1</sup>, Rafael Pinot<sup>1</sup>

---

# Table of Contents

- Problem Statement

  - General regression problem

- Fair regression and the problem of Wasserstein barycenter

- Statistical Guarantees

- Experiments



# Notations and General regression problem

## 1 Problem Statement


$$Y = f^*(X, S) + \xi \quad (1)$$

with  $\xi \in \mathbb{R}$  a centered r.v.

$(X, S) \sim \mathbb{P}_{X,S}$  on  $\mathbb{R}^d \times S$ , where  $S$  is the sensible attribute ;  $|S| < \infty$

$f^* : \mathbb{R}^d \times S \mapsto \mathbb{R}$  the regression function minimizing the squared risk.

$\forall f : \mathbb{R}^d \times S \mapsto \mathbb{R}$ ,  $\nu_{f|S}$  is the distribution of  $f(X, S)|S = s$ .

$$\text{i.e. } F_{\nu_{f|S}}(t) = \mathbb{P}(f(X, S) \leq t | S = s) \quad (2)$$

### Definition (Demographic Parity (DP))

$g : \mathbb{R}^d \times S \mapsto \mathbb{R}$  is fair if :

$$\forall s, s' \in S, \sup_{t \in \mathbb{R}} |\mathbb{P}(g(X, S) \leq t | S = s) - \mathbb{P}(g(X, S) \leq t | S = s')| = 0$$

In particular, we note that if  $g$  is fair, then  $\nu_{g|S} = \nu_g$

# Table of Contents

- ▶ Problem Statement
- ▶ Fair regression and the problem of Wasserstein barycenter
  - Characterization of fair optimal transport
  - General Form of the estimator
- ▶ Statistical Guarantees
- ▶ Experiments



# Wasserstein Distance

## 2 Fair regression and the problem of Wasserstein barycenter

### Definition (Wasserstein-2 distance)

For all  $\mu, \nu$  univariate probabilities measures, we define the Wasserstein-2 distance :

$$\mathcal{W}_2^2(\mu, \nu) = \inf_{\gamma \in \Gamma_{\mu, \nu}} \int |x - y|^2 d\gamma(x, y)$$

where  $\Gamma_{\mu, \nu}$  is the set of coupling measures on  $\mathbb{R} \times \mathbb{R}$

### Theorem (Monge map)

Let  $\mu$  and  $\nu$  be two univariate measures such that  $\nu$  has a density. And let  $X \sim \nu$ . Then there exists a mapping

$T : \mathbb{R} \mapsto \mathbb{R}$  such that  $\mathcal{W}_2^2(\mu, \nu) = \mathbb{E}(X - T(X))^2$ .

i.e.  $(X, T(X)) \sim \gamma^* \in \Gamma_{\mu, \nu}$  with  $\gamma^*$  and optimal coupling measure.

Also, we have :  $T = Q_\mu \circ F_\nu$

### Theorem (Inverse of the cdf)

Let  $\nu_1, \dots, \nu_{|S|}$  be  $|S|$  univariate probability measures admitting densities. Let  $p_1, \dots, p_{|S|} \geq 0$  s.t.  $\sum_{s=1}^{|S|} p_s = 1$ . Let us

denote :  $\nu^* \in \arg \min_{\nu} \sum_{s=1}^{|S|} p_s \mathcal{W}_2^2(\nu_s, \nu)$ .

Then :  $F_{\nu^*}^*(\cdot) = \left( \sum_{s=1}^{|S|} p_s Q_{\nu_s} \right)^{\leftarrow}(\cdot)$ .

### Theorem (Characterization of fair optimal prediction)

Let us assume  $\forall s \in S$  that the univariate measure  $\nu_{f^*|s}$  has a density and let  $p_s = \mathbb{P}(S = s)$ .  
Then :

$$\min_{g \text{ is fair}} \mathbb{E}[(f^*(X, S) - g(X, S))^2] = \min_v \sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|s}, \nu)$$

Moreover if  $g^*$  solves l.h.s and  $\nu^*$  the r.h.s, then  $\nu^* = \nu_g^*$  and :

$$g^*(x, s) = \left( \sum_{s' \in S} p_{s'} Q_{f^*|s'} \right) \circ F_{f^*|s}(f^*(x, s))$$

# Proof

## 2 Fair regression and the problem of Wasserstein barycenter

Let us show that

$$\min_{g \text{ fair}} \mathbb{E} [(f^*(X, S) - g(X, S))^2] = \min_v \sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|s}, \nu)$$

