

# Demographic parity constraint for algorithmic fairness

a statistical perspective

Statlearn'23  
Montpellier

Evgenii Chzhen  
CNRS, Univ. Paris-Saclay

# Today's plan

1. A **biased** intro to fairness and fairness zoology
2. **Demographic Parity** constraint and analogies
3. **Regression** with demographic parity constraint
4. Building estimators

# Fairness in ML: a major societal concern



PRODUCTS▼ CUSTOMERS▼ PRICING RESOURCES▼

REQUEST A DEMO



Talent Assessment | 16 Min Read

## How AI-based HR Chatbots are Simplifying Pre-screening

Source <https://www.mettl.com>

# Fairness in ML: a major societal concern

05-17-19

## **Schools are using software to help pick who gets in. What could go wrong?**

Admissions officers are increasingly turning to automation and AI with the hope of streamlining the application process and leveling the playing field.

Source <https://www.fastcompany.com>

# Fairness in ML: a major societal concern

SCIENCE ADVANCES | RESEARCH ARTICLE

---

## RESEARCH METHODS

# The accuracy, fairness, and limits of predicting recidivism

Julia Dressel and Hany Farid\*

Algorithms for predicting recidivism are commonly used to assess a criminal defendant's likelihood of committing a crime. These predictions are used in pretrial, parole, and sentencing decisions. Proponents of these systems argue that big data and advanced machine learning make these analyses more accurate and less biased than humans. We show, however, that the widely used commercial risk assessment software COMPAS is no more accurate or fair than predictions made by people with little or no criminal justice expertise. In addition, despite COMPAS's collection of 137 features, the same accuracy can be achieved with a simple linear predictor with only two features.

Copyright © 2018  
The Authors, some  
rights reserved;  
exclusive licensee  
American Association  
for the Advancement  
of Science. No claim to  
original U.S. Government  
Works. Distributed  
under a Creative  
Commons Attribution  
NonCommercial  
License 4.0 (CC BY-NC).

# EU regulation for AI



EUROPEAN COMMISSION

Brussels, 21.4.2021

COM(2021) 206 final

2021/0106(COD)

Proposal for a

**REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL**

**LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE (ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION LEGISLATIVE ACTS**

# Group fairness paradigm

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

Predictions:  $f : \mathcal{Z} \rightarrow \mathcal{Y}$

- ▶ Fairness through **awareness**:  $\mathcal{Z} = \mathcal{X} \times \mathcal{S}$  (disparate treatment)
- ▶ Fairness through **unawareness**:  $\mathcal{Z} = \mathcal{X}$  (legal reasons: regulations)

Risk:  $f \mapsto \mathcal{R}(f)$

- ▶ **classification**:  $\mathcal{R}(f) = \mathbb{P}(Y \neq f(\mathbf{Z}))$
- ▶ **regression**:  $\mathcal{R}(f) = \mathbb{E}(Y - f(\mathbf{Z}))^2$

Fairness criteria: **dichotomy** of prediction functions: which functions we call fair? There are a lot of definitions, maybe too many to parse.

# Popular definitions of fair classifiers

- ▶ Demographic Parity (DP) (Calders, Kamiran, and Pechenizkiy, 2009)

$$\mathbb{P}(f(\mathbf{Z}) = 1 \mid S = 0) = \mathbb{P}(f(\mathbf{Z}) = 1 \mid S = 1)$$

1. Prediction rate is the same for two groups
2. Random variable  $f(\mathbf{Z})$  is independent from  $S$
3. Only  $\mathbf{X}|S$  matters
4. Constant predictions satisfy DP



# Popular definitions of fair classifiers

- ▶ Demographic Parity (DP) (Calders, Kamiran, and Pechenizkiy, 2009)

$$\mathbb{P}(f(\mathbf{Z}) = 1 \mid S = 0) = \mathbb{P}(f(\mathbf{Z}) = 1 \mid S = 1)$$

1. Prediction rate is the same for two groups
  2. Random variable  $f(\mathbf{Z})$  is independent from  $S$
  3. Only  $\mathbf{X}|S$  matters
  4. Constant predictions satisfy DP
- ▶ Equalized Odds (Hardt, Price, and Srebro, 2016)
- $$\mathbb{P}(f(\mathbf{Z}) = y \mid Y = y, S = 0) = \mathbb{P}(f(\mathbf{Z}) = y \mid Y = y, S = 1) \quad \forall y \in \{0, 1\}$$
1. Equal True Positive and True Negative rates
  2. Requires more knowledge about the distribution
  3. Constant predictions satisfy Equalized Odds

# Popular definitions of fair classifiers

- ▶ Equal Opportunity (Hardt, Price, and Srebro, 2016)

$$\mathbb{P}(f(\mathbf{Z}) = 1 \mid Y = 1, S = 0) = \mathbb{P}(f(\mathbf{Z}) = 1 \mid Y = 1, S = 1)$$

