

Fairness and Discrimination in Insurance (Causal Inference meets Optimal Transport)

Arthur Charpentier¹

(with Emmanuel Flachaire², Ewen Gallic² François Hu¹ & Philipp Ratz¹)

Actuarial Science Workshop on 2023 SSC Annual Meeting, Ottawa, May 2023

¹ UQAM, Canada, ² Aix-Marseille University (AMSE), France

Preamble (extension of slide 7 in Valdez (2023) talk)

	CA	HI	GA	NC	NY	MA	PA	FL	TX	AL	ON	NB	NL	QC
Gender	X	X	•	X	•	X	X	•	•	•	•	X	X	•
Age	X	X	•	X*	•	X	•	•	•	•	*	•	X	•
Driving experience	•	X	•	•	•	•	•	•	•	•	•	•	•	•
Credit history	X	X	•	•	•	X	•*	•	•	X*	X	•*	X	•
Education	X	X	X	X	X	X	•	•	•	•	•	•	•	•
Occupation	X	X	X	•	X	X	•	•	•	•	•	•	•	•
Employment status	X	X	X	•	X	X	•	•	•	•	•	•	•	•
Marital status	•	X	•	•	•	X	•	•	•	•	•	•	•	•
Housing situation	X	X	•	•	•	X	•	•	•	X	X	•	•	•
Address/ZIP code	•	•	•	•	•	•	•	•	•	X	X	•	•	•
Insurance history	•	•	•	•	•	•	•	•	•	•	•	•	•	•

CA: Californie, HI: Hawaii, GA: Georgia, NC: Caroline du nord, NY: New York, MA: Massachusetts, PA: Pennsylvanie, FL: Floride, TX: Texas, AL: Alberta, ON: Ontario, NB: Nouveau-Brunswick, NL: Terre-Neuve-et-Labrador, QC: Québec

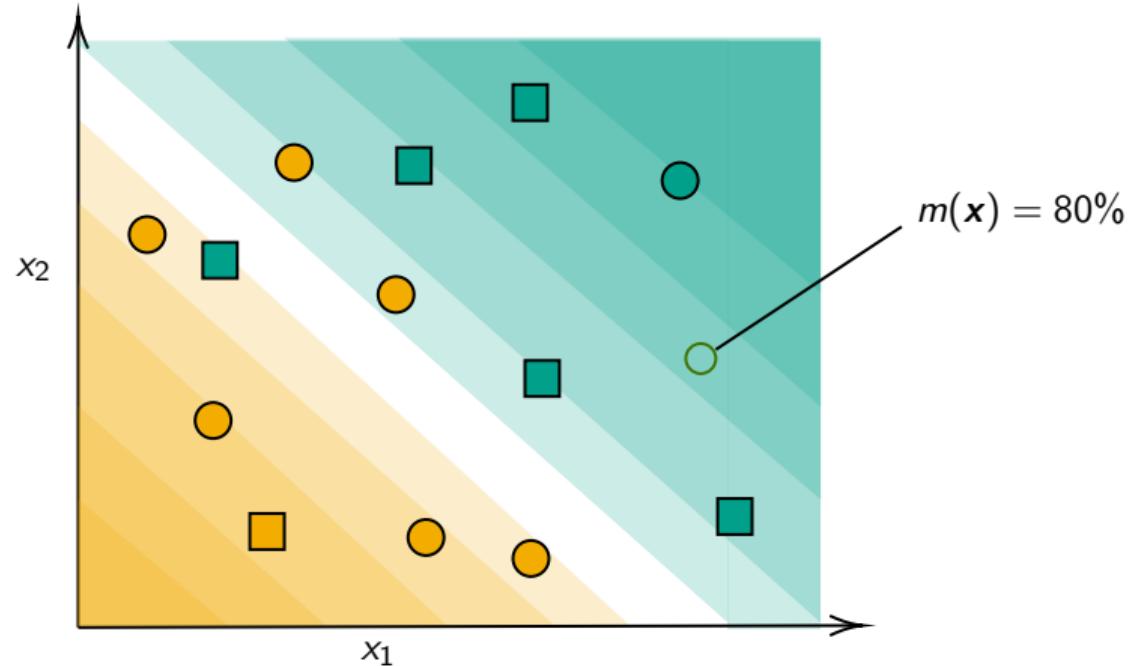
Motivation

- ▶ "at the core of insurance business lies discrimination between risky and non-risky insureds", Avraham (2017)
 - ▶ "Technology is neither good nor bad; nor is it neutral ", Kranzberg (1986)
 - ▶ "Machine learning won't give you anything like gender neutrality 'for free' that you didn't explicitly ask for ", Kearns and Roth (2019)

See Charpentier (2022, 2023a) for more details.

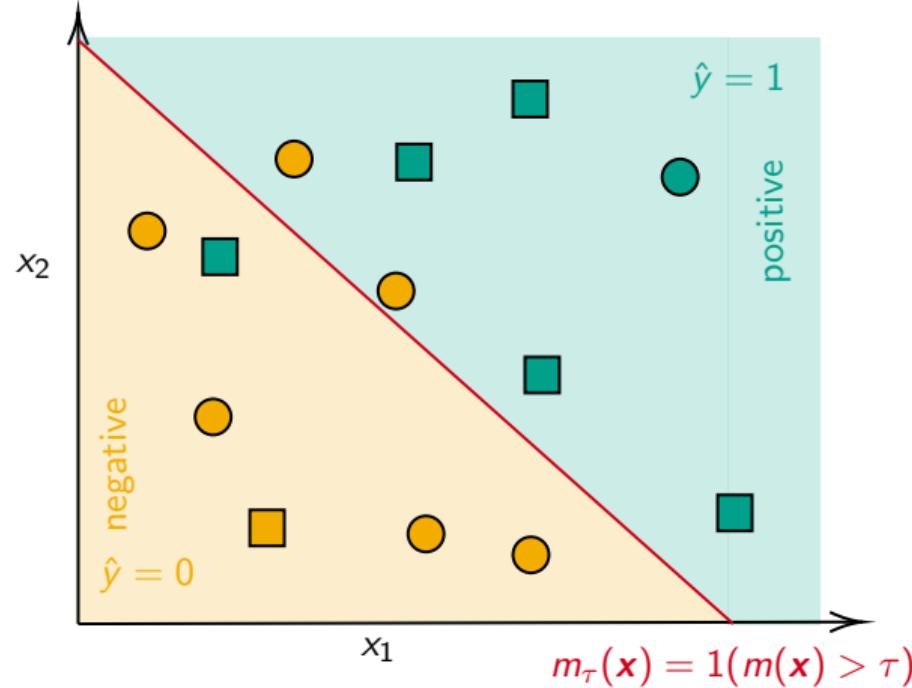


Notations



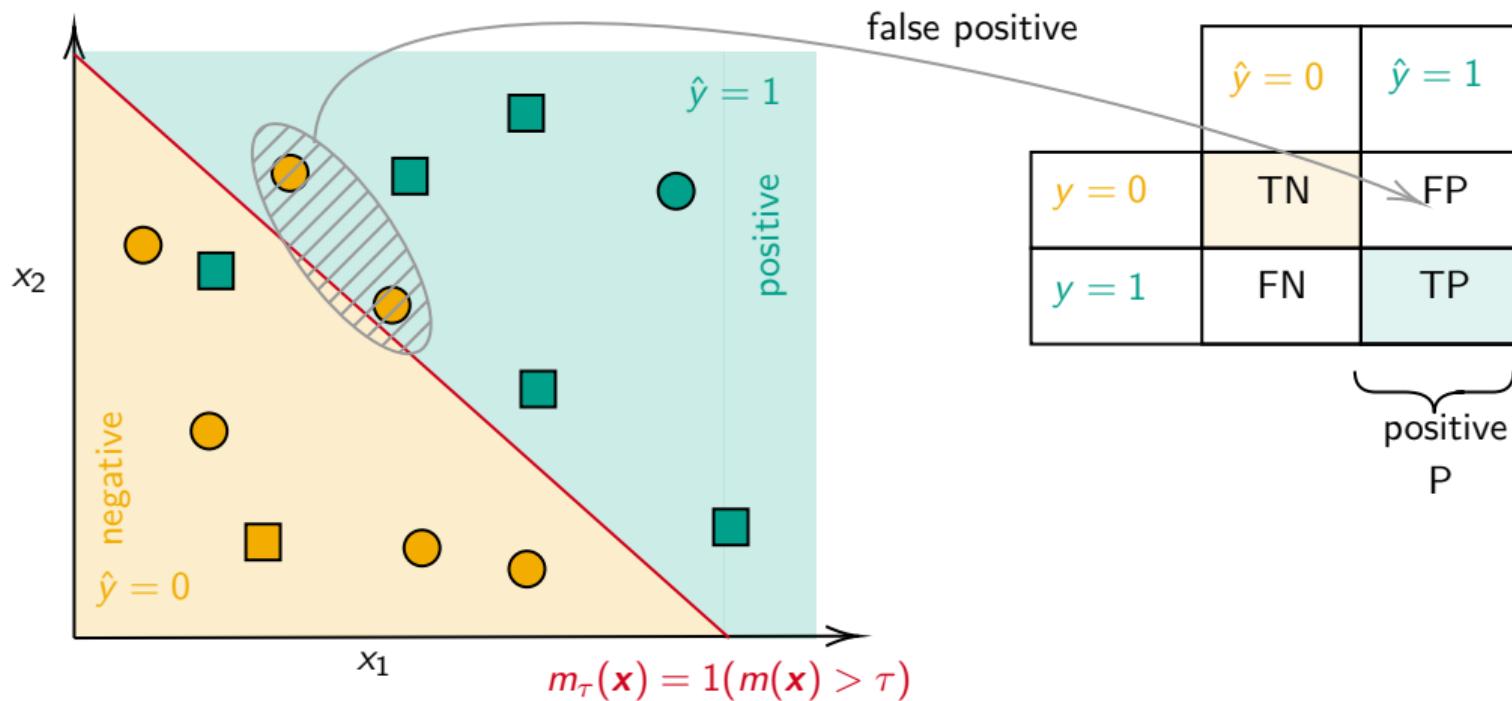
Dataset $\{(y_i, \mathbf{x}_i)\}$, $y_i \in \{0, 1\}$, score $m(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})$

Notations



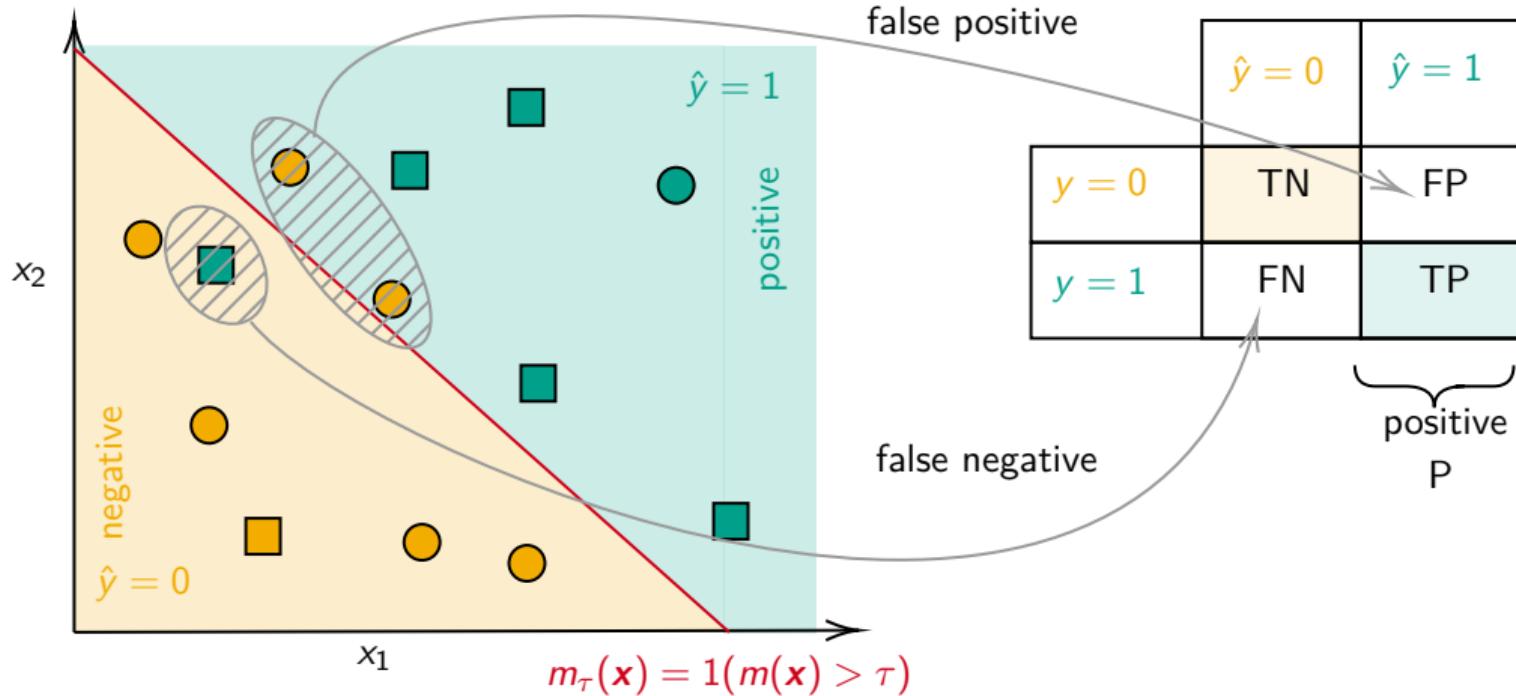
For $\tau \in (0, 1)$, define classifier $m_\tau(\mathbf{x}) = \mathbf{1}(m(\mathbf{x}) > \tau) \in \{0, 1\}$

Notations

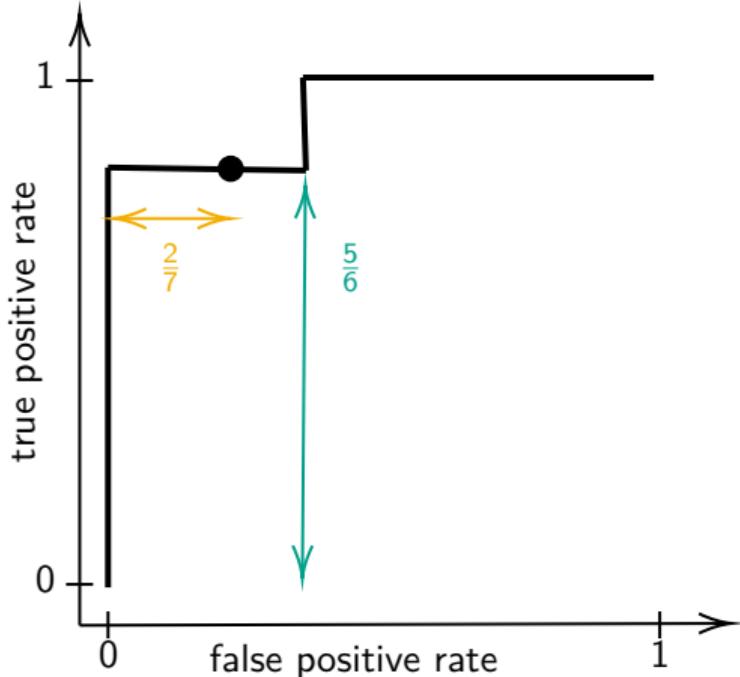
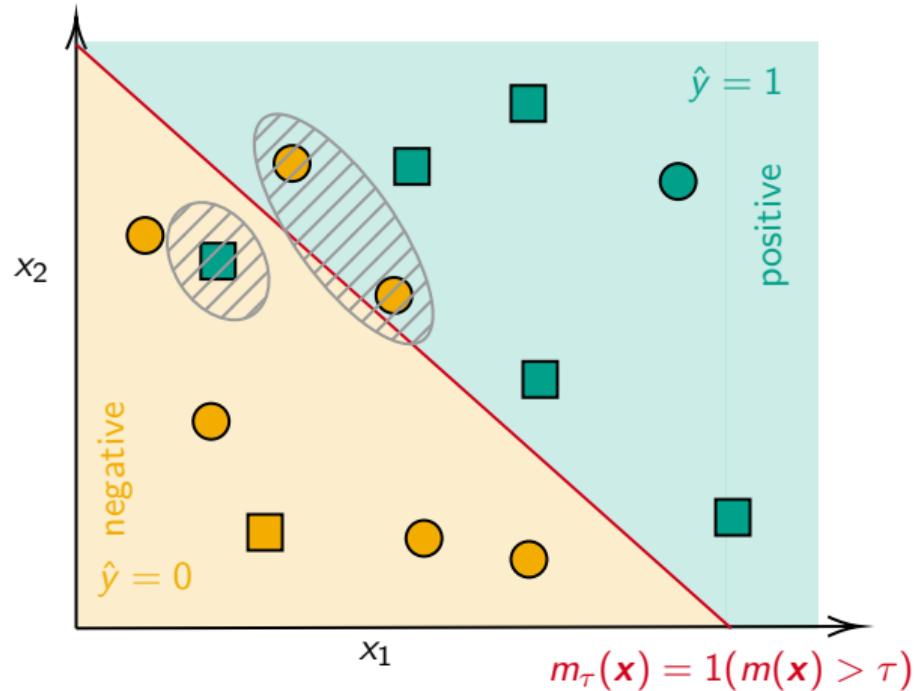


From m_τ consider the associated confusion matrix (negative/positive, true/false)

Notations

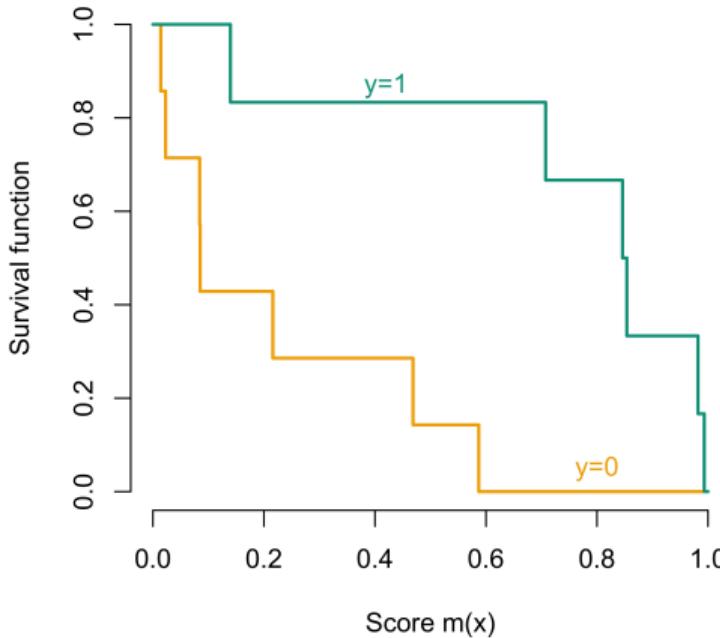
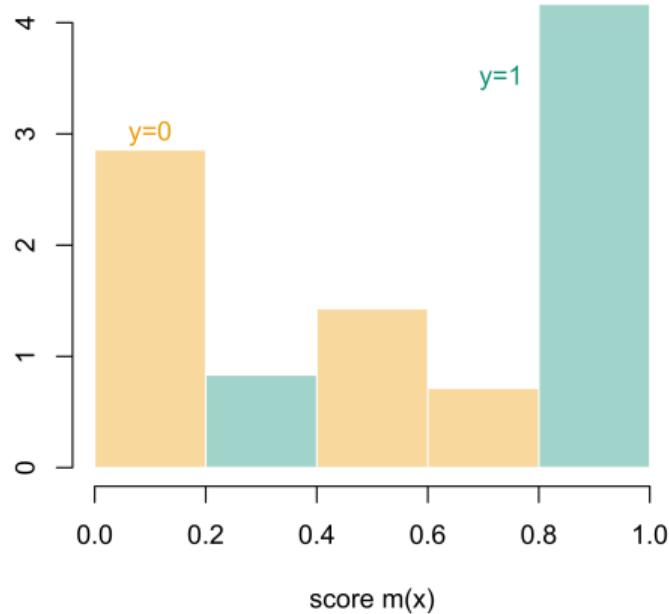


Notations



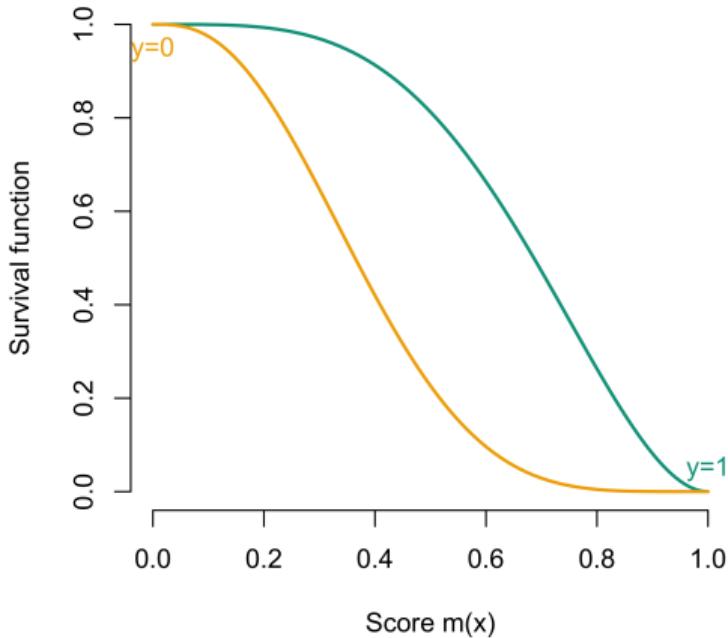
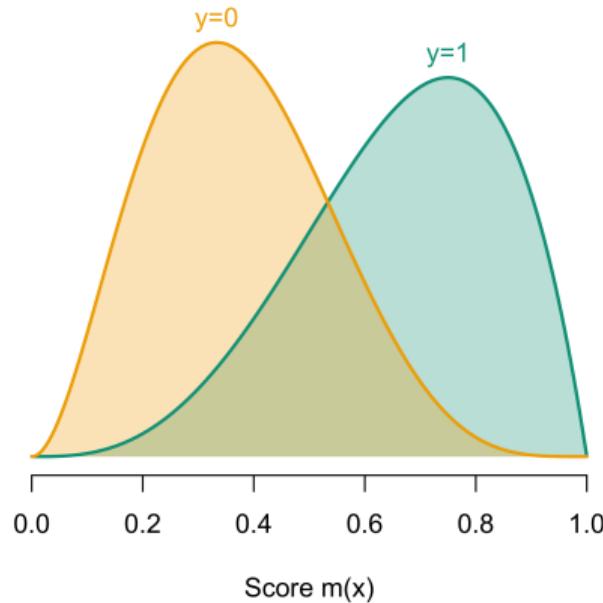
Changing τ lead to the ROC curve (TPR against FPR)

Notations



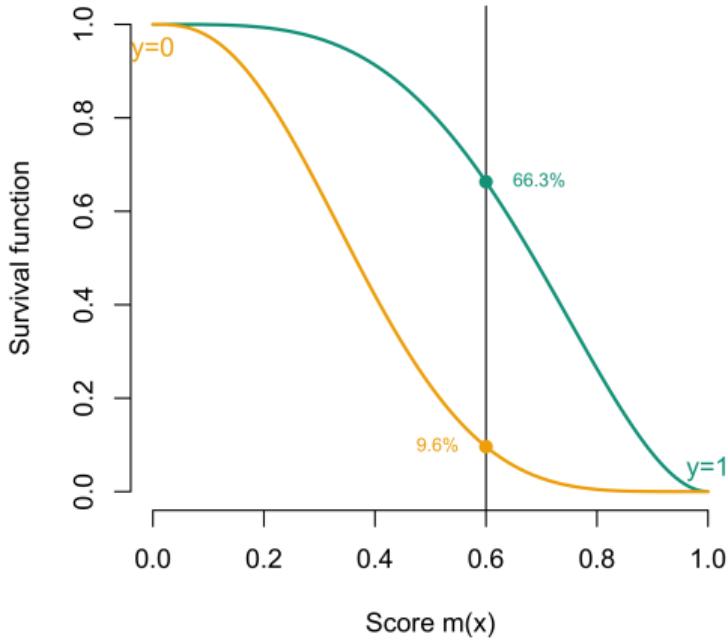
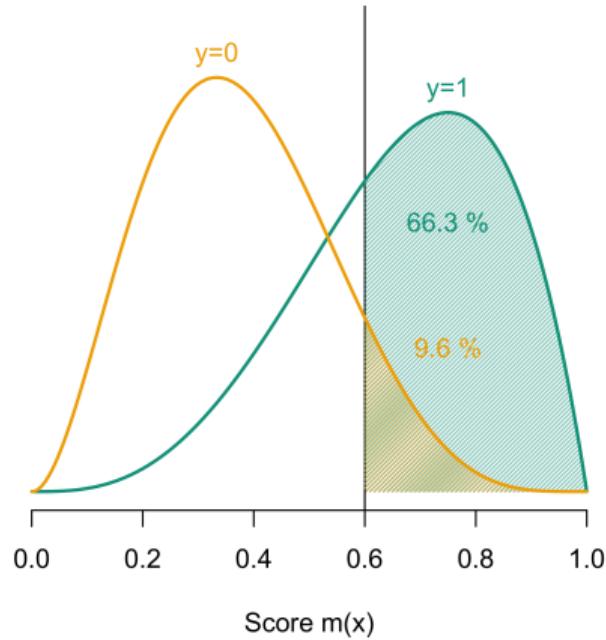
Consider the distribution of scores $m(x_i)$ when $y_i = 0$ and $y_i = 1$

Notations



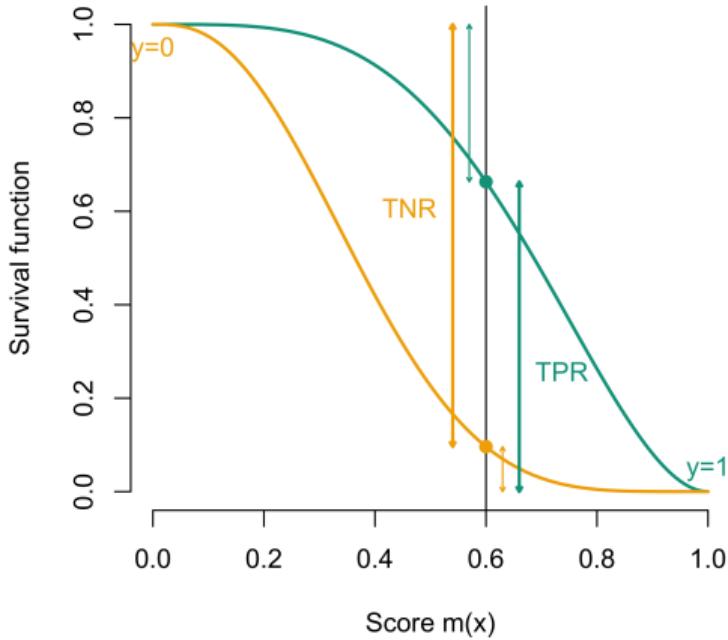
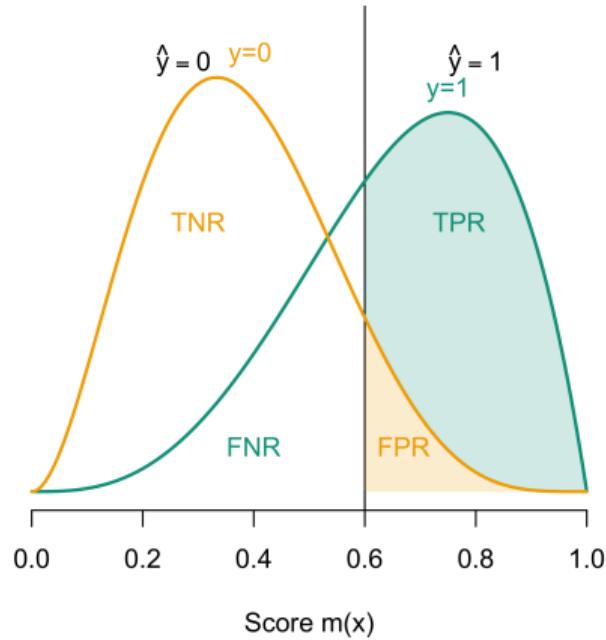
Consider the distribution of scores $m(x_i)$ when $y_i = 0$ and $y_i = 1$ (continuous version)

Notations



$\tau = 60\%$, $\text{FPR} \sim 9.6\%$ and $\text{TPR} \sim 66.3\%$ (continuous version)

Notations



$\tau = 60\%$, Given τ , one visualize FPR, TNR, TPR and FNR (continuous version)

Notations

$y \in \{0, 1\}$	variable of interest (binary)
$s \in \{0, 1\}$	protected variable (sensitive)
$\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$	'explanatory' variables
$m : \mathcal{X} \rightarrow [0, 1]$	score $m(\mathbf{x}) = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})$, e.g. $\text{logit}(\mathbf{x}^\top \boldsymbol{\beta})$
	$m(\mathbf{x}) = \mathbb{P}[Y = 1 \mathbf{X} = \mathbf{x}]$
$\tau \in (0, 1)$	threshold
$\hat{y} \in \{0, 1\}$	classifier $\hat{y} = m_\tau(\mathbf{x}) \mathbf{1}(m(\mathbf{x}) > \tau)$

Remark for people who prefer regression over classification:

instead of $\mathbb{P}[\hat{Y} = 1 | \dots]$ or $\mathbb{P}[Y = 1 | \dots]$, read

$\mathbb{E}[\hat{Y} \dots]$ or $\mathbb{E}[Y \dots]$	weak version
$\mathbb{P}[\hat{Y} \in A \dots]$ or $\mathbb{P}[Y \in A \dots] \quad \forall A \subset \mathcal{Y}$	strong version

Group fairness

At least 21 definitions of "fairness", Narayanan (2018).

