

# Exposing indirect discrimination

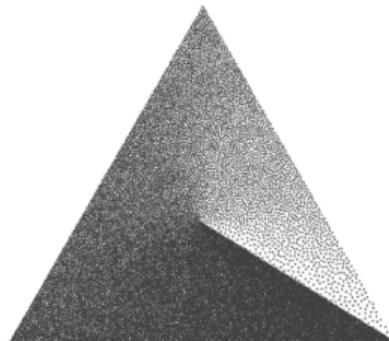
A scalable toolbox for fairness assessment of insurance rates

Presented by

**Olivier Côté**

Joint work with

**Prof. Marie-Pier Côté** (U. Laval) and  
**Prof. Arthur Charpentier** (UQÀM)



Séminaire étudiant UQAM

• Montréal

23 Juillet 2025

## Fairness and equity

1/53

The strategic goal of the Council of Europe in the field of anti-discrimination, diversity and inclusion is to ensure genuine equality and **full access to rights and opportunities for all members** of society.

2021 Report by the Secretary General of the Council of Europe entitled “State of democracy, human rights and the rule of law : A democratic renewal for Europe”

## Fairness and equity : no consensus

1/53

The strategic goal of the Council of Europe in the field of anti-discrimination, diversity and inclusion is to ensure genuine equality and **full access to rights and opportunities for all members** of society.

2021 Report by the Secretary General of the Council of Europe entitled “State of democracy, human rights and the rule of law : A democratic renewal for Europe”

The public release of these [Federal Equity Action] plans demonstrated immense public waste and **shameful discrimination**. That ends today. Americans deserve a government committed to **serving every person with equal dignity and respect** [...]

Executive order of The White House issued on January 20, 2025 entitled “Ending Radical And Wasteful Government DEI Programs And Preferencing”

## Differentiation at the core of insurance pricing

2/53

Insurance pricing relies on grouping policyholders by risk to set adequate premiums.

Modern predictive models excel at detecting statistical associations to differentiate between risks, but they can learn spurious correlations.

This raises concerns when sensitive variables may (intentionally or inadvertently) affect the **fairness of insurance pricing**.

# Notation

3/53

Age	Vehicle	Occupation	Gender	Religion	Credit	Claim
					800	
					700	
					650	
					435	

Europe

California

Ontario

Allowed variables  
 $X$

Prohibited variables  
 $D$   
(Collected)

Response  
 $Y$

## Some fair premiums

4/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project the allowed variables  $X$  in a space orthogonal to the prohibited  $D$  in a **pre-processing** step;

## Some fair premiums

4/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project the allowed variables  $X$  in a space orthogonal to the prohibited  $D$  in a **pre-processing** step;
- Lindholm et al. (2022) propose a **post-processing** of the premium in order to prevent  $D$  from indirectly influencing the calculation.

## Some fair premiums

4/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project the allowed variables  $X$  in a space orthogonal to the prohibited  $D$  in a **pre-processing** step;
- Lindholm et al. (2022) propose a **post-processing** of the premium in order to prevent  $D$  from indirectly influencing the calculation.
- 

These different premiums all aim to be *fair*. Again, no consensus.

# Fair premiums

5/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project allowed variables  $X$  orthogonally to prohibited  $D$  in **pre-processing**;

## Fair premiums

5/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project allowed variables  $X$  orthogonally to prohibited  $D$  in **pre-processing**;
- Lindholm et al. (2022) **post-process** premiums to limit  $D$ 's indirect influence.

## Fair premiums

5/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project allowed variables  $X$  orthogonally to prohibited  $D$  in **pre-processing**;
- Lindholm et al. (2022) **post-process** premiums to limit  $D$ 's indirect influence.
- Boucher and Pigeon (2024) find that **including telematics data** in  $X$  dilutes the significance of  $D$ .

Methods differ fundamentally but **all aim for fairness**.

# Fair premiums

5/53

Given a prohibited variable  $D$ ,

- Frees and Huang (2023) project allowed variables  $X$  orthogonally to prohibited  $D$  in **pre-processing**;
  - Lindholm et al. (2022) **post-process** premiums to limit  $D$ 's indirect influence.
  - Boucher and Pigeon (2024) find that **including telematics data** in  $X$  dilutes the significance of  $D$ .
- Methods differ fundamentally but **all aim for fairness**.
- Frees and Huang (2023) argue that **using variables related to  $D$**  constitutes indirect discrimination.
  - Lindholm et al. (2024) equate indirect and **proxy discrimination**, where formulas implicitly infer  $D$ .