1. Equal True Positive rates
2. If a person  $\mathbf{Z}$  is qualified ( $Y = 1$ ) then positive prediction ( $f(\mathbf{Z}) = 1$ ) is given with the same probability for any sensitive attribute

# Popular definitions of fair classifiers

- ▶ Equal Opportunity (Hardt, Price, and Srebro, 2016)

$$\mathbb{P}(f(\mathbf{Z}) = 1 \mid Y = 1, S = 0) = \mathbb{P}(f(\mathbf{Z}) = 1 \mid Y = 1, S = 1)$$

1. Equal True Positive rates
2. If a person  $\mathbf{Z}$  is qualified ( $Y = 1$ ) then positive prediction ( $f(\mathbf{Z}) = 1$ ) is given with the same probability for any sensitive attribute

- ▶ Test fairness (Chouldechova, 2017)

$$\mathbb{P}(Y = 1 \mid S = 0, f(\mathbf{Z}) = 1) = \mathbb{P}(Y = 1 \mid S = 1, f(\mathbf{Z}) = 1)$$

1.  $Y$  independent from  $S$  conditionally on  $f(\mathbf{Z}) = 1$ .
2. Closely related to group-wise calibration.

# Global view on group fairness constraints

Most of the definitions of fairness fall inside or try to reflect only 3 criteria

1.  $f(\mathbf{Z}) \perp\!\!\!\perp S$  - **independence** (DP, Statistical Parity)
2.  $(f(\mathbf{Z}) \perp\!\!\!\perp S) \mid Y$  - **separation** (Equal Odds, Equal Opportunity)
3.  $(Y \perp\!\!\!\perp S) \mid f(\mathbf{Z})$  - **sufficiency** (Test fairness)

**N.B.** Sometimes we consider a score function  $f(\mathbf{Z}) \in [0, 1]$ .

# Impossibilities for score functions

1.  $f(\mathbf{Z}) \perp\!\!\!\perp S$  - **independence** (DP, Statistical Parity)
2.  $(f(\mathbf{Z}) \perp\!\!\!\perp S) \mid Y$  - **separation** (Equal Odds, Equal Opportunity)
3.  $(Y \perp\!\!\!\perp S) \mid f(\mathbf{Z})$  - **sufficiency** (Test fairness)
  - ▶ If  $S$  and  $Y$  are not independent, then sufficiency and independence cannot both hold.
  - ▶ If  $Y \in \{0, 1\}$ ,  $S$  and  $Y$  are not independent,  $f(\mathbf{Z})$  is not independent from  $Y$ , then independence and separation cannot both hold.
  - ▶ If  $S$  and  $Y$  are not independent, and  $\mathbb{P}(Y = 1) \in (0, 1)$ , then separation and sufficiency cannot both hold.

# Impossibilities for score functions

1.  $f(\mathbf{Z}) \perp\!\!\!\perp S$  - **independence** (DP, Statistical Parity)
2.  $(f(\mathbf{Z}) \perp\!\!\!\perp S) \mid Y$  - **separation** (Equal Odds, Equal Opportunity)
3.  $(Y \perp\!\!\!\perp S) \mid f(\mathbf{Z})$  - **sufficiency** (Test fairness)
  - ▶ If  $S$  and  $Y$  are not independent, then sufficiency and independence cannot both hold.
  - ▶ If  $Y \in \{0, 1\}$ ,  $S$  and  $Y$  are not independent,  $f(\mathbf{Z})$  is not independent from  $Y$ , then independence and separation cannot both hold.
  - ▶ If  $S$  and  $Y$  are not independent, and  $\mathbb{P}(Y = 1) \in (0, 1)$ , then separation and sufficiency cannot both hold.

**A fact:** famous example of COMPAS nearly satisfied sufficiency, but failed to satisfy separation. Due to the latter propublica published an article that extremely influenced the field of algorithmic fairness (Chouldechova, 2017).

# Three (rough) types of methods: **pre-processing**

## Pre-processing – **Fair representation**

Find a feature representation  $\mathbf{Z} \mapsto \hat{\varphi}(\mathbf{Z})$  such that

$$\hat{\varphi}(\mathbf{Z}) \perp\!\!\!\perp S$$

then use any method on this representation.

Typically, (**unsupervised**) optimal fair representation is defined as

$$\varphi^* \in \arg \min \{ \mathbb{E}[d(\mathbf{X}, \varphi(\mathbf{Z}))] : \varphi(\mathbf{Z}) \perp\!\!\!\perp S \} .$$

# Three (rough) types of methods: pre-processing

## Pre-processing – Fair representation

Find a feature representation  $\mathbf{Z} \mapsto \hat{\varphi}(\mathbf{Z})$  such that

$$\hat{\varphi}(\mathbf{Z}) \perp\!\!\!\perp S$$

then use any method on this representation.

