Classification and Regression under Statistical Parity

Solenne Gaucher¹



Evgenii Chzhen²



Vincent Divol³



Nicolas Schreuder⁴

¹CREST ENSAE, IPP

²CNRS, LMO, Université Paris-Saclay

³CEREMADE, Université Paris Dauphine PSL

⁴CNRS, LIGM, Université Gustave Eiffel

- 1. Introduction
- 2. The awareness framework
 - (Short reminders on) fair regression
 - Fair classification
- 3. The unawareness framework
 - Fair regression
 - Fair classification

Introduction to statistical fairness

Observations: $(X, S, Y) \sim \mathbb{P}$ on $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$. feature sensitive attribute outcome Setting Regression $\mathcal{Y} = \mathbb{R}$ Outcome

Predictions	$\mathcal{F} \triangleq \{f : \mathcal{X} \times \mathcal{S} \to \mathbb{R}\}$
Risk	$\mathcal{R}^{sq} = \mathbb{E}\left[(Y - f(X, S))^2\right]$
Fairness criteria	f(X, S) ⊥⊥ S (Statistical Parity)

 $\text{Observations:} \ (\underbrace{\mathcal{X}}_{}, \underbrace{}_{}, \underbrace{}_{}_{}, \underbrace{}_{}_{}, \underbrace{}_{}_{}_{}, \underbrace{}_{}_{}_{}, \underbrace{}_{}_{}_{}_{}) \sim \mathbb{P} \ \text{on} \ \mathcal{X} \times \mathcal{S} \times \mathcal{Y}.$ feature sensitive attribute outcome Setting Regression Classification

0		
Outcome	$\mathcal{Y}=\mathbb{R}$	$\mathcal{Y} = \{0,1\}$
Predictions	$\mathcal{F} \triangleq \{f : \mathcal{X} \times \mathcal{S} \to \mathbb{R}\}$	$\mathcal{G} \triangleq \{g: \mathcal{X} \times \mathcal{S} \to \{0, 1\}\}$
Risk	$\mathcal{R}^{sq} = \mathbb{E}\left[(Y - f(X, S))^2\right]$	$\mathcal{R}^{0-1} = \mathbb{E}\left[Y \neq g(X, S)\right]$
Fairness criteria	f(X,S) ⊥⊥ S (Statistical Parity)	ℝ[g(X, S) S] ⊥⊥ S (Demographic Parity)

Observations : $(\underbrace{X}_{\text{feature sensitive attribute outcome}}, \underbrace{Y}_{\text{outcome}}) \sim \mathbb{P} \text{ on } \mathcal{X} \times \mathcal{S} \times \mathcal{Y}.$		
Setting	Regression	Classification
Outcome	$\mathcal{Y}=\mathbb{R}$	$\mathcal{Y} = \{0,1\}$
Predictions	$\mathcal{F} \triangleq \{f : \mathcal{X} \times \mathcal{S} \to \mathbb{R}\}$	$\mathcal{G} \triangleq \{g: \mathcal{X} \times \mathcal{S} \to \{0, 1\}\}$
Risk	$\mathcal{R}^{sq} = \mathbb{E}\left[(Y - f(X, S))^2\right]$	$\mathcal{R}^{0-1} = \mathbb{E}\left[Y \neq g(X, S)\right]$
Fairness criteria	f <mark>(X, S)</mark> ⊥⊥ S (Statistical Parity)	E[g(X, S) S] ⊥⊥ S (Demographic Parity)

This is the awareness framework... in the unawareness framework, f and g cannot depend on S.

