

# Fairness and discrimination in actuarial pricing

**Arthur Charpentier<sup>1</sup>, Laurence Barry<sup>2</sup>, Vincent Grari<sup>3</sup>,**  
**Lamprier Sylvain<sup>3</sup>, Detyniecki Marcin<sup>4</sup>**

<sup>1</sup> Univeristé du Québec à Montréal    <sup>2</sup> Chaire Pari    <sup>3</sup> Sorbonne Université    <sup>4</sup> AXA

CIRM - MLISTRAL - September 2022

## Agenda

- ▶ Charpentier (2022) Insurance: Discrimination, Biases and Fairness, *Institut Louis Bachelier* 
  - ▶ Grari et al. (2022) A fair pricing model via adversarial learning, *ArXiv:2202.12008*  
  - ▶ “*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)



## Summary

**Market value and market timing approach** This approach is based on the premise that stocks have the potential to rise or fall in value over time. It is based on the belief that the market is efficient and that it is impossible to consistently beat the market by predicting future price movements. Instead, the focus is on identifying undervalued stocks and holding them for the long term. This approach emphasizes the importance of diversification and rebalancing the portfolio periodically to reflect changes in the market environment.

involving a reduction and then re-enrichment. Hence, the measured increase in the number of individuals belonging to a single group participating in the same activity is often interpreted as reflecting a "loss" of individuals from the population.

Population dynamics can also be studied by examining the rates of change in the size of a population over time. The growth rate of a population is defined as the ratio of the increase in the size of the population to the total size of the population. This measure of growth is called the "rate of natural increase." The growth rate of a population is determined by the birth rate and the death rate. The birth rate is the number of live births per unit of time, and the death rate is the number of deaths per unit of time. The growth rate of a population is equal to the birth rate minus the death rate.

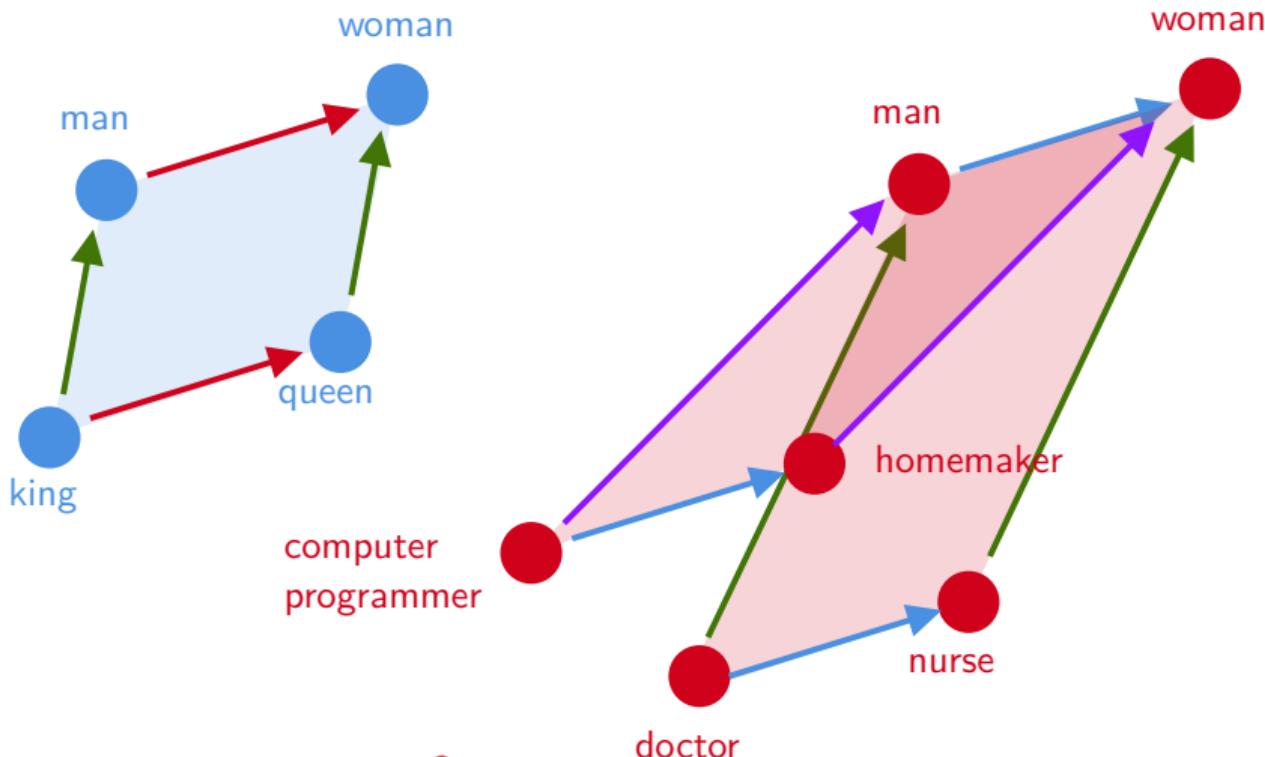
The growth rate of a population is also known as the "rate of natural increase." This measure of growth is called the "rate of natural increase." The growth rate of a population is determined by the birth rate and the death rate. The birth rate is the number of live births per unit of time, and the death rate is the number of deaths per unit of time. The growth rate of a population is equal to the birth rate minus the death rate.

It is a pleasure to receive these comments from our colleagues. We are grateful for the interest shown in our work and for the useful suggestions made. We have tried to incorporate some of the changes in the original version of the article. Thus, there are now more details on the methods used to estimate the parameters of the model. The discussion on the implications of the results has been modified to emphasize the fact that a reduction in the number of individuals in a population can lead to a reduction in the incidence of disease. This is particularly important in the case of a disease such as AIDS which is transmitted through sexual contact. The discussion on the implications of the results has also been modified to emphasize the fact that a reduction in the number of individuals in a population can lead to a reduction in the incidence of disease. This is particularly important in the case of a disease such as AIDS which is transmitted through sexual contact.

The project reflects the extensive experience of the two organizations in developing and applying the LSCM approach. The LSCM System: Theory and Application Guidebook is available at [www.lscm.org](http://www.lscm.org).