Typically, (unsupervised) optimal fair representation is defined as

$$\varphi^* \in \arg \min \{ \mathbb{E}[d(\mathbf{X}, \varphi(\mathbf{Z}))] : \varphi(\mathbf{Z}) \perp\!\!\!\perp S \} .$$

## Methods

- ▶ Linear models (Zemel et al., 2013)
- ▶ Kernel methods (Grünwälder and Khaleghi, 2021)
- ▶ GANs (Xu et al., 2018)



# Three (rough) types of methods: **in-processing**

Add the fairness **constraint into training**

$$f_{\mathcal{F}}^* \in \arg \min_{f \in \mathcal{F}} \{ \mathcal{R}(f) : f(\mathbf{Z}) \perp S \}$$

In-processing type method: Given data  $(\mathbf{X}_1, S_1, Y_1), \dots, (\mathbf{X}_n, S_n, Y_n)$  build an estimator  $\hat{f}$  as a solution

$$\min_{f \in \mathcal{F}} \left\{ \hat{\mathcal{R}}(f) + \lambda_0 \cdot \Omega_{\text{compl}}(f) + \lambda_1 \cdot \Omega_{\text{UNfairness}}(f) \right\}$$

# Three (rough) types of methods: **in-processing**

Add the fairness **constraint into training**

$$f_{\mathcal{F}}^* \in \arg \min_{f \in \mathcal{F}} \{ \mathcal{R}(f) : f(\mathbf{Z}) \perp S \}$$

In-processing type method: Given data  $(\mathbf{X}_1, S_1, Y_1), \dots, (\mathbf{X}_n, S_n, Y_n)$  build an estimator  $\hat{f}$  as a solution

$$\min_{f \in \mathcal{F}} \left\{ \hat{\mathcal{R}}(f) + \lambda_0 \cdot \Omega_{\text{compl}}(f) + \lambda_1 \cdot \Omega_{\text{UNfairness}}(f) \right\}$$

## Methods

- ▶ Regularized ERM methods (Oneto, Donini, and Pontil, 2019)
- ▶ MWU-type methods for minmax games (Agarwal et al., 2018)

## Three (rough) types of methods: **post-processing**

Given a **base algorithm**  $f$ , find a **transformation**

$$f \mapsto \hat{T}(f) ,$$

so that  $\hat{T}(f)$  satisfies your fairness constraint

## Three (rough) types of methods: **post-processing**

Given a **base algorithm**  $f$ , find a **transformation**

$$f \mapsto \hat{T}(f) ,$$

so that  $\hat{T}(f)$  satisfies your fairness constraint

Typical algorithm construction is based on the **connection** between

$$f_{\text{fair}}^* \in \arg \min_{f: \mathcal{Z} \rightarrow \mathcal{Y}} \{ \mathcal{R}(f) : f \text{ is fair} \} \quad \text{and} \quad f_{\text{Bayes}}^* \in \arg \min_{f: \mathcal{Z} \rightarrow \mathcal{Y}} \mathcal{R}(f)$$

Often we can show that

$$f_{\text{fair}}^* = T^*(f_{\text{Bayes}}^*) ,$$

treat the base algorithm  $f$  as if it were a Bayes and estimate  $T^*$

# Three (rough) types of methods: **post-processing**

Given a **base algorithm**  $f$ , find a **transformation**

$$f \mapsto \hat{T}(f) ,$$

so that  $\hat{T}(f)$  satisfies your fairness constraint

Typical algorithm construction is based on the **connection** between

$$f_{\text{fair}}^* \in \arg \min_{f: \mathcal{Z} \rightarrow \mathcal{Y}} \{\mathcal{R}(f) : f \text{ is fair}\} \quad \text{and} \quad f_{\text{Bayes}}^* \in \arg \min_{f: \mathcal{Z} \rightarrow \mathcal{Y}} \mathcal{R}(f)$$

Often we can show that

$$f_{\text{fair}}^* = T^*(f_{\text{Bayes}}^*) ,$$

treat the base algorithm  $f$  as if it were a Bayes and estimate  $T^*$

## Methods

- ▶ Threshold adjustments (Hardt, Price, and Srebro, 2016; Menon and Williamson, 2018; C. et al., 2019)
- ▶ Optimal transport based (C. et al., 2020; Le Gouic, Loubes, and Rigollet, 2020)

# What is the Demographic Parity constraint?

with C. Denis, S. Gaucher, M. Hebiri, L. Oneto, M. Pontil, and N. Schreuder

# Learning with Demographic Parity

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \mathcal{Y}$$

**Prediction:**  $f : \mathbb{R}^d \times \mathcal{S} \rightarrow \mathcal{Y}$

**Risk:**  $\mathcal{R}(f) = \mathbb{E}[(Y - f(\mathbf{X}, S))^2]$  or  $\mathcal{R}(f) = \mathbb{P}(Y \neq f(\mathbf{X}, S))$

**Demographic Parity** fairness

$$f(\mathbf{X}, S) \perp\!\!\!\perp S$$

Optimal **fair** prediction:

$$f_0^* \in \arg \min \{ \mathcal{R}(f) : f(\mathbf{X}, S) \perp\!\!\!\perp S \}$$