"*focus on equality of treatment among groups of people from criteria requiring equality of treatment among couples of similar individuals*", Castelnovo et al. (2022)

Fairness Through Unawareness

Kusner et al. (2017) We will speak of fairness through unawareness if the sensitive attribute s is not explicitly used in the decision function \hat{y} , i.e. neither in the construction of the score m , nor in the choice of the threshold level τ , allowing to pass from m to \hat{y} .

... obviously not sufficient.

Demographic Parity

"independence is equivalent to requiring the same positive prediction ratio across groups identified by the sensitive features. This form of independence is usually known as Demographic Parity (DP), statistical parity, or sometimes as group fairness", Castelnovo et al. (2022)

Demographic Parity

Corbett-Davies et al. (2017), Agarwal (2021) A decision function \hat{y} satisfies demographic parity if $\hat{Y} \perp\!\!\!\perp S$, i.e.

$$\mathbb{P}[\hat{Y} = y|S = 0] = \mathbb{P}[\hat{Y} = y|S = 1], \forall y \in \{0, 1\}$$

Disparate impact

Feldman et al. (2015)] A decision function \hat{Y} has a disparate impact, for a given threshold d , if,

$$\min\left\{\frac{\mathbb{P}[\hat{Y} = 1|S = 0]}{\mathbb{P}[\hat{Y} = 1|S = 1]}, \frac{\mathbb{P}[\hat{Y} = 1|S = 1]}{\mathbb{P}[\hat{Y} = 1|S = 0]}\right\} < d \text{ (usually 80%).}$$

Equalized Odds

Separation criteria: independence when $Y = 0$ or $Y = 1$

True positive equality, Equalized Odds

Hardt et al. (2016) We will speak of equality of opportunity, or parity of true positives, if

$$\mathbb{P}[\hat{Y} = 1|S = 0, Y = 1] = \mathbb{P}[\hat{Y} = 1|S = 1, Y = 1]$$

or equivalently $\text{TPR}_0 = \frac{\text{TP}_0}{\text{FN}_0 + \text{TP}_0} = \frac{\text{TP}_1}{\text{FN}_1 + \text{TP}_1} = \text{TPR}_1$.

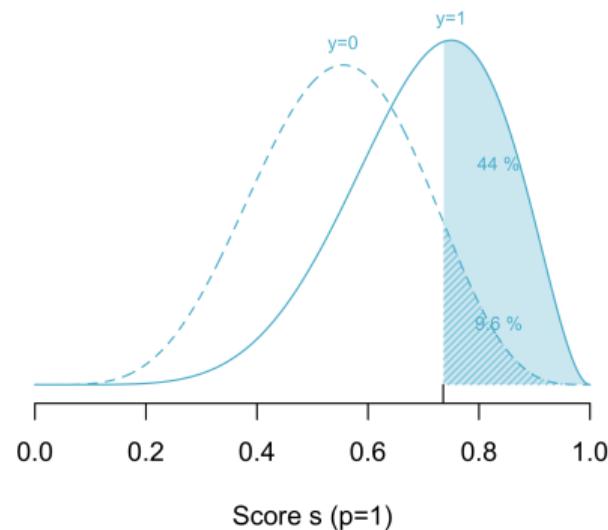
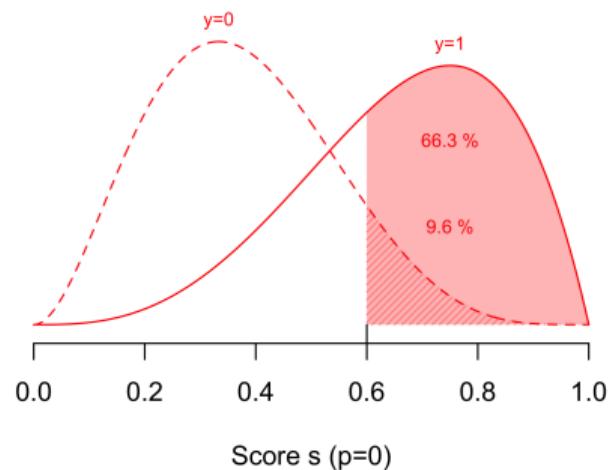
False positive equality

Hardt et al. (2016) We will speak of equality of false positives if

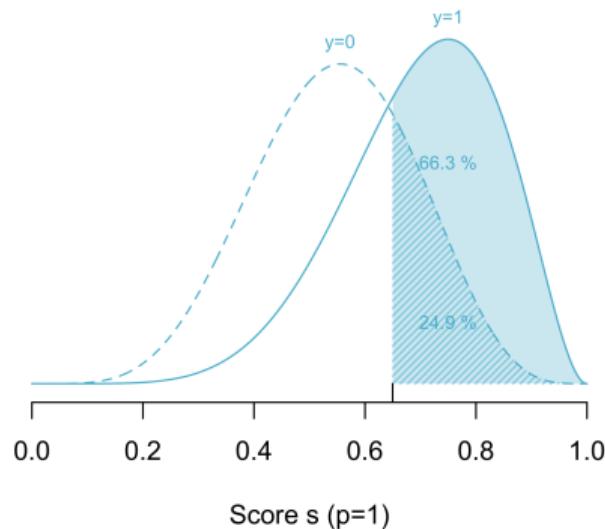
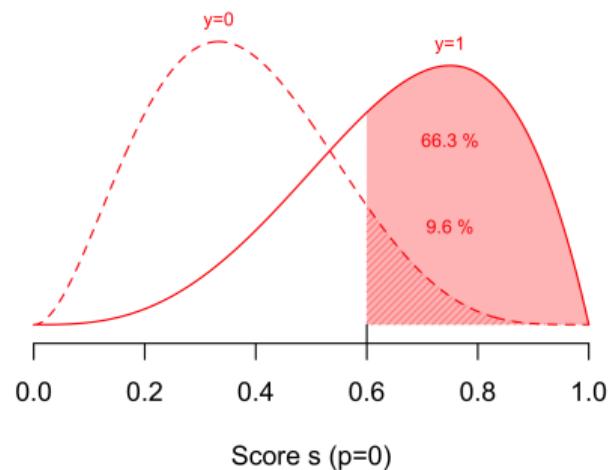
$$\mathbb{P}[\hat{Y} = 1|S = 0, Y = 0] = \mathbb{P}[\hat{Y} = 1|S = 1, Y = 0],$$

or equivalently $\text{FPR}_0 = \frac{\text{FP}_0}{\text{TN}_0 + \text{FP}_0} = \frac{\text{FP}_1}{\text{TN}_1 + \text{FP}_1} = \text{FPR}_1$.

Equalized Odds



Equalized Odds



Equalized Odds

Equalized opportunity

Hardt et al. (2016) The parity of false positives and true positives is called equality of opportunity,

$$\mathbb{P}[\hat{Y} = 1|S = 0, Y = y] = \mathbb{P}[\hat{Y} = 1|S = 1, Y = y], \forall y \in \{0, 1\}$$

in other words, $\hat{Y} \perp\!\!\!\perp S$ conditionally on Y .

One can also use any measure based on confusion matrices, such as ϕ , introduced by Matthews (1975),

ϕ -fairness

Chicco and Jurman (2020) We will have ϕ -fairness if $\phi_1 = \phi_0$, where ϕ_s denotes Matthews correlation coefficient for the s group,

$$\phi_s = \frac{\text{TP}_s \cdot \text{TN}_s - \text{FP}_s \cdot \text{FN}_s}{\sqrt{(\text{TP}_s + \text{FP}_s)(\text{TP}_s + \text{FN}_s) \cdot (\text{TN}_s + \text{FP}_s)(\text{TN}_s + \text{FN}_s)}}$$

Equalized Odds

All those measures are based on some choice of thresholds, but it is also possible to consider a global measures of calibration, such as the area under the curve,

AUC fairness

Borkan et al. (2019) We will have AUC fairness if $\text{AUC}_1 = \text{AUC}_0$, where AUC_s is the AUC for the s group.

We find a similar idea in Beutel et al. (2019). The problem with the AUC is that we can have identical AUCs, but very different underlying ROC curves. So, it can be interesting to consider a notion of fairness based on the ROC curves. As a reminder, we had defined the ROC curve as $t \mapsto \text{TPR} \circ \text{FPR}^{-1}(t)$.

Equality of ROC curves

Vogel et al. (2021) Let $\text{FRP}_s(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 0, S = s]$ and $\text{TPR}_s(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 1, S = s]$. Set $\Delta_{\text{TPR}}(t) = \text{TPR}_1 \circ \text{TPR}_0^{-1}(t) - t$ et $\Delta_{\text{FPR}}(t) = \text{FPR}_1 \circ \text{FPR}_0^{-1}(t) - t$. We will have an fairness of ROC curves if $\|\Delta_{\text{TPR}}\|_\infty = \|\Delta_{\text{FPR}}\|_\infty = 0$.

Equalized Odds

From Hardt et al. (2016)

$$\begin{cases} \text{True positive equality : } & \mathbb{P}[\hat{Y} = 1 | S = 0, Y = 1] = \mathbb{P}[\hat{Y} = 1 | S = 1, Y = 1] \\ \text{False positive equality : } & \mathbb{P}[\hat{Y} = 1 | S = 0, Y = 0] = \mathbb{P}[\hat{Y} = 1 | S = 1, Y = 0] \\ \text{False negative equality : } & \mathbb{P}[\hat{Y} = 0 | S = 0, Y = 1] = \mathbb{P}[\hat{Y} = 0 | S = 1, Y = 1] \end{cases}$$

but one can consider equality of ratios...

Equal treatment

Berk et al. (2021a) We have equality of treatment, the rate of false positives and false negatives are identical in the protected groups,

$$\frac{\mathbb{P}[\hat{Y} = 1 | S = 0, Y = 0]}{\mathbb{P}[\hat{Y} = 0 | S = 0, Y = 1]} = \frac{\mathbb{P}[\hat{Y} = 1 | S = 1, Y = 0]}{\mathbb{P}[\hat{Y} = 0 | S = 1, Y = 1]}$$

Calibration

Instead of \hat{Y} focus on $m(\mathbf{X})$ (if possible)

Class balance

Kleinberg et al. (2016) We will have class balance in the weak sense if

$$\mathbb{E}[m(\mathbf{X})|Y = y, S = 0] = \mathbb{E}[m(\mathbf{X})|Y = y, S = 1], \forall y \in \{0, 1\}$$

or in the strong sense if

$$\mathbb{P}[m(\mathbf{X}) \leq \mu | Y = y, S = 0] = \mathbb{P}[m(\mathbf{X}) \leq \mu | Y = y, S = 1], \forall \mu \in [0, 1], \forall y \in \{0, 1\}.$$

Calibration (or accuracy) parity

Kleinberg et al. (2016), Zafar et al. (2017)] We have calibration parity if

$$\mathbb{P}[Y = 1 | m(\mathbf{X}) = \mu, S = 0] = \mathbb{P}[Y = 1 | m(\mathbf{X}) = \mu, S = 1], \forall \mu \in [0, 1].$$

Calibration

We can go further by asking not only for parity, but also for a good calibration

Good calibration

Kleinberg et al. (2017) We have an fairness of good calibration if

$$\mathbb{P}[Y = 1|m(\mathbf{X}) = \mu, S = 0] = \mathbb{P}[Y = 1|m(\mathbf{X}) = \mu, S = 1] = \mu, \quad \forall \mu \in [0, 1].$$

Nice property, but usually never satisfied by ML models (without fairness issue),

$$\mathbb{E}[Y|m(\mathbf{X}) = \mu] \neq \mu, \quad \forall \mu$$

This “good calibration” property of the model m , also called “well-calibration” in Dawid (1982), and “autocalibration” in Van Calster et al. (2019), Krüger and Ziegel (2021) and Denuit et al. (2021) in the context of regression, i.e. $\mathbb{E}[Y|m(\mathbf{X}) = \mu] = \mu$, is a standard property in econometrics, in generalized linear models, but not in most machine learning algorithms.

freakonometrics

freakonometrics.hypotheses.org Arthur Charpentier, 2023

23 / 112

Non-reconstruction of the protected attribute

Kim (2017) If we cannot tell from the result (\mathbf{x} , $m(\mathbf{x})$, y and \hat{y}) whether the subject was a member of a protected group or not, we will talk about fairness by non-reconstruction of the protected attribute

$$\mathbb{P}[S = 0 | \mathbf{X}, m(\mathbf{X}), \hat{Y}, Y] = \mathbb{P}[S = 1 | \mathbf{X}, m(\mathbf{X}), \hat{Y}, Y].$$

Used in Zhang et al. (2018) to suggest a adversarial approach to mitigate discrimination

Wrap-Up on Group Fairness

<i>statistical parity</i> , Dwork et al. (2012)	$\mathbb{P}[\hat{Y} = 1 S = s] = \text{cst}, \forall s$	independence $\hat{Y} \perp\!\!\!\perp P$
<i>conditional stat. parity</i> , Corbett-Davies et al. (2017)	$\mathbb{P}[\hat{Y} = 1 S = s, X = x] = \text{cst}_x, \forall s, y$	
<i>equalized odds</i> , Hardt et al. (2016)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = y] = \text{cst}_y, \forall s, y$	separation
<i>equalized opportunity</i> , Hardt et al. (2016)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = 1] = \text{cst}, \forall s$	
<i>predictive equality</i> , Corbett-Davies et al. (2017)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = 0] = \text{cst}, \forall s$	$\hat{Y} \perp\!\!\!\perp S Y$
<i>balance (positive)</i> , Kleinberg et al. (2017)	$\mathbb{E}[m(\mathbf{X}) S = s, Y = 1] = \text{cst}, \forall s$	$S \perp\!\!\!\perp P Y$
<i>balance (negative)</i> , Kleinberg et al. (2017)	$\mathbb{E}[m(\mathbf{X}) S = s, Y = 0] = \text{cst}, \forall s$	
<i>conditional accuracy equality</i> , Berk et al. (2017)	$\mathbb{P}[Y = y S = s, \hat{Y} = y] = \text{cst}_y, \forall s, y$	sufficiency
<i>predictive parity</i> , Chouldechova (2017)	$\mathbb{P}[Y = 1 S = s, \hat{Y} = 1] = \text{cst}, \forall s$	
<i>calibration</i> , Chouldechova (2017)	$\mathbb{P}[Y = 1 S = s, m(\mathbf{X}) = m] = \text{cst}_s, \forall s, m$	$Y \perp\!\!\!\perp S \hat{Y}$
<i>well-calibration</i> , Chouldechova (2017)	$\mathbb{P}[Y = 1 S = s, m(\mathbf{X}) = m] = s, \forall s, m$	
<i>accuracy equality</i> , Berk et al. (2017)	$\mathbb{P}[\hat{Y} = Y m(\mathbf{X}) = m] = \text{cst}, \forall m$	
<i>treatment equality</i> , Berk et al. (2017)	$\frac{\text{FN}_s}{\text{FP}_s} = \text{cst}_s, \forall s$	

Individual fairness

"Individual fairness is embodied in the following principle: similar individuals should be given similar decisions. This principle deals with the comparison of single individuals rather than focusing on groups of people sharing some characteristics. On the other hand, group fairness starts from the idea that there are groups of people potentially suffering biases and unfair decisions, and thus tries to reach equality of treatment for groups instead of individuals", Castelnovo et al. (2022)

Lipschitz property

Duivesteijn and Feelders (2008), Luong et al. (2011) A decision function \hat{Y} satisfies the Lipschitz property if

$$d_y(\hat{y}_i, \hat{y}_j) \leq d_x(\mathbf{x}_i, \mathbf{x}_j), \quad \forall i, j = 1, \dots, n.$$

Two "close" individuals (in the sense of unprotected characteristics \mathbf{x}) must have the same forecast

Counterfactual fairness

"what would have been the decision if that individual had a different gender?"

Sensitive s	Outcome			Age	School	Height	Weight	
	s_i	y_i	$y_{i,S \leftarrow 1}^*$	$y_{i,S \leftarrow 0}^*$	$x_{1,i}$	$x_{2,i}$	$x_{3,i}$	$x_{4,i}$
1	0 – M	1	1	?	37	14	160	56
2	1 – F	1	?	1	28	12	156	54
3	0 – M	0	0	?	53	11	190	87

Counterfactual fairness

Kusner et al. (2017) If the prediction in the real world is the same as the prediction in the counterfactual world where the individual would have belonged to a different demographic group, we have counterfactual fairness, i.e.

$$\mathbb{P}[Y_{S \leftarrow 1}^* = 1 | \mathbf{X} = \mathbf{x}] - \mathbb{P}[Y_{S \leftarrow 0}^* = 1 | \mathbf{X} = \mathbf{x}] = 0, \quad \forall \mathbf{x}.$$

Counterfactual fairness

"what would have been the outcome if that individual had a different treatment?"

Treatment	Outcome			Age	School	Height	Weight	
	t_i	y_i	$y_{i,T \leftarrow 1}^*$	$y_{i,T \leftarrow 0}^*$	$x_{1,i}$	$x_{2,i}$	$x_{3,i}$	$x_{4,i}$
1	1	121	121	?	37	14	160	56
2	0	109	?	109	28	12	156	54
3	1	162	162	?	53	11	190	87

Counterfactual fairness

Kusner et al. (2017) If the prediction in the real world is the same as the prediction in the counterfactual world where the individual would have belonged to a different demographic group, we have counterfactual fairness, i.e.

$$\mathbb{P}[Y_{S \leftarrow 1}^* = 1 | \mathbf{X} = \mathbf{x}] - \mathbb{P}[Y_{S \leftarrow 0}^* = 1 | \mathbf{X} = \mathbf{x}] = 0, \quad \forall \mathbf{x}.$$

Counterfactual fairness

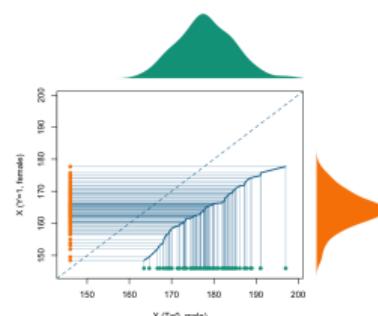
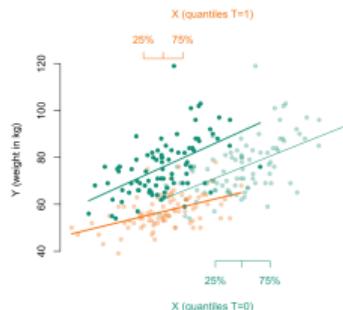
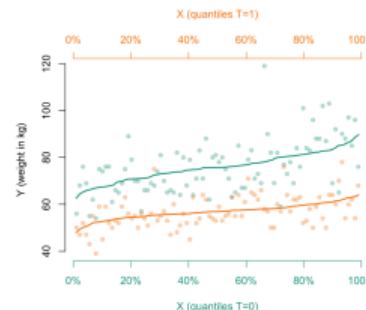
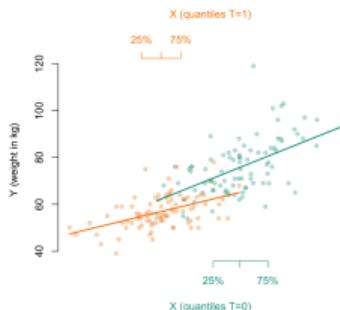
Related to the concept of **conditional average treatment**,

CATE

Hahn (1998), Heckman et al. (1998) Conditional ATE is CATE(x)

$$\mathbb{E}[Y_{T \leftarrow 1}^* - Y_{T \leftarrow 0}^* | X = x]$$

Unfortunately, standard estimates are *ceteris paribus*, see Charpentier et al. (2023) at ECONVN2023 for a *mutandis mutatis* estimate, using optimal transport...



Preamble: Ceteris Paribus & Mutatis Mutandis

Ceteris paribus sic stantibus

Ceteris paribus is a Latin phrase, meaning "*all other things being equal*" or "*other things held constant*".

Mutatis mutandis

Mutatis mutandis is a Latin phrase meaning "*with things changed that should be changed*" or "*once the necessary changes have been made*".

Consider a linear model $\hat{y} = m(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$, where $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$.

What happens if $x_1 \rightarrow x_1 + dx_1$?

- ▶ Ceteris paribus: $\hat{y} \rightarrow \hat{y} + \beta_1 dx_1$
- ▶ Mutatis mutandis: $\hat{y} \rightarrow \hat{y} + \beta_1 dx_1 + \beta_2 \frac{r\sigma_2}{\sigma_1} dx_1$

Causal inference

► Angrist and Pischke (2009, 2014)

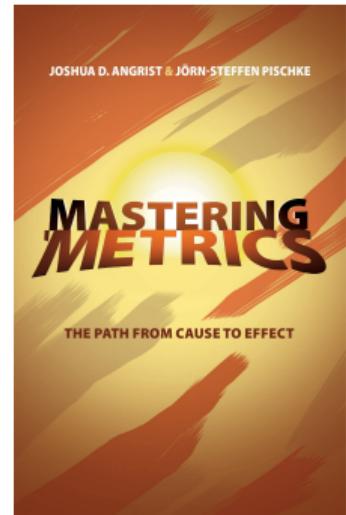
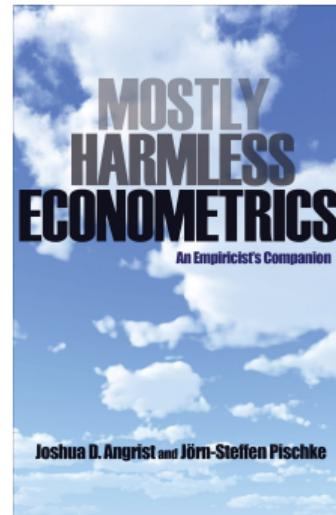
A naive comparison of averages by hospitalization status tells us something about potential outcomes, though not necessarily what we want to know. The comparison of average health conditional on hospitalization status is formally linked to the average causal effect by the equation below:

$$\underbrace{E[Y_i|D_i = 1] - E[Y_i|D_i = 0]}_{\text{Observed difference in average health}} = \underbrace{E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 1]}_{\text{average treatment effect on the treated}} + \underbrace{E[Y_{0i}|D_i = 1] - E[Y_{0i}|D_i = 0]}_{\text{selection bias}}$$

The term

$$E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 1] = E[Y_{1i} - Y_{0i}|D_i = 1]$$

is the *average causal effect of hospitalization on those who were hospitalized*. This term captures the average difference between the health of the hospitalized, $E[Y_{1i}|D_i = 1]$, and what would have happened to *them* had they not been hospitalized, $E[Y_{0i}|D_i = 1]$. The observed difference in health status however, adds to this causal effect a term called *selection bias*. This term is the difference in average Y_{0i} between those who



Causal inference

“Ladder of causation” from Pearl (2009)

3. Counterfactuals

(Imagining, “*what if I had done...*”)

2. Intervention

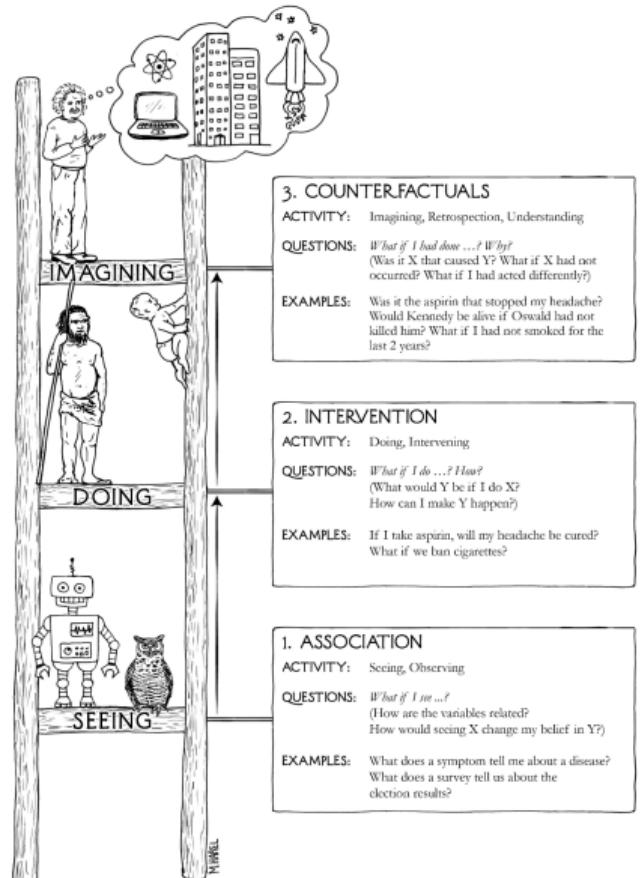
(Doing, “*what if I do...*”)

1. Association

(Seeing, “*what if I see...*”)

Picture source: Pearl and Mackenzie (2018)

What would be the impact of a treatment T on a variable of interest Y ?



