

An introduction to algorithmic fairness

+ a bit of fair regression with demographic parity constraint

Machine Learning in Montpellier
Theory & Practice

Evgenii Chzhen
CNRS, Univ. Paris-Saclay

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).

Motivating examples: Google translate

The screenshot shows a Google Translate interface. The source text in Hungarian is: "Ő szép. Ő okos. Ő olvas. Ő mosogat. Ő épít. Ő varr. Ő főz. Ő kutat. Ő gyereket nevel. Ő zenél. Ő takarító. Ő politikus. Ő sok pénzt keres. Ő süteményt süti. Ő professzor. Ő asszisztens." The target text in French is: "Elle est belle. Il est intelligent. Il lit. Elle lave la vaisselle. Il construit. Elle coud. Il enseigne. Elle cuisine. Il fait des recherches. Elle élève un enfant. Il joue de la musique. C'est une femme de ménage. C'est un politicien. Il gagne beaucoup d'argent. Elle prépare un gâteau. C'est un professeur. C'est une assistante." The interface includes language selection tabs at the top: ОПРЕДЕЛИТЬ ЯЗЫК, ВЕНГЕРСКИЙ, АНГЛИЙСКИЙ, ФРАНЦУЗСКИЙ, РУССКИЙ, АНГЛИЙСКИЙ, ФРАНЦУЗСКИЙ. Below the text boxes are various interaction icons like copy, edit, and share.

Source https://www.reddit.com/r/europe/comments/m9uphb/hungarian_has_no_gendered_pronouns_so_google/

Motivating examples: Google translate

The screenshot shows the Google Translate interface comparing Finnish (SUOMI) and English (ENGLANTI) translations. A blue callout box highlights a specific sentence from the English results: "She is taking care of the child." Below this, a circular note in Finnish discusses the complexity of the sentence's meaning.

Finnish (SUOMI):

- hän hoitaa lasta. hän lyö lasta. hän rakastaa vauvaa. hän vihaa lastaan.
- Hän hoitaa lasta.
- Hän lyö lasta.
- Hän rakastaa vauvaa.
- Hän vihaa lastaan.

English (ENGLANTI):

- She is taking care of the child.
- He hits the baby.
- She loves baby.
- He hates his child.

Note (Finnish):

Itse siis kääntäjällä näitä oon testannut. Eri selaimilla saattaa olla vähän eri näkymät. Aika paljon olisi tehtävää. Ei tupsu hommat eiliseen.

Source <https://kotiliesi.fi/ihmiset-ja-ilmiot/ilmiot/miksi-google-kaantajan-mukaan-mies-johtaa-ja-mies-tiskaa/>

Motivating examples: Twitter cropping

Fact: Twitter automatically **crops large** images in order to fit the size of an average mobile screen.

Original



Cropped



Motivating examples: Twitter cropping

Fact: Twitter automatically **crops large** images in order to fit the size of an **average mobile screen**.

Question: How will Twitter crop these two images??



Motivating examples: Twitter cropping

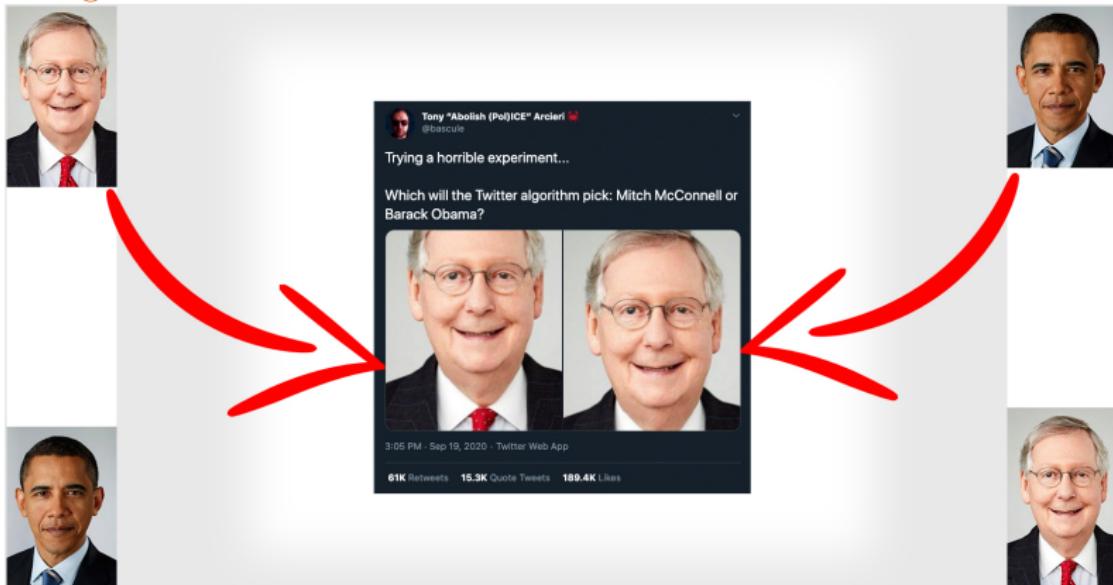
Fact: Twitter automatically **crops large** images in order to fit the size of an average mobile screen.