# Motivation

- Accuracy :  $\pi(\mathbf{x}) = \mathbb{E}_{\mathbb{P}}[Y|\mathbf{X} = \mathbf{x}]$  ( $\mathbb{P}$  historical probability) (is)
- Fairness :  $\pi^*(\mathbf{x}) = \mathbb{E}_{\mathbb{P}^*}[Y|\mathbf{X} = \mathbf{x}]$  ( $\mathbb{P}^*$  targeted probability) (ought, Hume (1739))



Anglais

a doctor, a nurse

Français

un médecin, une infirmière

Espagnol

Una doctora, una enfermera (feminin)

Un doctor, un enfermero (masculin)

freakonometrics

freakonometrics.hypotheses.org

# Actuarial Pricing I

- “at the core of insurance business lies discrimination between risky and non-risky insureds”, Avraham (2017)

$$\left\{ \begin{array}{l} \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d : \text{'explanatory' variables} \\ \mathbf{x}_c \in \mathcal{X}_c \subset \mathbb{R}^{d_c} : \text{car-vehicule 'explanatory' variables} \\ \mathbf{x}'_c = m_c(\mathbf{x}_c; \mathbf{x}) \in \mathbb{R} : \text{car-vehicle scoring} \\ \mathbf{x}_g \in \mathcal{X}_g \subset \mathbb{R}^{d_g} : \text{geographic 'explanatory' variables} \\ \mathbf{x}'_g = m_g(\mathbf{x}_g; \mathbf{x}) \in \mathbb{R} : \text{geographic scoring} \\ y \in \{0, 1\}, \mathbb{N} \text{ or } \mathbb{R}_+ : \text{variable of interest} \\ \hat{y} = m(\mathbf{x}, \mathbf{x}'_c, \mathbf{x}'_g) : \text{prediction (or score)} \end{array} \right.$$

Classically, a two-stage approach is considered to create a geographic score

- Fit a model  $\hat{m}$  to predict  $y$  based on  $\mathbf{x}$ , compute residuals  $\hat{\varepsilon}$
- Fit a model  $\hat{m}_g(\cdot; \mathbf{x})$  to predict  $\hat{\varepsilon}$  based on  $\mathbf{x}_g$
- define a score  $\mathbf{x}'_g = \hat{m}_g(\mathbf{x}_g; \mathbf{x})$  (that is function of  $\mathbf{x}$ , too)

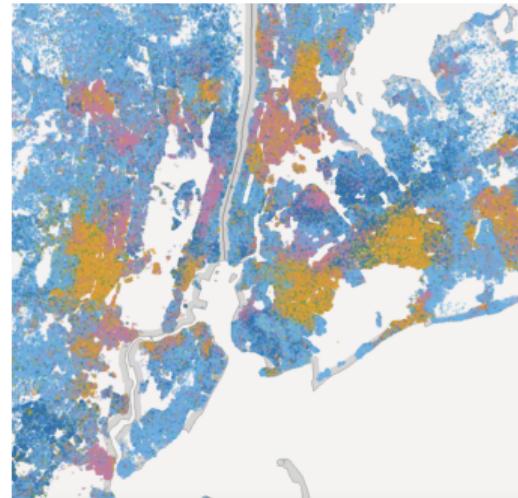
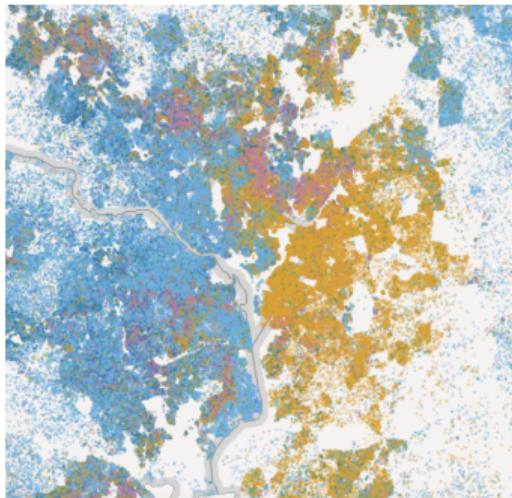
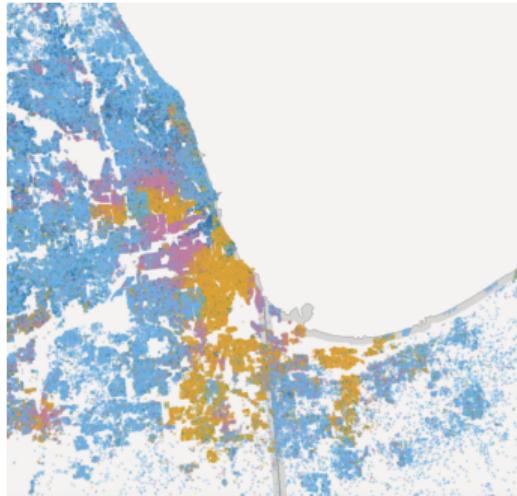
# Protected Attributes ? Motor insurance in North America

	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	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

# Protected Attributes ? Motor insurance in North America

Spatial information and racial bias (redlining)



## Defining Group Fairness when $y$ is binary I

$$\left\{ \begin{array}{l} \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d : \text{'explanatory' variables} \\ p \in \{0, 1\} : \text{protected variable or sensitive attribute} \\ \left\{ \begin{array}{l} \mathbf{x}_c \in \mathcal{X}_c \subset \mathbb{R}^{d_c} : \text{car-vehicle 'explanatory' variables} \\ \mathbf{x}'_c = m_c(\mathbf{x}_c; \mathbf{x}, p) \in \mathbb{R} : \text{car-vehicle scoring} \end{array} \right. \\ \left\{ \begin{array}{l} \mathbf{x}_g \in \mathcal{X}_g \subset \mathbb{R}^{d_g} : \text{geographic 'explanatory' variables} \\ \mathbf{x}'_g = m_g(\mathbf{x}_g; \mathbf{x}, p) \in \mathbb{R} : \text{geographic scoring} \end{array} \right. \\ y \in \{0, 1\} : \text{variable of interest (binary)} \\ s = m(\mathbf{x}, p, \mathbf{x}'_c, \mathbf{x}'_g) : \text{score} \\ \hat{y} = \mathbf{1}(s > \text{threshold}) : \text{prediction} \end{array} \right.$$

Fairness Through Unawareness, Kusner et al. (2017)

Protected attribute  $p$  is not explicitly used in decision function  $\hat{y}$ .

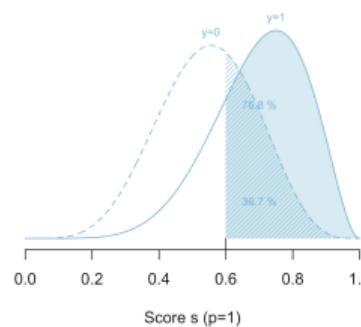
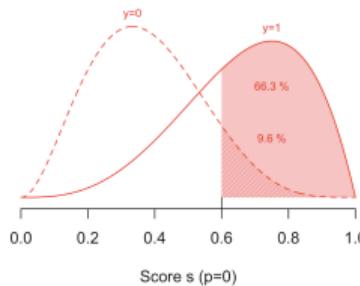
# Defining Group Fairness when $y$ is binary II

**Demographic Parity**, (Corbett-Davies et al. (2017), Agarwal (2021))

Decision function  $\hat{Y}$  satisfies demographic parity if  $\hat{Y} \perp\!\!\!\perp P$ , i.e.

$$\mathbb{P}[\hat{Y} = y | P = 0] = \mathbb{P}[\hat{Y} = y | P = 1], \forall y \text{ or } \mathbb{E}[\hat{Y} | P = 0] = \mathbb{E}[\hat{Y} | P = 1]$$