Causal inference

Gender	Treatment	Outcome (Weight)			Height	...
		t_i	y_i	$y_{i,T \leftarrow 0}^*$		
1	H	0	75	75	?	172
2	F	1	52	?	52	161
3	F	1	57	?	57	163
4	H	0	78	78	?	183

(different notations are used $y(1)$ and $y(0)$ in Imbens and Rubin (2015), y^1 and y^0 in Cunningham (2021), or $y_{t=1}$ and $y_{t=0}$ in Pearl and Mackenzie (2018))

ATE & SATE

$$\text{ATE} = \mathbb{E}[Y_{T \leftarrow 1}^* - Y_{T \leftarrow 0}^*] \text{ and } \text{SATE} = \frac{1}{n} \sum_{i=1}^n y_{i,T \leftarrow 1}^* - y_{i,T \leftarrow 0}^*$$

Causal inference

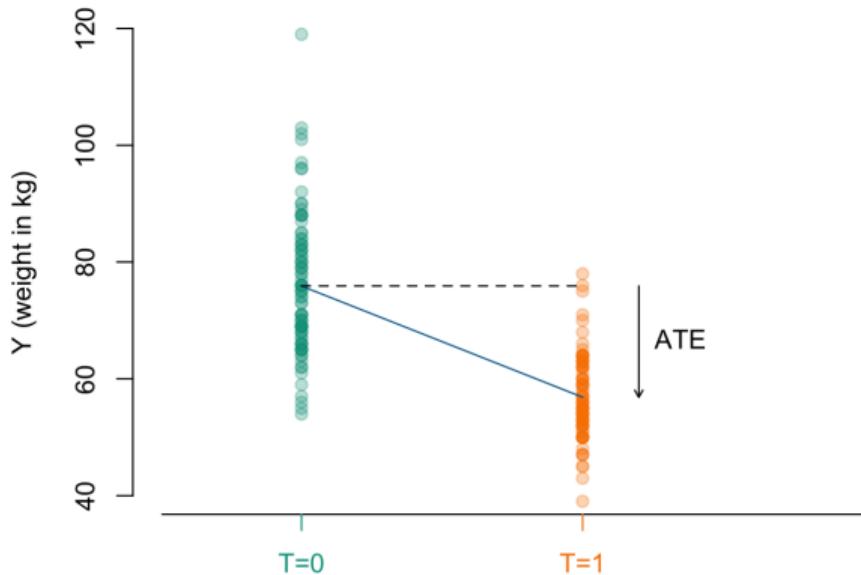
Let $y_i \in \mathcal{D}_0$ and $y_j \in \mathcal{D}_1$

"what would have been the weight of that person if that person had been a woman, and not a man? "

(too) simple sample estimate

$$\text{SATE} = \bar{y}_1 - \bar{y}_0$$

$$\text{SATE} = \frac{1}{n_1} \sum_{j \in \mathcal{D}_1} y_j - \frac{1}{n_0} \sum_{i \in \mathcal{D}_0} y_i$$



Causal inference

Consider a third variable x

CATE

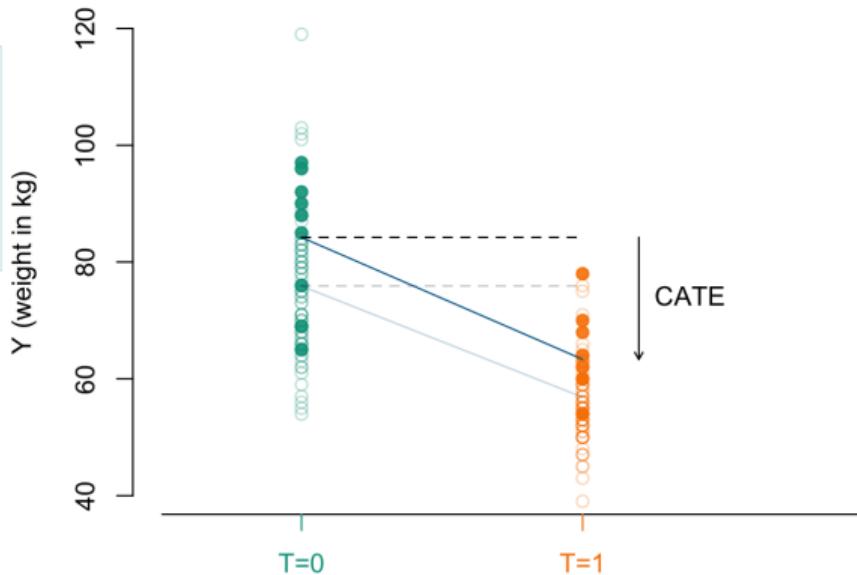
Conditional ATE is CATE(x)

$$\mathbb{E}[Y_{T \leftarrow 1}^* - Y_{T \leftarrow 0}^* | X = x]$$

A natural estimate would be

$$\text{SCATE} = \frac{1}{k} \sum_{x_j \in \mathcal{V}_{x,1}^k} y_j - \frac{1}{k} \sum_{x_i \in \mathcal{V}_{x,0}^k} y_i$$

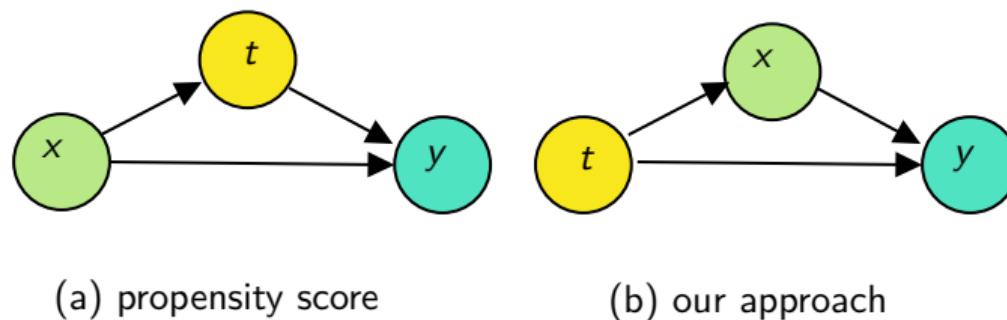
as in [Abrevaya et al. \(2015\)](#).



Causal inference

This $\text{CATE}(x)$ is a **ceteris paribus CATE** that does not take into account possible correlation between x , y and t .

Classical approach is based on the use of the propensity score, implying that x might influence the treatment t



Mutatis mutandis CATE

Given x "in the reference group" (0)

$$\text{CATE}(x) = \mathbb{E}[Y_{T \leftarrow 1}^* | x_{T \leftarrow 1}] - \mathbb{E}[Y_{T \leftarrow 0}^* | x]$$

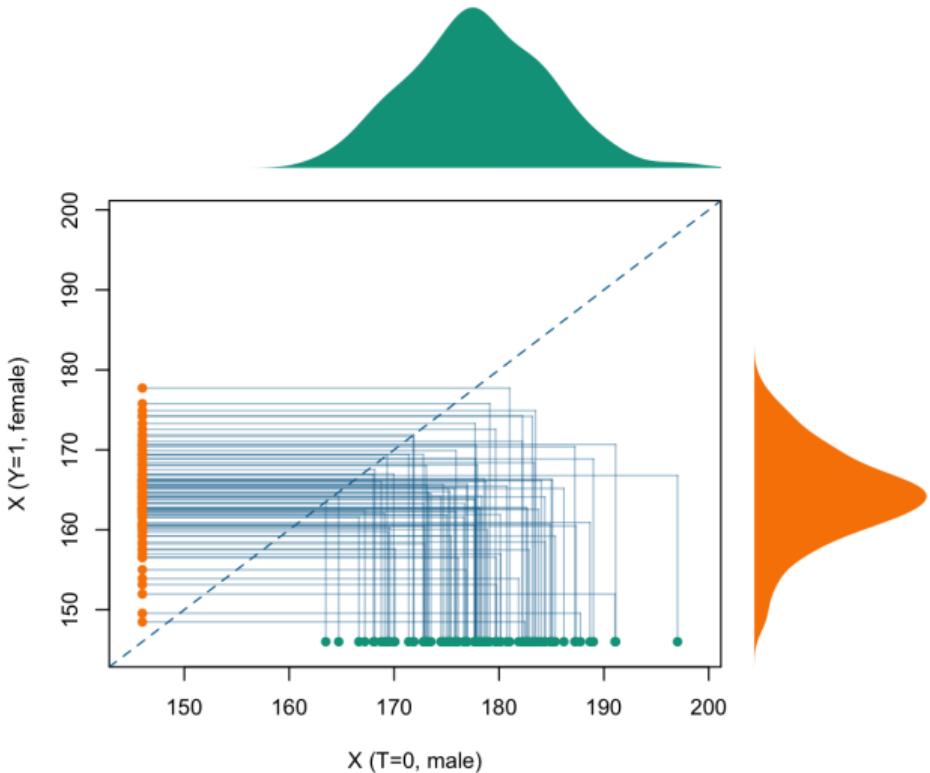
Matching

Between $x_i \in \mathcal{D}_0$ and $x_j \in \mathcal{D}_1$

$$j_i^* = \operatorname{argmin}_{j \in \mathcal{D}_1} \{d(x_i, x_j)\},$$

then remove observation from \mathcal{D}_1 .

Algorithm in Rubin (1973),
described in Stuart (2010) under
the name "1:1 nearest neighbor
matching", see Ho et al. (2007) or
Dehejia and Wahba (1999)
also called "Greedy Matching"



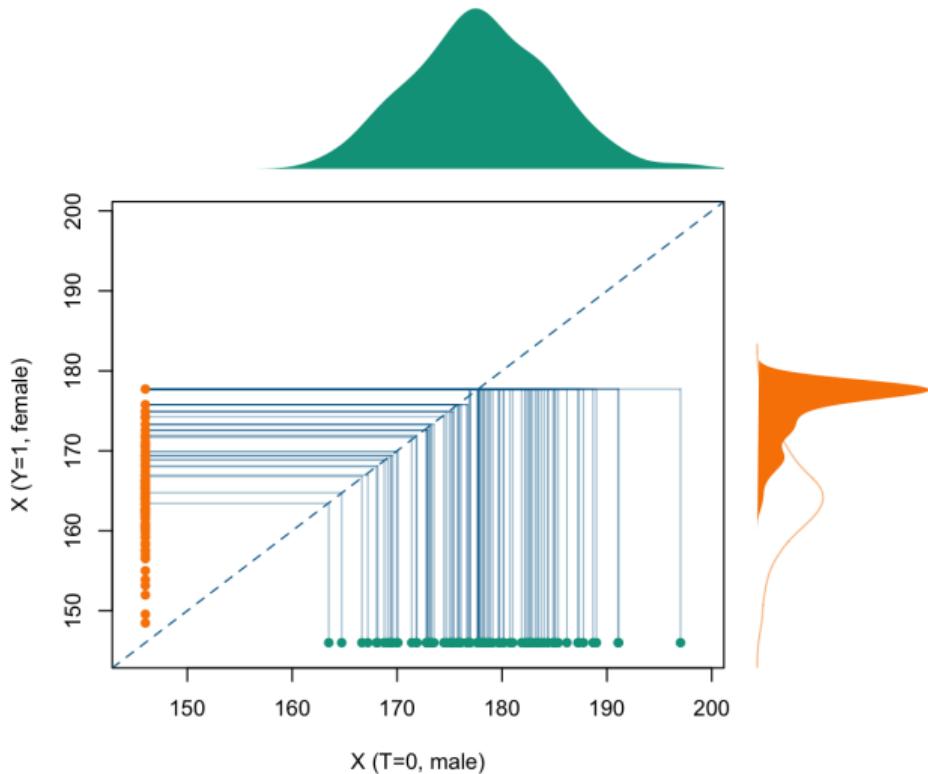
Matching

Between $x_i \in \mathcal{D}_0$ and $x_j \in \mathcal{D}_1$

$$j_i^* = \operatorname{argmin}_{j \in \mathcal{D}_1} \{d(x_i, x_j)\},$$

and **keeping** observation from \mathcal{D}_1
will distort distribution of x ,

$$\{x_j, j \in \mathcal{D}_1\} \stackrel{\mathcal{L}}{\neq} \{x_{j_i^*}, i \in \mathcal{D}_0\}$$

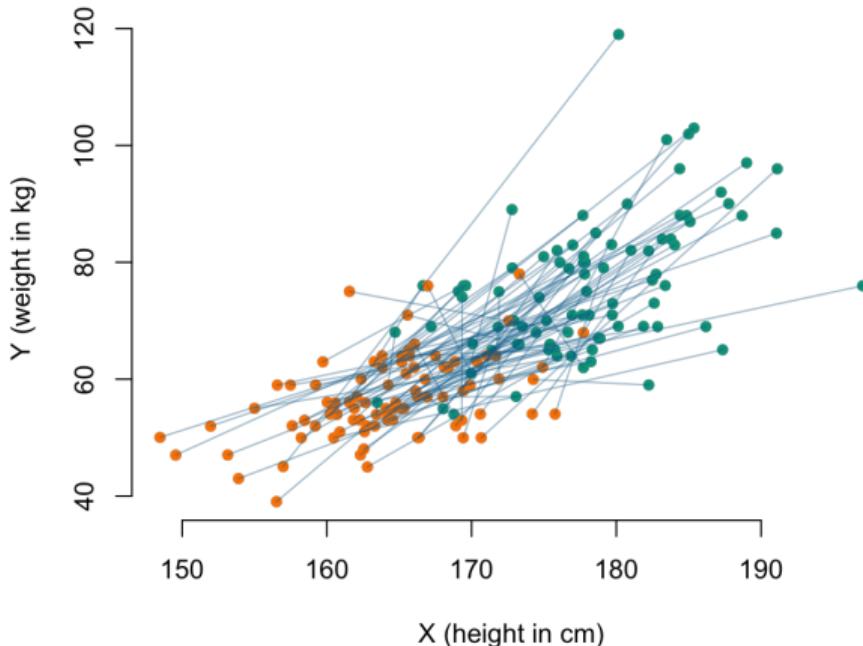


Matching

Done properly
(removing observation),
each individual (x_i, y_i)
in the control group (\mathcal{D}_0)
has counterfactual $(x_{j_i^*}, y_{j_i^*})$
in the treated group (\mathcal{D}_1)...

$$\text{SATE} = \frac{1}{n} \sum_{i=1}^n (y_{j_i^*} - y_i)$$

i.e. $\text{SATE} = \bar{y}_1 - \bar{y}_0$
(simply permute observations in
 \mathcal{D}_1)



Matching

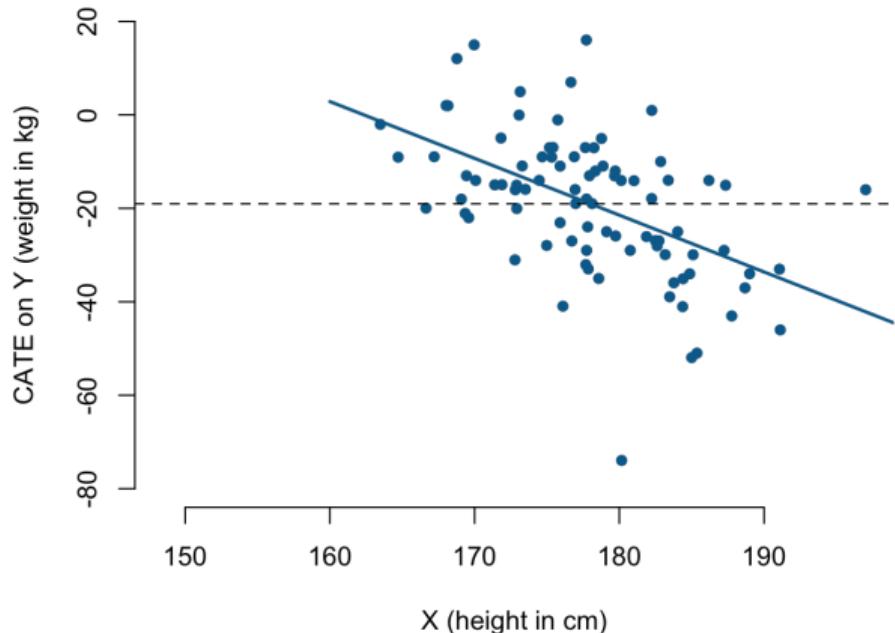
Let V_x^k denote the list of k nearest neighbors of x_i 's in \mathcal{D}_0 close to x ,

$$\text{SCATE}(x) = \frac{1}{k} \sum_{i \in V_x^k} (y_{j_i^*} - y_i)$$

Here scatter-plot \bullet of

$$\{(x_i, y_{j_i^*} - y_i)\}_{i=1, \dots, n}$$

and linear regression ——
Horizontal line --- is ATE



Matching

Algorithm 1 SATE, matching case (classical)

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n) shuffled, with indices i
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n), with indices j
 - 3: **for** $i = 1$ to n **do**
 - 4: $j_i^* = \underset{j:t_j=1}{\operatorname{argmin}}\{d(\mathbf{x}_i, \mathbf{x}_j)\}$ in \mathcal{D}_1 ,
 - 5: $d_i \leftarrow y_{j_i^*}^{(1)} - y_i^{(0)}$
 - 6: remove observation j_i^* from \mathcal{D}_1
 - 7: **end for**
 - 8: $\text{SATE} \leftarrow \frac{1}{n} \sum_{i=1}^n d_i$
-

freakonometrics

Matching

Algorithm 2 SCATE, matching case (classical)

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n) shuffled, with indices i
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n), with indices j
 - 3: $V_x^k \leftarrow$ list of k nearest neighbors of \mathbf{x}_i 's in \mathcal{D}_0 close to \mathbf{x}
 - 4: **for** $i = 1, 2, \dots, n$ **do**
 - 5: $j_i^* = \underset{j: t_j=1}{\operatorname{argmin}} \{d(\mathbf{x}_i, \mathbf{x}_j)\}$ in \mathcal{D}_1 ,
 - 6: $d_i \leftarrow y_{j_i^*}^{(1)} - y_i^{(0)}$
 - 7: remove observation j_i^* from \mathcal{D}_1
 - 8: **end for**
 - 9: $\text{SCATE}(\mathbf{x}) \leftarrow \frac{1}{k} \sum_{i \in V_x^k} d_i$
-

freakonometrics

freakonometrics.hypotheses.org Arthur Charpentier, 2023

42 / 112

Optimal Matching

C is the $n \times n$ matrix that quantifies the distance between individuals in the two groups, $C_{i,j} = d(\textcolor{teal}{x}_i, \textcolor{orange}{x}_j)^2 = (\textcolor{teal}{x}_i - \textcolor{orange}{x}_j)^2$, the optimal matching is solution of

$$\min_{P \in \mathcal{P}} \left\{ \langle P, C \rangle \right\} = \min_{P \in \mathcal{P}} \left\{ \sum_{i,j} P_{i,j} C_{i,j} \right\}, \quad (1)$$

where \mathcal{P} is the set of permutation matrices

$n \times n$ permutation matrix, P , with entries in $\{0, 1\}$, satisfying $P\mathbf{1}_n = \mathbf{1}_n$ and $P^{\star\top}\mathbf{1}_n = \mathbf{1}_n$, see [Bru Aldi \(2006\)](#).

Optimal Matching

Initial algorithm without the "no-replacement" rule (1:1), total cost 1.06

	7	8	9	10	11	12
1	0.41	0.55	0.22	0.64	0.04	0.25
2	0.28	0.24	0.73	0.22	0.64	0.80
3	0.28	0.47	0.32	0.52	0.16	0.37
4	0.28	0.62	0.81	0.25	0.64	0.85
5	0.41	0.37	0.89	0.25	0.81	0.97
6	0.66	0.76	0.21	0.89	0.22	0.14

	7	8	9	10	11	12	
1	1	.	$1 \rightarrow 11$
2	.	.	.	1	.	.	$2 \rightarrow 10$
3	1	.	$3 \rightarrow 11$
4	.	.	.	1	.	.	$4 \rightarrow 10$
5	.	.	.	1	.	.	$5 \rightarrow 10$
6	1	$6 \rightarrow 12$

Optimal Matching

Initial algorithm (not optimal), total cost 2.19

	7	8	9	10	11	12
1	0.41	0.55	0.22	0.64	0.04	0.25
2	0.28	0.24	0.73	0.22	0.64	0.80
3	0.28	0.47	0.32	0.52	0.16	0.37
4	0.28	0.62	0.81	0.25	0.64	0.85
5	0.41	0.37	0.89	0.25	0.81	0.97
6	0.66	0.76	0.21	0.89	0.22	0.14

	7	8	9	10	11	12	
1	1	.	.	.	1	.	$1 \leftrightarrow 11$
2	2	.	.	.	1	.	$2 \leftrightarrow 10$
3	3	1	$3 \leftrightarrow 7$
4	4	.	1	.	.	.	$4 \leftrightarrow 8$
5	5	.	.	1	.	.	$5 \leftrightarrow 9$
6	6	1	$6 \leftrightarrow 12$

Optimal Matching

Initial algorithm (not optimal), another initial shuffle, total cost 1.32

	7	8	9	10	11	12
1	0.41	0.55	0.22	0.64	0.04	0.25
2	0.28	0.24	0.73	0.22	0.64	0.80
3	0.28	0.47	0.32	0.52	0.16	0.37
4	0.28	0.62	0.81	0.25	0.64	0.85
5	0.41	0.37	0.89	0.25	0.81	0.97
6	0.66	0.76	0.21	0.89	0.22	0.14

	7	8	9	10	11	12	
5	1	.	.	1	.	.	$1 \leftrightarrow 9$
3	2	1	$2 \leftrightarrow 7$
4	3	1	$3 \leftrightarrow 11$
1	4	.	.	.	1	.	$4 \leftrightarrow 10$
2	5	.	1	.	.	.	$5 \leftrightarrow 8$
6	6	1	$6 \leftrightarrow 12$

Optimal Matching

Initial algorithm (optimal), total cost 1.27*

	7	8	9	10	11	12
1	0.41	0.55	0.22	0.64	0.04	0.25
2	0.28	0.24	0.73	0.22	0.64	0.80
3	0.28	0.47	0.32	0.52	0.16	0.37
4	0.28	0.62	0.81	0.25	0.64	0.85
5	0.41	0.37	0.89	0.25	0.81	0.97
6	0.66	0.76	0.21	0.89	0.22	0.14

	7	8	9	10	11	12
1	1	.
2	.	1
3	.	.	1	.	.	.
4	1
5	.	.	.	1	.	.
6	1

1 \leftrightarrow **11**

2 \leftrightarrow **8**

3 \leftrightarrow **9**

4 \leftrightarrow **7**

5 \leftrightarrow **10**

6 \leftrightarrow **12**

Optimal Matching

Algorithm 3 SATE, optimal matching case

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n) shuffled, with indices i
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n), with indices j
 - 3: $C \leftarrow$ matrix $n \times n$, $C_{i,j} = d(\mathbf{x}_i, \mathbf{x}_j)$ between points in \mathcal{D}_0 and \mathcal{D}_1
 - 4: $P^* \leftarrow$ solution of Problem (1)
 - 5: $SATE \leftarrow \frac{1}{n_0} \sum_{i=1}^n y_i^0 - P_i^{*\top} \mathbf{y}^1$
-

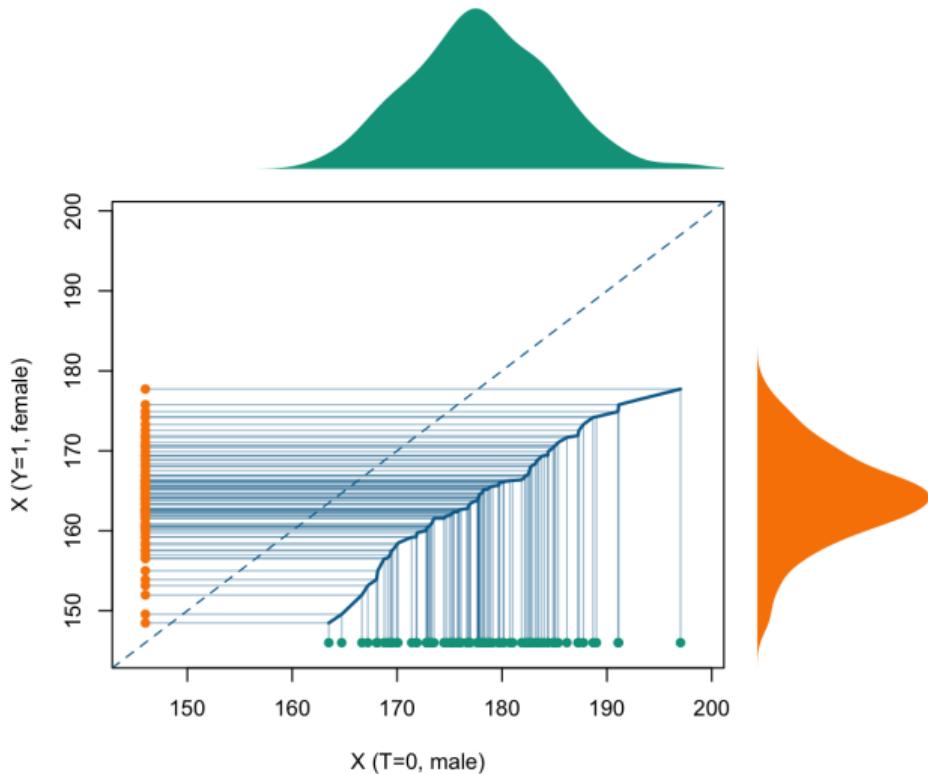
Actually Problem (1) has a simple explicit solution (based on ranks) and a nice geometric interpretation

Optimal Matching

$y_i \in \mathcal{D}_0$ and $y_j \in \mathcal{D}_1$

Let r_i denote the rank of $y_i \in \mathcal{D}_0$
and s_j denote the rank of $y_j \in \mathcal{D}_1$

Match $y_i \in \mathcal{D}_0$ with $y_{j_i^*} \in \mathcal{D}_1$ such
that $r_i = s_j$.

