# Review of 'Fair Regression with Wasserstein Barycenters'
## Chzhen et al., NeurIPS 2020

Presented by Kunwoong Kim

Department of Statistics, Seoul National University

September 30, 2025

- **Problem**: (Group-)fair regression
  We aim to find a function that minimizes the mean squared error under the demographic parity constraint.

- **Idea**: Alignment of predictions using Wasserstein barycenter.

- **Proposed method**: A post-processing algorithm for perfect fairness.

- Variables
  - $X \in \mathbb{R}^d$ : an input random vector
  - $Y \in \mathbb{R}$ : a real-valued output
  - $S \in \mathcal{S}$ : a sensitive attribute (e.g., $\mathcal{S} = \{0, 1\}$)
- Distributions
  - $\mathbb{P}$ : the joint distribution of $(X, S, Y)$.
  - $\mathbb{P}_{X,S}$ : the marginal distribution of $(X, S)$.
- Cumulative Distribution Function (CDF)
  For a given probability measure $\mu$, we denote $F_\mu$ as the CDF of $\mu$.
- Quantile Function
  For a given probability measure $\mu$, we denote $Q_\mu : [0, 1] \to \mathbb{R}$ as the quantile function of $\mu$. That is, $Q_\mu(t) = \inf\{y \in \mathbb{R} : F_\mu(y) > t\}$ for $t \in (0, 1]$.

- A standard regression model:
$$Y = f(X, S) + \eta,$$

  where $\eta \in \mathbb{R}$ is a centered random variable.

- Let $f^*$ be the true regression function such that
$$f^*(x, s) = \mathbb{E}\left(Y | X = x, S = s\right).$$

- Given $f$, denote $\nu_{f|s}$ as the conditional distribution of $f(X, S)|S = s$.
  The CDF of $\nu_{f|s}$ is given by
$$F_{\nu_{f|s}}(t) = \mathbb{P}(f(X, S) \leq t | S = s).$$

### Definition 1 ((Strong) demographic parity)

A prediction model $g : \mathbb{R}^d \times \mathcal{S} \to \mathbb{R}$ is fair if, for every $s, s' \in \mathcal{S}$

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}(g(X, S) \leq t | S = s) - \mathbb{P}(g(X, S) \leq t | S = s') \right| = 0. \tag{1}$$

- Strong demographic parity defined in this paper requires the Kolmogorov-Smirnov distance to be zero for all $s, s'$.

## Main results

<div style="border:1px solid">

### Theorem 2

Let $p_s := \mathbb{P}(S = s)$. Assume that $\nu_{f^*|s}$ has a density for each $s \in \mathcal{S}$. Then, we have

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

Moreover, if $g^*$ and $\nu^*$ solve the left-hand-side and the right-hand-side of Equation (2) respectively, then $\nu^* = \nu_{g^*}$ and

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

</div>

- Implication: We can obtain an optimal fair regression model by: sequentially doing (i) quantile matching and (ii) transforming to barycenter.

In other words, the optimal fair prediction model $g^*$ is a transformation of $f^*$ defined by

$$g^*(x,s) = p_s f^*(x,s) + (1-p_s)t^*(x,s),$$

where $t^*$ is a correction so that the quantile of $f^*(X,s)$ is the same as the quantile of $f^*(X,s')$ for $s \neq s'$.
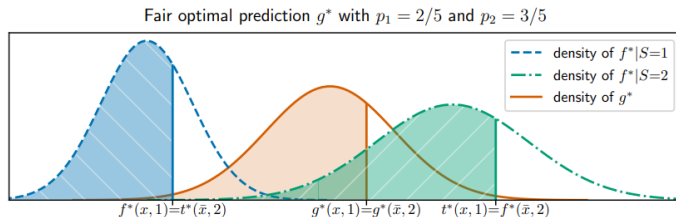


Figure 1: For a new point $(x,1)$, the value $t^*(x,1)$ is chosen such that the shaded `Green Area` (//) $= \mathbb{P}(f^*(X,S) \leq t^*(x,1)|S=2)$ equals to the shaded `Blue Area` (\\) $= \mathbb{P}(f^*(X,S) \leq f^*(x,1)|S=1)$. The final prediction $g^*(x,1)$ is a convex combination of $f^*(x,1)$ and $t^*(x,1)$. The same is done for $(\bar{x},2)$.

## Main results

- Let $\mathcal{D}_n := \{(x_i, s_i, y_i)\}_{i=1}^n$ be a given dataset. Let $\mathcal{D}_n^s := \{(x_i, s_i, y_i) \in \mathcal{D}_n\}_{i:s_i=s}$ be a subset of $\mathcal{D}_n$ conditional on $s$ and let $n_s = |\mathcal{D}_n^s|$.