We can compare  $s(\mathbf{X})$  conditional on  $Y$ , but also on  $P$



# Defining Group Fairness when $y$ is binary III

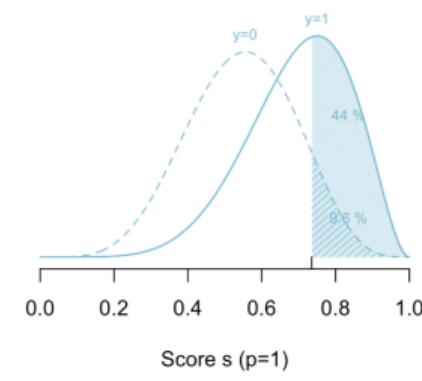
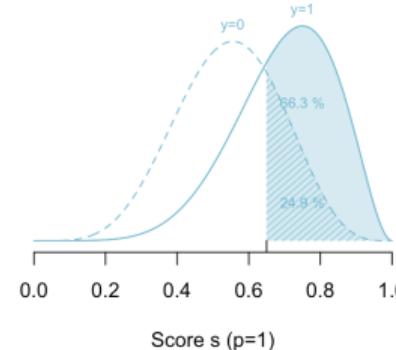
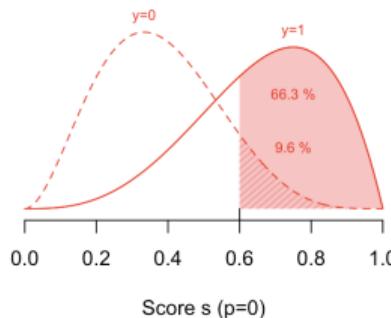
**Equal Opportunity**, Hardt et al. (2016)

True positive parity

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

or false positive parity

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



# Defining Group Fairness when $y$ is binary IV

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

## Dependence measures I

Group fairness is characterized by independence or conditional independence properties. Given two random variables  $U$  and  $V$ ,

$$C(U, V) = \begin{cases} \text{corr}[U, V] & \text{Pearson's correlation} \\ \text{corr}[F_U(U), F_V(V)] & \text{Spearman's rank correlation} \\ \tau[U, V] & \text{Kendall's tau} \end{cases}$$

that can be extended to conditional measures, as [Lawrance \(1976\)](#), since

$$\text{corr}[U, V] = \mathbb{E}[UV] \text{ when } \begin{cases} \mathbb{E}[U] = \mathbb{E}[V] = 0 \\ \mathbb{E}[U^2] = \mathbb{E}[V^2] = 1 \end{cases}$$

$$\begin{cases} \textbf{Demographic Parity} : \hat{Y} \perp\!\!\!\perp P \implies C(\hat{Y}, P) = 0 \\ \textbf{Equalized Odds} : \hat{Y} \perp\!\!\!\perp P|Y \implies C(\hat{Y}, P|Y = y) = 0, \forall y \end{cases}$$

## Dependence measures II

Hirschfeld (1935), Gebelein (1941) and Rényi (1959)

$$HGR(U, V) = \max \{ \text{corr}[f(U), g(V)] \} = \max_{f \in \mathcal{S}_U, g \in \mathcal{S}_V} \{ \mathbb{E}[f(U)g(V)] \}$$

where  $\mathcal{S}_U = \{f : \mathcal{U} \rightarrow \mathbb{R} : \mathbb{E}[f(U)] = 0 \text{ and } \mathbb{E}[f(U)^2] = 1\}$  and similarly  $\mathcal{S}_V$ .  
One can also consider a conditional version,

$$HGR(U, V|Z) = \max_{f \in \mathcal{S}_{U|Z}, g \in \mathcal{S}_{V|Z}} \{ \mathbb{E}[f(U)g(V)|Z] \}$$

where  $\mathcal{S}_{U|Z} = \{f : \mathcal{U} \rightarrow \mathbb{R} : \mathbb{E}[f(U)|Z] = 0 \text{ and } \mathbb{E}[f(U)^2|Z] = 1\}$ .

This measure can be used to characterize independence,

$$U \perp\!\!\!\perp V \iff HGR(U, V) = 0,$$

and if  $(U, V)$  is a Gaussian vector,  $HGR(U, V) = |\text{corr}(U, V)|$ .

## Dependence measures III

Thus, this measure can be used to characterize fairness,

$$\begin{cases} \text{Demographic Parity} : \hat{Y} \perp\!\!\!\perp P \iff HGR(\hat{Y}, P) = 0 \\ \text{Equalized Odds} : \hat{Y} \perp\!\!\!\perp P|Y \iff HGR(\hat{Y}, P|Y = y) = 0, \forall y \end{cases}$$

$HGR$  can be difficult to estimate, but one can use some Neural Networks

$$HGR_{NN}(U, V) = \max_{\omega_f, \omega_g} \left\{ \mathbb{E}[f_{\omega_f}(U)g_{\omega_g}(V)] \right\}$$

See also [Breiman and Friedman \(1985\)](#) on the estimation of this maximal correlation, in the context of regression

## Dependence measures IV

More generally,  $\mathbf{V}$  can be a vector on  $\mathcal{V} \subset \mathbb{R}^k$ , then

$$HGR(U, \mathbf{V}) = \max_{h: \mathcal{V} \rightarrow \mathbb{R}} \{ HGR[U, h(\mathbf{V})] \} = \max_{f \in \mathcal{S}_U, g \in \mathcal{S}_{\mathcal{V}}} \{ \mathbb{E}[f(U)g(\mathbf{V})] \}$$

where  $\mathcal{S}_{\mathcal{V}} = \{g : \mathcal{V} \rightarrow \mathbb{R} : \mathbb{E}[g(\mathbf{V})] = 0 \text{ and } \mathbb{E}[g(\mathbf{V})^2] = 1\}$ . A conditional version exists, and one can estimate that measure using a neural network,

$$HGR_{NN}(U, \mathbf{V}) = \max_{\omega_f, \omega_g} \{ \mathbb{E}[f_{\omega_f}(U)g_{\omega_g}(\mathbf{V})] \}$$

$$\begin{cases} \text{Demographic Parity} : \hat{Y} \perp\!\!\!\perp \mathbf{P} & \iff HGR(\hat{Y}, \mathbf{P}) = 0 \\ \text{Equalized Odds} : \hat{Y} \perp\!\!\!\perp \mathbf{P} | Y & \iff HGR(\hat{Y}, \mathbf{P} | Y = y) = 0, \forall y \end{cases}$$

# Adversarial Approach I

In a classical ML or econometric pricing model, solve

$$\operatorname{argmin}_{\theta} \{\mathcal{L}(\hat{y}, y)\}, \text{ where } \mathcal{L}(\hat{y}, y) = \sum_{i=1}^n \ell(\hat{y}_i, y_i) \text{ and } \hat{y} = h_{\theta}(x)$$

either related to some loss, or some log-likelihood,

To avoid over-fit: penalize complexity (penalty  $\mathcal{P}$ )

$$\operatorname{argmin}_{\theta} \{\mathcal{L}(h_{\theta}(x), y) + \lambda \mathcal{P}(h_{\theta})\}$$