Observations: ($(\underbrace{X}_{\text{feature sensitive attribute outcome}}, \underbrace{Y}_{\text{feature sensitive attribute outcome}}) \sim$	\mathbb{P} on $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$.
Setting	Regression	Classification
Outcome	$\mathcal{Y}=\mathbb{R}$	$\mathcal{Y} = \{0,1\}$
Predictions	$\mathcal{F} \triangleq \{f : \mathcal{Z} \to \mathbb{R}\}$	$\mathcal{G} \triangleq \{g: \mathbf{\mathcal{Z}} \to \{0,1\}\}$
Risk	$\mathcal{R}^{sq} = \mathbb{E}\left[(Y - f(Z))^2\right]$	$\mathcal{R}^{0-1} = \mathbb{E}\left[Y \neq g(\mathbf{Z})\right]$
Fairness criteria	f(Z) ⊥⊥ S (Statistical Parity)	E[g(Z) S] ⊥⊥ S (Demographic Parity)

Awareness framework : $\mathcal{Z} = \mathcal{X} \times \mathcal{S}$, Z = (X, S)

Unawareness framework : $\mathcal{Z} = \mathcal{X}$, Z = X

Without fairness constraint, the Bayes classifier solution to

 $\underset{g \in \mathcal{G}}{\text{minimize}} \quad \mathbb{P}\left[Y \neq g(Z)\right]$

is given by

$$g^{Bayes}(z) = \mathbb{1}\left\{\eta(z) \geq \frac{1}{2}\right\}$$

where $\eta(Z) \triangleq \mathbb{E}[Y|Z]$ is the solution to

minimize
$$\mathbb{E}\left[(Y-f(Z))^2\right]$$
.

This relationship between classification and regression can be used to design and study classifiers [Yang, 1999], [Massart and Nédélec, 2006], [Audibert and Tsybakov, 2007], [Biau et al., 2008]

Consider the risks

$$\mathcal{R}_{\tau}(g) = \tau \mathbb{P}[Y = 0, g(Z) = 1] + (1 - \tau)\mathbb{P}[Y = 1, g(Z) = 0].$$

We have

$$\mathcal{R}_{\tau}(g) = (1-\tau)\mathbb{E}[Y] + \mathbb{E}\left[g(Z)(\tau - \eta(Z))\right].$$

 \implies The Bayes classifier g_{τ}^{Bayes} is given by

$$g_{\tau}^{\text{Bayes}}(z) = \mathbb{1}\left\{\eta(z) \geq \tau\right\}.$$

There is an equivalence between solving the regression problem and solving the classification problem for all τ .

minimize $\mathbb{E}\left[(Y - f(Z))^2\right]$ such that $f(Z) \perp S$.

We can define $g: z \mapsto \mathbb{1} \{ f^*(z) \ge \tau \}.$

- ▶ Is g optimal for some threshold τ ?
- Does g verify Demographic Parity?

minimize $\mathbb{E}\left[(Y - f(Z))^2\right]$ such that $f(Z) \perp S$.

We can define $g: z \mapsto \mathbb{1} \{ f^*(z) \ge \tau \}$.

- Is g optimal for some threshold τ ?
- Does g verify Demographic Parity? Yes... it verifies Strong Demographic Parity

minimize $\mathbb{E}\left[(Y - f(Z))^2\right]$ such that $f(Z) \perp S$.

We can define $g: z \mapsto \mathbb{1} \{ f^*(z) \ge \tau \}$.

- ▶ Is g optimal for some threshold τ ?
- Does g verify Demographic Parity? Yes... it verifies Strong Demographic Parity

Strong Demographic parity [JPSJC19]: A classifier *g* verifies **Strong Demographic parity** if $g(z) = \mathbb{1} \{ f(z) \ge \tau \}$ for some threshold τ , and *f* verifies Statistical Parity.

minimize $\mathbb{E}\left[(Y - f(Z))^2\right]$ such that $f(Z) \perp S$.

We can define $g: z \mapsto \mathbb{1} \{ f^*(z) \ge \tau \}$.

- ▶ Is g optimal for some threshold τ ? It depends...
- Does g verify Demographic Parity? Yes... it verifies Strong Demographic Parity

Strong Demographic parity [JPSJC19]: A classifier *g* verifies **Strong Demographic parity** if $g(z) = \mathbb{1} \{ f(z) \ge \tau \}$ for some threshold τ , and *f* verifies Statistical Parity.