Twitter's response: (https://blog.twitter.com/en_us/topics/product/2020/transparency-image-cropping.html)

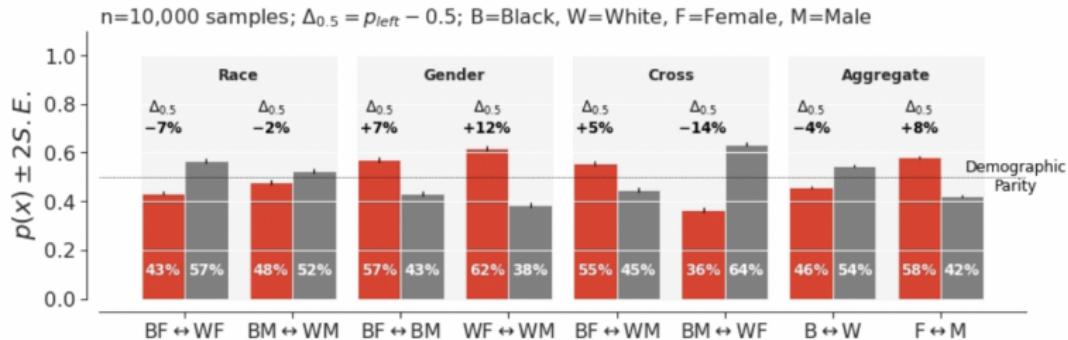
Motivating examples: Twitter cropping

Fact: Twitter automatically **crops large** images in order to fit the size of an average mobile screen.



Twitter's response: (https://blog.twitter.com/en_us/topics/product/2020/transparency-image-cropping.html)

Motivating examples: Twitter cropping



Source:

https://blog.twitter.com/engineering/en_us/topics/insights/2021/sharing-learnings-about-our-image-cropping-algorithm

More details in associated paper (Yee, Tantipongpipat, and Mishra, 2021)

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**

Today's plan: a **biased** intro to fairness

1. Fairness zoology with the emphasize on group fairness
2. Three types of approaches
3. Regression with demographic parity constraint—how to build post-processing algorithms?
4. Some open questions (on the board and if the time permits)

Individual fairness paradigm

“*treat like cases as like*” (\leq Aristotel)

“*Ensure that similar individuals are treated similarly*” (Dwork et al., 2012)

We observe $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$. Individual fairness *often* considers randomized predictions $f : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$

Individual fairness paradigm

“treat like cases as like” (\leq Aristotel)

“Ensure that *similar individuals* are *treated similarly*” (Dwork et al., 2012)

We observe $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$. Individual fairness *often* considers randomized predictions $f : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$

1. **Similarity of predictions:** $D : \Delta(\mathcal{Y}) \times \Delta(\mathcal{Y}) \rightarrow \mathbb{R}_+$
2. **Similarity of individuals:** $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$

A prediction $f : \mathcal{X} \rightarrow \Delta(\{0, 1\})$ is called *perfectly* (D, d) -individually fair if $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$

$$D(f(\mathbf{x}_1), f(\mathbf{x}_2)) \leq d(\mathbf{x}_1, \mathbf{x}_2)$$

see (Rothblum and Yona, 2018) for relaxations

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. DP (not differential privacy!) cares only about $\mathbf{X}|S$.
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. DP (not differential privacy!) cares only about $\mathbf{X}|S$.
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 $f(\mathbf{Z}) \in [0, 1]$. Above notions applied in this case ensure that any threshold will result in fair classification : incurs higher drop in accuracy; used in regression.

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](#)).