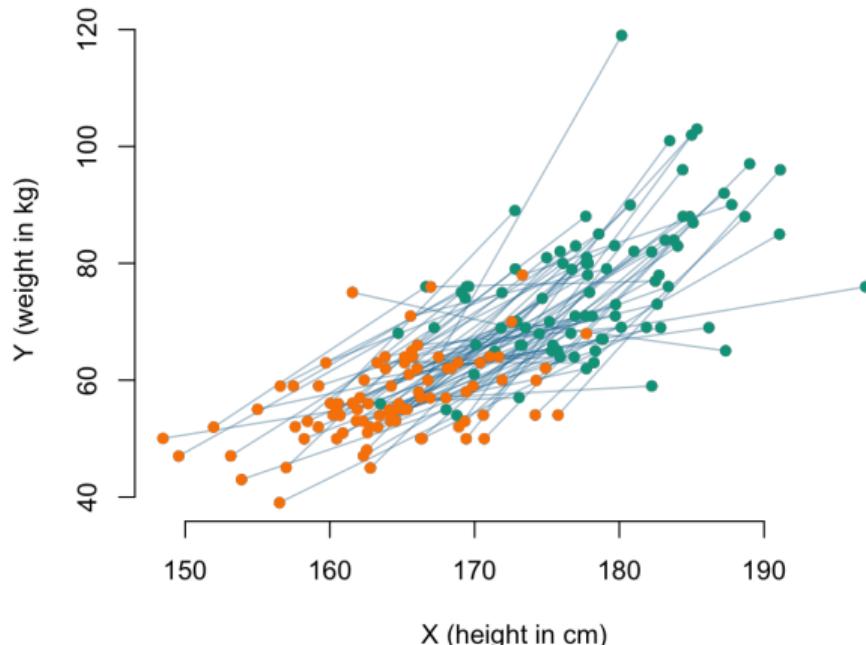


Optimal Matching

Each individual (x_i, y_i)
in the control group
has counterfactual $(x_{j_i^*}, y_{j_i^*})$
in the treated group...

$$\text{SATE} = \frac{1}{n} \sum_{i=1}^n (y_{j_i^*} - y_i)$$

i.e. $= \bar{y}_1 - \bar{y}_0$
(again, simple permutation)



Optimal Matching

Let V_x^k denote the list of k nearest neighbors of x_i 's in \mathcal{D}_0 close to x ,

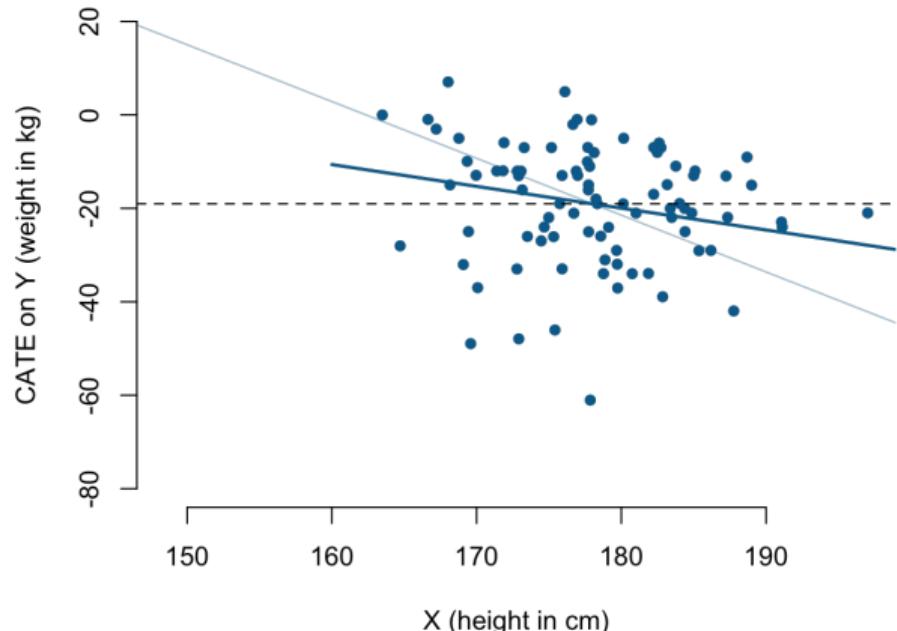
$$\text{SCATE}(x) = \frac{1}{k} \sum_{i \in V_x^k} (y_{j_i^*} - y_i)$$

Here scatter-plot • of

$$\{(x_i, y_{j_i^*} - y_i)\}_{i=1, \dots, n}$$

and linear regression —————

Horizontal line —— is ATE
(same as the previous one)



Optimal Coupling

$r_i^{(1)}$ denote the rank of $x_i^{(1)}$ in the treated dataset $\{x_1^{(1)}, \dots, x_n^{(1)}\}$. The procedure then becomes simply a matching based on ranks, in the sense that j_i^* satisfies $r_{j_i^*}^{(1)} = r_i^{(0)}$, as discussed in Chapter 2 of Santambrogio (2015).

In a very general setting, if $\mathbf{a}_0 \in \mathbb{R}_+^{n_0}$ and $\mathbf{a}_1 \in \mathbb{R}_+^{n_1}$ satisfy $\mathbf{a}_0^\top \mathbf{1}_{n_0} = \mathbf{a}_1^\top \mathbf{1}_{n_1}$ (identical sums), define

$$U(\mathbf{a}_0, \mathbf{a}_1) = \{M \in \mathbb{R}_+^{n_0 \times n_1} : M\mathbf{1}_{n_1} = \mathbf{a}_0 \text{ and } M^\top \mathbf{1}_{n_0} = \mathbf{a}_1\}.$$

This set of matrices is a convex polytope (see Brualdi (2006)).

In our case, let U_{n_0, n_1} denote $U\left(\mathbf{1}_0, \frac{n_0}{n_1}\mathbf{1}_1\right)$

$$P^* \in \underset{P \in U_{n_0, n_1}}{\operatorname{argmin}} \left\{ \langle P, C \rangle \right\} \text{ or } \underset{P \in U_{n_0, n_1}}{\operatorname{argmin}} \left\{ \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} P_{i,j} C_{i,j} \right\}. \quad (2)$$

Optimal Coupling

	7	8	9	10	11	12	13	14	15	16
1	0.41	0.55	0.22	0.64	0.04	0.25	0.24	0.77	0.74	0.55
2	0.28	0.24	0.73	0.22	0.64	0.80	0.76	0.76	0.12	0.10
3	0.28	0.47	0.32	0.52	0.16	0.37	0.27	0.68	0.63	0.45
4	0.28	0.62	0.81	0.25	0.64	0.85	0.58	0.32	0.51	0.48
5	0.41	0.37	0.89	0.25	0.81	0.97	0.91	0.81	0.05	0.25
6	0.66	0.76	0.21	0.89	0.22	0.14	0.33	0.96	0.99	0.79

	7	8	9	10	11	12	13	14	15	16
1	.	.	1/5	.	3/5	.	1/5	.	.	.
2	.	2/5	3/5
3	3/5	2/5	.	.	.
4	.	.	.	2/5	.	.	.	3/5	.	.
5	.	1/5	.	1/5	3/5	.
6	.	.	2/5	.	.	3/5

Optimal Coupling

Algorithm 4 SATE, optimal coupling case

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n_0) shuffled, with indices i
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n_1), with indices j
 - 3: $C \leftarrow$ matrix $n_0 \times n_1$, $C_{i,j} = d(\mathbf{x}_i, \mathbf{x}_j)$ between points in \mathcal{D}_0 and \mathcal{D}_1
 - 4: $P^* \leftarrow$ solution of Problem (4)
 - 5: $\text{SATE} \leftarrow \frac{1}{n_0} \sum_{i=1}^n y_i^0 - P_i^{*\top} \mathbf{y}^1$
-

Optimal Coupling

Algorithm 5 SCATE, optimal coupling case

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n_0) shuffled, with indices i
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n_1), with indices j
 - 3: $V_x^k \leftarrow$ list of k nearest neighbors of \mathbf{x}_i 's in \mathcal{D}_0 close to \mathbf{x}
 - 4: $C \leftarrow$ matrix $n_0 \times n_1$, $C_{i,j} = d(\mathbf{x}_i, \mathbf{x}_j)$ between points in \mathcal{D}_0 and \mathcal{D}_1
 - 5: $P^* \leftarrow$ solution of Problem (4)
 - 6: $\text{SCATE}(\mathbf{x}) \leftarrow \frac{1}{k} \sum_{i \in V_x^k} y_i^0 - P_i^{*\top} \mathbf{y}^1$
-

Quantile CATE

If $X_0 \sim F_0$, then $X_1 = \mathcal{T}(X_0) \sim F_1$, where $\mathcal{T} : x_0 \mapsto x_1 = F_1^{-1} \circ F_0(x_0)$.

Mutatis mutandis Quantile CATE

$$\text{QCATE}(u) = \mathbb{E}[Y_{T \leftarrow 1}^* | X = F_1^{-1}(u)] - \mathbb{E}[Y_{T \leftarrow 0}^* | X = F_0^{-1}(u)], \quad u \in (0, 1),$$

where F_t is the cumulative distribution function of X , conditional on $T = t$

Mutatis mutandis CATE

$$\mathbb{E}[Y_{T \leftarrow 1}^* | X = \mathcal{T}(x)] - \mathbb{E}[Y_{T \leftarrow 0}^* | X = x], \quad \mathcal{T} = F_1^{-1} \circ F_0$$

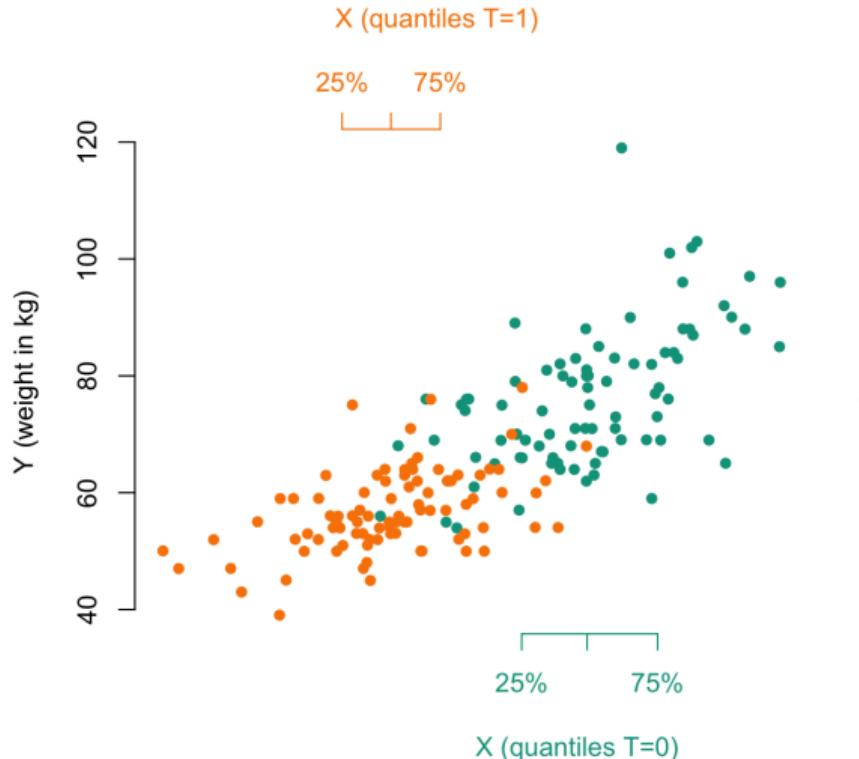
where x is considered with respect to the control group.

Quantile CATE

$(x_i, y_i) \in \mathcal{D}_0$ and $(x_j, y_j) \in \mathcal{D}_1$

Instead of x scale, visualize

$\begin{cases} \text{top : } & \text{probability, } x_i \in \mathcal{D}_0 \\ \text{bottom : } & \text{probability, } x_j \in \mathcal{D}_1 \end{cases}$

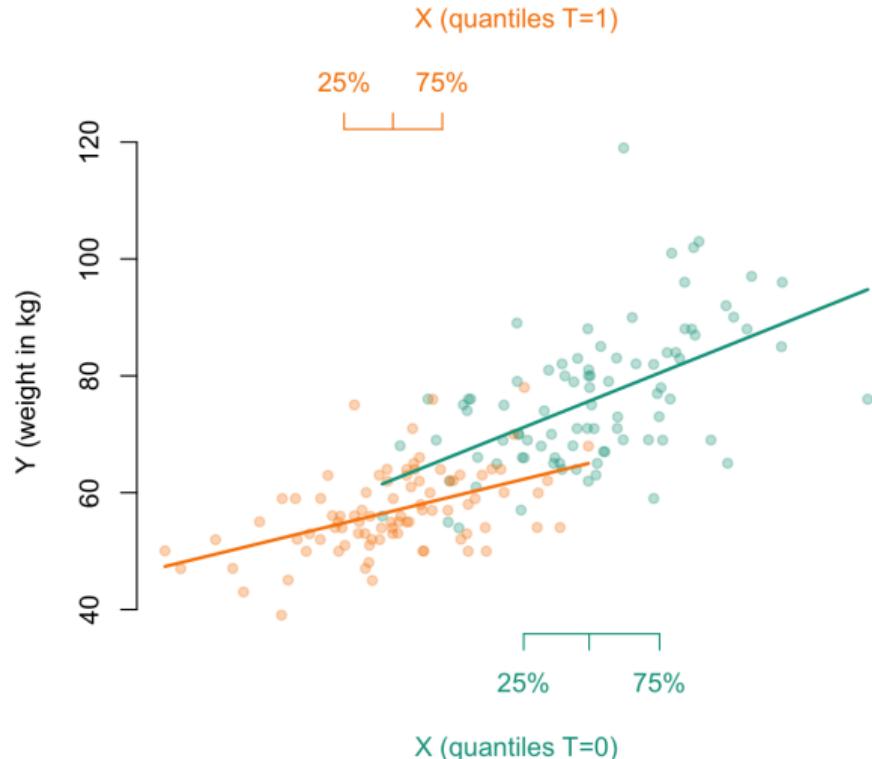


Quantile CATE

$$\begin{cases} \text{top : probability, } x_i \in \mathcal{D}_0 \\ \text{bottom : probability, } x_j \in \mathcal{D}_1 \end{cases}$$

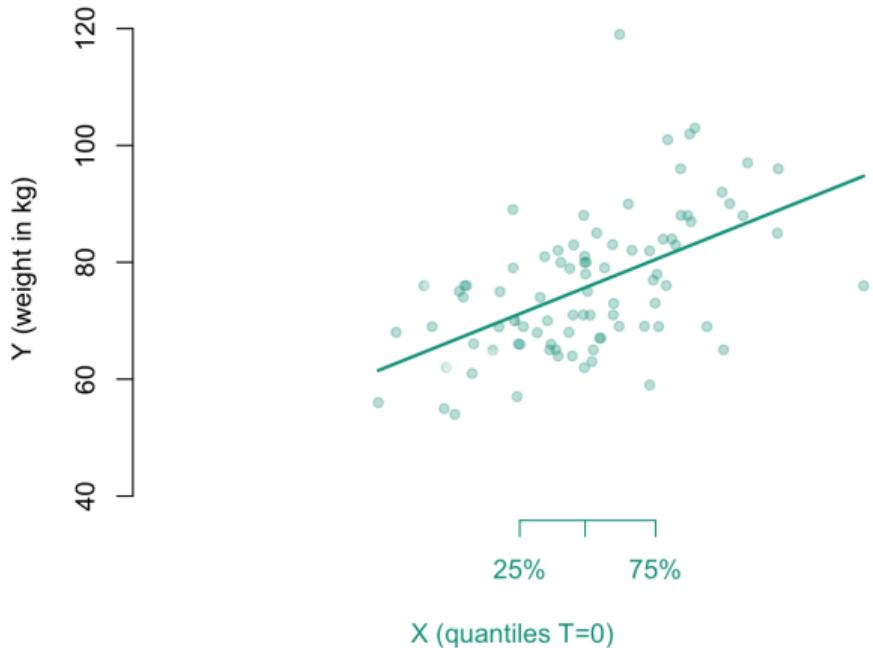
Fit two models m_0 and m_1

$$\begin{cases} m_0(x) = \mathbb{E}[Y|X = x, T = 0] \\ m_1(x) = \mathbb{E}[Y|X = x], T = 1 \end{cases}$$



Quantile CATE

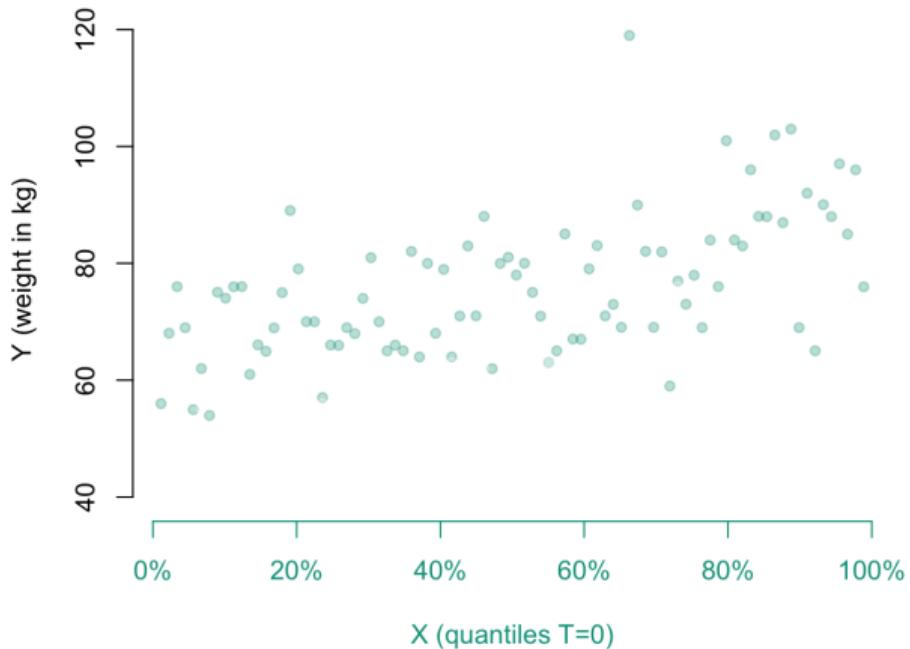
Instead of (x_i, y_i) (in \mathcal{D}_0)



Quantile CATE

Plot $(F_0(x_i), y_i)$ (in \mathcal{D}_0)

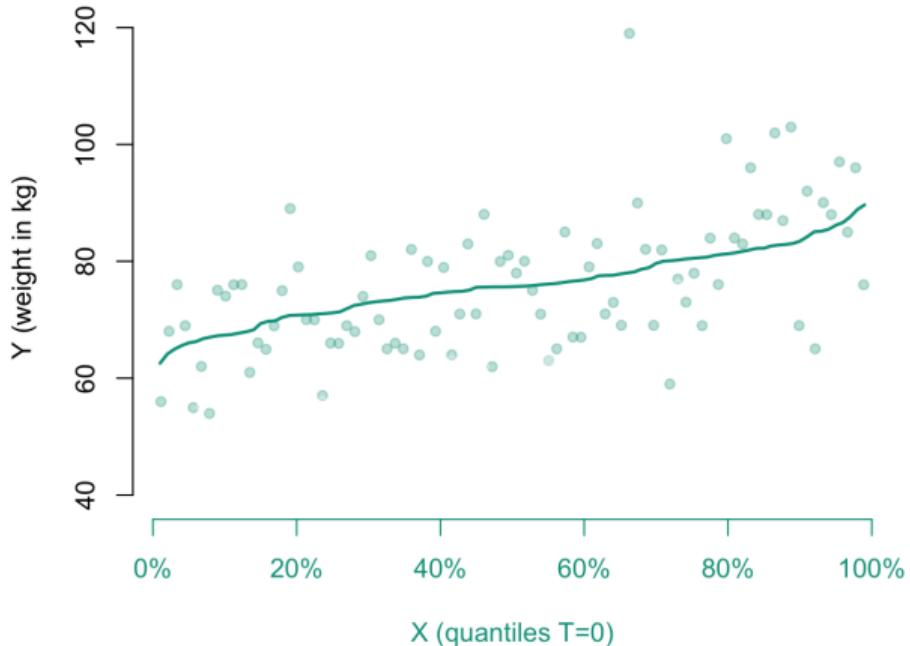
Note: $F_0(x_i) \propto r_i$



Quantile CATE

$$\mu_0(u) = \mathbb{E}[Y|X = F_0^{-1}(u), T = 0]$$

i.e. $\mu_0(u) = m_0(F_0^{-1}(u))$



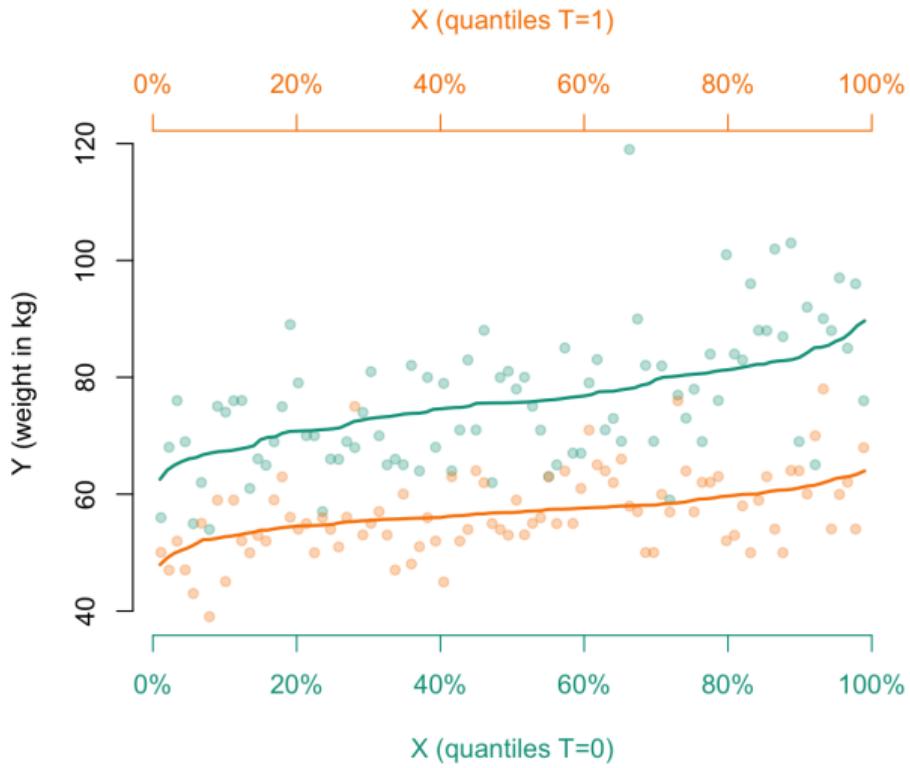
Quantile CATE

$$\mu_0(u) = \mathbb{E}[Y|X = F_0^{-1}(u), T = 0]$$

i.e. $\mu_0(u) = m_0(F_0^{-1}(u))$

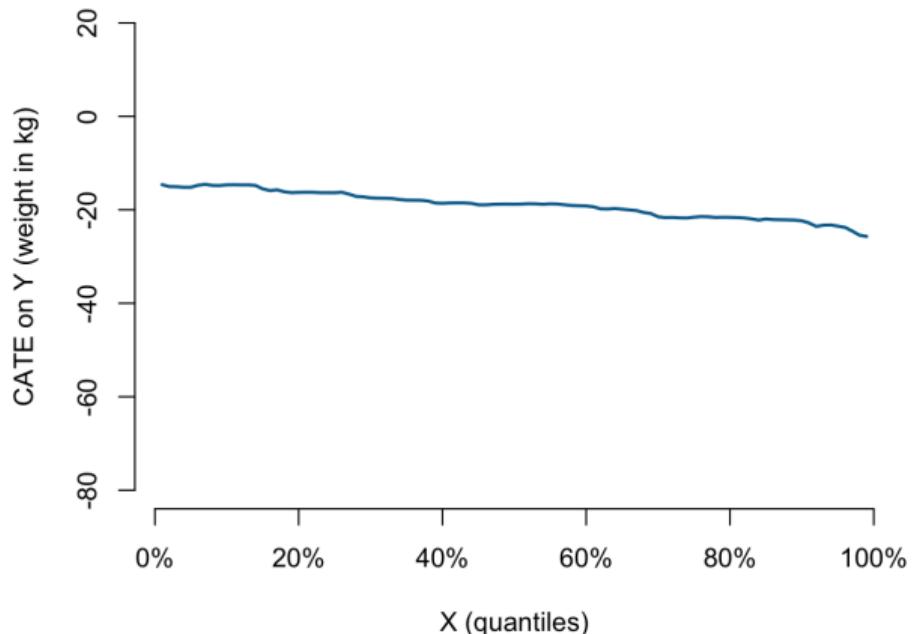
$$\mu_1(u) = \mathbb{E}[Y|X = F_1^{-1}(u), T = 0]$$

i.e. $\mu_1(u) = m_1(F_1^{-1}(u))$



Quantile CATE

Thus QCATE(u)
= $\mathbb{E}[Y_{T \leftarrow 1}^* | X = F_1^{-1}(u)]$
- $\mathbb{E}[Y_{T \leftarrow 0}^* | X = F_0^{-1}(u)]$
for $u \in (0, 1)$.



Quantile CATE

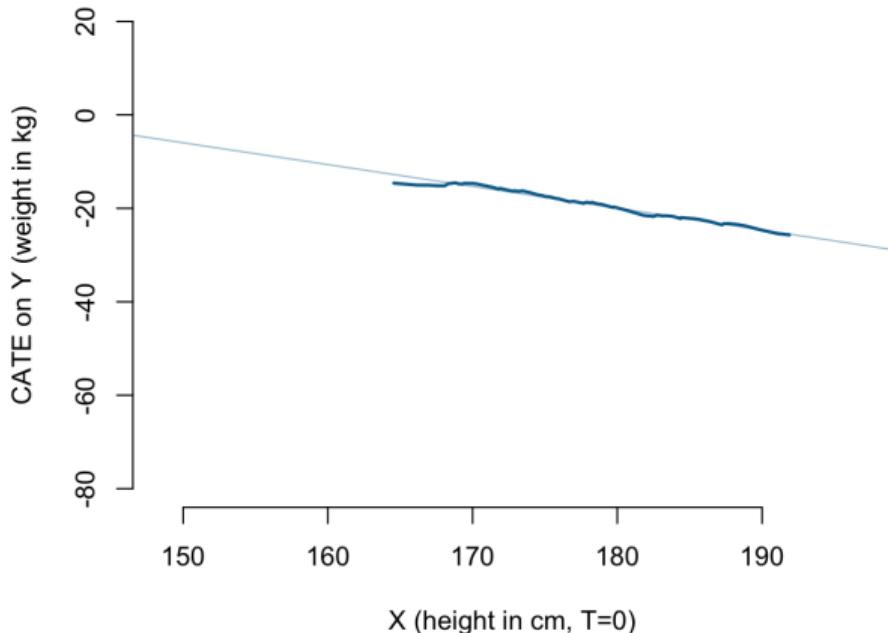
Mutatis mutandis CATE

$$= \mathbb{E}[Y_{T \leftarrow 1}^* | X = \mathcal{T}(x)]$$

$$- \mathbb{E}[Y_{T \leftarrow 0}^* | X = x],$$

$$\text{where } \mathcal{T} = F_1^{-1} \circ F_0$$

and x is considered with respect to
the control group.