## Adversarial Approach II

Inspired by Goodfellow et al. (2018) (but also Bechavod and Ligett (2017) or Cho et al. (2020)), to avoid un-fairness: penalize according to  $HGR(\hat{y}, p)$  (for demographic parity),

$$\operatorname{argmin}_{\theta, \omega_f, \omega_g} \left\{ \mathcal{L}(h_{\theta}(\mathbf{x}), y) + \lambda HGR_{\omega_f, \omega_g}(h_{\theta}(\mathbf{x}), p) \right\}$$

i.e.

$$\operatorname{argmin}_{\theta} \left\{ \max_{\omega_f, \omega_g} \left\{ \mathcal{L}(h_{\theta}(\mathbf{X}), Y) + \lambda \mathbb{E}_{(\mathbf{X}, S) \sim \mathcal{D}} (\hat{f}_{\omega_f}(h_{\theta}(\mathbf{X})) \hat{g}_{\omega_g}(P)) \right\} \right\}$$

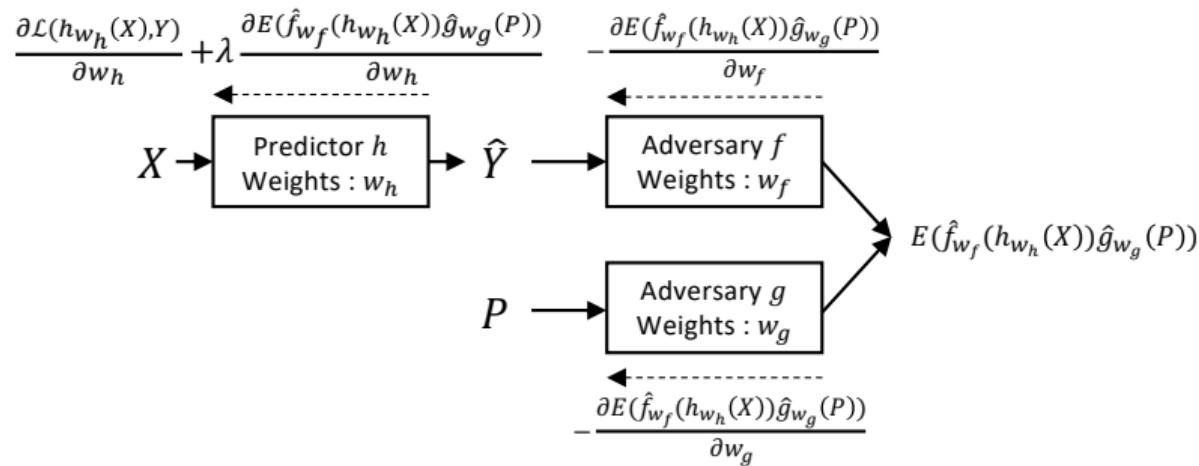
or  $HGR(\hat{y}, p|y)$  (for equalized odds), i.e. when  $y \in \{0, 1\}$

$$\begin{aligned} \operatorname{argmin}_{\theta} \left\{ \max_{\omega_{f0}, \omega_{g0}, \omega_{f1}, \omega_{g1}} \left\{ \mathcal{L}(h_{\theta}(\mathbf{X}), Y) + \lambda_0 \mathbb{E}_{(\mathbf{X}, P) \sim \mathcal{D}_0} (\hat{f}_{\omega_{f0}}(h_{\theta}(\mathbf{X})) \hat{g}_{\omega_{g0}}(P)) \right. \right. \\ \left. \left. + \lambda_1 \mathbb{E}_{(\mathbf{X}, P) \sim \mathcal{D}_1} (\hat{f}_{\omega_{f1}}(h_{\theta}(\mathbf{X})) \hat{g}_{\omega_{g1}}(P)) \right\} \right\} \end{aligned}$$

or, more generally when  $y \in \Omega_Y$  (e.g.  $\{0, 1, 2, 3+\}$ ), if  $k = \#\Omega_y$

# Adversarial Approach III

$$\operatorname{argmin}_{\theta} \left\{ \max_{\omega_{f0}, \omega_{g0}, \dots, \omega_{fk}, \omega_{gk}} \left\{ \mathcal{L}(h_{\theta}(\mathbf{X}), Y) + \sum_{y \in \Omega_y} \lambda_y \mathbb{E}_{(\mathbf{X}, P) \sim \mathcal{D}_y} (\hat{f}_{\omega_f} (h_{\theta}(\mathbf{X})) \hat{g}_{\omega_g} (P)) \right\} \right\}$$



## Dealing with high dimension I

- ▶ geographic / spatial information,  $\mathbf{X}_g$
- ▶ car type / make / model,  $\mathbf{X}_c$
- ▶ other classical ratemaking variables,  $\mathbf{X}$  (non protected)

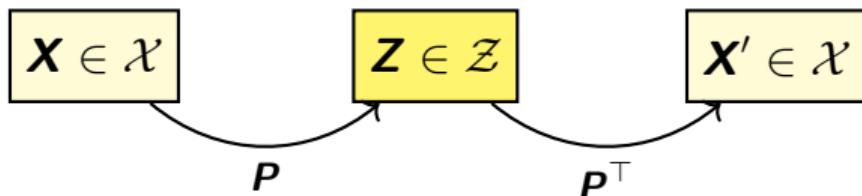
Some features can be in high dimension, natural solution would be PCA or autoencoders (see [Shi and Shi \(2021\)](#) about feature embedding in high dimension).

## Dealing with high dimension II

Principal Component Analysis (PCA)

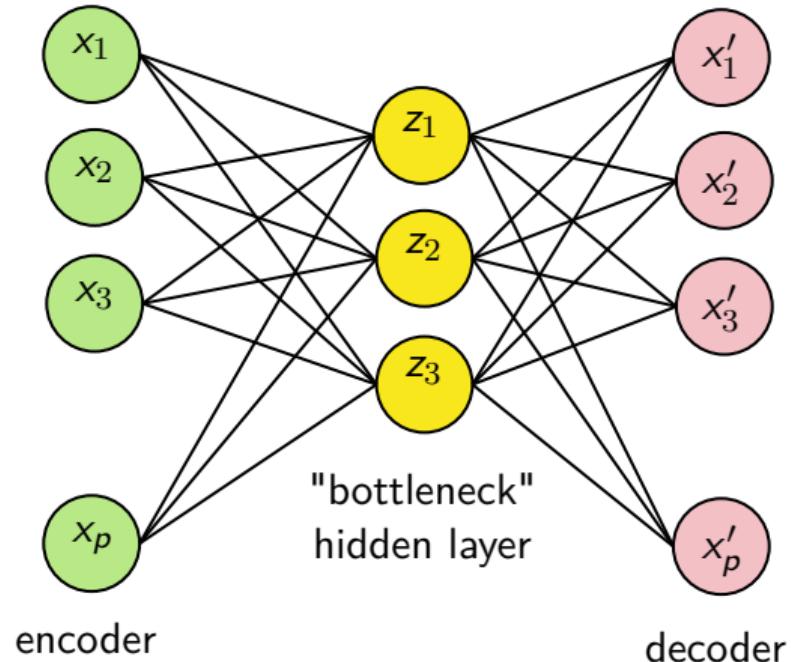
$$\min_{P \in \Pi} \{ \|X - P^T P X\|_F^2 \} \text{ s.t. } \text{rank}(P) = k$$

where  $\Pi$  is the set of projection matrices.