# Our goals

1. Understand a relation between **regression** and **classification** under the **Demographic Parity constraint**
2. Understand a relation between **constraint** and **unconstraint** (Bayes optimal) problems
3. Try to explain the notion of Demographic Parity in a **simple** language
4. Figure out an **estimation** strategy and get some **bounds** on risk and unfairness



# Classical classification-regression link

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

$$g^* \in \arg \min_{g: \mathcal{X} \times \mathcal{S} \rightarrow \{0, 1\}} \mathbb{P}(Y \neq g(\mathbf{X}, S)) \quad f^* \in \arg \min_{f: \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}[(Y - f(\mathbf{X}, S))^2]$$

# Classical classification-regression link

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

$$g^* \in \arg \min_{g: \mathcal{X} \times \mathcal{S} \rightarrow \{0, 1\}} \mathbb{P}(Y \neq g(\mathbf{X}, S)) \quad f^* \in \arg \min_{f: \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}[(Y - f(\mathbf{X}, S))^2]$$

---

---

**A folklore result**

---

---

$$f^*(\mathbf{X}, S) = \mathbb{E}[Y \mid \mathbf{X}, S] \quad g^*(\mathbf{X}, S) = \mathbf{1}\{f^*(\mathbf{X}, S) \geq 1/2\}$$

---

---

present in **every** ML/Stat book

**N.B.** Simple to prove, but very useful in theory and in practice.

# Classification-regression link under DP

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathcal{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

Can we expect the **same result** under the Demographic parity constraint?

There is really **no reason** for such a relation...

# Classification-regression link under DP

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

Can we expect the **same result** under the Demographic parity constraint?

There is really **no reason** for such a relation... Indeed, if

$$g_0^* \in \arg \min_{g: \mathcal{X} \times \mathcal{S} \rightarrow \{0,1\}} \{ \mathbb{P}(Y \neq g(\mathbf{X}, S)) : g(\mathbf{X}, S) \perp\!\!\!\perp S \}$$

$$f_0^* \in \arg \min_{f: \mathcal{X} \times \mathcal{S} \rightarrow \{0,1\}} \{ \mathbb{E}[(Y - f(\mathbf{X}, S))^2] : f(\mathbf{X}, S) \perp\!\!\!\perp S \}$$

are such that

$$g_0^*(\mathbf{X}, S) = \mathbb{1}\{f_0^*(\mathbf{X}, S) \geq 1/2\} ,$$

then  $g_0^*$  is “**much fairer**” than we expect— $f_0^*$  is fair at **every threshold**, while  $g_0^*$  needs to be fair **only at one of them**.

# Classification-regression link under DP

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

$$g_0^* \in \arg \min_{g: \mathcal{X} \times \mathcal{S} \rightarrow \{0, 1\}} \{ \mathbb{P}(Y \neq g(\mathbf{X}, S)) : g(\mathbf{X}, S) \perp\!\!\!\perp S \}$$

$$f^* \in \arg \min_{f: \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}} \mathbb{E}[(Y - f(\mathbf{X}, S))^2]$$

---

---

## Lemma

---

---

$$g_0^*(\mathbf{X}, S) = \mathbb{1} \left\{ f^*(\mathbf{X}, S) \geq \frac{1}{2} + \frac{\lambda_s^*}{2w_s} \right\}$$

where  $w_s = \mathbb{P}(S = s)$  and

$$(\lambda_1^*, \dots, \lambda_K^*) \in \arg \min_{(\lambda_1, \dots, \lambda_K) \in \mathbb{R}^K} \left\{ \mathbb{E} \left| 2f^*(\mathbf{X}, S) - 1 - \frac{\lambda_S}{w_S} \right| : \sum_{s \in \mathcal{S}} \lambda_s = 0 \right\}$$

---

---

(Menon and Williamson, 2018; Gaucher, Schreuder, and C., 2023)

# Classification-regression link under DP

$$\underbrace{(\text{feature})}_{\mathbf{X}}, \underbrace{(\text{sensitive attribute})}_{\mathcal{S}}, \underbrace{(\text{signal})}_{Y} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \{0, 1\}$$

Nevertheless

$$g_0^* \in \arg \min_{g: \mathcal{X} \times \mathcal{S} \rightarrow \{0, 1\}} \{\mathbb{P}(Y \neq g(\mathbf{X}, S)) : g(\mathbf{X}, S) \perp\!\!\!\perp S\}$$

$$f_0^* \in \arg \min_{f: \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}} \{\mathbb{E}[(Y - f(\mathbf{X}, S))^2] : f(\mathbf{X}, S) \perp\!\!\!\perp S\}$$

---

---

**Lemma**

---

---

$$g^*(\mathbf{X}, S) = \mathbb{1}\{f_0^*(\mathbf{X}, S) \geq 1/2\}$$

$$f_0^*(\mathbf{X}, S) = ??$$

---

---

(Gaucher, Schreuder, and C., 2023)