Unclear fairness terms impede progress.

# Goals

6 / 53

- 1 Understand the behaviour of fair pricing methodologies.
  - A Fair price to pay : exploiting causal graphs for fairness in insurance, *Journal of Risk and Insurance*

# Goals

6 / 53

- 1 Understand the behaviour of fair pricing methodologies.
  - ▶ A Fair price to pay : exploiting causal graphs for fairness in insurance, *Journal of Risk and Insurance*
- 2 Quantify the **confounding bias** on real data.
  - ▶ A scalable toolbox for exposing indirect discrimination in insurance rates, *Forthcoming*

# Goals

6 / 53

- 1 Understand the behaviour of fair pricing methodologies.
  - ▶ A Fair price to pay : exploiting causal graphs for fairness in insurance, *Journal of Risk and Insurance*
- 2 Quantify the **confounding bias** on real data.
  - ▶ A scalable toolbox for exposing indirect discrimination in insurance rates, *Forthcoming*

# Goals

6 / 53

- 1 Understand the behaviour of fair pricing methodologies.
  - ▶ A Fair price to pay : exploiting causal graphs for fairness in insurance, *Journal of Risk and Insurance*
- 2 Quantify the **confounding bias** on real data.
  - ▶ A scalable toolbox for exposing indirect discrimination in insurance rates, *Forthcoming*



This joint work with Marie-Pier Côté and Arthur Charpentier is supported by a Canadian insurance company.

# Exploiting causal graphs for fairness

---

## 1 Exploiting causal graphs for fairness

- Sources of bias
- Detailing the causal graph
- Defining direct and indirect discrimination
- Two examples of fair methodologies
- Families of fair premiums
- Familles précise, inconsciente et consciente

## 2 A scalable toolbox for exposing indirect discrimination

# Causal graphs for insurance

7 / 53

Actuarial predictive models **are based on proxies** rather than risk factors  
(Embrechts and Wüthrich, 2022).

# Causal graphs for insurance

Actuarial predictive models **are based on proxies** rather than risk factors  
(Embrechts and Wüthrich, 2022).

But, **causality is key** for fairness (Makhlof et al., 2024) :

## Causal graphs for insurance

Actuarial predictive models **are based on proxies** rather than risk factors (Embrechts and Wüthrich, 2022).

But, **causality is key** for fairness (Makhlof et al., 2024) :

- Araiza Iturria et al. (2024) relate discrimination-free premiums to causal inference techniques.
- Fernandes Machado et al. (2024) explore fairness through counterfactuals.

## Causal graphs for insurance

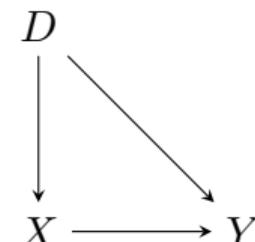
7 / 53

Actuarial predictive models **are based on proxies** rather than risk factors (Embrechts and Wüthrich, 2022).

But, **causality is key** for fairness (Makhlof et al., 2024) :

- Araiza Iturria et al. (2024) relate discrimination-free premiums to causal inference techniques.
- Fernandes Machado et al. (2024) explore fairness through counterfactuals.

Both are **supported by causal graphs** (DAG).



## Motivation

8 / 53

- The causal graph that relates  $X$ ,  $D$  and  $Y$ 
  - ▶ has properties that are **not fully exploited**.
  - ▶ is too simple,
  - ▶ does not relate the premium (score) with the other variables,

## Motivation

- The causal graph that relates  $X$ ,  $D$  and  $Y$ 
  - ▶ has properties that are **not fully exploited**.
  - ▶ is too simple,
  - ▶ does not relate the premium (score) with the other variables,
- There is **confusion** around fairness terms.
  - ▶ Indirect discrimination lacks a consensual definition.
  - ▶ The interplay between fairness criteria and indirect discrimination is ambiguous.

## Motivation

- The causal graph that relates  $X$ ,  $D$  and  $Y$ 
  - ▶ has properties that are **not fully exploited**.
  - ▶ is too simple,
  - ▶ does not relate the premium (score) with the other variables,
- There is **confusion** around fairness terms.
  - ▶ Indirect discrimination lacks a consensual definition.
  - ▶ The interplay between fairness criteria and indirect discrimination is ambiguous.
- Many fairness methodologies exist and do not match. The *Pre-In-Post* categorization **does not clarify the expected fairness properties**.

# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$		
Conditioned	$A \leftarrow [C] \rightarrow B$		

# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$		
Conditioned		$A \leftarrow C \rightarrow B$	



# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$ 	$A \rightarrow C \rightarrow B$	
Conditioned	$A \leftarrow [C] \rightarrow B$	$A \rightarrow [C] \rightarrow B$	

# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$	$A \rightarrow C \rightarrow B$	
Conditioned	$A \leftarrow [C] \rightarrow B$	$A \rightarrow [C] \rightarrow B$	

# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$	$A \rightarrow C \rightarrow B$	$A \rightarrow C \leftarrow B$
Conditioned	$A \leftarrow [C] \rightarrow B$	$A \rightarrow [C] \rightarrow B$	$A \rightarrow [C] \leftarrow B$

# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$	$A \rightarrow C \rightarrow B$	$A \rightarrow C \leftarrow B$
Conditioned	$A \leftarrow [C] \rightarrow B$	$A \rightarrow [C] \rightarrow B$	$A \rightarrow [C] \leftarrow B$

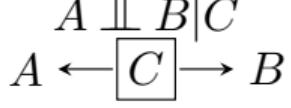
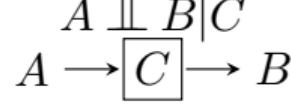
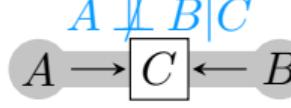
# Dependence implications of a causal graph

9 / 53

Status of C	Causal type of $C$		
	Confounder	Mediator	Collider
Unconditioned	$A \leftarrow C \rightarrow B$	$A \rightarrow C \rightarrow B$	$A \rightarrow C \leftarrow B$
Conditioned	$A \leftarrow [C] \rightarrow B$	$A \rightarrow [C] \rightarrow B$	$A \rightarrow [C] \leftarrow B$

# Dependence implications of a causal graph

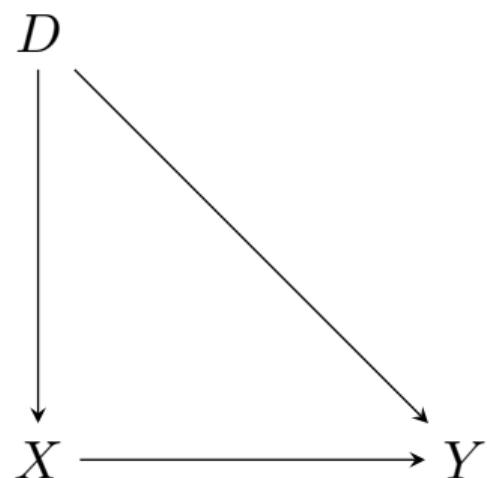
9 / 53

	Causal type of $C$		
	Confounder	Mediator	Collider
Status of $C$	$A \perp\!\!\!\perp B$ 	$A \perp\!\!\!\perp B$ 	$A \perp\!\!\!\perp B$ $A \rightarrow C \leftarrow B$
	$A \perp\!\!\!\perp B   C$ 	$A \perp\!\!\!\perp B   C$ 	$A \perp\!\!\!\perp B   C$ 

Causal relationships have **dependence implications**.

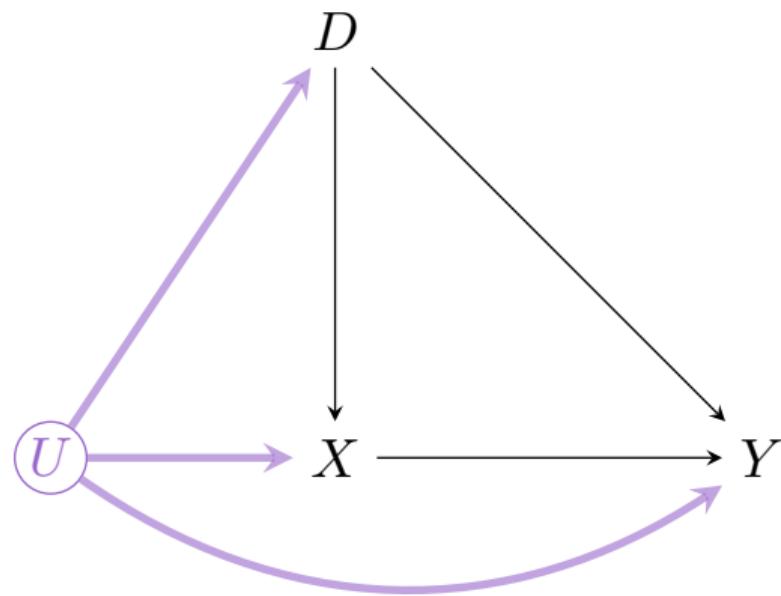
# Potential bias : unobserved confounders