The awareness framework

Fair regression

Assumption (A1)

 $\mu_s = Law(\eta(X, S)|S = s)$ are continuous and have finite second moments.

Theorem (Chzhen et al., 2020, Le Gouic et al., 2020)

Let *f*^{*} be the solution to

minimize	$\mathbb{E}\left[\left(Y-f(X,S)\right)^{2}\right]$
such that	$f(X,S) \perp S.$

Under Assumption (A1),

$$f^*(\mathbf{X},\mathbf{S}) = F_{\overline{\mu}}^{-1} \odot F_{\mu_{\mathbf{S}}}(\eta(\mathbf{X},\mathbf{S}))$$

where $F_{\mu_s}(t) = \mathbb{P}[\eta(X, S) \le t | S = s]$, and $F_{\overline{\mu}}^{-1}(\epsilon) = \sum_s p_s F_{\mu_s}^{-1}(\epsilon)$.

Fair classification

First results

Assumption (A2)

 μ_s = Law($\eta(X, S)|S = s$) are continuous and supported on an interval.

Theorem (G., Schreuder and Chzhen, 2023)

Let g* be the solution to

minimize	$\mathbb{P}\left[Y\neq g(X,S)\right]$
such that	g(X, S) ⊥⊥ S.

Under Assumption (A2), g* can be expressed as

$$g^*(x,s) = \mathbb{1}\left\{f^*(x,s) \ge \frac{1}{2}\right\}.$$

We look for classifiers $g(x,s) = \mathbb{1} \{ \eta(x,s) \ge \kappa_s \}$.



We look for classifiers $g(x,s) = \mathbb{1} \{ \eta(x,s) \ge \kappa_s \}$.



We look for classifiers $g(x,s) = \mathbb{1} \{\eta(x,s) \ge \kappa_s\}$.



We consider general performance measures

$$\mathcal{U}_{n,d}(g) \triangleq \frac{n_0 + n_1 \mathbb{P}\left[g(X,S) = 1, Y = 1\right] + n_2 \mathbb{P}\left[g(X,S) = 1\right]}{d_0 + d_1 \mathbb{P}\left[g(X,S) = 1, Y = 1\right] + d_2 \mathbb{P}\left[g(X,S) = 1\right]}$$

where n, d can depend on $\mathbb{P}[Y = 1]$ (but not on g).

Accuracy U_{n,d}(g) = -ℙ[g(X, S) ≠ Y].
 F_b-score U_{n,d}(g) = (1+b²)ℙ[g(X,S)=1,Y=1]/b²ℙ[Y=1]+ℙ[g(X,S)=1].
 Jaccard Index U_{n,d}(g) = <u>ℙ[g(X,S)=1,Y=0]+ℙ[g(X,S)=1].</u>
 AM Measure U_{n,d}(g) = 1/2 (ℙ[g(X,S) = 1|Y = 1] + ℙ[g(X,S) = 0|Y = 0]).

Theorem (G., Schreuder and Chzhen, 2023 - Informal)

Let $g_{n,d}^*$ be the solution to

 $\begin{array}{ll} \mbox{maximize} & \mathcal{U}_{n,d}(g) \\ \mbox{such that} & g(X,S) \perp L S. \end{array}$

Under Assumption (A2) and assumptions on n and d, $g_{n,d}^*$ can be expressed as

$$g_{n,d}^*(x,s) = \mathbb{1} \{ f^*(x,s) \ge \theta_{n,d}^* \}.$$

Assumption (A3)

$$\begin{split} & 1. \ d_0 + (d_1 + (d_2)_+)_+ \geq 0 \\ & 2. \ \begin{cases} d_2 n_1 > n_2 d_1 \ and \ d_0 n_1 - n_0 d_1 \geq (n_0 d_2 - d_0 n_2)_+ \\ \frac{n_0 d_2 - d_0 n_2}{n_2 d_1 - d_2 n_1} \leq \mathbb{P} \left[Y = 1 \right] \\ & 3. \ \begin{cases} d_2 n_1 = n_2 d_1 \ and \ n_1 d_0 > d_1 n_0 \\ \frac{d_0 n_2 - n_0 d_2}{n_2 d_1 - d_0 n_1} \in [0, 1] \end{cases} \end{split}$$