**N.B.** It remains to understand the regression case

# Regression + Demographic Parity

$$\underbrace{(\text{feature}, \text{sensitive attribute}, \text{signal})}_{\mathbf{X} \times \mathcal{S} \times \mathbb{R}} \sim \mathbb{P} \text{ on } \mathbb{R}^d \times \underbrace{\mathcal{S}}_{=\{1, \dots, K\}} \times \mathbb{R}$$

**Prediction:**  $f : \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}$

**Risk:**  $\mathcal{R}(f) = \mathbb{E}[(f^*(\mathbf{X}, S) - f(\mathbf{X}, S))^2]$  where  $f^*(\mathbf{X}, S) = \mathbb{E}[Y | \mathbf{X}, S]$

**Demographic Parity** fairness

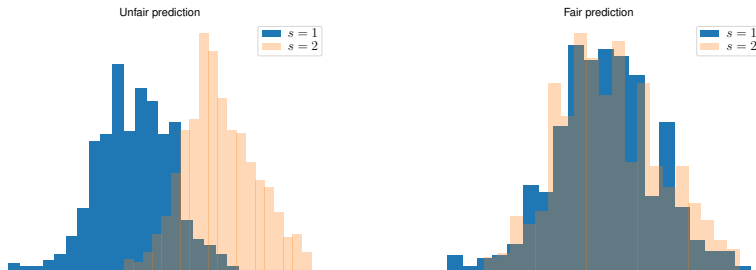
$$f(\mathbf{X}, S) \perp\!\!\!\perp S$$

Optimal **fair** prediction:

$$f_0^* \in \arg \min \{ \mathcal{R}(f) : f(\mathbf{X}, S) \perp\!\!\!\perp S \}$$

# An illustration and main assumption

$$f(\mathbf{X}, S) \perp\!\!\!\perp S$$



---

---

## Assumption (A)

---

---

The group-wise prediction distributions  $\text{Law}(f^*(\mathbf{X}, S) \mid S = s)$  have **finite second moment** and are **non-atomic** for  $s$  in  $\mathcal{S}$ .

---

---



# Optimal transport and the Wasserstein-2 metric

Define, for  $\mu, \nu \in \mathcal{P}_2(\mathbb{R})$ ,

$$W_2^2(\mu, \nu) := \inf \left\{ \mathbb{E}_{(X,Y)} (\mathbf{X} - \mathbf{Y})^2 : \mathbf{X} \sim \mu, \mathbf{Y} \sim \nu \right\}.$$

- ▶ Metric on  $\mathcal{P}_2(\mathbb{R}^d)$
- ▶ Optimal  $T_{\mu \rightarrow \nu}^* \equiv F_{\nu}^{-1} \circ F_{\mu}$
- ▶ Nice interpretations

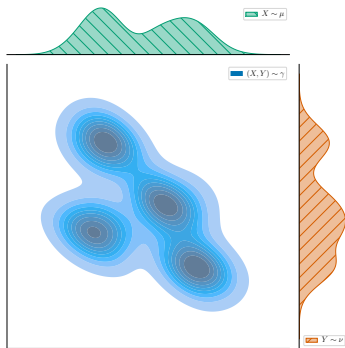


Figure: Transport plan illustration

# Reminder: post-processing

Optimal fair:  $f_0^* \in \arg \min_{f: \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}} \{\mathcal{R}(f) : f(\mathbf{X}, S) \perp\!\!\!\perp S\}$

Bayes optimal:  $f^* \in \arg \min_{f: \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}} \mathcal{R}(f)$

Question: is there a link between  $f_0^*$  and  $f^*$ ?

More precisely, can we show that

$$f_0^* \equiv T \circ f^* ?$$

# Main insight

Optimal fair:  $f_0^* \in \arg \min_{f: \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}} \{\mathcal{R}(f) : f(\mathbf{X}, S) \perp\!\!\!\perp S\}$

Bayes optimal:  $f^* \in \arg \min_{f: \mathbb{R}^d \times \mathcal{S} \rightarrow \mathbb{R}} \mathcal{R}(f)$

Question: is there a link between  $f_0^*$  and  $f^*$ ?

---

---

## Theorem

---

---

Set  $w_s = \mathbb{P}(S=s)$ . Let Assumption (A) be satisfied, then

$$\text{Law}(f_0^*(\mathbf{X}, S)) = \arg \min_{\nu \in \mathcal{P}_2(\mathbb{R})} \underbrace{\sum_{s \in \mathcal{S}} w_s \mathcal{W}_2^2 \left( \text{Law}(f^*(\mathbf{X}, S) \mid S = s), \nu \right)}_{\text{Wasserstein barycenter problem}},$$

$$f_0^*(\mathbf{x}, 1) = w_1 f^*(\mathbf{x}, 1) + w_2 T_{1 \rightarrow 2}^* \circ f^*(\mathbf{x}, 1), \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

$T_{1 \rightarrow 2}^*$  – optimal transport map from  $\text{Law}(f^* \mid S = 1)$  to  $\text{Law}(f^* \mid S = 2)$ .