This yields the following proper definition

Mutatis mutandis Sample CATE

Consider two models, $\hat{m}_0(x)$ and $\hat{m}_1(x)$, that estimate, respectively, $\mathbb{E}[Y|X = x, T = 0]$ and $\mathbb{E}[Y|X = x, T = 1]$,

$$\text{SCATE}(x) = \hat{m}_1(\hat{\mathcal{T}}(x)) - \hat{m}_0(x)$$

where $\hat{\mathcal{T}}(x) = \hat{F}_1^{-1} \circ \hat{F}_0(x)$, with \hat{F}_0 and \hat{F}_1 denoting the empirical distribution functions of x conditional on $t = 0$ and $t = 1$, respectively.

Gaussian x

Note that a simple parametric transformation can be obtained, based on the assumption that X conditional on T is Gaussian. More precisely, if

$$\textcolor{teal}{X}_0 \stackrel{\mathcal{L}}{=} X|t=0 \sim \mathcal{N}(\mu_0, \Sigma_0) \text{ and } \textcolor{orange}{X}_1 \stackrel{\mathcal{L}}{=} X|t=1 \sim \mathcal{N}(\mu_1, \Sigma_1)$$

$$\underbrace{\mu_1 + \sigma_1 \cdot \frac{\textcolor{teal}{X}_0 - \mu_0}{\sigma_0}}_{\mathcal{T}(\textcolor{teal}{X}_0)} \stackrel{\mathcal{L}}{=} \textcolor{orange}{X}_1$$

Quite naturally, we can use a Gaussian approximation for \mathcal{T} ,

Mutatis Mutandis Gaussian SCATE

Mutatis mutandis Gaussian CATE

Consider two models two models, $\hat{m}_0(x)$ and $\hat{m}_1(x)$, that estimate, respectively, $\mathbb{E}[Y|X = x, T = 0]$ and $\mathbb{E}[Y|X = x, T = 1]$,

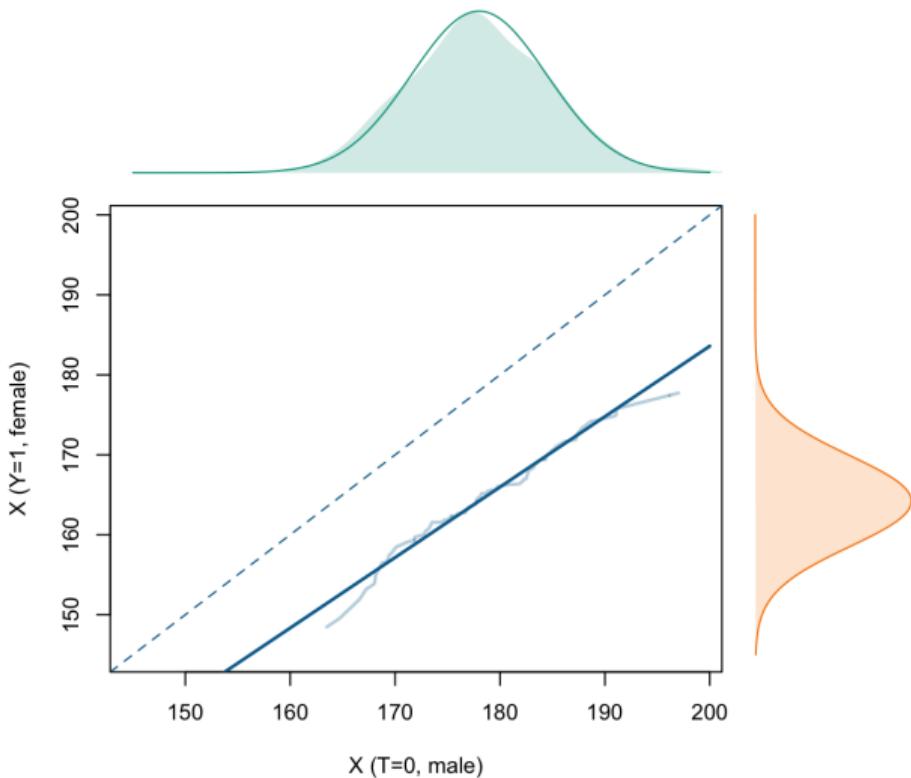
$$\text{SCATE}(x) = \hat{m}_1(\hat{T}_{\mathcal{N}}(x)) - \hat{m}_0(x)$$

where $\hat{T}_{\mathcal{N}}(x) = \bar{x}_1 + s_1 s_0^{-1}(x - \bar{x}_0)$, \bar{x}_0 and \bar{x}_1 being respectively the averages of x in the two sub-populations, and s_0 and s_1 the sample standard deviations.

Mutatis Mutandis Gaussian SCATE

As previously,
consider $x_i \in \mathcal{D}_0$ and $x_j \in \mathcal{D}_1$

Suppose $X|T=0 \sim \mathcal{N}(\mu_0, \sigma_0^2)$
and $X|T=1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$



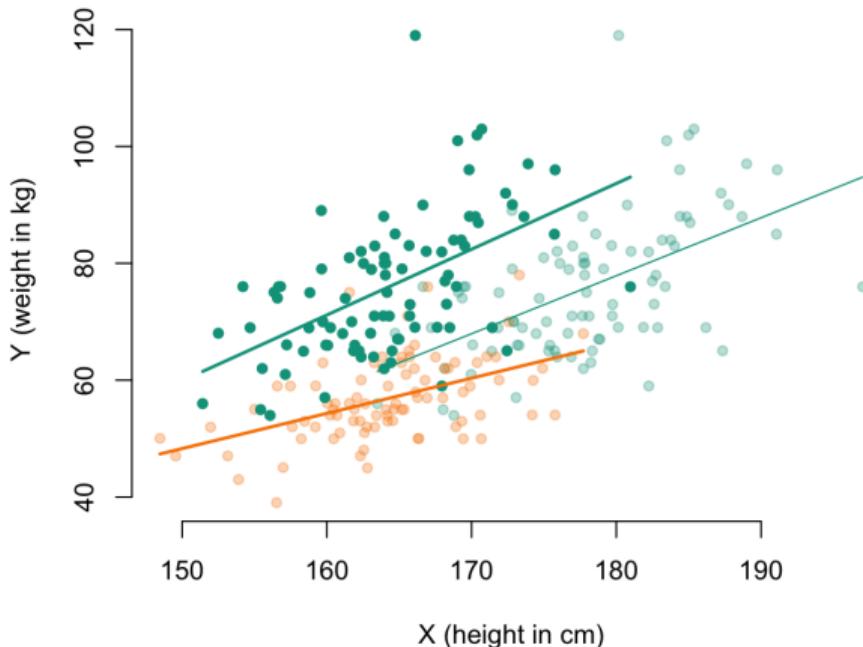
Mutatis Mutandis Gaussian SCATE

Previously, we had a matching,

$$i \in \mathcal{D}_0 \leftrightarrow j \in \mathcal{D}_1$$

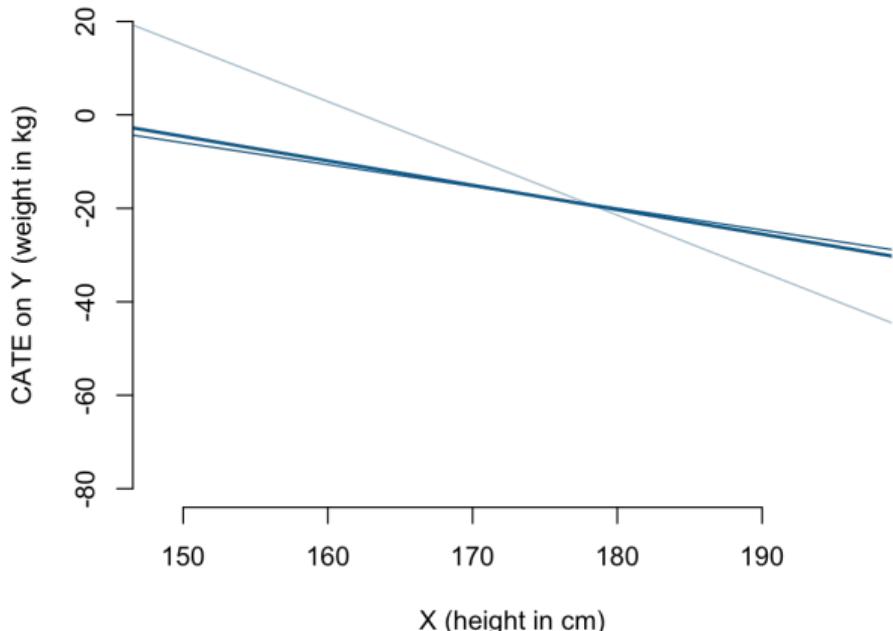
We have here an explicit mapping

$$\mathcal{T}$$



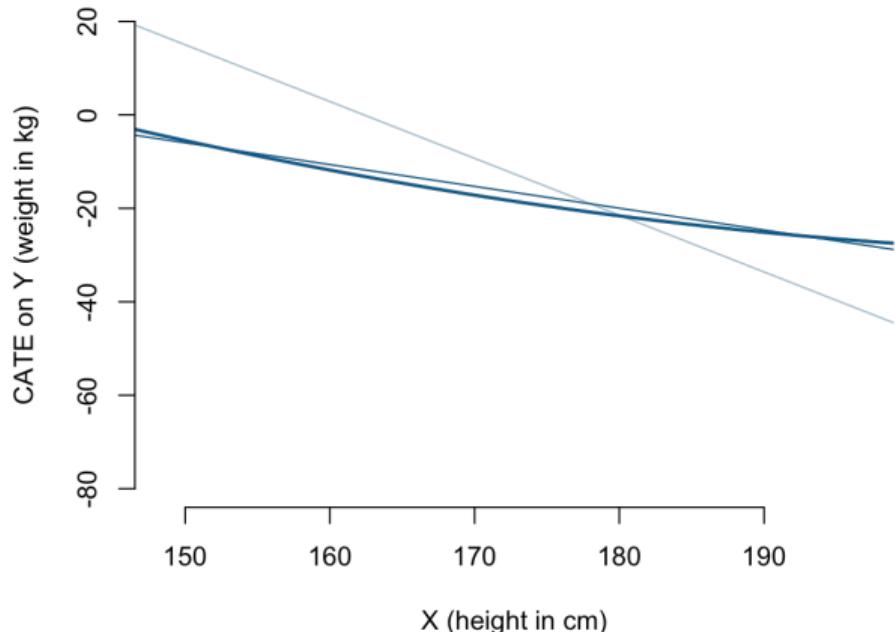
Mutatis Mutandis Gaussian SCATE

We can actually plot
 $x \mapsto \hat{m}_1(\hat{T}_N(x)) - \hat{m}_0(x)$
(and not only estimate it
from matched samples)
Here \hat{m}_0 and \hat{m}_1 are linear

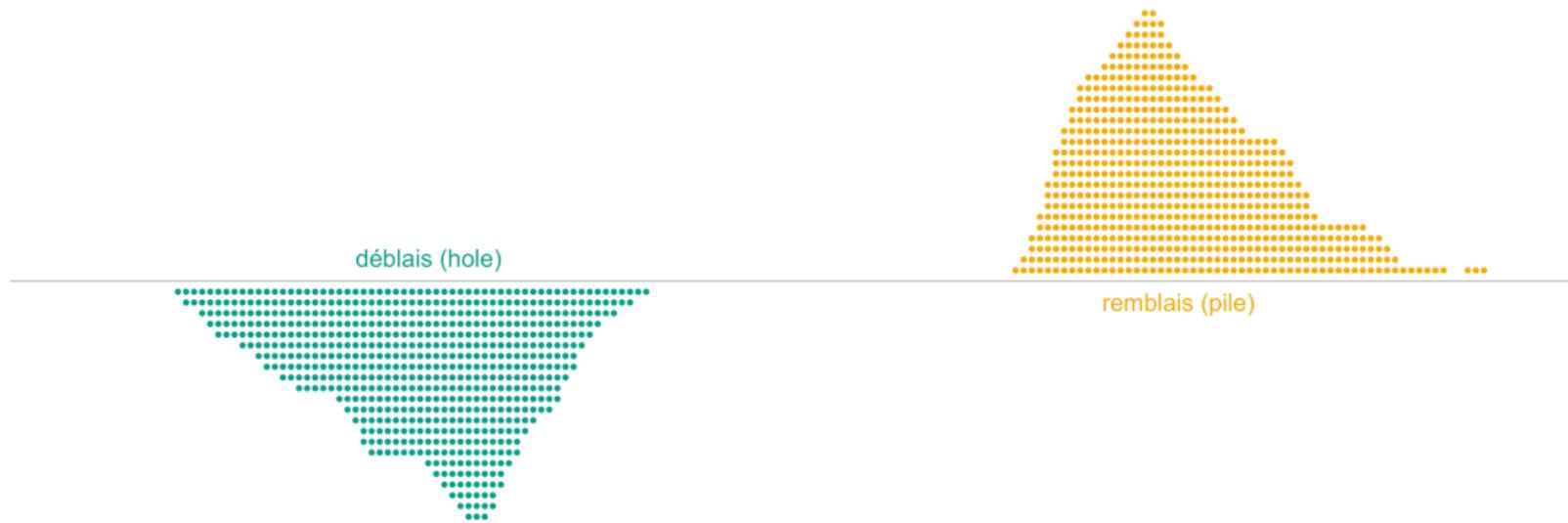


Mutatis Mutandis Gaussian SCATE

We can actually plot
 $x \mapsto \hat{m}_1(\hat{T}_N(x)) - \hat{m}_0(x)$
(and not only estimate it
from matched samples)
Here \hat{m}_0 and \hat{m}_1 are non-linear



Optimal Transport



Monge (1781), *Mémoire sur la Théorie des Déblais et des Remblais*

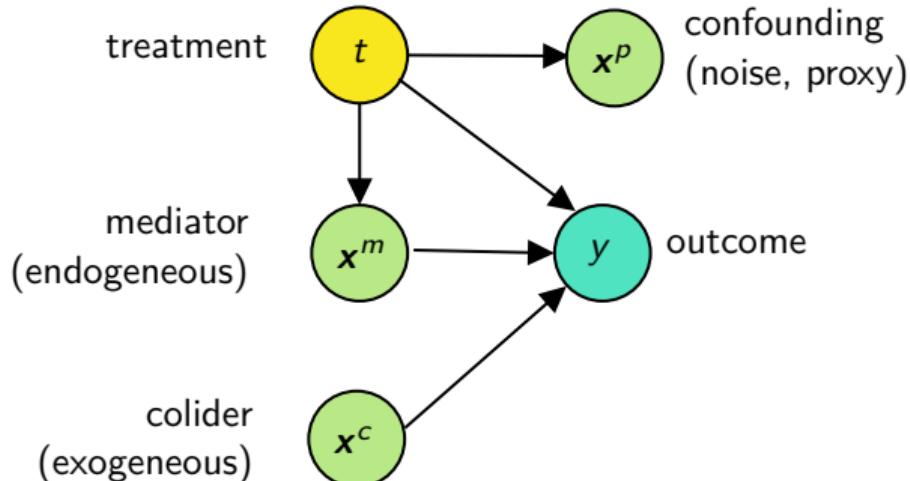
CATE(x)

We can consider a multivariate x

But we should not be transporting all components of x

Mutatis Mutandis CATE would be
 $m_1(\mathcal{T}(x^m), x^c, x^p) - m_0(x^m, x^c, x^p)$

- ▶ x^p is correlated with y
not causal, so not in $m(\cdot)$
- ▶ x^c is not influenced by t
- ▶ x^m is influenced by t
it should be transported



See [Charpentier et al. \(2023\)](#), detailed example with Gaussian Structural Causal Model.

CATE(\mathbf{x})

Given $\mathcal{T} : \mathbb{R}^k \rightarrow \mathbb{R}^k$, define the “*push-forward*” measure,

$$\mathbb{P}_1(A) = \mathcal{T}_\# \mathbb{P}_0(A) = \mathbb{P}_0(\mathcal{T}^{-1}(A)), \quad \forall A \subset \mathbb{R}^k.$$

An optimal transport \mathcal{T}^* (in Brenier's sense, from [Brenier \(1991\)](#), see [Villani \(2009\)](#) or [Galichon \(2016\)](#)) from \mathbb{P}_0 towards \mathbb{P}_1 will be solution of

$$\mathcal{T}^* \in \operatorname{arginf}_{\mathcal{T}: \mathcal{T}_\# \mathbb{P}_0 = \mathbb{P}_1} \left\{ \int_{\mathbb{R}^k} \gamma(\mathbf{x} - \mathcal{T}(\mathbf{x})) d\mathbb{P}_0(\mathbf{x}) \right\},$$

that we can write

$$\min_{\nu} \int_{\mathbb{R}^k \times \mathbb{R}^k} \underbrace{\gamma(\mathbf{x}, \mathbf{y})}_{=C} \underbrace{\nu(d\mathbf{x}, d\mathbf{y})}_{=P}, \text{ where } \nu \text{ is a coupling with margins } \mathbb{P}_0 \text{ and } \mathbb{P}_1,$$

for some cost function $\gamma : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}_+$.

Matching & Coupling

Previous problems were actually defined in any dimension. For matching, solve

$$\operatorname{argmin}_{P \in \mathcal{P}} \left\{ \langle P, C \rangle \right\} \text{ or } \operatorname{argmin}_{P \in \mathcal{P}} \left\{ \sum_{i,j} P_{i,j} C_{i,j} \right\}, \quad (3)$$

and, for the coupling case

$$P^* \in \operatorname{argmin}_{P \in U_{n_0, n_1}} \left\{ \langle P, C \rangle \right\} \text{ or } \operatorname{argmin}_{P \in U_{n_0, n_1}} \left\{ \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} P_{i,j} C_{i,j} \right\}. \quad (4)$$



Gaussian x

In the case where $\mathbf{X}|t=1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathbf{X}|t=0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$, there is an explicit expression for the optimal transport, which is simply an affine map (see [Villani \(2003\)](#) for more details). In the univariate case, $x_1 = \mathcal{T}_{\mathcal{N}}^*(x_0) = \mu_1 + \frac{\sigma_1}{\sigma_0}(x_0 - \mu_0)$, while in the multivariate case, an analogous expression can be derived:

$$x_1 = \mathcal{T}_{\mathcal{N}}^*(x_0) = \boldsymbol{\mu}_1 + \mathbf{A}(x_0 - \boldsymbol{\mu}_0),$$

where \mathbf{A} is a symmetric positive matrix that satisfies $\mathbf{A}\boldsymbol{\Sigma}_0\mathbf{A} = \boldsymbol{\Sigma}_1$, which has a unique solution given by $\mathbf{A} = \boldsymbol{\Sigma}_0^{-1/2}(\boldsymbol{\Sigma}_0^{1/2}\boldsymbol{\Sigma}_1\boldsymbol{\Sigma}_0^{1/2})^{1/2}\boldsymbol{\Sigma}_0^{-1/2}$, where $\mathbf{M}^{1/2}$ is the square root of the square (symmetric) positive matrix \mathbf{M} based on the Schur decomposition ($\mathbf{M}^{1/2}$ is a positive symmetric matrix), as described in [Higham \(2008\)](#).

Gaussian SCATE

Mutatis mutandis Gaussian CATE

Consider two models two models, $\hat{m}_0(\mathbf{x})$ and $\hat{m}_1(\mathbf{x})$, that estimate, respectively, $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}, T = 0]$ and $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}, T = 1]$,

$$\text{SCATE}(\mathbf{x}) = \hat{m}_1(\hat{T}_{\mathcal{N}}(\mathbf{x})) - \hat{m}_0(\mathbf{x})$$

where $\hat{T}_{\mathcal{N}}(\mathbf{x}) = \bar{\mathbf{x}}_1 + \hat{\mathbf{A}}(\mathbf{x} - \bar{\mathbf{x}}_0)$, with $\bar{\mathbf{x}}_0$ and $\bar{\mathbf{x}}_1$ being, respectively, the averages of \mathbf{x} in the two sub-populations, and $\hat{\mathbf{A}} = \hat{\Sigma}_0^{-1/2} (\hat{\Sigma}_0^{1/2} \hat{\Sigma}_1 \hat{\Sigma}_0^{1/2})^{1/2} \hat{\Sigma}_0^{-1/2}$ where $\hat{\Sigma}_0$ and $\hat{\Sigma}_1$ denote the sample variance.

Gaussian SCATE

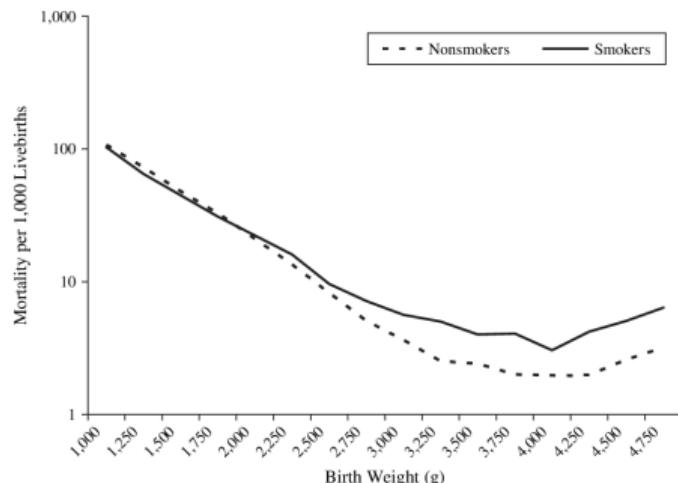
Algorithm 6 SCATE, Gaussian transport

Require: dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i, t_i)\}$

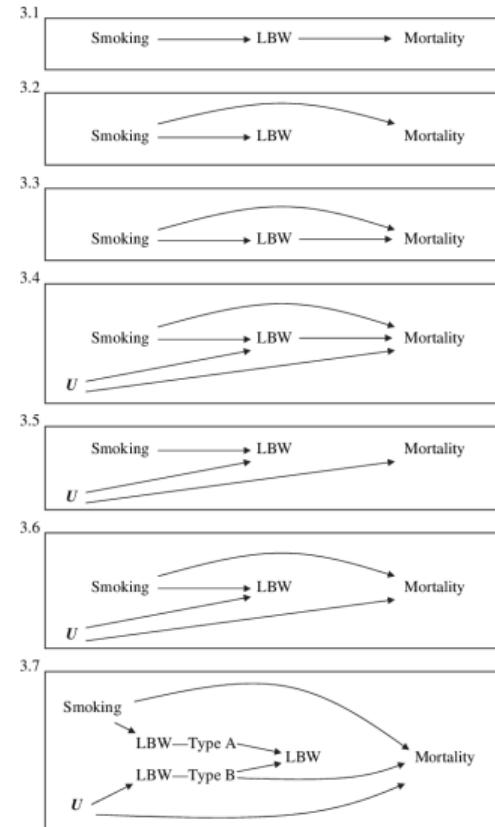
- 1: $\mathcal{D}_0 \leftarrow$ subset of \mathcal{D} when $t = 0$ (size n_0)
 - 2: $\mathcal{D}_1 \leftarrow$ subset of \mathcal{D} when $t = 1$ (size n_1)
 - 3: $\hat{m}_0 \leftarrow$ model to predict y based on \mathbf{x} , trained on \mathcal{D}_0
 - 4: $\hat{m}_1 \leftarrow$ model to predict y based on \mathbf{x} , trained on \mathcal{D}_1
 - 5: estimate moments of \mathbf{x}_t 's $\hat{\mu}_0$, $\hat{\mu}_1$, $\hat{\Sigma}_0$ and $\hat{\Sigma}_1$,
 - 6: $\hat{\mathbf{A}} \leftarrow \hat{\Sigma}_0^{-1/2} (\hat{\Sigma}_0^{1/2} \hat{\Sigma}_1 \hat{\Sigma}_0^{1/2})^{1/2} \hat{\Sigma}_0^{-1/2}$
 - 7: $\tilde{\mathbf{x}} \leftarrow \hat{\mu}_1 + \hat{\mathbf{A}}(\mathbf{x} - \hat{\mu}_0)$
 - 8: $\text{SCATE}_{\mathcal{N}}(\mathbf{x}) \leftarrow \hat{m}_1(\tilde{\mathbf{x}}) - \hat{m}_0(\mathbf{x})$
-

Application on newborn infant deliveries (natural, or not)

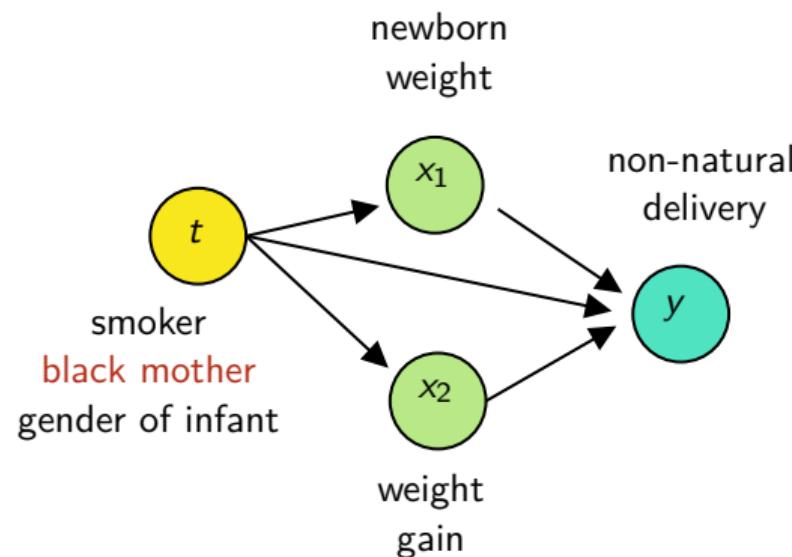
“Low birth weight paradox”, introduced by Wilcox (1993, 2001): Low birth weight (LBW) of babies x is strongly associated with increased neonatal mortality y . LBW infants born to mothers who smoke $t = 1$ usually have lower mortality rates than LBW infants born to nonsmoking mothers $t = 0$.