Theorem (G., Schreuder and Chzhen, 2023)

 $\textit{Under} (A2) + (A3)(1.) + (A3)(2. \textit{ or } 3.), g^*_{n,d}(x,s) = \mathbb{1} \left\{ f^*(x,s) \geq \theta^*_{n,d} \right\}, \textit{ where }$

▶ if (A3)(2.) holds,
$$\theta_{n,d}^*$$
 solves

$$\mathbb{E}\left[(f^{*}(X,S) - \theta)_{+}\right] = \theta \cdot \left\{\frac{n_{0}d_{1} - d_{0}n_{1}}{n_{2}d_{1} - d_{2}n_{1}}\right\} + \left\{\frac{n_{0}d_{2} - d_{0}n_{2}}{n_{2}d_{1} - d_{2}n_{1}}\right\}$$

► if (A3)(3.) holds,
$$\theta^*_{(n,d)} = \frac{d_0n_2 - n_0d_2}{n_0d_1 - d_0n_1}$$

Algorithm:

- 1. Estimate η
- 2. Estimate *f**
- 3. Estimate n and d
- 4. Estimate threshold $\theta_{n,d}^*$ (explicit formula or fixed-point equation)
- 5. Use double-plug-in estimator

$$\widehat{g}_{n,d}(x,s) = \mathbb{1}\left\{\widehat{f}^*(x,s) \geq \widehat{\theta}_{\widehat{n},\widehat{d}}\right\}$$

In the awareness framework, DP-fair optimal classification with performance measure $\mathcal{U}_{n,d}$

▶ is given by
$$g_{n,d}^*(x,s) = \mathbb{1}\left\{f^*(x,s) \ge \theta_{n,d}^*\right\}$$

- has desirable properties :
 - $\cdot\,$ does no harm to the protected group
 - preserves rational ordering
 - preserves monotonicity

In the awareness framework, DP-fair optimal classification with performance measure $\mathcal{U}_{n,d}$

▶ is given by
$$g_{n,d}^*(x,s) = \mathbb{1}\left\{f^*(x,s) \ge \theta_{n,d}^*\right\}$$

has desirable properties :

- · does no harm to the protected group
- preserves rational ordering
- preserves monotonicity

(Informal) A family of classifiers g_{θ}^{*} preserves monotonicity if $\mu\text{-almost-surely,}$

 $\{\mathbb{P}\left[g_{\theta}^*(X,S)=1\right] > \mathbb{P}\left[g_{\theta'}^*(X,S)=1\right]\} \implies \{g_{\theta}^*(x,s) \ge g_{\theta'}^*(x,s)\}.$

The unawareness framework

Fair regression

Problem: We want to solve

minimize	$\mathbb{E}\left[(Y-f(X))^2\right]$
such that	<i>f</i> (X) ⊥⊥ S.

In the following, we assume $S = \{1, 2\}$.

Notations: $X \sim \mu$, $X|S = 1 \sim \mu_1$, $X|S = 2 \sim \mu_2$.

Jordan decomposition We write $\mu_1 - \mu_2 = (\mu_+ - \mu_-)$, where

- \blacktriangleright μ_+ and μ_- are positive measures
- ▶ $\mathcal{X}_+ \triangleq \operatorname{supp}(\mu_+)$ and $\mathcal{X}_- \triangleq \operatorname{supp}(\mu_-)$ are disjoint.

Jordan decomposition We write $\mu_1 - \mu_2 = M(\mu_+ - \mu_-)$, where

- μ_+ and μ_- are probability measures
- ▶ $\mathcal{X}_+ \triangleq \operatorname{supp}(\mu_+)$ and $\mathcal{X}_- \triangleq \operatorname{supp}(\mu_-)$ are disjoint.