Taken from Chapter 2 of ([Barocas, Hardt, and Narayanan, 2019](#))

propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

Fairness notions based on the risk

- Equalized risks

$$\mathbb{E}[\ell(f(\mathbf{Z}), Y) \mid S = s] = \mathbb{E}[\ell(f(\mathbf{Z}), Y)]$$

- Minmax fairness: a fair prediction function should minimize

$$\max_{s \in \mathcal{S}} \mathbb{E}[\ell(f(\mathbf{Z}), Y) \mid S = s]$$

- Group-wise no-regret: a fair prediction function should satisfy no-regret by group

$$\max_{s \in \mathcal{S}} \left\{ \mathbb{E}[\ell(f(\mathbf{Z}), Y) \mid S = s] - \min_{f \in \mathcal{F}} \mathbb{E}[\ell(f(\mathbf{Z}), Y) \mid S = s] \right\} = o(1)$$

- Group-wise calibration:

$$\mathbb{E}[Y \mid S = s, f(\mathbf{Z})] = f(\mathbf{Z})$$

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 (Grünewälder and Khaleghi, 2021)
- ▶ GANs (Xu et al., 2018)

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

Add the fairness constraint into training

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

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

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

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; C. et al., 2019)
- ▶ Optimal transport based (C. et al., 2020)
- ▶ Conformal predictions (Romano et al., 2019)

Regression with Demographic Parity

an example of a post-processing scheme

joint works with C. Denis, M. Hebiri, L. Oneto, M. Pontil, and N. Schreuder

Regression + Demographic Parity

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

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

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

Demographic Parity fairness

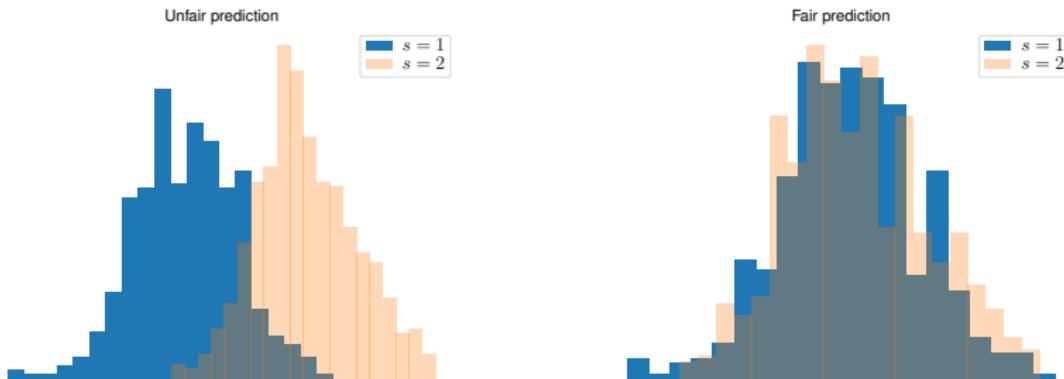
$$f(\mathbf{X}, S) \perp\!\!\!\perp S$$

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 $\text{Law}(f^*(\mathbf{X}, S) | S = s)$ have **finite second moment** and are **non-atomic** for any 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

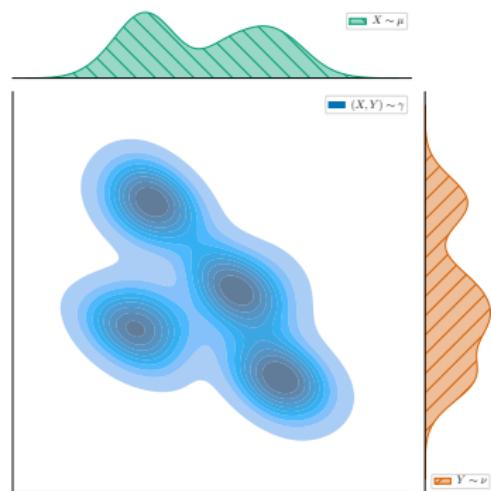


Figure: Transport plan illustration

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 (informal with $\mathcal{S} = \{1, 2\}$)