$$\min \|X - X'\|^2 = \min \|X - P^T P X\|^2$$

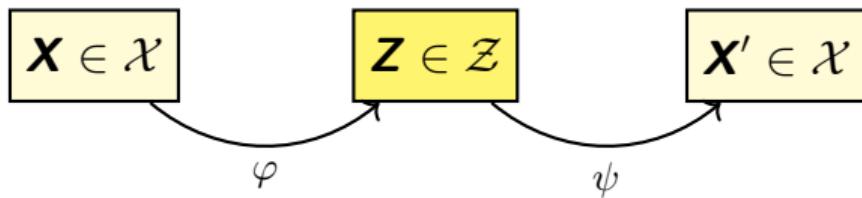
$$= \min \sum_{i=1}^n (\mathbf{P}^T \mathbf{P} x_i - x_i)^T (\mathbf{P}^T \mathbf{P} x_i - x_i)$$



# Dealing with high dimension III

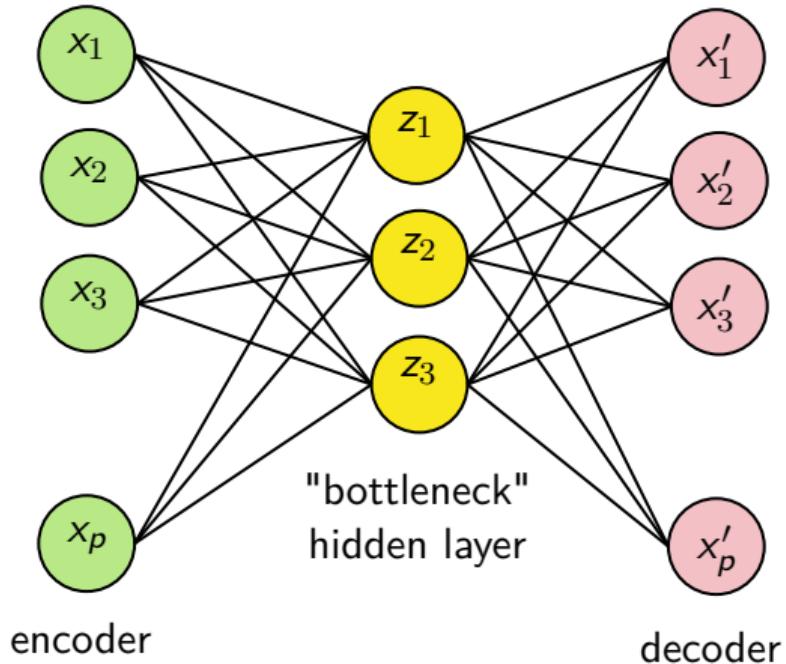
Autoencoder

$$\min_{\psi} \{ \| \mathbf{X} - \psi \circ \varphi \mathbf{X} \|_F^2 \}$$



$$\min \| \mathbf{X} - \mathbf{X}' \|^2 = \min \| \mathbf{X} - \psi \circ \varphi(\mathbf{X}) \|^2$$

$$\min \sum_{i=1}^n (\psi \circ \varphi(\mathbf{x}_i) - \mathbf{x}_i)^\top (\psi \circ \varphi(\mathbf{x}_i) - \mathbf{x}_i)$$