Jordan decomposition We write $\mu_1 - \mu_2 = M(\mu_+ - \mu_-)$, where

- μ_+ and μ_- are probability measures
- ▶ $\mathcal{X}_+ \triangleq \operatorname{supp}(\mu_+)$ and $\mathcal{X}_- \triangleq \operatorname{supp}(\mu_-)$ are disjoint.

Lemma (Chzhen and Schreuder, 2020)

A function $f : \mathcal{X} \to \mathbb{R}$ verifies Statistical Parity if and only if $f \sharp \mu_+ = f \sharp \mu_-$.

Jordan decomposition We write $\mu_1 - \mu_2 = M(\mu_+ - \mu_-)$, where

- μ_+ and μ_- are probability measures
- ▶ $\mathcal{X}_+ \triangleq \operatorname{supp}(\mu_+)$ and $\mathcal{X}_- \triangleq \operatorname{supp}(\mu_-)$ are disjoint.

Lemma (Chzhen and Schreuder, 2020)

A function $f: \mathcal{X} \to \mathbb{R}$ verifies Statistical Parity if and only if $f \sharp \mu_+ = f \sharp \mu_-$.

Consequence We can look for functions

$$f(x) = \begin{cases} f_+(x) & \text{if } x \in \mathcal{X}_+ \\ f_-(x) & \text{if } x \in \mathcal{X}_- \\ \eta(x) & \text{if } x \in \mathcal{X}_= \triangleq \{x : \mu_1(x) = \mu_2(x)\} \\ \text{such that } f_+ \sharp \mu_+ = f_- \sharp \mu_-. \end{cases}$$

$$\mathcal{R}^{sq}(f) = \int_{\mathcal{X}_{+}} c\left(\left(\eta(x), \Delta(x)\right), f_{+}(x)\right) d\mu_{+}(x) + \int_{\mathcal{X}_{-}} c\left(\left(\eta(x), \Delta(x)\right), f_{-}(x)\right) d\mu_{-}(x) + cste$$

where $c: (\mathbb{R} \times \mathbb{R}^*_+) \times \mathbb{R} \to \mathbb{R}_+$ is a cost, and $\Delta(x) = \begin{cases} \frac{d\mu_+}{d\mu}(x) \text{ if } x \in \mathcal{X}_+\\ \frac{d\mu_-}{d\mu}(x) \text{ if } x \in \mathcal{X}_- \end{cases}$.

$$\mathcal{R}^{sq}(f) = \int_{\mathcal{X}_+} c\left(\left(\eta(x), \Delta(x)\right), f_+(x)\right) d\mu_+(x) + \int_{\mathcal{X}_-} c\left(\left(\eta(x), \Delta(x)\right), f_-(x)\right) d\mu_-(x) + cste$$

where $c: \left(\mathbb{R} \times \mathbb{R}^*_+\right) \times \mathbb{R} \to \mathbb{R}_+$ is a cost, and $\Delta(x) = \begin{cases} \frac{d\mu_+}{d\mu}(x) & \text{if } x \in \mathcal{X}_+\\ \frac{d\mu_-}{d\mu}(x) & \text{if } x \in \mathcal{X}_- \end{cases}$. $\blacktriangleright f_{\pm}$ should depend only on $\left(\eta(x), \Delta(x)\right): f_{\pm}(x) = \tilde{f}_{\pm}\left(\eta(x), \Delta(x)\right).$

$$\mathcal{R}^{sq}(f) = \int_{\mathcal{X}_+} c\left(\left(\eta(x), \Delta(x)\right), f_+(x)\right) d\mu_+(x) + \int_{\mathcal{X}_-} c\left(\left(\eta(x), \Delta(x)\right), f_-(x)\right) d\mu_-(x) + cste$$