See Hernández-Díaz et al. (2006) and Wilcox (2006)



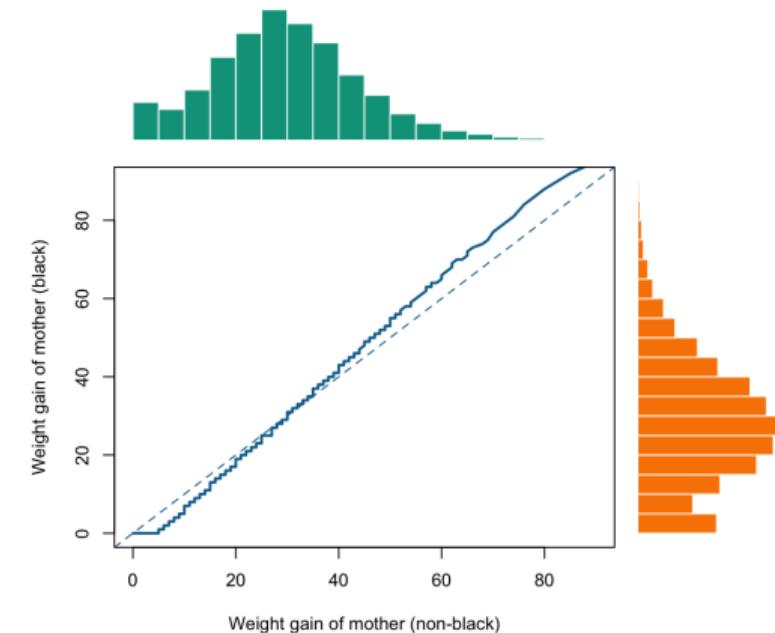
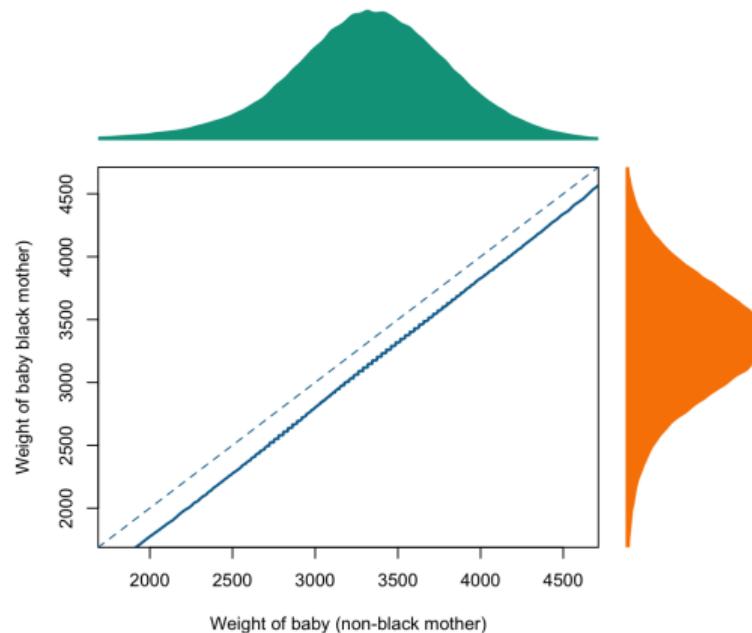
Application on newborn infant deliveries (natural, or not)



See [Charpentier et al. \(2023\)](#) for more details. Here y is a binary variable.

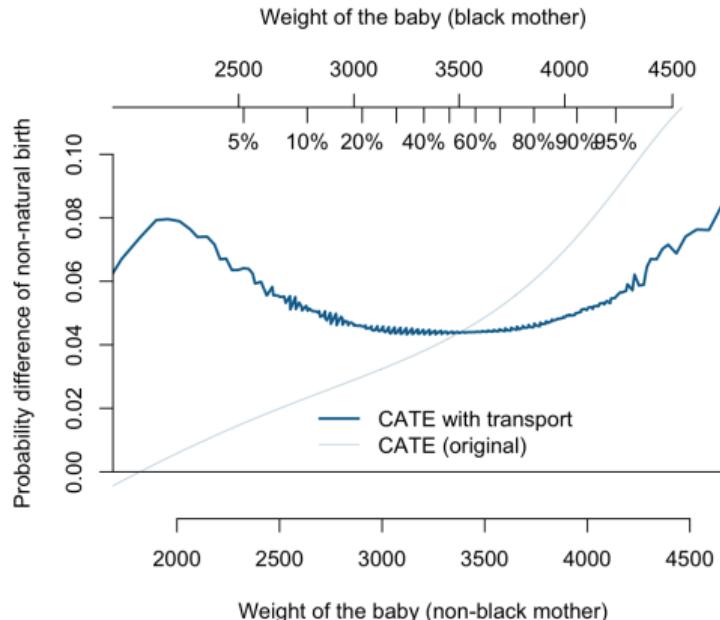
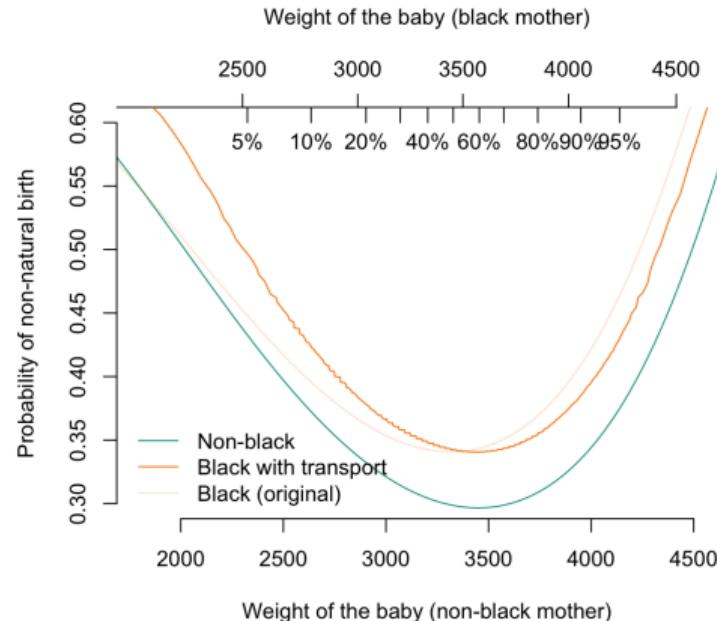
Application on newborn infant deliveries (natural, or not)

$x_1 \leftrightarrow x_1$ (newborn weight) and $x_2 \leftrightarrow x_2$ (weight gain of the mother)



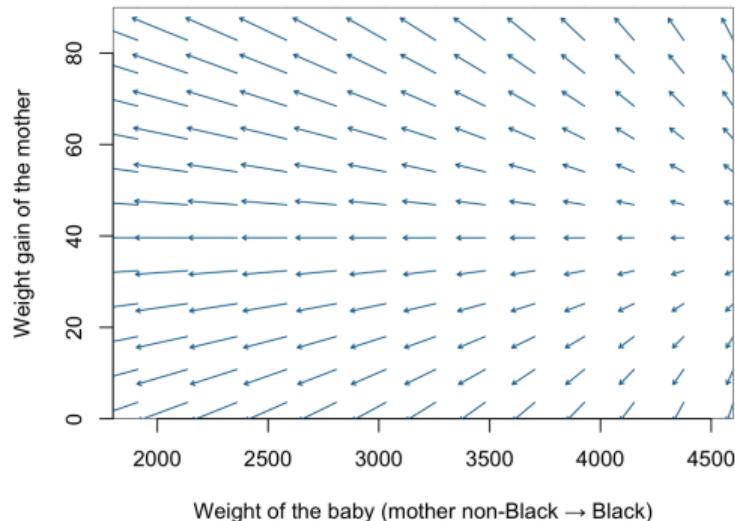
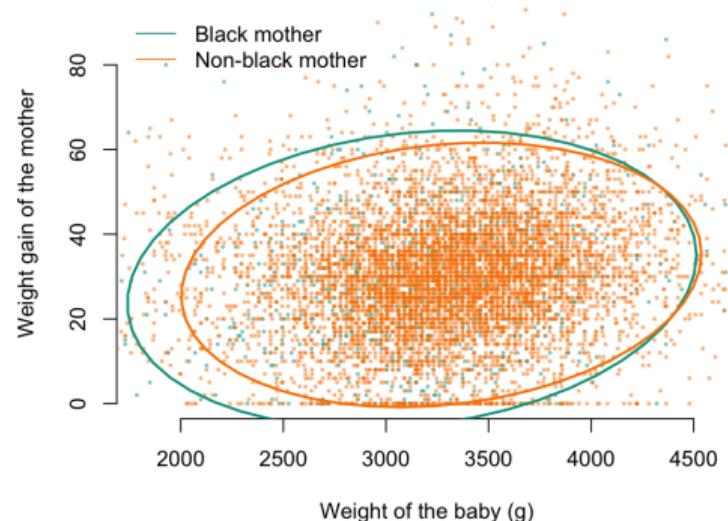
Application on newborn infant deliveries (natural, or not)

Ceteribus Paris vs. Mutatis Mutandis CATE(x_1)



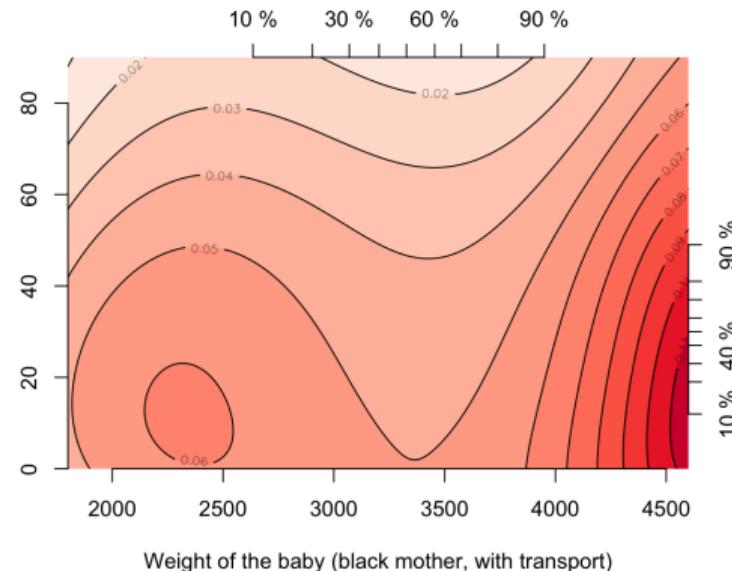
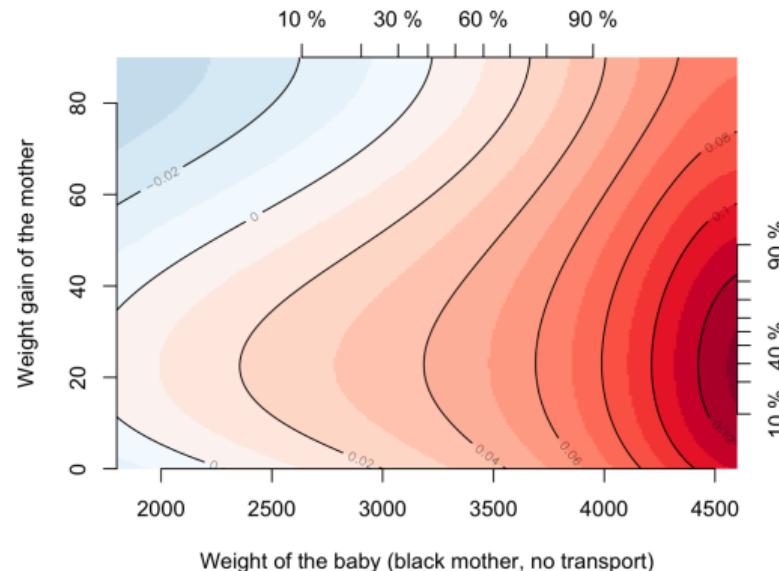
Application on newborn infant deliveries (natural, or not)

$(x_1, x_2) \leftrightarrow (x_1, x_2)$ (newborn weight, weight gain of the mother)



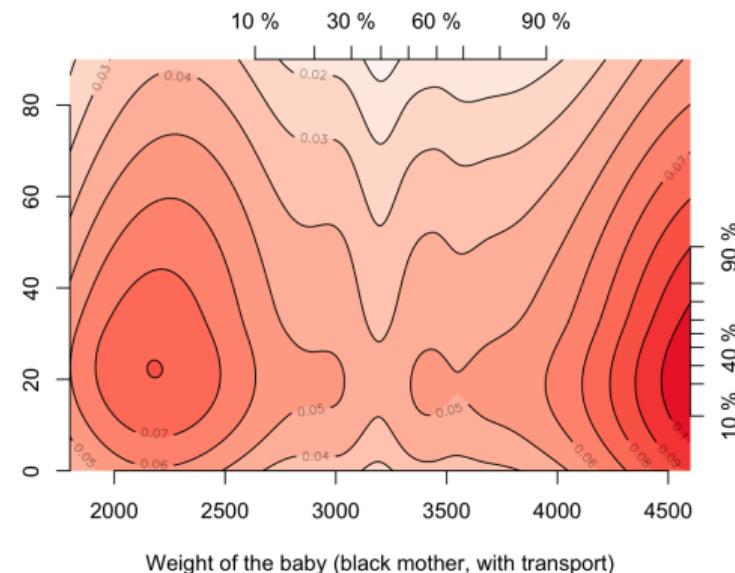
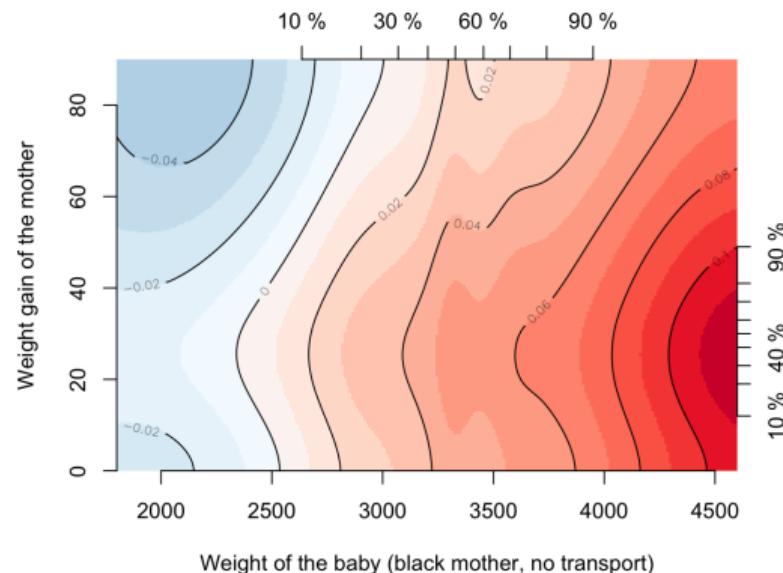
Application on newborn infant deliveries (natural, or not)

Ceteribus Paris vs. Mutatis Mutandis CATE(x_1, x_2)



Application on newborn infant deliveries (natural, or not)

Ceteribus Paris vs. Mutatis Mutandis CATE(x_1, x_2)

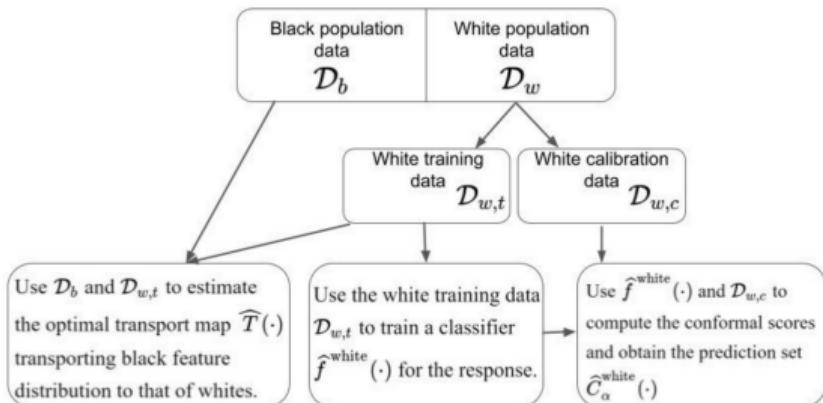


Sidenote (1)

Improving Fairness in Criminal Justice Algorithmic Risk Assessments Using Optimal Transport and Conformal Prediction Sets*

Richard A. Berk
University of Pennsylvania
Arun Kumar Kuchibhotla
Carnegie Mellon University
Eric Tchetgen Tchetgen
University of Pennsylvania

► Berk et al. (2021b)



Abstract

In the United States and elsewhere, risk assessment algorithms are being used to help inform criminal justice decision-makers. A common intent is to forecast an offender's "future dangerousness." Such algorithms have been correctly criticized for potential unfairness, and there is an active cottage industry trying to make repairs. In this paper, we use counterfactual reasoning to consider the prospects for improved fairness when members of a disadvantaged class are treated by a risk algorithm as if they are members of an advantaged class. We combine a machine learning classifier trained in a novel manner with an optimal transport adjustment for the relevant joint probability distributions, which together provide a constructive response to claims of bias-in-bias-out. A key distinction is made between fairness claims that are empirically testable and fairness claims that are not. We then use confusion tables and conformal prediction sets to evaluate achieved fairness for estimated risk. Our data are a random sample of 300,000 offenders at their arraignments for a large metropolitan area in the United States during which decisions to release or detain are made. We show that substantial improvement in fairness can be achieved consistent with a Pareto improvement for legally protected classes.

*Cary Coglianese and Sandra Mayson provided many insightful suggestions for legal conceptions of fairness and the prospect for criminal justice reform. Emanuele Candès offered several very instructive insights when commenting on this work at the Stanford/Berkeley Online Causal Inference Seminar. We also received very helpful feedback from a group of researchers at MIT and Harvard who work on causal inference. In that regard, a special thanks go to Devavrat Shah. Thanks also go to three thoughtful reviewers.

Sidenote (2)

► Hallin et al. (2021)

Theorem 2.1 (McCann 1985) Let P_1 and P_2 denote two distributions in \mathcal{P}_d . Then, (i) the class of functions

$$\nabla\Psi_{P_1;P_2} := \{\nabla\psi \mid \psi : \mathbb{R}^d \rightarrow \mathbb{R} \text{ convex, lower semi-continuous, and} \quad (2.1) \\ \text{such that } \nabla\psi \# P_1 = P_2\}$$

is not empty; (ii) if $\nabla\psi'$ and $\nabla\psi''$ are two elements of $\nabla\Psi_{P_1;P_2}$, they coincide P_1 -a.s.;⁹ (iii) if P_1 and P_2 have finite moments of order two, any element of $\nabla\Psi_{P_1;P_2}$ is an optimal quadratic transport pushing P_1 forward to P_2 .

Definition 2.2 Call $F_\pm := \nabla\phi$ the *center-outward distribution function* of $P \in \mathcal{P}_d$.

Definition 2.3 Call *empirical center-outward distribution function* any of the mappings $F_\pm^{(n)} : (\mathbf{Z}_1^{(n)}, \dots, \mathbf{Z}_n^{(n)}) \mapsto (F_\pm^{(n)}(\mathbf{Z}_1^{(n)}), \dots, F_\pm^{(n)}(\mathbf{Z}_n^{(n)})) =: F_\pm^{(n)}(\mathbf{Z}^{(n)})$ satisfying

$$\sum_{i=1}^n \|\mathbf{Z}_i^{(n)} - F_\pm^{(n)}(\mathbf{Z}_i^{(n)})\|^2 = \min_{T \in \mathcal{T}} \sum_{i=1}^n \|\mathbf{Z}_i^{(n)} - T(\mathbf{Z}_i^{(n)})\|^2 \quad (2.10)$$

or, equivalently,

$$\sum_{i=1}^n \|\mathbf{Z}_i^{(n)} - F_\pm^{(n)}(\mathbf{Z}_i^{(n)})\|^2 = \min_{\pi} \sum_{i=1}^n \|\mathbf{Z}_{\pi(i)}^{(n)} - F_\pm^{(n)}(\mathbf{Z}_i^{(n)})\|^2 \quad (2.10)$$

where the set $\{F_\pm^{(n)}(\mathbf{Z}_i^{(n)}) \mid i = 1, \dots, n\}$ consists of the n points of the augmented grid and π ranges over the $n!$ possible permutations of $\{1, 2, \dots, n\}$.

The Annals of Statistics
2021, Vol. 49, No. 2, 1139–1165
<https://doi.org/10.1214/20-AOS1996>
© Institute of Mathematical Statistics, 2021

DISTRIBUTION AND QUANTILE FUNCTIONS, RANKS AND SIGNS IN DIMENSION d : A MEASURE TRANSPORTATION APPROACH

BY MARC HALLIN¹, EUSTASIO DEL BARRO^{2,*},
JUAN CUESTA-ALBERTOS³ AND CARLOS MATRÁN²

¹ECARES and Département de Mathématique, Université libre de Bruxelles, Belgium, mhallin@ulb.ac.be

²IMUVA and Departamento de Estadística e Investigación Operativa, Universidad de Valladolid, Spain,
eustasio.delbarrio@uv.es, carlos.matran@uv.es

³Departamento de Matemáticas, Estadística y Computación, Facultad de Ciencias, Universidad de Cantabria, Spain,
juan.cuesta@unican.es

Unlike the real line, the real space \mathbb{R}^d , for $d \geq 2$, is not canonically ordered. As a consequence, such fundamental univariate concepts as quantile and distribution functions and their empirical counterparts, involving ranks and signs, do not canonically extend to the multivariate context. Palliating that lack of a canonical ordering has been an open problem for more than half a century, generating an abundant literature and motivating, among others, the development of statistical depth and copula-based methods. We show that, unlike the many definitions proposed in the literature, the measure transportation-based ranks and signs introduced in Chernozhukov, Galichon, Hallin and Henry (*Ann. Statist.* **45** (2017) 223–256) enjoy all the properties that make univariate ranks a successful tool for semiparametric inference. Related with those ranks, we propose a new *center-outward* definition of multivariate distribution and quantile functions, along with their empirical counterparts, for which we establish a Glivenko–Cantelli result. Our approach is based on McCann (*Duke Math. J.* **80** (1995) 309–323) and our results do not require any moment assumptions. The resulting ranks and signs are shown to be strictly distribution-free and essentially maximal ancillary in the sense of Basu (*Sankhyā* **21** (1959) 247–256) which, in semiparametric models involving noise with unspecified density, can be interpreted as a finite-sample form of semiparametric efficiency. Although constituting a sufficient summary of the sample, empirical center-outward distribution functions are defined at observed values only. A continuous extension to the entire d -dimensional space, yielding smooth empirical quantile contours and sign curves while preserving the essential monotonicity and Glivenko–Cantelli features of the concept, is provided. A numerical study of the resulting empirical quantile contours is conducted.

1. Introduction. Unlike the real line, the real space \mathbb{R}^d , for $d \geq 2$, is not canonically ordered. As a consequence, such fundamental concepts as quantile and distribution functions, which are strongly related to the ordering of the observation space, and their empirical counterparts—ranks and empirical quantiles—playing, in dimension $d = 1$, a fundamental role in statistical inference, do not canonically extend to dimension $d \geq 2$.

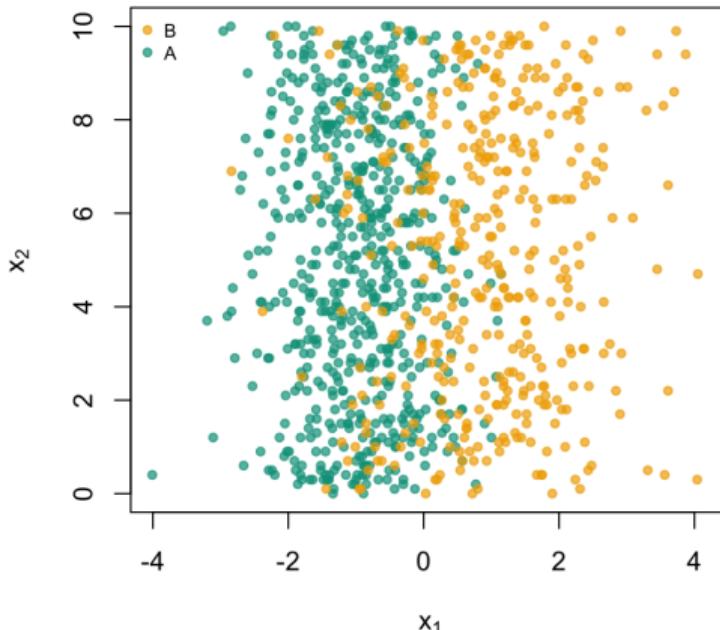
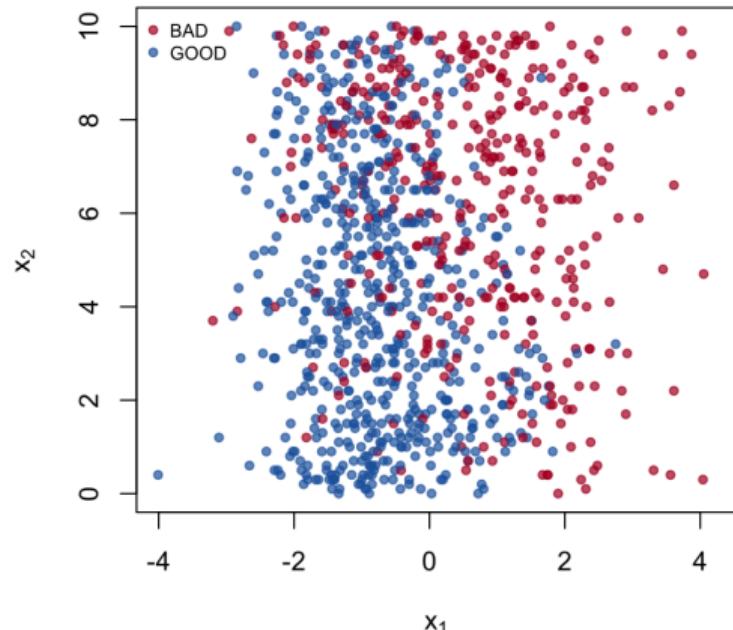
Of course, a classical concept of distribution function—the familiar one, based on marginal orderings—does exist. That concept, from a probabilistic point of view, does the job of characterizing the underlying distribution. However, the corresponding quantile function does not mean much (see, e.g., Genest and Rivest (2001)), and the corresponding empirical versions

Received May 2019; revised June 2020.

MSC2020 subject classifications. Primary 62G30; secondary 62B05.