---

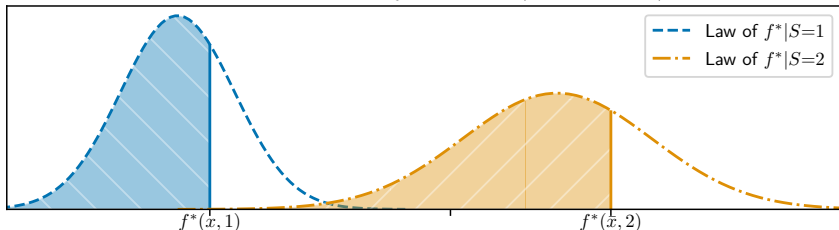
---

(C. et al., 2020)

# Interpretation for $\mathcal{S} = \{1, 2\}$

Fair optimal:  $f_0^*(\mathbf{x}, 1) = w_1 f^*(\mathbf{x}, 1) + w_2 F_{f^*|S=2}^{-1} \circ F_{f^*|S=1} \circ f^*(\mathbf{x}, 1)$

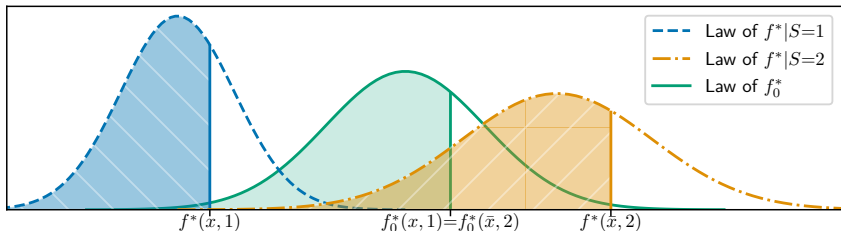
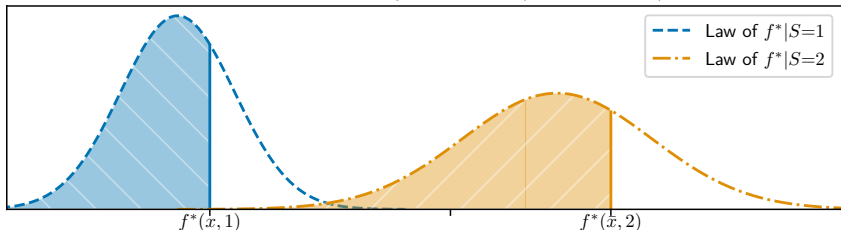
Fair optimal prediction  $f_0^*$  with  $w_1 = 2/5$  and  $w_2 = 3/5$



# Interpretation for $\mathcal{S} = \{1, 2\}$

Fair optimal:  $f_0^*(x, 1) = w_1 f^*(x, 1) + w_2 F_{f^*|S=2}^{-1} \circ F_{f^*|S=1} \circ f^*(x, 1)$

Fair optimal prediction  $f_0^*$  with  $w_1 = 2/5$  and  $w_2 = 3/5$



# Generic post-processing estimator ( $\mathcal{S} = \{1, 2\}$ )

**Fair optimal:**  $f_0^*(\mathbf{x}, 1) = w_1 f^*(\mathbf{x}, 1) + w_2 T_{1 \rightarrow 2}^* \circ f^*(\mathbf{x}, 1)$

- ▶ **Base estimator:**  $\hat{f} : \mathbb{R}^d \times \{1, 2\} \rightarrow \mathbb{R}$  trained independently from the following data.
- ▶ **Unlabeled data:**  $\forall s \in \mathcal{S}$  we observe  $\mathbf{X}_1^s, \dots, \mathbf{X}_{N_s}^s \stackrel{i.i.d.}{\sim} \mathbb{P}_{\mathbf{X}|S=s}$

**Meta algo:**

1. estimate  $w_s$  if needed
2. estimate transport maps  $T_{1 \rightarrow 2}^*$  and  $T_{2 \rightarrow 1}^*$  using **unlabeled data** and **base estimator**

# Generic post-processing estimator ( $\mathcal{S} = \{1, 2\}$ )

**Fair optimal:**  $f_0^*(\mathbf{x}, 1) = w_1 f^*(\mathbf{x}, 1) + w_2 T_{1 \rightarrow 2}^* \circ f^*(\mathbf{x}, 1)$

- ▶ **Base estimator:**  $\hat{f} : \mathbb{R}^d \times \{1, 2\} \rightarrow \mathbb{R}$  trained independently from the following data.
- ▶ **Unlabeled data:**  $\forall s \in \mathcal{S}$  we observe  $\mathbf{X}_1^s, \dots, \mathbf{X}_{N_s}^s \stackrel{i.i.d.}{\sim} \mathbb{P}_{\mathbf{X}|S=s}$

**Meta algo:**

1. estimate  $w_s$  if needed
2. estimate transport maps  $T_{1 \rightarrow 2}^*$  and  $T_{2 \rightarrow 1}^*$  using **unlabeled data** and **base estimator**