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

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

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

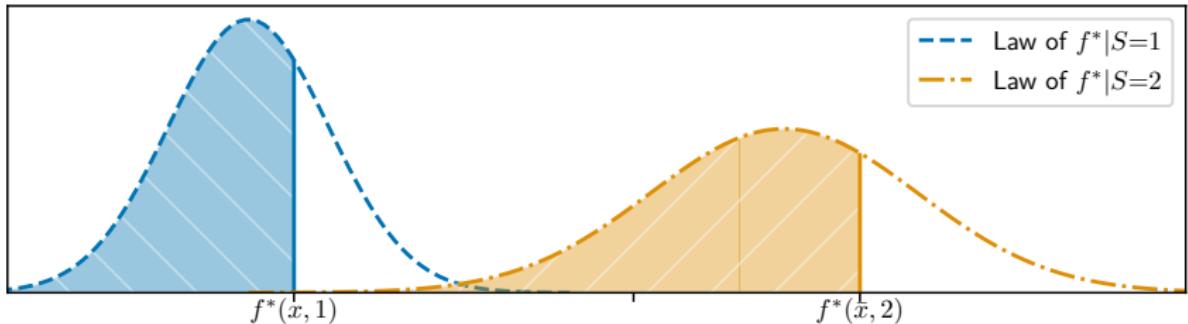
$T_{1 \rightarrow 2}^*$ – optimal transport map from $\text{Law}(f^* \mid S = 1)$ to $\text{Law}(f^* \mid S = 2)$.

(C. et al., 2020; Le Gouic, Loubes, and Rigollet, 2020)

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

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

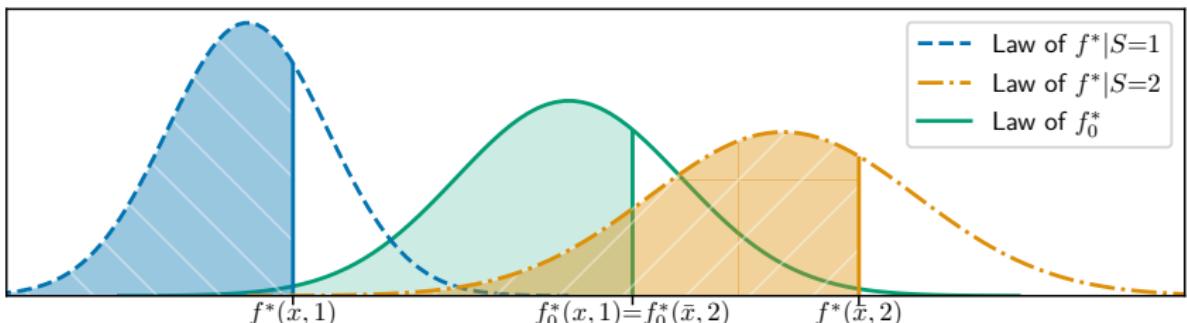
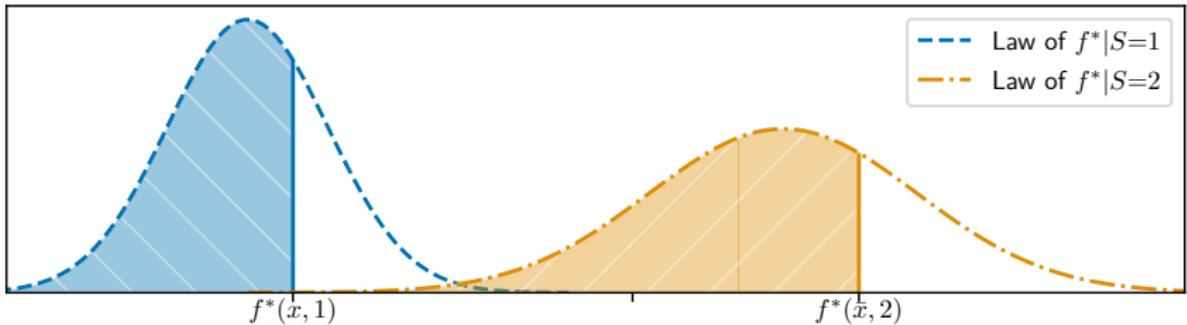
Fair optimal prediction f_0^* with $w_1 = 2/5$ and $w_2 = 3/5$



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

- **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 (informal)

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

$$\mathbf{E} \left[\sup_{t \in \mathbb{R}} \left| \mathbf{P}(\hat{f}_0(\mathbf{X}, S) \leq t \mid S=1, \mathcal{D}) - \mathbf{P}(\hat{f}_0(\mathbf{X}, S) \leq t \mid S=2, \mathcal{D}) \right| \right] \lesssim \frac{1}{\sqrt{N_1 \wedge N_2}}$$

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

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

(C. et al., 2020)

Additional assumptions: $(f^*(\mathbf{X}, S) \mid S = s)$ admits density which is **upper** and **lower** bounded (leading constant for the risk rate depends on this upper/lower bound)

How to measure unfairness ?