where
$$c: \left(\mathbb{R} \times \mathbb{R}^*_+\right) \times \mathbb{R} \to \mathbb{R}_+$$
 is a cost, and $\Delta(x) = \begin{cases} \frac{d\mu_+}{d\mu}(x) & \text{if } x \in \mathcal{X}_+\\ \frac{d\mu_-}{d\mu}(x) & \text{if } x \in \mathcal{X}_- \end{cases}$
 f_{\pm} should depend only on $(\eta(x), \Delta(x)): f_{\pm}(x) = \tilde{f}_{\pm}(\eta(x), \Delta(x)).$
 $Let \ \tilde{\mu}_{\pm} = Law((\eta(X), \Delta(X)) \mid X \sim \mu_{\pm}). \ \tilde{f}_{\pm}$ should solve:
minimize $\int c((\eta, \Delta), \tilde{f}_+(\eta, \Delta)) d\tilde{\mu}_+(\eta, \Delta) + \int c((\eta, \Delta), \tilde{f}_-(\eta, \Delta)) d\tilde{\mu}_-(\eta, \Delta)$
such that $\tilde{f}_+ \sharp \tilde{\mu}_+ = \tilde{f}_- \sharp \tilde{\mu}_-$

► Let
$$\tilde{\mu}_{\pm} = \text{Law}\Big((\eta(X), \Delta(X)) \mid X \sim \mu_{\pm}\Big)$$
. \tilde{f}_{\pm} should solve:
minimize $\int c\Big((\eta, \Delta), \tilde{f}_{\pm}(\eta, \Delta)\Big) d\tilde{\mu}_{\pm}(\eta, \Delta) + \int c\Big((\eta, \Delta), \tilde{f}_{\pm}(\eta, \Delta)\Big) d\tilde{\mu}_{\pm}(\eta, \Delta)$
such that $\tilde{f}_{\pm} \sharp \tilde{\mu}_{\pm} = \tilde{f}_{\pm} \sharp \tilde{\mu}_{\pm}$

► Let
$$\tilde{\mu}_{\pm} = \text{Law}\left(\left(\eta(X), \Delta(X)\right) \mid X \sim \mu_{\pm}\right)$$
. \tilde{f}_{\pm} should solve:
minimize $\int c\left((\eta, \Delta), \tilde{f}_{+}(\eta, \Delta)\right) d\tilde{\mu}_{+}(\eta, \Delta) + \int c\left((\eta, \Delta), \tilde{f}_{-}(\eta, \Delta)\right) d\tilde{\mu}_{-}(\eta, \Delta)$
such that $\tilde{f}_{+} \sharp \tilde{\mu}_{+} = \tilde{f}_{-} \sharp \tilde{\mu}_{-}$

 $\blacktriangleright \tilde{f}_+ \sharp \tilde{\mu}_+ = \tilde{f}_- \sharp \tilde{\mu}_-$ should solve the barycenter problem:

$$\min_{\nu} OT_{c}(\tilde{\mu}_{+},\nu) + OT_{c}(\tilde{\mu}_{-},\nu)$$

where

$$OT_{c}(\tilde{\mu},\nu) = \inf_{\gamma \in \Pi(\tilde{\mu},\nu)} \int c((\eta,\Delta),y) d\gamma((\eta,\Delta),y).$$

Assumption (A4)

The measures $\tilde{\mu}_+$ and $\tilde{\mu}_-$ are continuous, and the interior of their support has measure 1.