Let us denote  $\bar{g}$  a minimizer of the l.h.s., and  $\nu_{\bar{g}}$  its distribution. As  $\nu_{f^*|s}$  admits a density, we can use the precedent theorem : for each  $s \in S$ , there exists  $T_s = Q_{\nu_{\bar{g}}} \circ F_{f^*|s}$  such that with  $\tilde{g} = T_s = Q_{\nu_{\bar{g}}} \circ F_{f^*|s} \circ f^*$ :

$$\sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|s}, \nu_{\bar{g}}) = \mathbb{E} [(f^*(X, S) - \tilde{g}(X, S))^2] \quad (3)$$

We can write:

$$\mathbb{P}(\tilde{g}(X, S) \leq t) = \sum_{s \in S} p_s \mathbb{P}(Q_{\nu_{\bar{g}}} \circ F_{f^*|s} \circ f^*(X, s) \leq t) \quad (4)$$

$$\begin{aligned} &= \sum_{s \in S} p_s \mathbb{P}(f^*(X, s) \leq Q_{f^*|s} \circ F_{\nu_{\bar{g}}}(t)) \\ &= \sum_{s \in S} p_s F_{f^*}(X, s) (Q_{f^*|s} \circ F_{\nu_{\bar{g}}}(t)) \\ &= \sum_{s \in S} p_s F_{\nu_{\bar{g}}}(t) \end{aligned} \quad (5)$$

# Proof

## 2 Fair regression and the problem of Wasserstein barycenter

$$\mathbb{P}(\tilde{g}(X, S) \leq t) = \sum_{s \in S} p_s F_{\nu_{\tilde{g}}}(t) \quad (5)$$

And since  $\tilde{g}$  is fair by definition, this implies that  $\tilde{g}$  is also fair.

We now define  $T^* = Q_{\nu^*} \circ F_{f^*|S}$  as the optimal transport map between  $\nu_{f^*|S}$  and

$\nu^* = \operatorname{argmin}_{\nu} (\sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|S}, \nu))$ , and we define  $g^* = T^* \circ f^*$ .

By the result of the Monge map theorem, we get:

$$\min_{\nu} \sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|S}, \nu) = \mathbb{E}[(f^*(X, S) - g^*(X, S))^2]$$

$g^*$  minimizing the LHS, we note that  $\nu^*$  is independent of  $s$  by construction, so  $g^*$  is DP. Thus we can bound the LHS by taking the minimum over the DP-fair set:

$$\min_{\nu} \sum_{s \in S} p_s \mathcal{W}_2^2(\nu_{f^*|S}, \nu) \geq \min_{g \text{ fair}} \mathbb{E}[(f^*(X, S) - g(X, S))^2] \quad (6)$$

But, by optimality of  $\tilde{g}$ , we have  $\mathbb{E}[(f^*(X, S) - \tilde{g}(X, S))^2] \geq \mathbb{E}[(f^*(X, S) - g^*(X, S))^2]$ .

But we also have that:

$$\begin{aligned} \mathcal{W}_2^2(\nu_{f^*|S}, \nu_{\tilde{g}}) &\leq \mathbb{E}[(f^*(X, S) - g^*(X, S))^2 | S = s] \\ \implies \min_{\nu} \sum_{s \in S} \mathcal{W}_2^2(\nu_{f^*|S}, \nu_{\tilde{g}}) &\leq \min_{g \text{ is fair}} \mathbb{E}[(f^*(X, S) - g(X, S))^2] \end{aligned}$$

This gives the equality, we can then set  $\tilde{g} = g^*$ , as  $g^*$  is fair.



# General Form of the estimator

## 2 Fair regression and the problem of Wasserstein barycenter

Given a base estimator  $\hat{f}$  of  $f^*$ , we define the final estimator  $\hat{g}$  of  $g^*$  as :

$$\hat{g}(x, s) = \left( \sum_{s' \in S} \hat{p}_{s'} \hat{Q}_{\hat{f}|s'} \right) \circ \hat{F}_{\hat{f}|s}(\hat{f}(x, s) + \varepsilon)$$

where  $\varepsilon \sim \mathcal{U}([- \sigma, \sigma])$ , independant from the others r.v.

With  $\hat{p}_s$  the empirical frequency of  $S = s$  evaluated on  $\{S_i\}_{i=1 \dots N} \stackrel{iid}{\sim} \mathbb{P}_S$

With  $\hat{F}_{f|s}$  and  $\hat{Q}_{f|s}$  the empirical CDF and quantile function of  $(f(X, S) + \varepsilon)|S = s$  based on  $\{f(X_i^s, S) + \varepsilon_{is}\}_{i \in I_1^s}$  and  $\{f(X_i^s, S) + \varepsilon_{is}\}_{i \in I_0^s}$  respectively.