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

- ▶ **Problem:** too stiff — either **fair** or **unfair**.
- ▶ **Question:** how to quantify unfairness *i.e.*, violation of DP?
- ▶ **Question:** how to trade accuracy for fairness?

Popular measure is based on KS distance ([Agarwal, Dudik, and Wu, 2019](#); [Oneto, Donini, and Pontil, 2019](#))

$$\mathcal{U}_{\text{KS}}(f) := \sum_{s \in \mathcal{S}} \text{KS}(\text{Law}(f(\mathbf{X}, S) \mid S = s), \text{Law}(f(\mathbf{X}, S)))$$

How to measure unfairness ?

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

- ▶ **Problem:** too stiff — either **fair** or **unfair**.
- ▶ **Question:** how to quantify unfairness *i.e.*, violation of DP?
- ▶ **Question:** how to trade accuracy for fairness?

Popular measure is based on KS distance ([Agarwal, Dudik, and Wu, 2019](#); [Oneto, Donini, and Pontil, 2019](#))

$$\mathcal{U}_{\text{KS}}(f) := \sum_{s \in \mathcal{S}} \text{KS}(\text{Law}(f(\mathbf{X}, S) \mid S = s), \text{Law}(f(\mathbf{X}, S)))$$

We consider: $\mathcal{U}(f) = \min_{\nu} \sum_{s \in \mathcal{S}} w_s W_2^2(\text{Law}(f(\mathbf{X}, S) \mid S = s), \nu)$

From previous result: $\mathcal{R}(f_0^*) = \mathcal{U}(f^*)$

Improving unfairness oracles

α -Relative Improvement $f_\alpha^* \in \arg \min \left\{ \mathcal{R}(f) : \boxed{\mathcal{U}(f) \leq \alpha \mathcal{U}(f^*)} \right\}$

- ▶ f_α^* – $1/\alpha$ times fairer than f^* .
- ▶ f_0^* – optimal DP fair prediction.
- ▶ $f_1^* \equiv f^*$ – Bayes optimal prediction.

Improving unfairness oracles

α -Relative Improvement $f_\alpha^* \in \arg \min \left\{ \mathcal{R}(f) : \boxed{\mathcal{U}(f) \leq \alpha \mathcal{U}(f^*)} \right\}$

- ▶ f_α^* – $1/\alpha$ times fairer than f^* .
- ▶ f_0^* – optimal DP fair prediction.
- ▶ $f_1^* \equiv f^*$ – Bayes optimal prediction.

Theorem

Under Assumption (A), for all $\alpha \in [0, 1]$ it holds that

$$f_\alpha^* \equiv \sqrt{\alpha} f_1^* + (1 - \sqrt{\alpha}) f_0^*$$

$$\alpha\text{-RI} \equiv \sqrt{\alpha} \cdot \text{Bayes optimal} + (1 - \sqrt{\alpha}) \cdot \text{Fair optimal}$$

(C. and Schreuder, 2020)

N.B. We can use previous algorithm to estimate f_0^* and *any* standard algorithm for estimation of f^*

Idea of the proof

Goal: $\min_{f: \mathcal{Z} \rightarrow \mathbb{R}} \left\{ \sum_{s=1}^K w_s \mathbb{E}[(f(\mathbf{X}, S) - f^*(\mathbf{X}, S))^2 \mid S = s] : \mathcal{U}(f) \leq \alpha \mathcal{U}(f^*) \right\}$

LB: $\sum_{s=1}^K w_s W_2^2(\text{Law}(f(\mathbf{X}, S) \mid S = s), \text{Law}(f^*(\mathbf{X}, S) \mid S = s))$

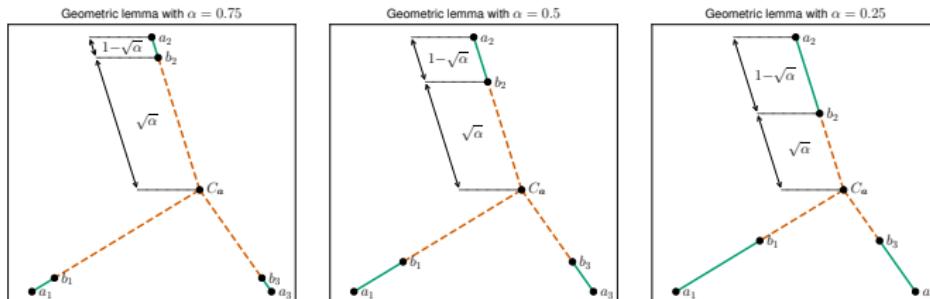
Idea of the proof

Goal: $\min_{f: \mathcal{Z} \rightarrow \mathbb{R}} \left\{ \sum_{s=1}^K w_s \mathbb{E}[(f(\mathbf{X}, S) - f^*(\mathbf{X}, S))^2 \mid S = s] : \mathcal{U}(f) \leq \alpha \mathcal{U}(f^*) \right\}$