**Put together:** 3.  $\hat{f}_0(\mathbf{x}, 1) = w_1 \hat{f}(\mathbf{x}, 1) + w_2 \hat{T}_{1 \rightarrow 2} \circ \hat{f}(\mathbf{x}, 1)$

# Theoretical guarantees

---

---

## Theorem

---

---

For **any** joint distribution  $\mathbb{P}$  of  $(\mathbf{X}, S, Y)$ , **any** base estimator  $\hat{f}$  it holds that

$$\hat{f}_0(\mathbf{X}, S) \perp\!\!\!\perp S$$

Under **additional assumptions** on  $\mathbb{P}$  we have

$$\mathbf{E}\|\hat{f}_0 - f_0^*\|_1 \lesssim \underbrace{\mathbf{E}\|\hat{f} - f^*\|_1}_{\text{quality of base estimator}} \bigvee \underbrace{\sum_{s \in \mathcal{S}} w_s N_s^{-1/2}}_{\text{transport estimation}}$$

---

---

(C. and Schreuder, 2022)

**Additional assumptions:**  $(f^*(\mathbf{X}, S) \mid S = s)$  admits density which is **upper** and **lower** bounded



# How did we get exact independence and a cute lemma from conformal prediction theory

---

---

## Lemma for “smoothed ranks”

---

---

Let  $\mathbf{V} = (V, V_1, \dots, V_n)$  be *i.i.d.* real valued random variables and let  $U$  be distributed uniformly on  $(0, 1)$  and independent of  $\mathbf{V}$ . Let

$$F(U, V_1, \dots, V_n, x) = \frac{1}{n+1} \left( \sum_{i=1}^n \mathbb{1}\{V_i < x\} + U \cdot \left( 1 + \sum_{i=1}^n \mathbb{1}\{V_i = x\} \right) \right) .$$

Then,  $F(U, V_1, \dots, V_n, V)$  is distributed uniformly on  $(0, 1)$ .

---

---

V. Vovk and A. Gammerman

**N.B.** No assumptions on the distribution of the data, to compare with rank statistics.

# How did we get exact independence and a cute lemma from conformal prediction theory

---

---

## Lemma for “smoothed ranks”

---

---

Let  $\mathbf{V} = (V, V_1, \dots, V_n)$  be *i.i.d.* real valued random variables and let  $U$  be distributed uniformly on  $(0, 1)$  and independent of  $\mathbf{V}$ .

$$F(U, V_1, \dots, V_n, V) \sim \text{Unif}(0, 1)$$

---

---

The optimal fair prediction can be expressed as

$$f_0^*(\mathbf{x}, s) = Q \circ (F_s(f^*(\mathbf{x}, s))) \text{ ,}$$

where  $Q$  is a **monotone** and  $F_s$  is the CDF of  $\text{Law}(f^*(\mathbf{X}, S) \mid S = s)$ .

**Idea.** Use the above lemma for estimation of  $F_s(f^*(\mathbf{x}, s))$  as it always produces uniform distributions on  $(0, 1)$  (conditionally on  $S = s$ )

# Conclusions

1. **Group fairness** – enforce some independence criterion

$$f(\mathbf{Z}) \perp\!\!\!\perp S, \quad (f(\mathbf{Z}) \perp\!\!\!\perp S) \mid Y, \quad (Y \perp\!\!\!\perp S) \mid f(\mathbf{Z})$$

2. Demographic parity **preserves** classical classification-regression

$$g_0^* = \mathbb{1}\{f_0^* \geq 1/2\}$$

3. Regression with demographic parity ( $f(\mathbf{Z}) \perp\!\!\!\perp S$ ) can be characterized by **Wasserstein barycenter** problem
4. Demographic parity simply **matches ranks** of individuals from different groups

# Thank you for your attention! Questions?

## PROHIBITED ARTIFICIAL INTELLIGENCE PRACTICES

### Article 5

1. The following artificial intelligence practices shall be prohibited:
  - (a) the placing on the market, putting into service or use of an AI system that deploys subliminal techniques beyond a person's consciousness in order to materially distort a person's behaviour in a manner that causes or is likely to cause that person or another person physical or psychological harm;
  - (b) the placing on the market, putting into service or use of an AI system that exploits any of the vulnerabilities of a specific group of persons due to their age, physical or mental disability, in order to materially distort the behaviour of a person pertaining to that group in a manner that causes or is likely to cause that person or another person physical or psychological harm;
  - (c) the placing on the market, putting into service or use of AI systems by public authorities or on their behalf for the evaluation or classification of the trustworthiness of natural persons over a certain period of time based on their social behaviour or known or predicted personal or personality characteristics, with the social score leading to either or both of the following:
    - (i) detrimental or unfavourable treatment of certain natural persons or whole groups thereof in social contexts which are unrelated to the contexts in which the data was originally generated or collected;
    - (ii) detrimental or unfavourable treatment of certain natural persons or whole groups thereof that is unjustified or disproportionate to their social behaviour or its gravity;