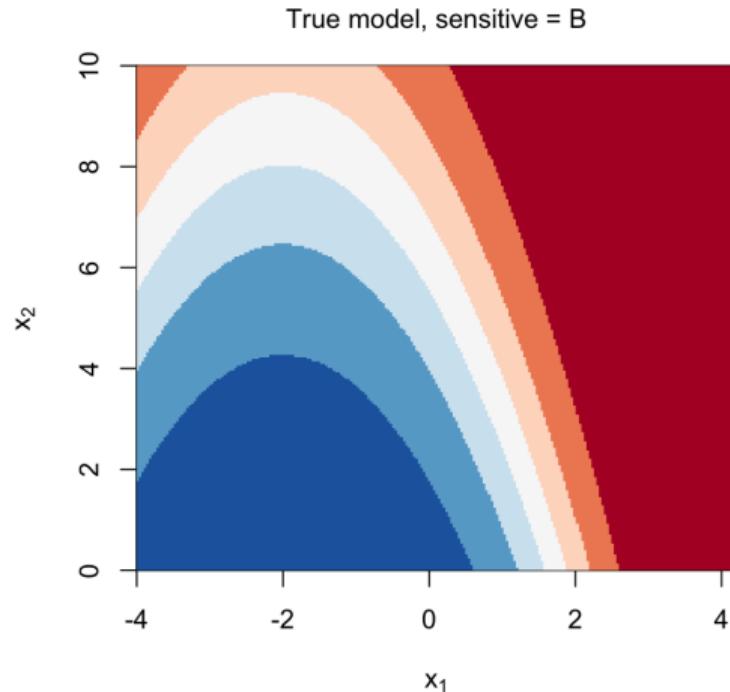
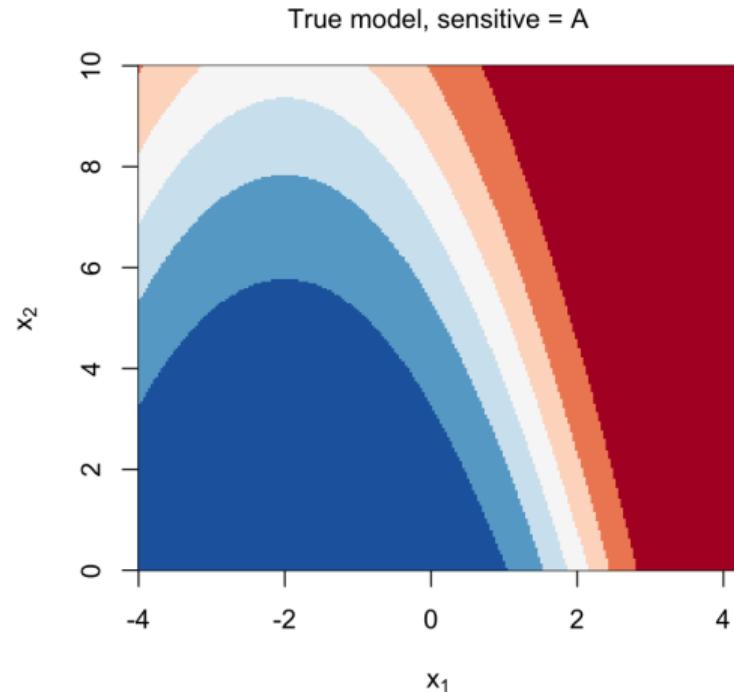
Key words and phrases. Multivariate distribution function, multivariate quantiles, multivariate ranks, multivariate signs, Glivenko–Cantelli theorem, Basu theorem, distribution-freeness, ancillarity, cyclical monotonicity.

Application on a toy example, from Charpentier (2023a)



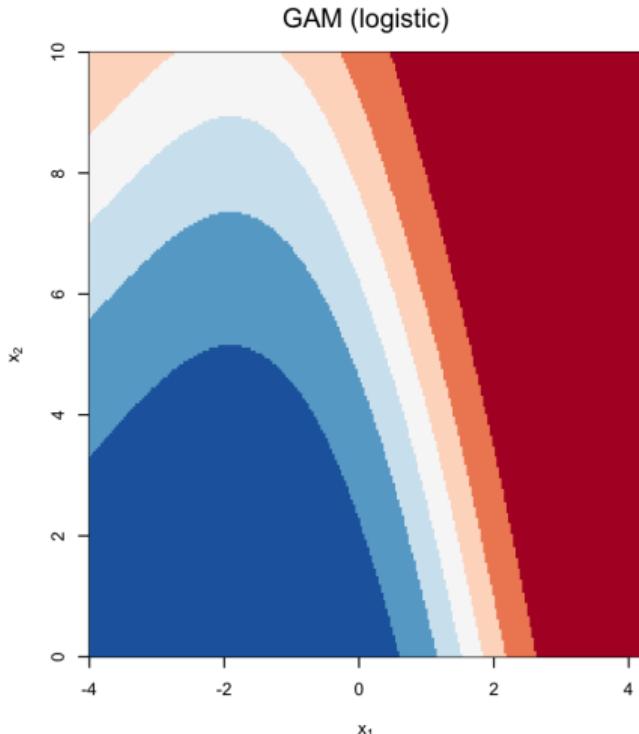
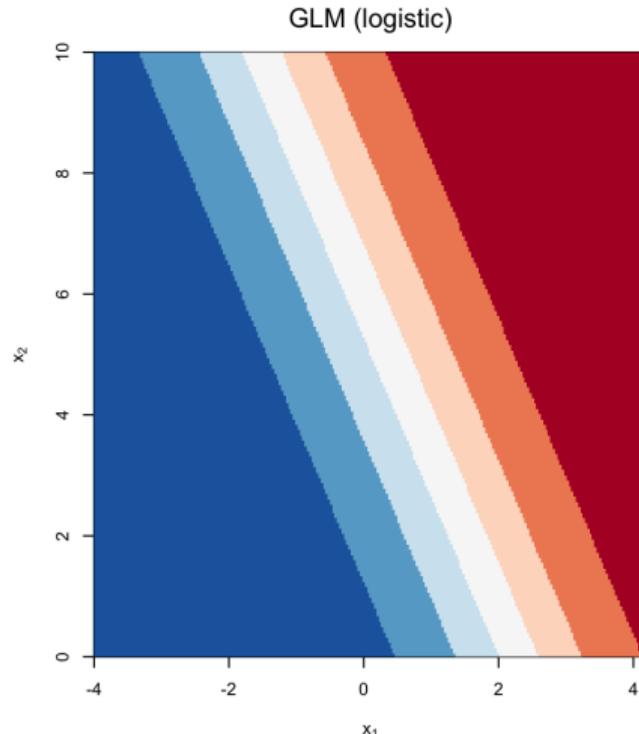
Scatterplot $(x_{1,i}, x_{2,i})$, $y_i \in \{0, 1\}$ (left) $s_i \in \{A, B\}$ (right)

Application on a toy example, from Charpentier (2023a)



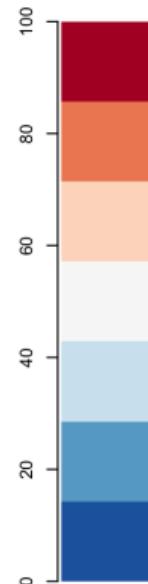
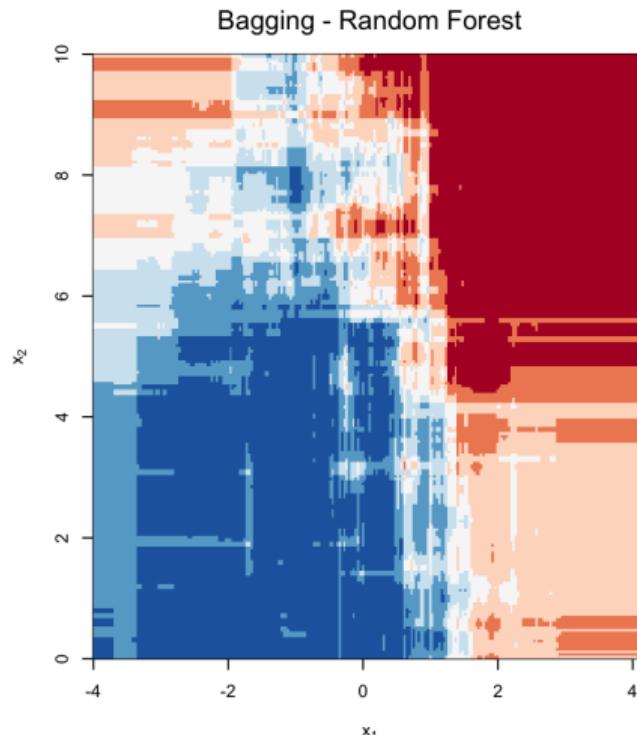
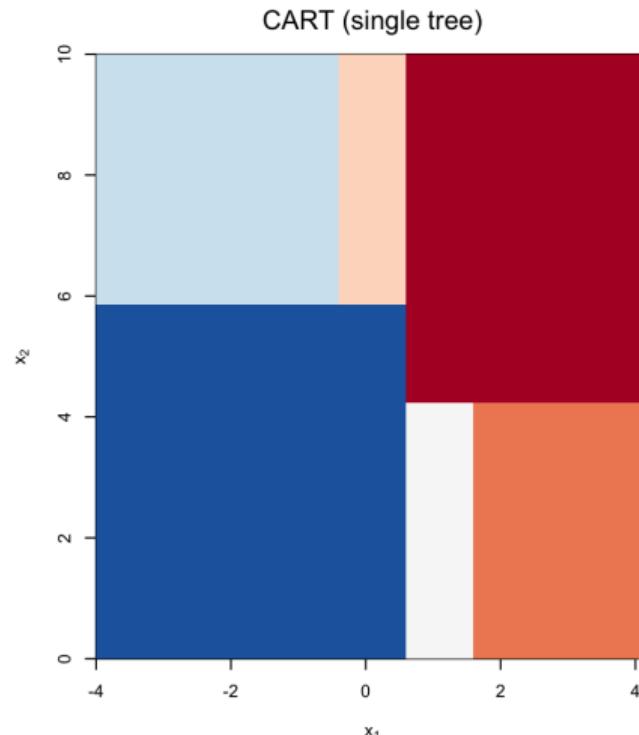
$\mu(\mathbf{x}, s) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}, S = s]$, with $\mu(\mathbf{x}, s = A)$ (left) $\mu(\mathbf{x}, s = B)$ (right)

Application on a toy example, from Charpentier (2023a)



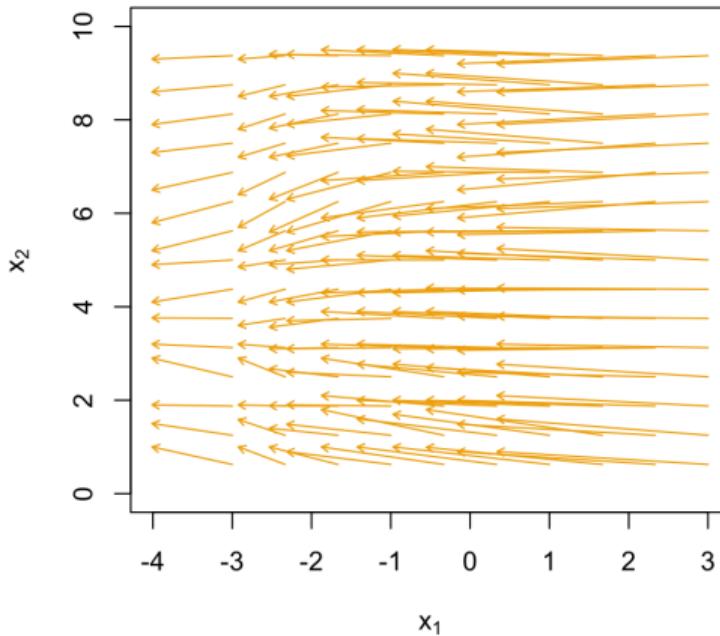
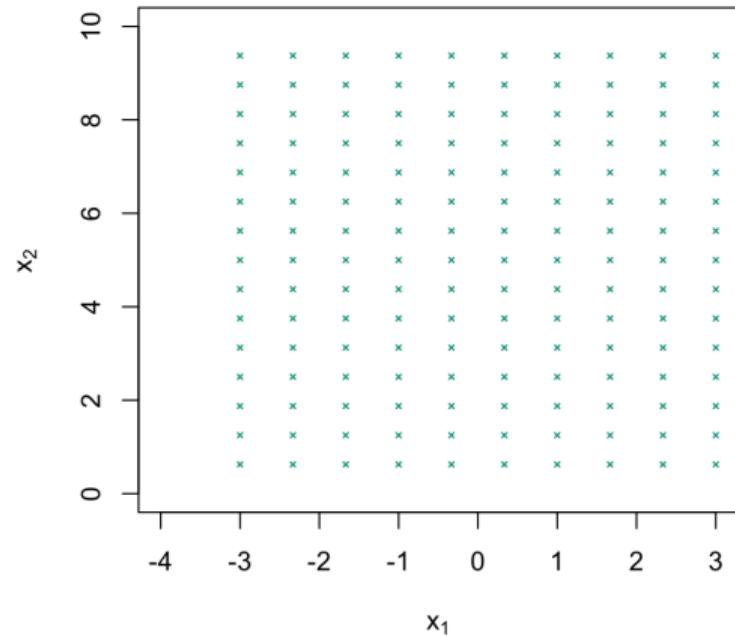
Blind fitted models $\hat{m}(\mathbf{x})$, logistic regression (left) GAM (right)

Application on a toy example, from Charpentier (2023a)



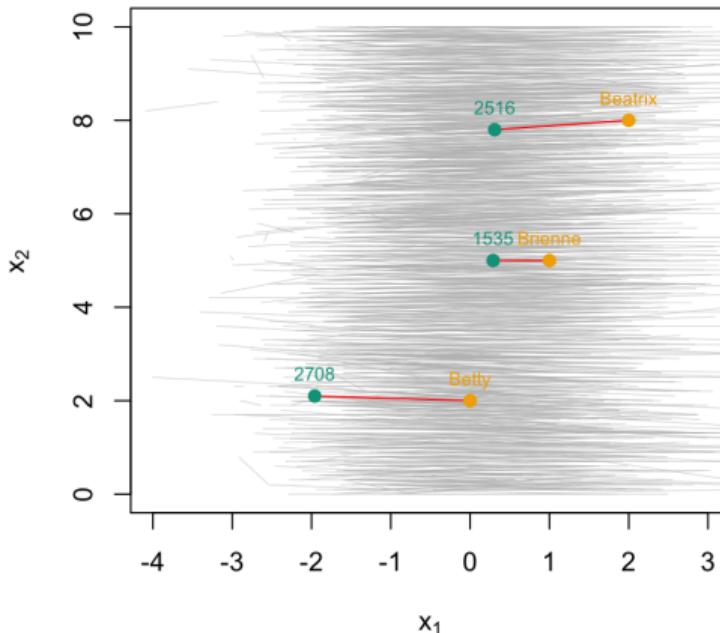
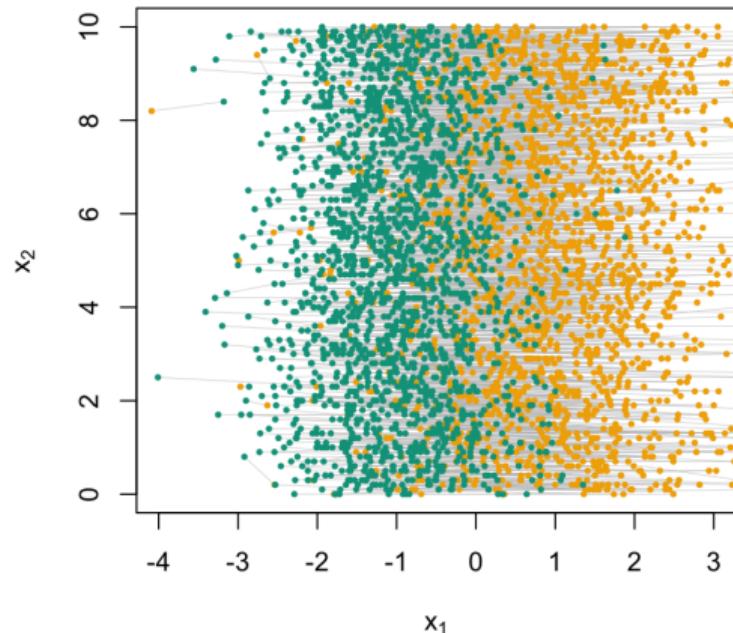
Blind fitted models $\hat{m}(x)$, classification tree (left) random forest (right)

Application on a toy example, from Charpentier (2023a)



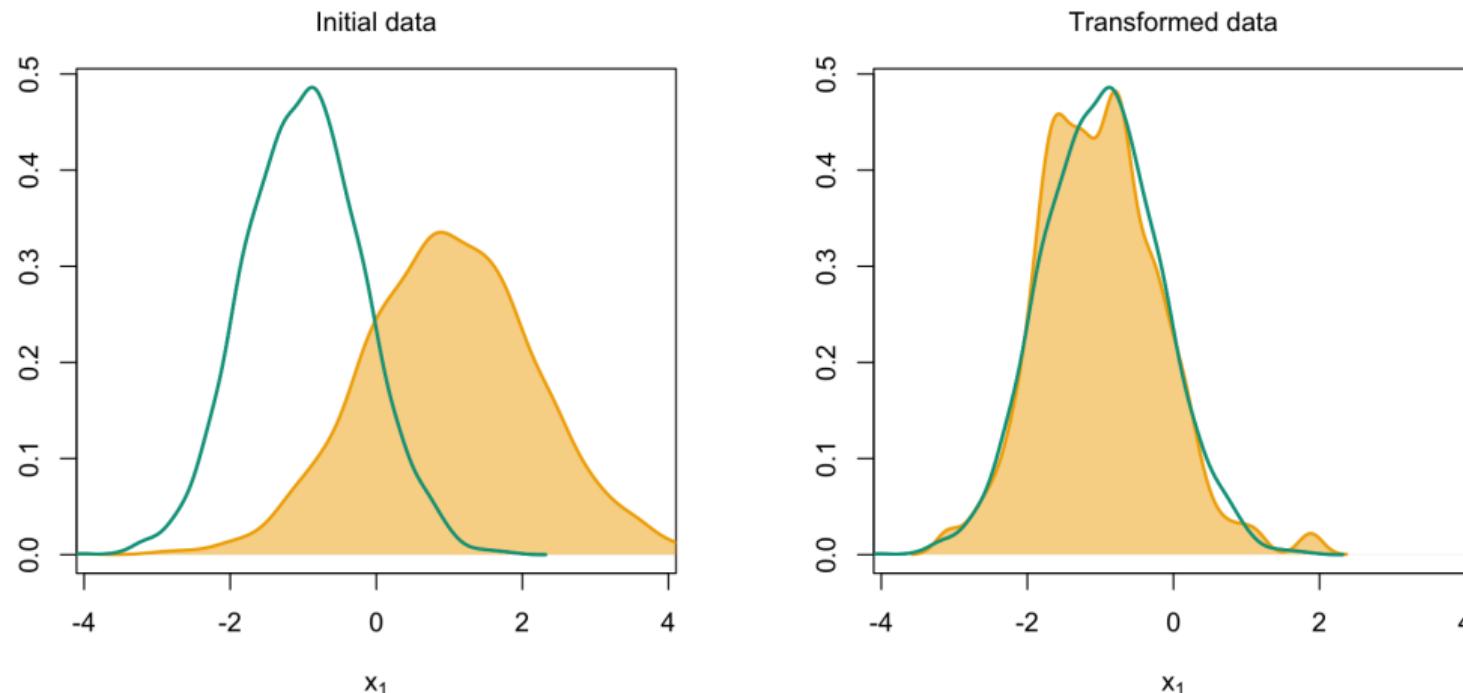
Counterfactual optimal transport, $\mathbf{x} = (x_1, x_2) \rightarrow \mathcal{T}(x_1, x_2)$, A (left) and B (right)

Application on a toy example, from Charpentier (2023a)



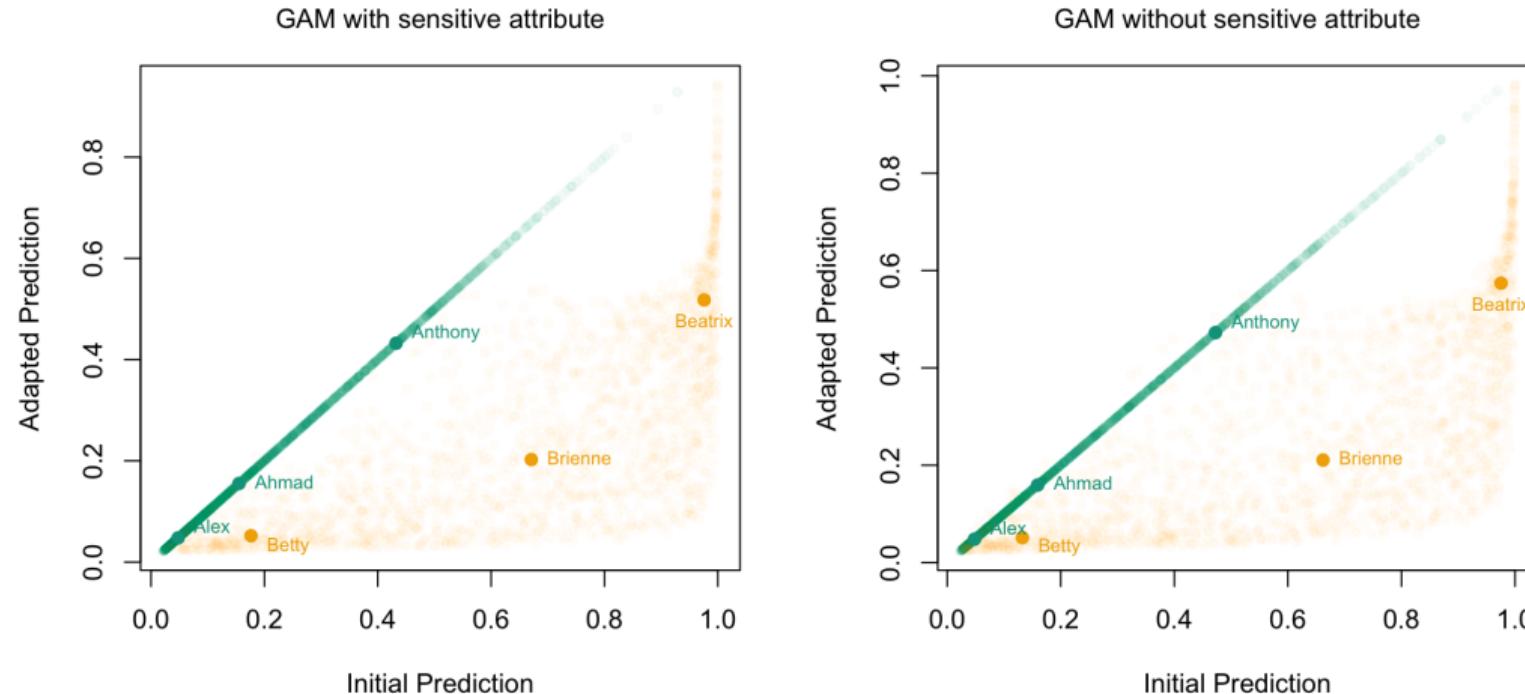
Optimal matching between individuals in group B and A.

Application on a toy example, from Charpentier (2023a)



Distribution of $x_1|s = A$, with the $x_1|s = B$ (left) and $\mathcal{T}(x_1)|s = B$ (right)

Application on a toy example, from Charpentier (2023a)



Scatterplot $\hat{m}(\mathcal{T}(x_i), s_i)$ against $\hat{m}(x_i, s_i)$ (left) $\hat{m}(\mathcal{T}(x_i))$ against $\hat{m}(x_i)$ (right)

Mitigating Discrimination

This transport can be used to quantify discrimination ("what would have been the prediction \hat{y} for that person if that person had been a man, and not a woman?") but cannot be used to mitigate discrimination ([Grari et al. \(2022\)](#) or [Hu et al. \(2023a,b\)](#))

We can therefore define some sort of average measure, solution of

$$\mathbb{P}^* = \operatorname{argmin}_{\mathbb{Q}} \left\{ \sum_{s \in S} \omega_s d(\mathbb{Q}, \mathbb{P}_s)^2 \right\},$$

for some distance (or divergence) d , as in [Nielsen and Boltz \(2011\)](#).

[Jeffreys \(1946\)](#) consider the empirical case of "*averaging histograms*" (and not theoretical measures \mathbb{P}_i), extended in [Nielsen and Nock \(2009\)](#) as the [Nielsen \(2013\)](#) as "*generalized Kullback–Leibler centroid*".

An alternative (see [Aguech and Carlier \(2011\)](#)), is to use the Wasserstein distance, to define "*Wassertein barycenter*". As shown in [Santambrogio \(2015\)](#), if one of the measures \mathbb{P}_i is absolutely continuous, the minimization problem has a unique solution.