LB: $\sum_{s=1}^K w_s W_2^2(\text{Law}(f(\mathbf{X}, S) \mid S = s), \text{Law}(f^*(\mathbf{X}, S) \mid S = s))$

New problem

$\min_{\mathbf{b} \in \mathcal{P}_2^K(\mathbb{R})} \left\{ \sum_{s=1}^K w_s W_2^2(b_s, a_s) : \sum_{s=1}^K w_s W_2^2(b_s, C_{\mathbf{b}}) \leq \alpha \sum_{s=1}^K w_s W_2^2(a_s, C_{\mathbf{a}}) \right\}$



Risk/fairness trade-off

α -Relative Improvement $f_\alpha^* \in \arg \min \left\{ \mathcal{R}(f) : \boxed{\mathcal{U}(f) \leq \alpha \mathcal{U}(f^*)} \right\}$

Proposition

Under Assumption (A), for all $\alpha \in [0, 1]$ it holds that

$$\mathcal{R}(f_\alpha^*) = (1 - \sqrt{\alpha})^2 \boxed{\mathcal{U}(f^*)} \quad \text{and} \quad \mathcal{U}(f_\alpha^*) = \alpha \boxed{\mathcal{U}(f^*)}$$

(C. and Schreuder, 2020)

Risk/fairness trade-off

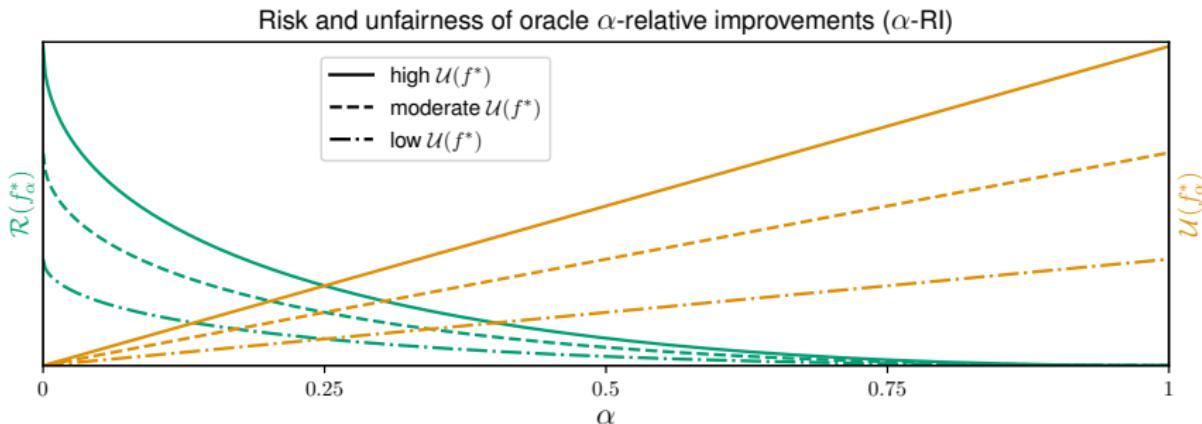
α -Relative Improvement $f_\alpha^* \in \arg \min \left\{ \mathcal{R}(f) : \boxed{\mathcal{U}(f) \leq \alpha \mathcal{U}(f^*)} \right\}$