- Let $\hat{F}_{f|s}$ and $\hat{Q}_{f|s}$ be the empirical CDF and empirical quantile function for a given $f$, respectively. Let $\hat{f}$ be a given prediction model (e.g., empirical risk minimizer) and define

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

where $\epsilon \sim \mathsf{Unif}([-\sigma, \sigma])$.

- Assume that (i) $\nu_{f^*|s}$ admits a bounded density for each $s \in \mathcal{S}$ and (ii) there exists a positive constant $c$ and a sequence $b_n$ such that $\mathbb{E}|f^*(X,S) - \hat{f}(X,S)| \leq c b_n^{-1/2}$.

### Theorem 3

Set $\sigma \lesssim \min_{s \in \mathcal{S}} n_s^{-1/2} \wedge b_n^{-1/2}$. Then, we have

$$\mathbb{E}|g^*(X,S) - \hat{g}(X,S)| \lesssim b_n^{-1/2} \bigvee \left( \sum_{s \in \mathcal{S}} p_s n_s^{-1/2} \right) \bigvee \sqrt{\frac{|\mathcal{S}|}{n}}. \tag{3}$$

- Performance measures
  - Prediction

$$\mathsf{MSE}(g) = \frac{1}{n} \sum_{(x_i, s_i, y_i) \in \mathcal{D}_n} (y_i - g(x_i, s_i))^2$$

  - Fairness

$$\mathsf{KS}(g) = \max_{s, s' \in \mathcal{S}} \sup_{t \in \mathbb{R}} \left| \frac{1}{n_s} \sum_{(x_i, s_i, y_i) \in \mathcal{D}_n^s} \mathbb{I}(g(x_i, s_i) \leq t) - \frac{1}{n_{s'}} \sum_{(x_i, s_i, y_i) \in \mathcal{D}_n^{s'}} \mathbb{I}(g(x_i, s_i) \leq t) \right| \quad (4)$$

| Method | CRIME | | LAW | | NLSY | | STUD | | UNIV | |
|---|---|---|---|---|---|---|---|---|---|---|
| | MSE | KS | MSE | KS | MSE | KS | MSE | KS | MSE | KS |
| RLS | .033±.003 | .55±.06 | .107±.010 | .15±.02 | .153±.016 | .73±.07 | 4.77±.49 | .50±.05 | 2.24±.22 | .14±.01 |
| RLS+Berk | .037±.004 | .16±.02 | .121±.013 | .10±.01 | .189±.019 | .49±.05 | 5.28±.57 | .32±.03 | 2.43±.23 | .05±.01 |
| RLS+Oneto | .037±.004 | .14±.01 | .112±.012 | .07±.01 | .156±.016 | .50±.05 | 5.02±.54 | .23±.02 | 2.44±.26 | .05±.01 |
| RLS+Ours | .041±.004 | .12±.01 | .141±.014 | .02±.01 | .203±.019 | .09±.01 | 5.62±.52 | .04±.01 | 2.98±.32 | .02±.01 |
| KRLS | .024±.003 | .52±.05 | .040±.004 | .09±.01 | .061±.006 | .58±.06 | 3.85±.36 | .47±.05 | 1.43±.15 | .10±.01 |
| KRLS+Oneto | .028±.003 | .19±.02 | .046±.004 | .05±.01 | .066±.007 | .06±.01 | 4.07±.39 | .18±.02 | 1.46±.13 | .04±.01 |
| KRLS+Perez | .033±.003 | .25±.02 | .048±.005 | .04±.01 | .065±.007 | .08±.01 | 3.97±.38 | .14±.02 | 1.50±.15 | .06±.01 |
| KRLS+Ours | .034±.004 | .09±.01 | .056±.005 | .01±.01 | .081±.008 | .03±.01 | 4.46±.43 | .03±.01 | 1.71±.16 | .02±.01 |
| RF | .020±.002 | .45±.04 | .046±.005 | .11±.01 | .055±.006 | .55±.06 | 3.59±.39 | .45±.05 | 1.31±.13 | .10±.01 |
| RF+Raff | .030±.003 | .21±.02 | .058±.006 | .06±.01 | .066±.006 | .08±.01 | 4.28±.40 | .09±.01 | 1.38±.12 | .02±.01 |
| RF+Agar | .029±.003 | .13±.01 | .050±.005 | .04±.01 | .065±.006 | .07±.01 | 3.87±.41 | .07±.01 | 1.40±.13 | .02±.01 |
| RF+Ours | .033±.003 | .08±.01 | .064±.006 | .02±.01 | .070±.007 | .03±.01 | 4.18±.38 | .02±.01 | 1.49±.14 | .01±.01 |

Table 1: Results for all the datasets and all the methods concerning MSE and KS.

- Performs well for various datasets and models.
- MSE is slightly larger than the baselines, while KS is slightly lower than the baselines.

# Thank you