Theorem (Divol and G., 2024)

Under Assumptions (A1) and (A4), the solution $\nu^{\rm bar}$ of the barycenter problem

$$\min_{\nu} OT_{c}(\tilde{\mu}_{+},\nu) + OT_{c}(\tilde{\mu}_{-},\nu)$$

is such that $OT_c(\tilde{\mu}_+, \nu^{bar})$ and $OT_c(\tilde{\mu}_-, \nu^{bar})$ are solved by transport maps \tilde{f}^*_+ and \tilde{f}^*_- . Moreover, the optimal fair prediction is given by

$$f(x) = \begin{cases} \tilde{f}_+^*(\eta(x), \Delta(x)) & \text{if } x \in \mathcal{X}_+ \\ \tilde{f}_-^*(\eta(x), \Delta(x)) & \text{if } x \in \mathcal{X}_- \\ \eta(x) & \text{else.} \end{cases}$$

Remarks:

- Under Assumptions (A1) and (A4), the optimal prediction is deterministic.
- ▶ The optimal prediction tries to guess the sensitive attribute.
- Unless η verifies Statistical Parity, the optimal fair prediction does not verify rational ordering within group.

Fair classification

We want to solve

minimize $\mathcal{R}_{\tau}(g) \triangleq \tau \mathbb{P}[Y = 0, g(X) = 1] + (1 - \tau)\mathbb{P}[Y = 1, g(X) = 0]$ such that $g(X) \perp S$.

Remark:

$$\mathcal{R}_{\tau}(g) = (1-\tau)\mathbb{E}[Y] + \mathbb{E}\left[g(X)(\tau - \eta(X))\right].$$

 \implies The Bayes classifier g_{τ}^{Bayes} is given by

$$\mathcal{G}_{\tau}^{\text{Bayes}}(\mathbf{X}) = \mathbb{1}\left\{\eta(\mathbf{X}) \geq \tau\right\}.$$

Family of risks \mathcal{R}_{τ} with corresponding thresholds τ .

Lemma (Chzhen and Schreuder, 2020)

A function $f: \mathcal{X} \to \mathbb{R}$ verifies Statistical Parity if and only if $f \sharp \mu_+ = f \sharp \mu_-.$

▶ g verifies Demographic Parity if and only if

 $\mathbb{E}_{\mu_+}\left[g(X)\right] = \mathbb{E}_{\mu_-}\left[g(X)\right].$

Lemma (Chzhen and Schreuder, 2020)

A function $f: \mathcal{X} \to \mathbb{R}$ verifies Statistical Parity if and only if $f \sharp \mu_+ = f \sharp \mu_-$.

▶ g verifies Demographic Parity if and only if

 $\mathbb{E}_{\mu_+}\left[g(X)\right] = \mathbb{E}_{\mu_-}\left[g(X)\right].$

 $\blacktriangleright \text{ We can look for } g(x) = \begin{cases} g_+(x) & \text{if } x \in \mathcal{X}_+ \\ g_-(x) & \text{if } x \in \mathcal{X}_- \\ \mathbbm{1} \{\eta(x) \ge \tau\} & \text{else.} \end{cases}$

Lemma (Chzhen and Schreuder, 2020)

A function $f: \mathcal{X} \to \mathbb{R}$ verifies Statistical Parity if and only if $f \sharp \mu_+ = f \sharp \mu_-$.

▶ g verifies Demographic Parity if and only if

 $\mathbb{E}_{\mu_+}\left[g(X)\right] = \mathbb{E}_{\mu_-}\left[g(X)\right].$

$$\blacktriangleright \text{ We can look for } g(x) = \begin{cases} g_+(x) & \text{if } x \in \mathcal{X}_+ \\ g_-(x) & \text{if } x \in \mathcal{X}_- \\ \mathbbm{1} \{\eta(x) \ge \tau\} & \text{else.} \end{cases}$$

Decomposition + change of measure:

$$\mathcal{R}_{\tau}(g) = \mathbb{E}_{\mu_{+}}\left[\left(\frac{\tau - \eta(X)}{\Delta(X)}\right)g_{+}(X)\right] + \mathbb{E}_{\mu_{-}}\left[\left(\frac{\tau - \eta(X)}{\Delta(X)}\right)g_{-}(X)\right] + cste.$$

We should choose $g_{\pm}(x) = \mathbb{1}\left\{\frac{\eta(X) - \tau}{\Delta(X)} \ge \kappa_{\pm}\right\}$.