Proposition

Under Assumption (A), for all $\alpha \in [0, 1]$ it holds that

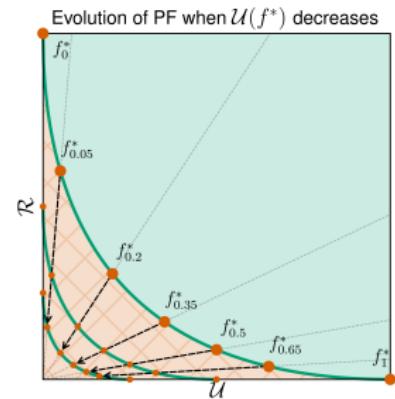
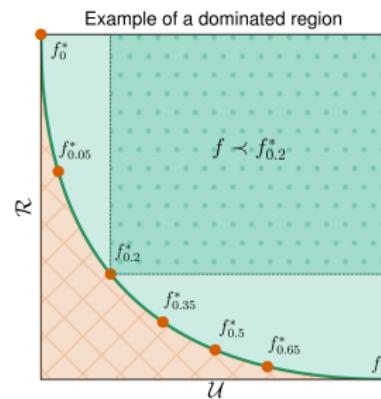
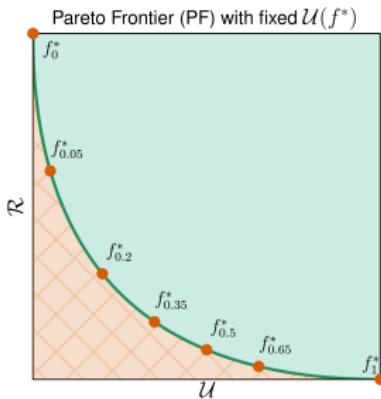
$$\mathcal{R}(f_\alpha^*) = (1 - \sqrt{\alpha})^2 \boxed{\mathcal{U}(f^*)} \quad \text{and} \quad \mathcal{U}(f_\alpha^*) = \alpha \boxed{\mathcal{U}(f^*)}$$

(C. and Schreuder, 2020)



Pareto efficiency

- ▶ Multi-objective optimization: $\min_{f: \mathcal{Z} \rightarrow \mathbb{R}} (\mathcal{U}(f), \mathcal{R}(f))$.
- ▶ Each prediction f defines a point $(\mathcal{U}(f), \mathcal{R}(f))$.
- ▶ f is **dominated** by f' iff $\mathcal{R}(f') \leq \mathcal{R}(f)$ and $\mathcal{U}(f') \leq \mathcal{U}(f)$.



Minimax statistical framework

Data: $(\mathbf{X}_1, S_1, Y_1), \dots, (\mathbf{X}_n, S_n, Y_n) \stackrel{i.i.d.}{\sim} \mathbf{P}_{(f^*, \boldsymbol{\theta})}$, $(f^*, \boldsymbol{\theta}) \in \mathcal{F} \times \Theta$

Given $\alpha \in [0, 1]$ and $t > 0$, the goal of the statistician is to construct an estimator \hat{f} , which simultaneously satisfies

1. Uniform fairness guarantee:

$$\forall (f^*, \boldsymbol{\theta}) \in \mathcal{F} \times \Theta \quad \mathbf{P}_{(f^*, \boldsymbol{\theta})} \left(\mathcal{U}(\hat{f}) \leq \alpha \mathcal{U}(f^*) \right) \geq 1 - t ,$$

2. Uniform risk guarantee:

$$\forall (f^*, \boldsymbol{\theta}) \in \mathcal{F} \times \Theta \quad \mathbf{P}_{(f^*, \boldsymbol{\theta})} \left(\mathcal{R}(\hat{f}) \leq r_{n, \alpha, f^*}(\mathcal{F}, \Theta, t) \right) \geq 1 - t .$$

Problem-dependent lower bound

For $t \in (0, 1)$, let $\delta_n(\mathcal{F}, \Theta, t)$ be a sequence that verifies

$$\inf_{\hat{f}} \sup_{(f^*, \boldsymbol{\theta}) \in \mathcal{F} \times \Theta} \mathbf{P}_{(f^*, \boldsymbol{\theta})} \left(\mathcal{R}(\hat{f}) \geq \delta_n(\mathcal{F}, \Theta, t) \right) \geq t$$

Theorem

Any estimator \hat{f} satisfying

$$\inf_{(f^*, \boldsymbol{\theta}) \in \mathcal{F} \times \Theta} \mathbf{P}_{(f^*, \boldsymbol{\theta})} \left(\mathcal{U}(\hat{f}) \leq \alpha \mathcal{U}(f^*) \right) \geq 1 - t'$$

verifies

$$\sup_{\substack{f^* \in \mathcal{F} \\ \boldsymbol{\theta} \in \Theta}} \mathbf{P}_{(f^*, \boldsymbol{\theta})} \left(\mathcal{R}^{1/2}(\hat{f}) \geq \delta_n^{1/2}(\mathcal{F}, \Theta, t) \vee \underbrace{(1 - \sqrt{\alpha}) \mathcal{U}^{1/2}(f^*)}_{= \mathcal{R}^{1/2}(f_\alpha^*)} \right) \geq t \wedge (1 - t')$$