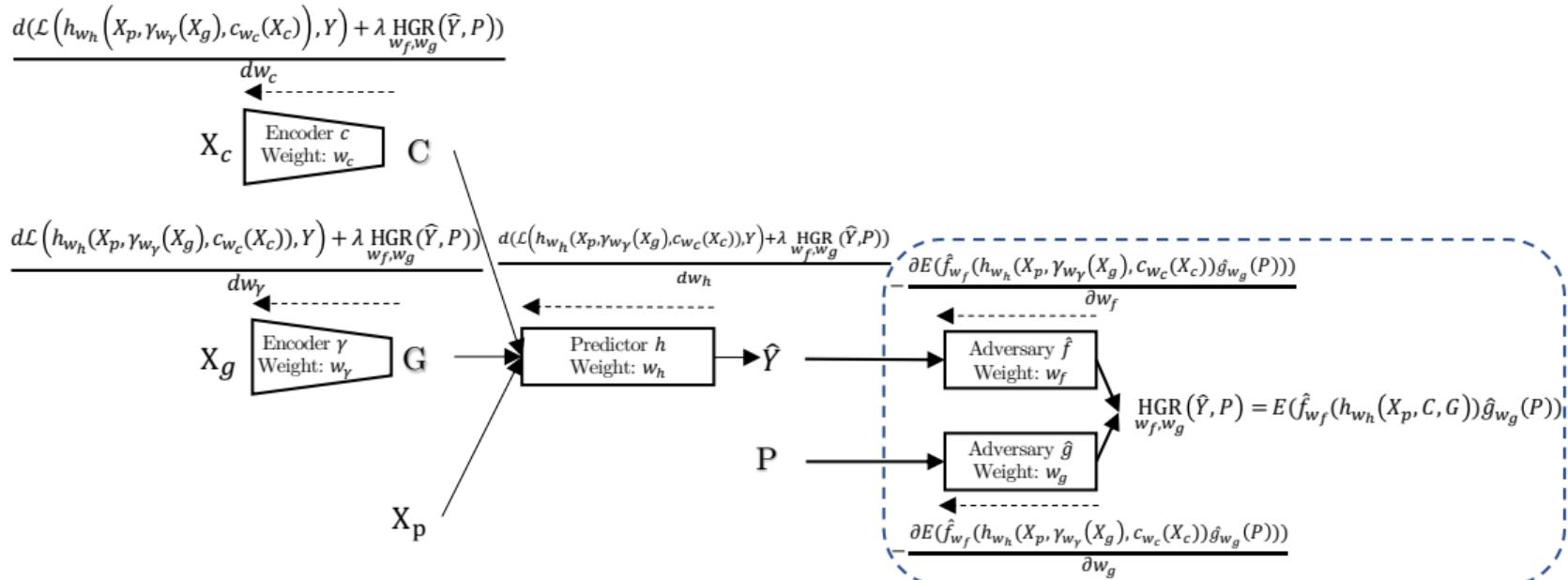


encoder

"bottleneck"  
hidden layer

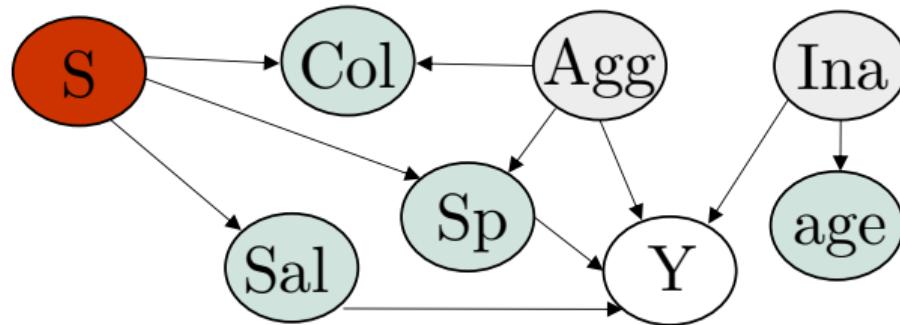
decoder

# Dealing with high dimension IV

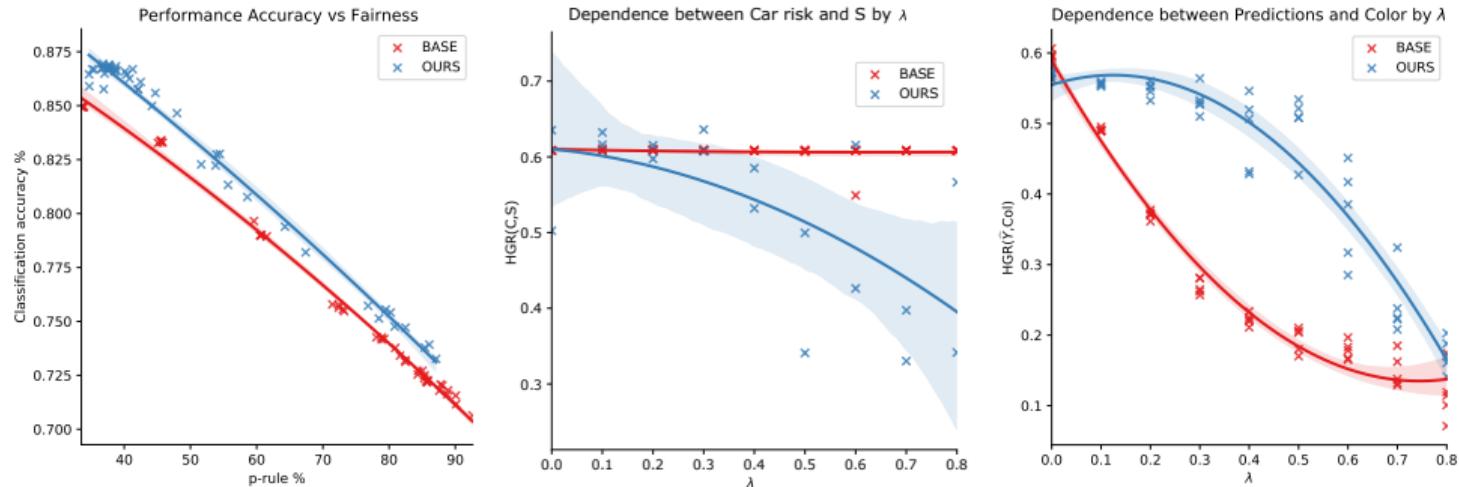


## Application on synthetic data I

- ▶ S: sensitive / protected (gender)
- ▶ Col: color of the car
- ▶ Sp: maximum speed of the car
- ▶ Sal: average salary of the policyholders area
- ▶ Age: age of the driver
- ▶ Ina: inattention
- ▶ Agg: aggressivity
- ▶ Y: total cost



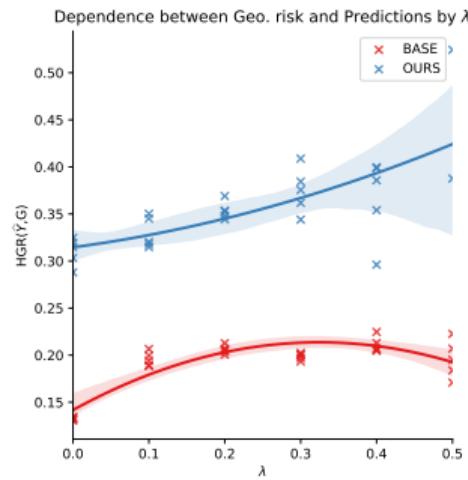
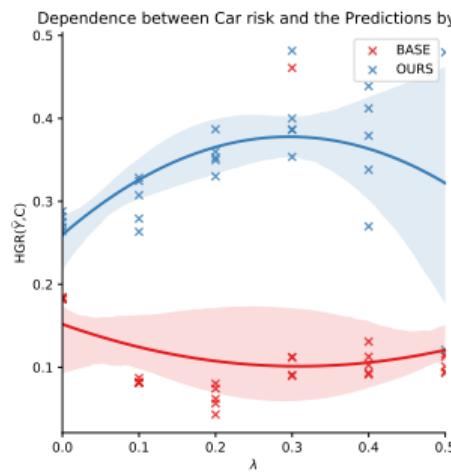
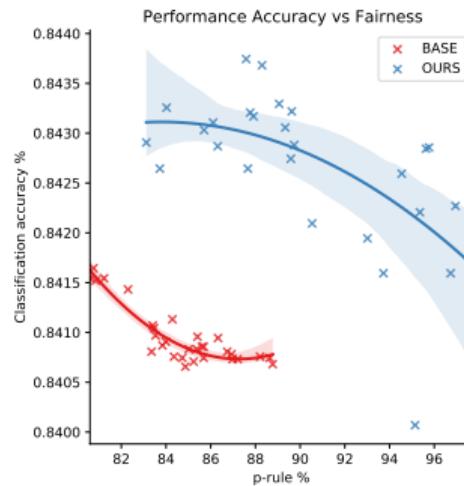
## Application on synthetic data II



- ▶  $\lambda$ : fairness penalty
- ▶  $p$ -rule:  $\min \left\{ \frac{\mathbb{P}(\hat{Y} = 1|P = 1)}{\mathbb{P}(\hat{Y} = 1|P = 0)}, \frac{\mathbb{P}(\hat{Y} = 1|P = 0)}{\mathbb{P}(\hat{Y} = 1|P = 1)} \right\}$

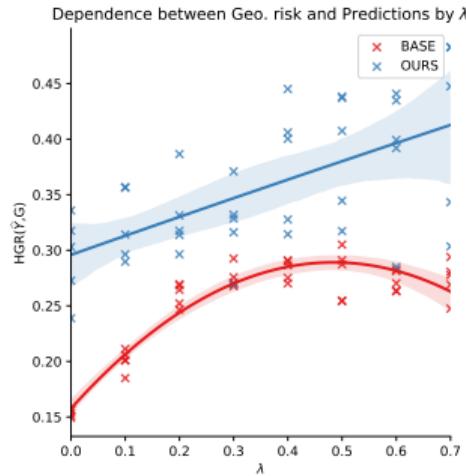
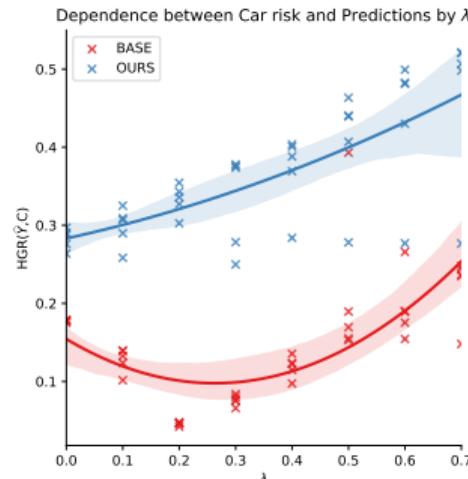
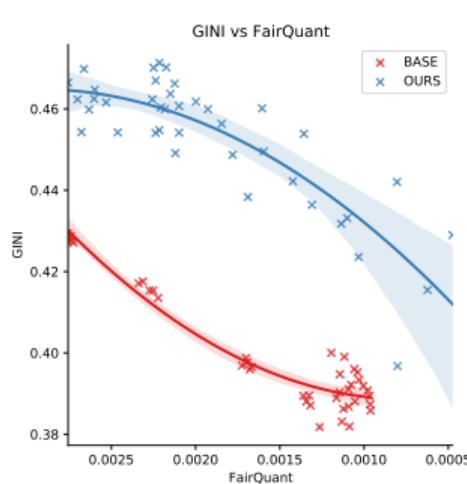
# Application on real data (pricing game 2015) I