They are defined as :

$$\hat{F}_{f|s} = F_{\hat{v}_{f|s}^1} \quad \text{and} \quad \hat{Q}_{f|s} = Q_{\hat{v}_{f|s}^0}$$

Here,  $\hat{v}_{f|s}^0$  and  $\hat{v}_{f|s}^1$  are estimators of  $v_{f|s}$ . Fixing  $I_0^s$  and  $I_1^s$  as an equal partition of  $[N_s]$ , these estimators read as :

$$\hat{v}_{f|s}^0 = \frac{1}{|I_0^s|} \sum_{i \in I_0^s} \delta(f(X_i^s, s) + \varepsilon_{is} - \cdot) \quad \text{and} \quad \hat{v}_{f|s}^1 = \frac{1}{|I_1^s|} \sum_{i \in I_1^s} \delta(f(X_i^s, s) + \varepsilon_{is} - \cdot) \quad \text{with} \quad \varepsilon_{is} \stackrel{iid}{\sim} \mathcal{U}([- \sigma, \sigma])$$

# Table of Contents

- ▶ Problem Statement
- ▶ Fair regression and the problem of Wasserstein barycenter
- ▶ **Statistical Guarantees**
- ▶ Experiments



### Theorem (Fairness guarantee)

For any  $\mathbb{P}$  on  $(X, S, Y)$ , for any  $\hat{f}$  constructed on labeled data, and for any  $s, s' \in S$ , we have fairness guarantee in expectation over the data :

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(\hat{g}(X, S) \leq t | S = s) - \mathbb{P}(\hat{g}(X, S) \leq t | S = s')| \leq \frac{2}{\min(N_s, N_{s'}) + 2} \mathbb{1}_{N_s \neq N_{s'}}$$

We also have the bound over the expected (over data) violation of the fairness definition (with  $\mathcal{D}$  the union of all the datasets):

$$\mathbb{E} \left[ \sup_{t \in \mathbb{R}} |\mathbb{P}(\hat{g}(X, S) \leq t | S = s, \mathcal{D}) - \mathbb{P}(\hat{g}(X, S) \leq t | S = s', \mathcal{D})| \right] \leq \frac{6}{\sqrt{\min(N_s, N_{s'}) + 1}}$$

**PROOF IDEAS:**  $\rightarrow$  For the first result, we use that

$$\forall s, s' \in S, \sup_{t \in \mathbb{R}} |\mathbb{P}(\hat{g}(X^s, s) \leq t) - \mathbb{P}(\hat{g}(X^{s'}, s') \leq t)| \leq \sup_{t \in (0, 1)} |\mathbb{P}(\hat{F}_{\hat{f}|s}(\hat{f}(X^s, s) + \varepsilon) \leq t) - \mathbb{P}(\hat{F}_{\hat{f}|s'}(\hat{f}(X^{s'}, s') + \varepsilon) \leq t)|$$

Then, fixing  $t \in (0, 1)$  and  $k_s(t) \in \{1, \dots, |I_1^s|\}$  such that  $t \in \left[ \frac{k_s(t)-1}{|I_1^s|}, \frac{k_s(t)}{|I_1^s|} \right)$ , and using the fact that conditionally on labeled data, the random variables

$$\sum_{i \in I_1^s} \mathbb{1}_{\hat{f}(X_i^s) + \varepsilon_{is} \leq \hat{f}(x, s) + \varepsilon} \sim \mathcal{U}(\{0, \dots, |I_1^s|\}), \text{ we have: } \mathbb{P}(\hat{F}_{\hat{f}|s}(\hat{f}(X^s, s) + \varepsilon) \leq t) = \frac{k_s(t)}{|I_1^s| + 1}$$

$$\text{We use the same argument for } s' \text{ to get: } \sup_{t \in \mathbb{R}} |\mathbb{P}(\hat{g}(X^s, s) \leq t) - \mathbb{P}(\hat{g}(X^{s'}, s') \leq t)| \leq \sup_{t \in (0, 1)} \left| \frac{k_s(t)}{(N_s/2+1)} - \frac{k_{s'}(t)}{(N_{s'}/2+1)} \right|$$

$\rightarrow$  For the second result, the beginning of the demonstration is the same, then we use the triangle inequality to finally get the bound

$$2\mathbb{E}\|F_{\hat{f}|s} - \hat{F}_{\hat{f}|s}\|_{\infty} + 2\mathbb{E}\|F_{\hat{f}|s'} - \hat{F}_{\hat{f}|s'}\|_{\infty}, \text{ where } F_{\hat{f}|s} = \mathbb{P}(\hat{f}(X^s, s) + \varepsilon \leq t | \mathcal{D}). \text{ Then, we conclude applying DKL inequality, conditionally on } \mathcal{L}.$$