Conclusions

1. Individual fairness – predict with Lipschitz functions

$$D(f(\mathbf{x}), f(\mathbf{x}')) \leq d(\mathbf{x}, \mathbf{x}')$$

2. Group fairness – enforce some independence criterion

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

3. Regression with demographic parity ($f(\mathbf{Z}) \perp\!\!\!\perp S$) can be characterized by Wasserstein barycenter problem

$$\mathcal{R}(f_0^*) = \mathcal{U}(f^*)$$

4. Risk/fairness trade-off can be characterized explicitly for introduced notion of unfairness

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;

Causal fairness

- ▶ Causal fairness aims to **identify** sources of unfairness
- ▶ The relations between attributes (X, S) and their influence on outcome is **modeled by structural equations**
- ▶ These structural equations **capture the influence** of sensitive attributes
- ▶ The objective is to **remove all discriminatory influences**

The notions of causal fairness heavily rely on the causal model. The accuracy of this model is critical.

Definition

A prediction is **counterfactually fair** if, in the causal graph, it does not depend on a **descendant of the sensitive attribute**.

N.B. I know nothing about causality. (questions are not allowed ;))

Bibliography I

- Agarwal, A. et al. (2018). “A reductions approach to fair classification”. In: *arXiv preprint arXiv:1803.02453*.
- Agarwal, Alekh, Miroslav Dudik, and Zhiwei Steven Wu (2019). “Fair regression: Quantitative definitions and reduction-based algorithms”. In: *International Conference on Machine Learning*. PMLR, pp. 120–129.
- Barocas, Solon, Moritz Hardt, and Arvind Narayanan (2019). *Fairness and Machine Learning*. <http://www.fairmlbook.org>. fairmlbook.org.
- C., E. and N. Schreuder (2020). “A minimax framework for quantifying risk-fairness trade-off in regression”. In: *arXiv preprint arXiv:2007.14265*.
- 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.

Bibliography II

- Dwork, Cynthia et al. (2012). “Fairness through awareness”. In: *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226.
- 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.
- Oneto, L., M. Donini, and M. Pontil (2019). “General Fair Empirical Risk Minimization”. In: *arXiv preprint arXiv:1901.10080*.
- Romano, Yaniv et al. (2019). “With malice towards none: Assessing uncertainty via equalized coverage”. In: *arXiv preprint arXiv:1908.05428*.

Bibliography III

- Rothblum, Guy N and Gal Yona (2018). “Probably Approximately Metric-Fair Learning”. In: *arXiv e-prints*, arXiv–1803.
- 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.
- Yee, K., U. Tantipongpipat, and S. Mishra (2021). *Image Cropping on Twitter: Fairness Metrics, their Limitations, and the Importance of Representation, Design, and Agency*. arXiv: 2105.08667 [cs.CY].
- 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*.
- Agarwal, Alekh, Miroslav Dudik, and Zhiwei Steven Wu (2019). “Fair regression: Quantitative definitions and reduction-based algorithms”. In: *International Conference on Machine Learning*. PMLR, pp. 120–129.
- Barocas, Solon, Moritz Hardt, and Arvind Narayanan (2019). *Fairness and Machine Learning*. <http://www.fairmlbook.org>. fairmlbook.org.
- C., E. and N. Schreuder (2020). “A minimax framework for quantifying risk-fairness trade-off in regression”. In: *arXiv preprint arXiv:2007.14265*.

Bibliography IV

- 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.
- Dwork, Cynthia et al. (2012). “Fairness through awareness”. In: *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226.
- 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*.

Bibliography V

- 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.
- Oneto, L., M. Donini, and M. Pontil (2019). “General Fair Empirical Risk Minimization”. In: *arXiv preprint arXiv:1901.10080*.
- Romano, Yaniv et al. (2019). “With malice towards none: Assessing uncertainty via equalized coverage”. In: *arXiv preprint arXiv:1908.05428*.
- Rothblum, Guy N and Gal Yona (2018). “Probably Approximately Metric-Fair Learning”. In: *arXiv e-prints*, arXiv–1803.
- 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.
- Yee, K., U. Tantipongpipat, and S. Mishra (2021). *Image Cropping on Twitter: Fairness Metrics, their Limitations, and the Importance of Representation, Design, and Agency*. arXiv: 2105.08667 [cs.CY].

Bibliography VI

Zemel, R. et al. (2013). “Learning fair representations”. In: *International Conference on Machine Learning*.