# Bibliography I

- Agarwal, A. et al. (2018). “A reductions approach to fair classification”. In: *arXiv preprint arXiv:1803.02453*.
- Barocas, Solon, Moritz Hardt, and Arvind Narayanan (2019). *Fairness and Machine Learning*. <http://www.fairmlbook.org>. fairmlbook.org.
- C., E. and N. Schreuder (2022). “A minimax framework for quantifying risk-fairness trade-off in regression”. In: *Annals of Statistics*.
- C., E et al. (2019). “Leveraging Labeled and Unlabeled Data for Consistent Fair Binary Classification”. Submitted to NeurIPS19.
- (2020). “Fair Regression with Wasserstein Barycenters”. In: *NeurIPS 2020*.
- Calders, T., F. Kamiran, and M. Pechenizkiy (2009). “Building classifiers with independency constraints”. In: *IEEE international conference on Data mining*.
- Chouldechova, Alexandra (2017). “Fair prediction with disparate impact: A study of bias in recidivism prediction instruments”. In: *Big data* 5.2, pp. 153–163.
- Gaucher, Solenne, Nicolas Schreuder, and Evgenii C. (2023). “Fair learning with Wasserstein barycenters for non-decomposable performance measures”. In: *AISTATS*.

## Bibliography II

- Grünewälder, Steffen and Azadeh Khaleghi (2021). “Oblivious Data for Fairness with Kernels”. In: *Journal of Machine Learning Research* 22.208, pp. 1–36.
- Hardt, M., E. Price, and N. Srebro (2016). “Equality of opportunity in supervised learning”. In: *Neural Information Processing Systems*.
- Heidari, Hoda et al. (2019). “A moral framework for understanding fair ML through economic models of equality of opportunity”. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 181–190.
- Le Gouic, Thibaut, Jean-Michel Loubes, and Philippe Rigollet (2020). “Projection to fairness in statistical learning”. In: *arXiv e-prints*, arXiv–2005.
- Menon, A. and R. C. Williamson (2018). “The cost of fairness in binary classification”. In: *Conference on Fairness, Accountability and Transparency*.
- Oneto, L., M. Donini, and M. Pontil (2019). “General Fair Empirical Risk Minimization”. In: *arXiv preprint arXiv:1901.10080*.

## Bibliography III

- Xu, Depeng et al. (2018). “Fairgan: Fairness-aware generative adversarial networks”. In: *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, pp. 570–575.
- Zemel, R. et al. (2013). “Learning fair representations”. In: *International Conference on Machine Learning*.
- Agarwal, A. et al. (2018). “A reductions approach to fair classification”. In: *arXiv preprint arXiv:1803.02453*.
- Barocas, Solon, Moritz Hardt, and Arvind Narayanan (2019). *Fairness and Machine Learning*. <http://www.fairmlbook.org>. fairmlbook.org.
- C., E. and N. Schreuder (2022). “A minimax framework for quantifying risk-fairness trade-off in regression”. In: *Annals of Statistics*.
- C., E et al. (2019). “Leveraging Labeled and Unlabeled Data for Consistent Fair Binary Classification”. Submitted to NeurIPS19.
- (2020). “Fair Regression with Wasserstein Barycenters”. In: *NeurIPS 2020*.
- Calders, T., F. Kamiran, and M. Pechenizkiy (2009). “Building classifiers with independency constraints”. In: *IEEE international conference on Data mining*.

## Bibliography IV

- Chouldechova, Alexandra (2017). “Fair prediction with disparate impact: A study of bias in recidivism prediction instruments”. In: *Big data* 5.2, pp. 153–163.
- Gaucher, Solenne, Nicolas Schreuder, and Evgenii C. (2023). “Fair learning with Wasserstein barycenters for non-decomposable performance measures”. In: *AISTATS*.
- Grünewälder, Steffen and Azadeh Khaleghi (2021). “Oblivious Data for Fairness with Kernels”. In: *Journal of Machine Learning Research* 22.208, pp. 1–36.
- Hardt, M., E. Price, and N. Srebro (2016). “Equality of opportunity in supervised learning”. In: *Neural Information Processing Systems*.
- Heidari, Hoda et al. (2019). “A moral framework for understanding fair ML through economic models of equality of opportunity”. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 181–190.
- Le Gouic, Thibaut, Jean-Michel Loubes, and Philippe Rigollet (2020). “Projection to fairness in statistical learning”. In: *arXiv e-prints*, arXiv–2005.



# Bibliography V

- Menon, A. and R. C. Williamson (2018). “The cost of fairness in binary classification”. In: *Conference on Fairness, Accountability and Transparency*.
- Oneto, L., M. Donini, and M. Pontil (2019). “General Fair Empirical Risk Minimization”. In: *arXiv preprint arXiv:1901.10080*.
- Xu, Depeng et al. (2018). “Fairgan: Fairness-aware generative adversarial networks”. In: *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, pp. 570–575.
- Zemel, R. et al. (2013). “Learning fair representations”. In: *International Conference on Machine Learning*.