# Estimation guarantee

## 3 Statistical Guarantees

Let us make the following assumptions :

- For any  $s \in \mathcal{S}$ ,  $\nu_{f^*|s}$  admits a density  $q_s$  such that  $0 < \underline{\lambda}_s \leq q_s \leq \bar{\lambda}_s$ .
- There exist  $c$  and  $C$  positive and independent of  $n, N, N_1, \dots, N_{|\mathcal{S}|}$  and  $b_n$  a positive sequence such that for all  $\delta > 0$ , we have :

$$\mathbb{P}(|f^*(x, s) - \hat{f}(x, s)| \geq \delta) \leq \exp(-Cb_n\delta^2)$$

### Theorem (Estimation guarantee)

Under these assumptions, let us set  $\sigma \lesssim \min_{s \in \mathcal{S}} \{ \max(\frac{1}{\sqrt{N_s}}, \frac{1}{\sqrt{bn}}) \}$ , then  $\hat{g}$  verifies :

$$\mathbb{E}[|g^*(X, S) - \hat{g}(X, S)|] \lesssim \min \left\{ \frac{1}{\sqrt{b_n}}; \sum_{s \in \mathcal{S}} p_s \frac{1}{\sqrt{N_s}}; \sqrt{\frac{|\mathcal{S}|}{N}} \right\}$$

# Table of Contents

- ▶ Problem Statement
- ▶ Fair regression and the problem of Wasserstein barycenter
- ▶ Statistical Guarantees
- ▶ Experiments

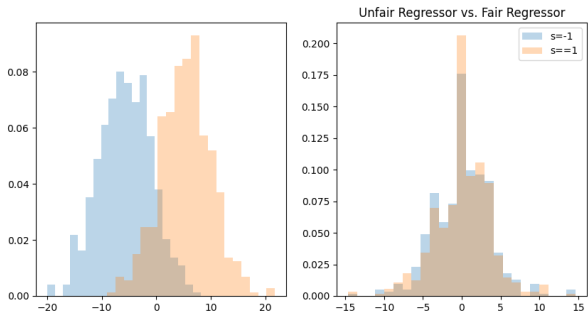


# Experiments

## 4 Experiments

Gaussian distributions:

- Binary:

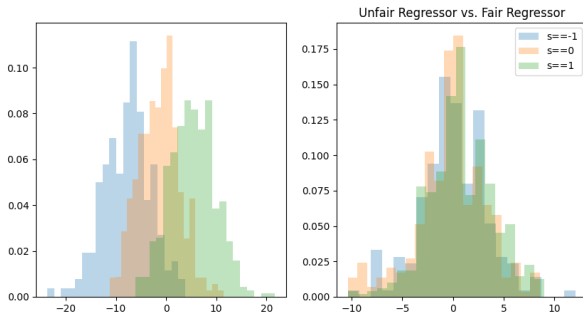


# Experiments

## 4 Experiments

Gaussian distributions:

- Three values:

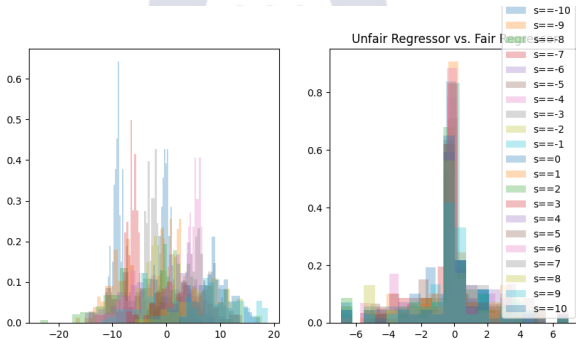


# Experiments

## 4 Experiments

Gaussian distributions:

- More values (21):



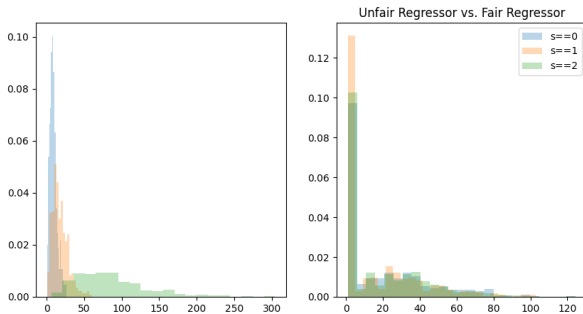


# Experiments

## 4 Experiments

Exponential distributions:

- Three values:





## **Fair Regression with Wasserstein Barycenters**

Thank you for listening !  
Any Questions ?