Theorem (G., Schreuder and Chzhen, 2023)

Let g_{τ}^* be the classifier minimizing \mathcal{R}_{τ} under the Demographic parity constraint. Under Assumption (A4), g_{τ}^* can be expressed as

$$g_{\tau}^{*}(x) = \begin{cases} \mathbbm{1}\left\{\frac{\eta(x)-\tau}{\Delta(x)} \ge \kappa(\tau)\right\} & \text{if } x \in \mathcal{X}_{+} \\ \mathbbm{1}\left\{\frac{\eta(x)-\tau}{\Delta(x)} \ge -\kappa(\tau)\right\} & \text{if } x \in \mathcal{X}_{-} \\ \mathbbm{1}\left\{\eta(x) \ge \tau\right\} & \text{else} \end{cases}$$

where $\kappa(\tau)$ is such that

$$\mathbb{P}_{\mu_+}\left[\frac{\eta(x)-\tau}{\Delta(x)} \geq \kappa(\tau)\right] = \mathbb{E}_{\mu_-}\left[\frac{\eta(x)-\tau}{\Delta(x)} \geq -\kappa(\tau)\right].$$

Remarks: In the unawareness framework, DP-fair optimal classification with risk measure \mathcal{R}_{τ}

- tries to guess the sensitive attribute;
- harms some individuals from protected group;
- does not preserves rational ordering

Remarks: In the unawareness framework, DP-fair optimal classification with risk measure \mathcal{R}_{τ}

- tries to guess the sensitive attribute;
- harms some individuals from protected group;
- does not preserves rational ordering
- may not preserves monotonicity

Remarks: In the unawareness framework, DP-fair optimal classification with risk measure \mathcal{R}_{τ}

- tries to guess the sensitive attribute;
- harms some individuals from protected group;
- does not preserves rational ordering
- may not preserves monotonicity

Monotonicity would imply that μ -almost surely,

$$\tau \ge \tau' \implies \begin{cases} \left\{ x \in \mathcal{X}_{+} : \underbrace{\frac{\eta(x) - \tau}{\Delta(x)} \ge \kappa(\tau)}{g_{\tau}^{*}(x) = 1} \right\} \subset \left\{ x \in \mathcal{X}_{+} : \underbrace{\frac{\eta(x) - \tau'}{\Delta(x)} \ge \kappa(\tau')}{g_{\tau'}^{*}(x) = 1} \right\} \\ \left\{ x \in \mathcal{X}_{-} : \underbrace{\frac{\eta(x) - \tau}{\Delta(x)} \ge -\kappa(\tau)}{g_{\tau}^{*}(x) = 1} \right\} \subset \left\{ x \in \mathcal{X}_{-} : \underbrace{\frac{\eta(x) - \tau'}{\Delta(x)} \ge -\kappa(\tau')}{g_{\tau'}^{*}(x) = 1} \right\}.\end{cases}$$

Original question: Is $\mathbb{1} \{ f^*(x) \ge \tau \}$ optimal?

- ▶ Without monotonicity, no! (The optimal classifier cannot be of the form $\mathbb{1} \{f(x) \ge \tau\}$.)
- (Divol and G., 2024 Informal) With monotonicity, yes!
- Both behaviours can be observed.

Conclusion

In the awareness framework:

- Optimal fair classifier verifies desirable properties : does no harm to the protected group, preserves rational ordering, preserves monotonicity.
- ▶ For general performance measures, $g^*(x,s) = \mathbb{1} \{ f^*(x,s) \ge \theta \}$

In the unawareness framework:

- Both the optimal fair classifier and the optimal fair regression function rely on guessing the sensitive attribute;
- Neither of them preserve rational ordering, fair classification harms indivual from the protected group.
- ▶ If monotonicity is verified, $g_{\tau}^*(x) = \mathbb{1} \{ f^*(x) \ge \tau \}$
- ▶ If monotonicity is not verified, $g_{\tau}^*(x) \neq \mathbb{1} \{ f^*(x) \geq \tau \}$.