$y \in \{0, 1\}$  (claim occurrence), and  $p$  is the (binary) gender



# Application on real data (pricing game 2015) II

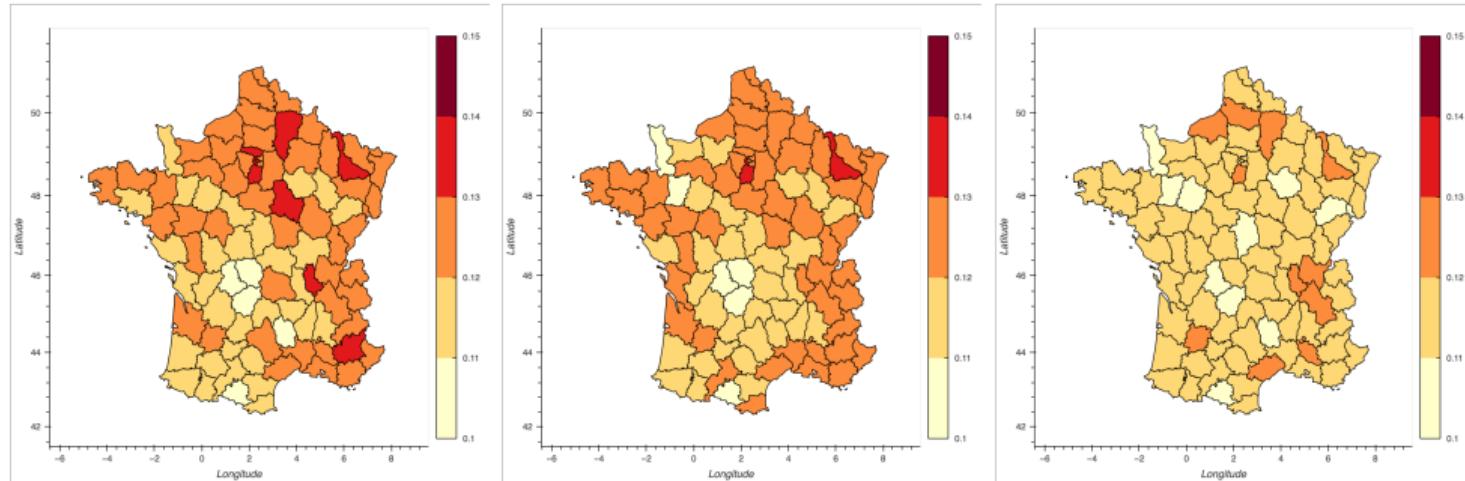
$y \in \{0, 1, 2+\}$  (claim frequency), and  $p$  is the (binary) gender



see Grari et al. (2022) for more examples (including the case where  $y \in \mathbb{R}^+$ )

# Application on real data (pricing game 2015) III

$y \in \{0, 1, 2+\}$  (claim frequency), and  $p$  is spatial information (redlining)



## From correlation to causality I

- ▶ “*classifying projection methods as using demographic/actuarial models or non-demographic/causal models*”  
Keilman (2003) and Hudson (2007)
- ▶ “*Article 5(2) allowed Member States to Permit proportionate differences in individuals premiums and benefits where the use of sex is a determining factor in the assessment of risk based on relevant and accurate actuarial and statistical data.*”  
Thiery and Van Schoubroeck (2006) and Schmeiser et al. (2014)
- ▶ “*Two judges on the Supreme Court dissented in the Zurich case. In their view, an insurer must not only prove a statistical correlation between a particular group and higher risk, but a causal connection*”  
Gomery et al. (2011)

# From correlation to causality II

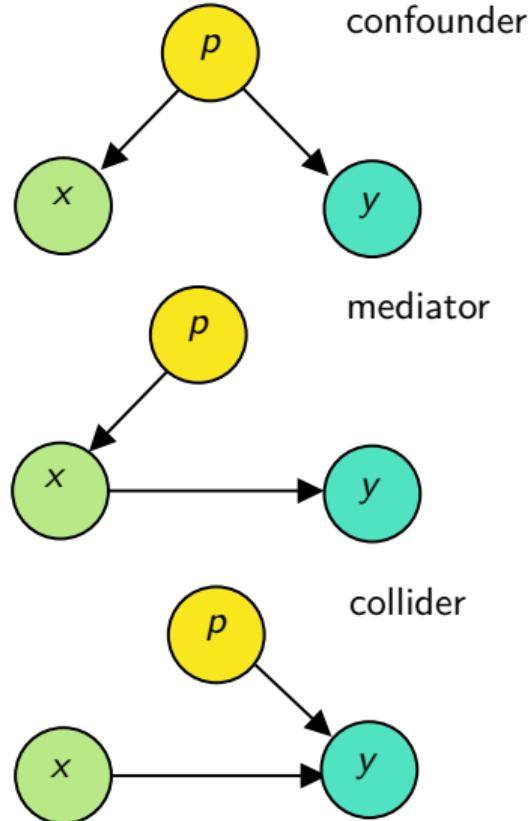
The screenshot shows the CBC News homepage with a search bar and sign-in options. Below is a news article from Calgary:

**Alberta man changes gender on government IDs for cheaper car insurance**

He says he saved almost \$1,100

Reid Southwick · CBC News · Posted: Jul 20, 2018 1:24 PM MT | Last Updated: July 26, 2018

- ▶ DAGs are important
- ▶ Looking for a **counterfactual**



## From correlation to causality III

Consider some distances  $D$  on  $\{0,1\} \times \{0,1\}$  or  $[0,1] \times [0,1]$ , and  $d$  on  $\mathbb{R}^p \times \mathbb{R}^p$ ,

**Lipschitz property**, Duivesteijn and Feelders (2008)

$$D(\hat{y}_i, \hat{y}_j) \text{ or } D(s_i, s_j) \leq d(\mathbf{x}_i, \mathbf{x}_j), \quad \forall i, j = 1, \dots, n.$$