Mitigating Discrimination

Given a reference measure, say \mathbb{P}_1 , it is possible to write the barycenter as the “*average push-forward*” transformation of \mathbb{P}_1 : if $\mathbb{P}_s = \mathcal{T}_{\#}^{1 \rightarrow s} \mathbb{P}_1$ (with the convention that $\mathcal{T}_{\#}^{1 \rightarrow 1}$ is the identity),

$$\mathbb{P}^* = \left(\sum_{s \in \mathcal{S}} \omega_s \mathcal{T}^{1 \rightarrow s} \right)_{\#} \mathbb{P}_1.$$

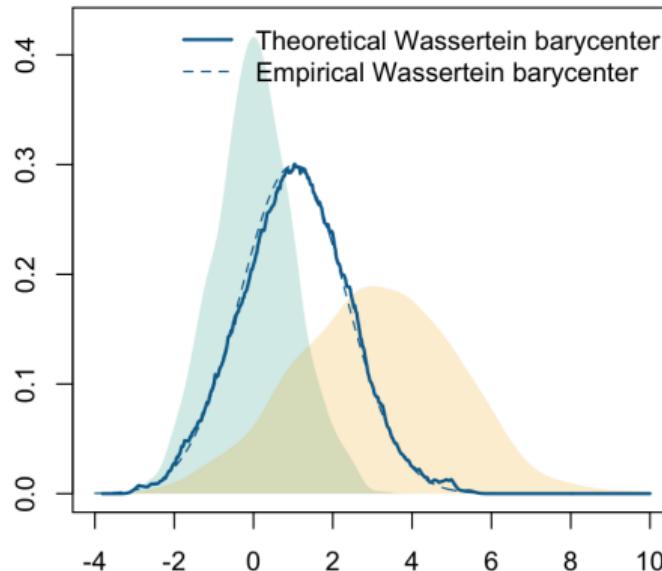
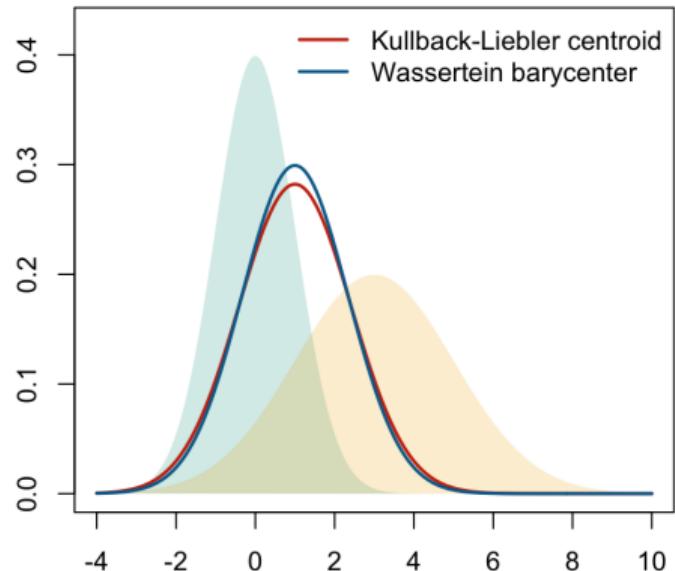
In the Gaussian case, Wasserstein barycenter is

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*), \text{ where } \boldsymbol{\mu}^* = \sum_{s \in \mathcal{S}} \omega_s \boldsymbol{\mu}_s, \boldsymbol{\Sigma}^* = \sum_{s \in \mathcal{S}} \omega_s (\boldsymbol{\Sigma}^{*1/2} \boldsymbol{\Sigma}_s \boldsymbol{\Sigma}^{*1/2})^{1/2}.$$

Jeffrey-Kullback–Leibler-centroid of those distribution would be

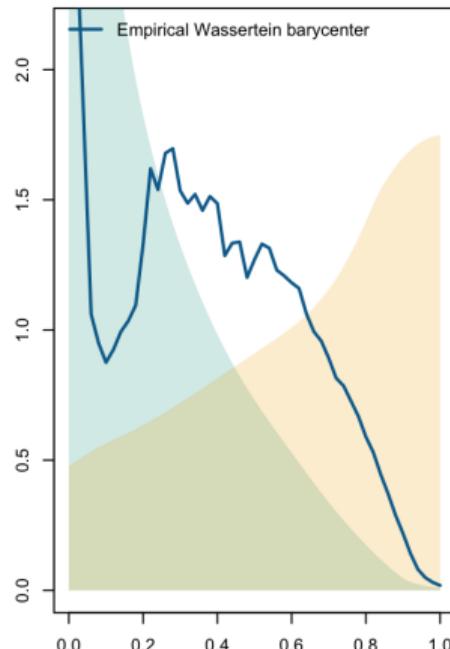
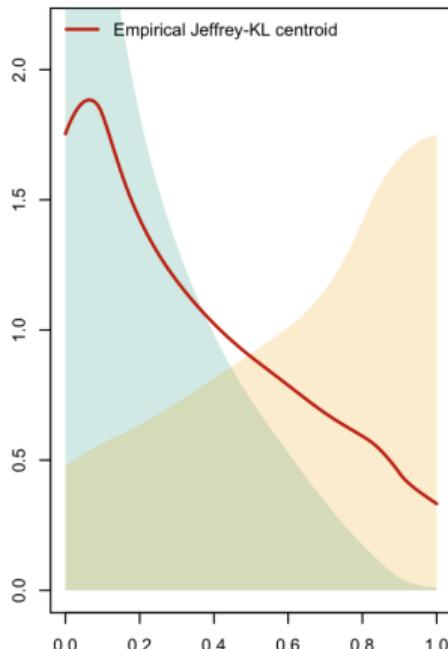
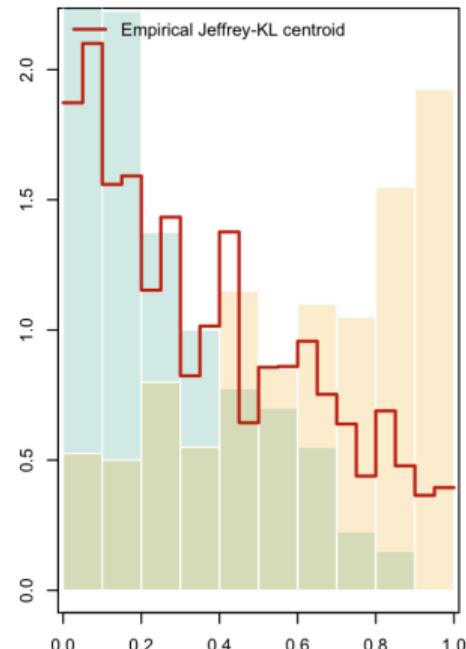
$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*), \text{ where } \boldsymbol{\mu}^* = \sum_{s \in \mathcal{S}} \omega_s \boldsymbol{\mu}_s \text{ and } \boldsymbol{\Sigma}^* = \sum_{s \in \mathcal{S}} \omega_s \boldsymbol{\Sigma}_s.$$

Mitigating Discrimination



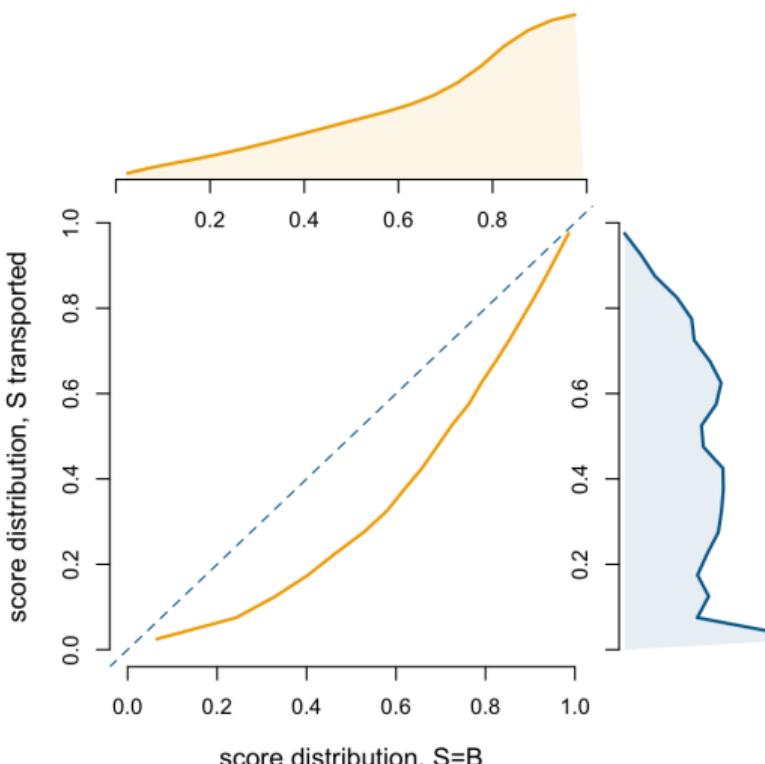
Centroid / barycenter of two Gaussian distributions

Mitigating Discrimination

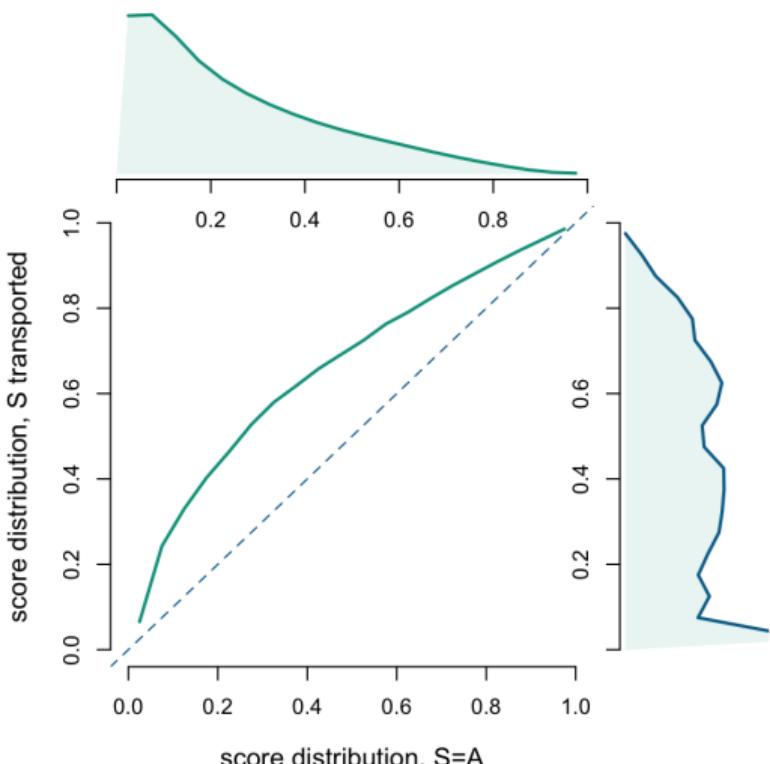


Centroid / barycenter of scoring functions $\hat{m}(x_i)$'s, $s_i \in \{A, B\}$

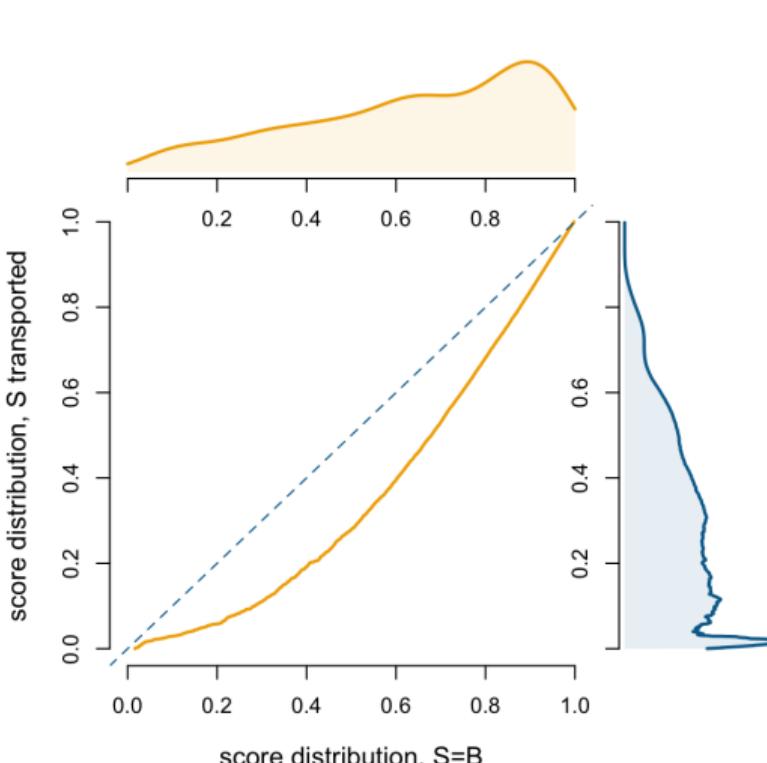
Mitigating Discrimination



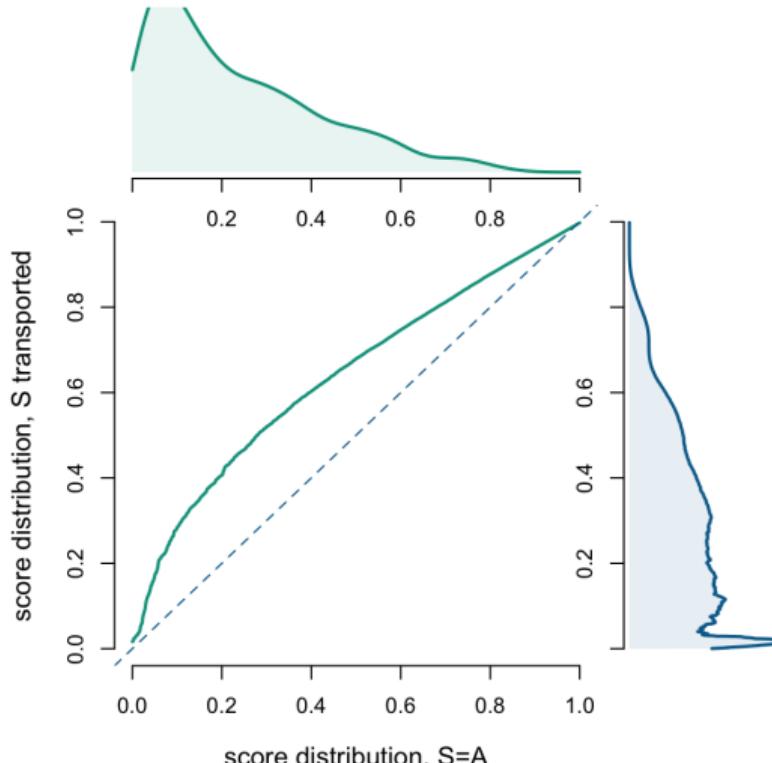
Centroid / barycenter of scoring functions $\hat{m}(x_i)$'s, $s_i \in \{A, B\}$ and barycenter



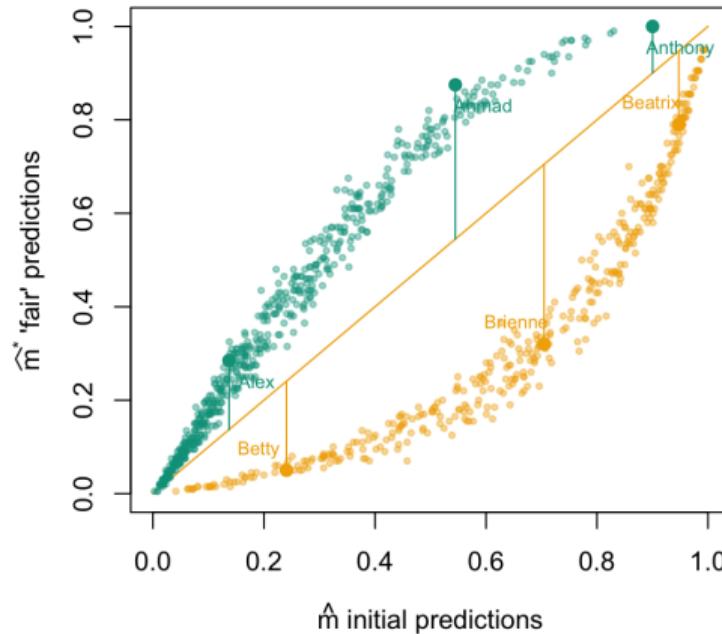
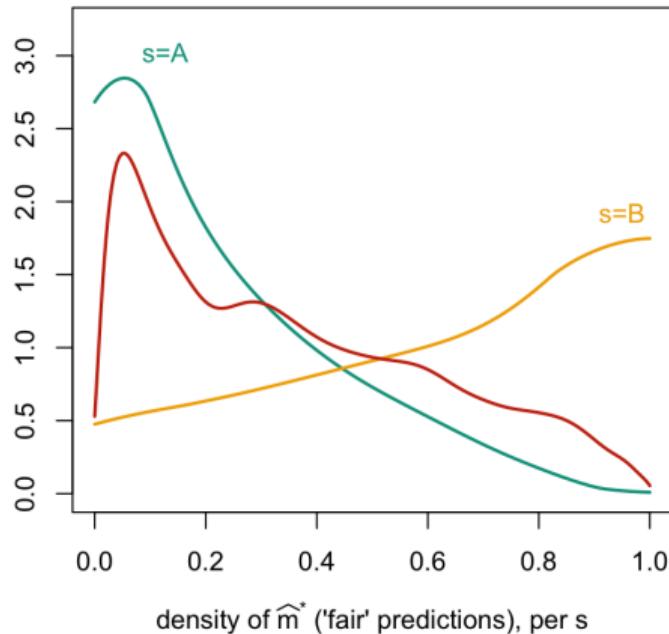
Mitigating Discrimination



Centroid / barycenter of scoring functions $\hat{m}(x_i)$'s, $s_i \in \{A, B\}$ and barycenter



Mitigating Discrimination

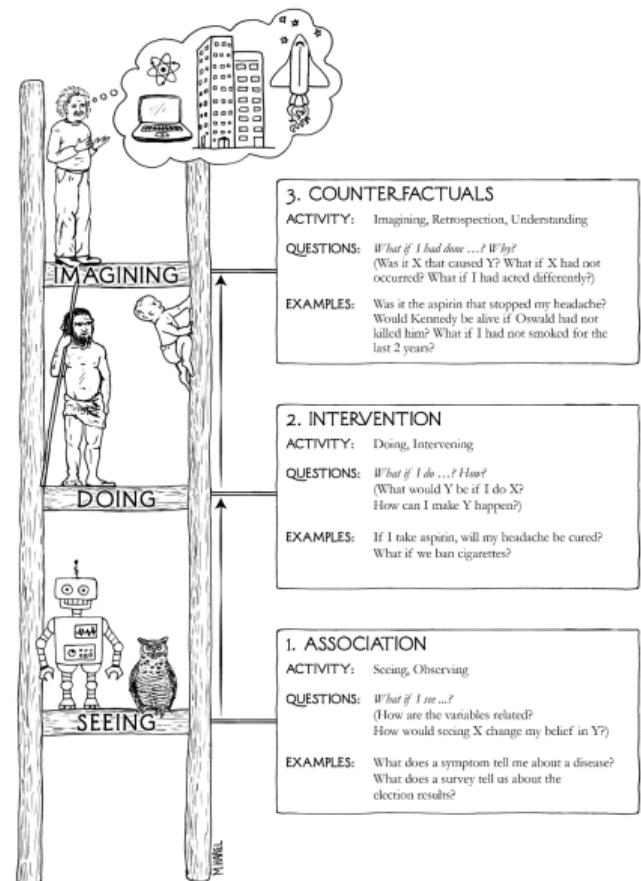


Distribution of transformed scores in the two groups $s_i \in \{A, B\}$ and barycenter

Take-away

- ▶ quantifying discrimination is hot topic
- ▶ connections with optimal transport are promising, including individual fairness
- ▶ more generally to get proper counterfactuals simple in dimension 1 (quantiles)
less natural in higher dimension
(see [Hallin et al. \(2021\)](#))
- ▶ see [Hu et al. \(2023a,b\)](#)
- ▶ see chapter 11 and 14, [Charpentier \(2023a\)](#)

Picture source: [Pearl and Mackenzie \(2018\)](#)



References

- Abrevaya, J., Hsu, Y.-C., and Lieli, R. P. (2015). Estimating conditional average treatment effects. *Journal of Business & Economic Statistics*, 33(4):485–505.
- Agarwal, S. (2021). Trade-offs between fairness and interpretability in machine learning. In *IJCAI 2021 Workshop on AI for Social Good*.
- Aguech, M. and Carlier, G. (2011). Barycenters in the wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924.
- Angrist, J. D. and Pischke, J.-S. (2009). *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.
- Angrist, J. D. and Pischke, J.-S. (2014). *Mastering'metrics: The path from cause to effect*. Princeton university press.
- Avraham, R. (2017). Discrimination and insurance. In Lippert-Rasmussen, K., editor, *Handbook of the Ethics of Discrimination*, pages 335–347. Routledge.
- Berk, R., Heidari, H., Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., Neel, S., and Roth, A. (2017). A convex framework for fair regression. *arXiv*, 1706.02409.
- Berk, R., Heidari, H., Jabbari, S., Kearns, M., and Roth, A. (2021a). Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 50(1):3–44.

References

- Berk, R. A., Kuchibhotla, A. K., and Tchetgen, E. T. (2021b). Improving fairness in criminal justice algorithmic risk assessments using optimal transport and conformal prediction sets. *arXiv*, 2111.09211.
- Beutel, A., Chen, J., Doshi, T., Qian, H., Wei, L., Wu, Y., Heldt, L., Zhao, Z., Hong, L., Chi, E. H., et al. (2019). Fairness in recommendation ranking through pairwise comparisons. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2212–2220.
- Borkan, D., Dixon, L., Sorensen, J., Thain, N., and Vasserman, L. (2019). Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pages 491–500.
- Brenier, Y. (1991). Polar factorization and monotone rearrangement of vector-valued functions. *Communications on pure and applied mathematics*, 44(4):375–417.
- Bruacli, R. A. (2006). *Combinatorial matrix classes*, volume 13. Cambridge University Press.
- Castelnovo, A., Crupi, R., Greco, G., Regoli, D., Penco, I. G., and Cosentini, A. C. (2022). A clarification of the nuances in the fairness metrics landscape. *Scientific Reports*, 12(1):1–21.
- Charpentier, A. (2022). *Insurance: biases, discrimination and fairness*. Institut Louis Bachelier.
- Charpentier, A. (2023a). *Insurance: biases, discrimination and fairness*. Springer Verlag.

References

- Charpentier, A. (2023b). Quantifying fairness and discrimination in predictive models. In Kreinovich, V., SriboonchiNa, S., and Yamaka, W., editors, *Machine Learning for Econometrics and Related Topics*. Springer Verlag.
- Charpentier, A., Flachaire, E., and Gallic, E. (2023). Causal inference with optimal transport. In Thach, N. N., Kreinovich, V., Ha, D. T., and Trung, N. D., editors, *Optimal Transport Statistics for Economics and Related Topics*. Springer Verlag.
- Chicco, D. and Jurman, G. (2020). The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. *BMC genomics*, 21(1):1–13.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163.
- Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., and Huq, A. (2017). Algorithmic decision making and the cost of fairness. *arXiv*, 1701.08230.
- Cunningham, S. (2021). *Causal inference*. Yale University Press.
- Dawid, A. P. (1982). The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610.
- de Lara, L., González-Sanz, A., Asher, N., Risser, L., and Loubes, J.-M. (2023). Transport-based counterfactual models. 2108.13025.

References

- Dehejia, R. H. and Wahba, S. (1999). Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American Statistical Association*, 94(448):1053–1062.
- Denuit, M., Charpentier, A., and Trufin, J. (2021). Autocalibration and tweedie-dominance for insurance pricing with machine learning. *Insurance: Mathematics & Economics*.
- Duivesteijn, W. and Feelders, A. (2008). Nearest neighbour classification with monotonicity constraints. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 301–316. Springer.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. (2012). Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226.
- Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. (2015). Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 259–268.
- Galichon, A. (2016). *Optimal transport methods in economics*. Princeton University Press.
- Grari, V., Charpentier, A., Lamprier, S., and Detyniecki, M. (2022). A fair pricing model via adversarial learning. *ArXiv*, 2202.12008.

References

- Hahn, J. (1998). On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, pages 315–331.
- Hallin, M., Del Barrio, E., Cuesta-Albertos, J., and Matrán, C. (2021). Distribution and quantile functions, ranks and signs in dimension d: A measure transportation approach. *The Annals of Statistics*, 49(2):1139–1165.
- Hardt, M., Price, E., and Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29:3315–3323.
- Heckman, J. J., Ichimura, H., and Todd, P. (1998). Matching as an econometric evaluation estimator. *The Review of Economic Studies*, 65(2):261–294.
- Hernández-Díaz, S., Schisterman, E. F., and Hernán, M. A. (2006). The birth weight “paradox” uncovered? *American journal of epidemiology*, 164(11):1115–1120.
- Higham, N. J. (2008). *Functions of matrices: theory and computation*. Society for Industrial and Applied Mathematics.
- Ho, D. E., Imai, K., King, G., and Stuart, E. A. (2007). Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Political analysis*, 15(3):199–236.
- Hu, F., Ratz, P., and Charpentier, A. (2023a). Fairness in multi-task learning via wasserstein barycenters. *submitted*.

References

- Hu, F., Ratz, P., and Charpentier, A. (2023b). Multivariate fair learning with statistical guarantees. *submitted*.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007):453–461.
- Kearns, M. and Roth, A. (2019). *The ethical algorithm: The science of socially aware algorithm design*. Oxford University Press.
- Kim, P. T. (2017). Auditing algorithms for discrimination. *University of Pennsylvania Law Review*, 166:189.
- Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. (2017). Human Decisions and Machine Predictions. *The Quarterly Journal of Economics*, 133(1):237–293.
- Kleinberg, J., Mullainathan, S., and Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *arXiv*, 1609.05807.
- Kranzberg, M. (1986). Technology and history:" kranzberg's laws". *Technology and culture*, 27(3):544–560.

References

- Krüger, F. and Ziegel, J. F. (2021). Generic conditions for forecast dominance. *Journal of Business & Economic Statistics*, 39(4):972–983.
- Kusner, M. J., Loftus, J., Russell, C., and Silva, R. (2017). Counterfactual fairness. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30*, pages 4066–4076. NIPS.
- Luong, B. T., Ruggieri, S., and Turini, F. (2011). k -nn as an implementation of situation testing for discrimination discovery and prevention. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 502–510.
- Matthews, B. W. (1975). Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451.
- Monge, G. (1781). Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des Sciences de Paris*.
- Narayanan, A. (2018). Translation tutorial: 21 fairness definitions and their politics. In *Proc. Conf. Fairness Accountability Transp.*, New York, USA, volume 1170, page 3.
- Nielsen, F. (2013). Jeffreys centroids: A closed-form expression for positive histograms and a guaranteed tight approximation for frequency histograms. *IEEE Signal Processing Letters*, 20(7):657–660.

References

- Nielsen, F. and Boltz, S. (2011). The burbea-rao and bhattacharyya centroids. *IEEE Transactions on Information Theory*, 57(8):5455–5466.
- Nielsen, F. and Nock, R. (2009). Sided and symmetrized bregman centroids. *IEEE transactions on Information Theory*, 55(6):2882–2904.
- Pearl, J. (2009). *Causality*. Cambridge university press.
- Pearl, J. and Mackenzie, D. (2018). *The book of why: the new science of cause and effect*. Basic books.
- Rubin, D. B. (1973). Matching to remove bias in observational studies. *Biometrics*, pages 159–183.
- Santambrogio, F. (2015). Optimal transport for applied mathematicians.
- Stuart, E. A. (2010). Matching methods for causal inference: A review and a look forward. *Statistical Science*, 25(1):1.
- Van Calster, B., McLernon, D. J., Van Smeden, M., Wynants, L., and Steyerberg, E. W. (2019). Calibration: the achilles heel of predictive analytics. *BMC medicine*, 17(1):1–7.
- Villani, C. (2003). *Topics in optimal transportation*, volume 58. American Mathematical Society.
- Villani, C. (2009). *Optimal transport: old and new*, volume 338. Springer Berlin, Heidelberg.

References

- Vogel, R., Bellet, A., Clémén, S., et al. (2021). Learning fair scoring functions: Bipartite ranking under roc-based fairness constraints. In *International Conference on Artificial Intelligence and Statistics*, pages 784–792. PMLR.
- Wilcox, A. J. (1993). Birth weight and perinatal mortality: the effect of maternal smoking. *American journal of epidemiology*, 137(10):1098–1104.
- Wilcox, A. J. (2001). On the importance—and the unimportance—of birthweight. *International journal of epidemiology*, 30(6):1233–1241.
- Wilcox, A. J. (2006). Invited commentary: the perils of birth weight—a lesson from directed acyclic graphs. *American Journal of Epidemiology*, 164(11):1121–1123.
- Zafar, M. B., Valera, I., Rodriguez, M. G., and Gummadi, K. P. (2017). Fairness constraints: Mechanisms for fair classification. *arXiv*, 1507.05259.
- Zhang, B. H., Lemoine, B., and Mitchell, M. (2018). Mitigating unwanted biases with adversarial learning. *arXiv*, 1801.07593.