**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 equity, i.e.

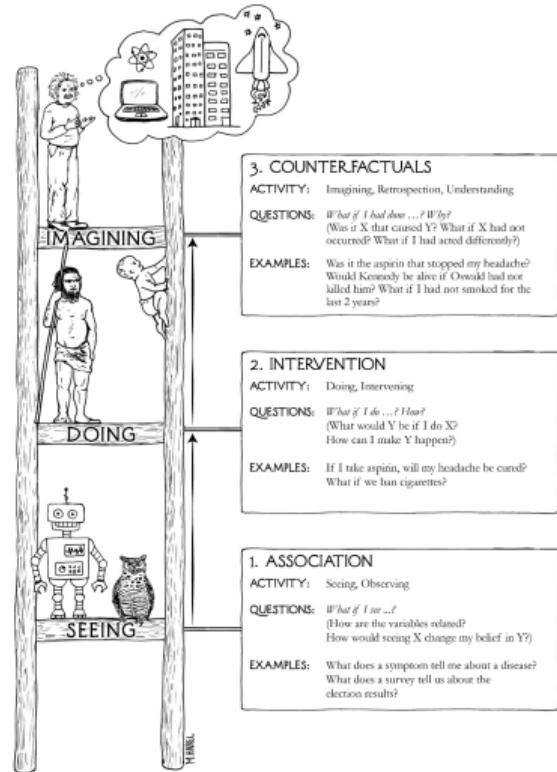
$$\mathbb{P}[Y_{P \leftarrow p}^* = y | \mathbf{X} = \mathbf{x}] = \mathbb{P}[Y_{P \leftarrow p'}^* = y | \mathbf{X} = \mathbf{x}], \quad \forall p', \mathbf{x}, y.$$

# From correlation to causality IV

- ▶ counterfactuals  
*(what if I had done...? )*
- ▶ intervention
- ▶ association  
*(what if I see...? )*

what would have happened if this person had had treatment 1 instead of treatment 0 ?

(picture Pearl & Mackenzie (2018))



# From correlation to causality V

Causal inference litterature,

- ▶  $t$  some binary treatment ( $t \in \{0, 1\}$ )
- ▶  $x$  some covariates
- ▶  $y$  denote the observed outcome,  $y_{i,T \leftarrow 1}^*$  and  $y_{i,T \leftarrow 0}^*$  the potential outcomes

	treatment	outcome		age	gender	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	F	160	56
2	0	109	?	109	28	F	156	54
3	1	162	162	?	53	M	190	87

There will be a significant impact of treatment  $t$  on  $y$  if  $y_{T \leftarrow 0}^* \neq y_{T \leftarrow 1}^*$  (see [Rubin \(1974\)](#), [Hernán and Robins \(2010\)](#) or [Imai \(2018\)](#)).

The causal effect for individual  $i$  is  $\tau_i = y_{i,T \leftarrow 1}^* - y_{i,T \leftarrow 0}^*$

## From correlation to causality VI

One can define the sample average treatment effect (SATE)

$$\text{SATE} = \frac{1}{n} \sum_{i=1}^n y_{i,T \leftarrow 1}^* - y_{i,T \leftarrow 0}^*$$

the average treatment effect (ATE)

$$\tau = \text{ATE} = \mathbb{E}[Y_{i,T \leftarrow 1}^* - Y_{i,T \leftarrow 0}^*]$$

and, for possibly heterogeneous effects, conditional average treatment effect (CATE)

$$\tau(\mathbf{x}) = \text{CATE}(\mathbf{x}) = \mathbb{E}[Y_{i,T \leftarrow 1}^* - Y_{i,T \leftarrow 0}^* | \mathbf{X} = \mathbf{x}]$$

## References I

- ,
- Agarwal, S. (2021). Trade-offs between fairness and interpretability in machine learning. In *IJCAI 2021 Workshop on AI for Social Good*.
- Avraham, R. (2017). Discrimination and insurance. In Lippert-Rasmussen, K., editor, *Handbook of the Ethics of Discrimination*, pages 335–347. Routledge.
- Barry, L. and Charpentier, A. (2022). The Fairness of Machine Learning in Insurance: New Rags for an Old Man? . *ArXiv*.
- Bechavod, Y. and Ligett, K. (2017). Penalizing unfairness in binary classification. *arXiv preprint arXiv:1707.00044*.
- 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.
- Breiman, L. and Friedman, J. H. (1985). Estimating optimal transformations for multiple regression and correlation. *Journal of the American statistical Association*, 80(391):580–598.
- Charpentier, A. (2022a). *Assurance: biais, discrimination et équité*. Institut Louis Bachelier.
- Charpentier, A. (2022b). *Insurance: biases, discrimination and fairness*. Institut Louis Bachelier.

## References II

- Cho, J., Hwang, G., and Suh, C. (2020). A fair classifier using mutual information. In *2020 IEEE International Symposium on Information Theory (ISIT)*, pages 2521–2526. IEEE.
- 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.
- 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.
- Gebelein, H. (1941). Das statistische problem der korrelation als variations- und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung. *ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 21(6):364–379.
- Gomery, S., Renault, O., and John, N. (2011). Gender neutral. *Canadian Underwriter*.
- Goodfellow, I., McDaniel, P., and Papernot, N. (2018). Making machine learning robust against adversarial inputs. *Communications of the ACM*, 61(7):56–66.

## References III

- Grari, V., Charpentier, A., Lamprier, S., and Detyniecki, M. (2022). A fair pricing model via adversarial learning. *ArXiv*, 2202.12008.
- Hardt, M., Price, E., and Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29:3315–3323.
- Hernán, M. A. and Robins, J. M. (2010). Causal inference.
- Hirschfeld, H. O. (1935). A connection between correlation and contingency. *Mathematical Proceedings of the Cambridge Philosophical Society*, 31(4):520524.
- Hudson, R. (2007). Mortality projections and unisex pricing of annuities in the uk. *Journal of Financial Regulation and Compliance*.
- Hume, D. (1739). *A Treatise of Human Nature*. Cambridge University Press Archive.
- Imai, K. (2018). *Quantitative social science: an introduction*. Princeton University Press.
- Kearns, M. and Roth, A. (2019). *The ethical algorithm: The science of socially aware algorithm design*. Oxford University Press.
- Keilman, N. (2003). Types of models for projecting mortality. *Perspectives on mortality forecasting*, 1:19–27.

## References IV

- 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.
- Kranzberg, M. (1986). Technology and history:" kranzberg's laws". *Technology and culture*, 27(3):544–560.
- 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.
- Lawrance, A. (1976). On conditional and partial correlation. *The American Statistician*, 30(3):146–149.
- Rényi, A. (1959). On measures of dependence. *Acta mathematica hungarica*, 10(3-4):441–451.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688.
- Schmeiser, H., Störmer, T., and Wagner, J. (2014). Unisex insurance pricing: consumers perception and market implications. *The Geneva Papers on Risk and Insurance-Issues and Practice*, 39(2):322–350.
- Shi, P. and Shi, K. (2021). Nonlife insurance risk classification using categorical embedding. *SSRN*, 3777526.
- Thiery, Y. and Van Schouwbroeck, C. (2006). Fairness and equality in insurance classification. *The Geneva Papers on Risk and Insurance-Issues and Practice*, 31(2):190–211.