10/53



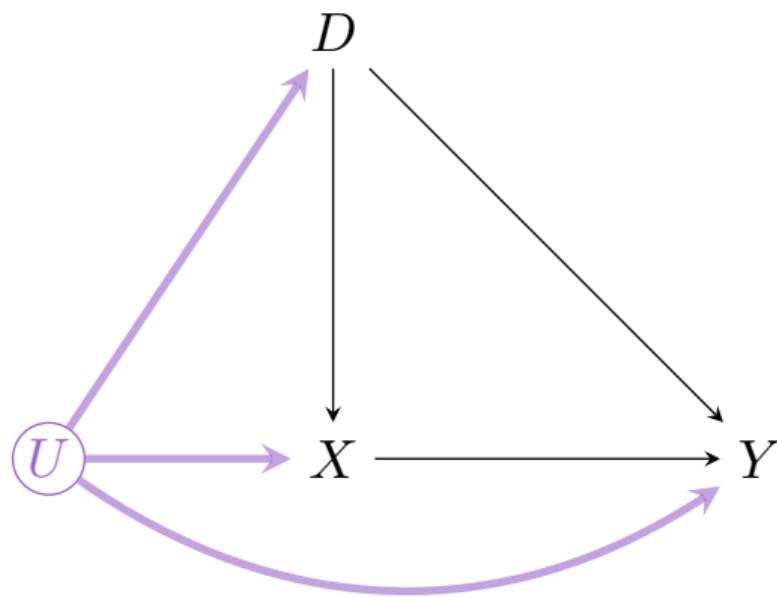
## Potential bias : unobserved confounders

10/53



## Potential bias : unobserved confounders

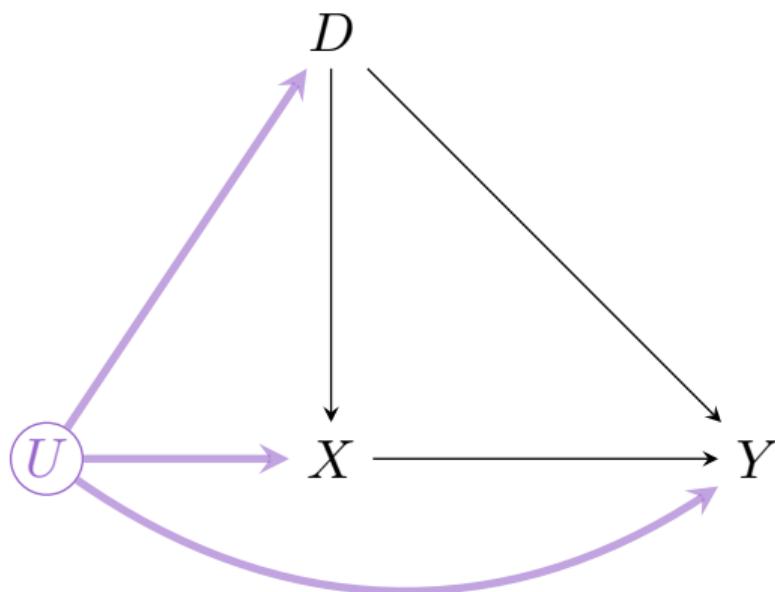
10 / 53



- **Omitted variable bias** arises because  $D \leftarrow U \rightarrow X$  is indistinguishable from  $D \rightarrow X$ .

## Potential bias : unobserved confounders

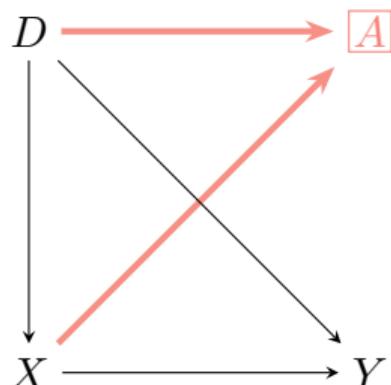
10 / 53



- **Omitted variable bias** arises because  $D \leftarrow U \rightarrow X$  is indistinguishable from  $D \rightarrow X$ .
- Detecting unobserved confounders relies on expert judgment, making these confounders **unverifiable**.

# Portfolio selection bias

11/53



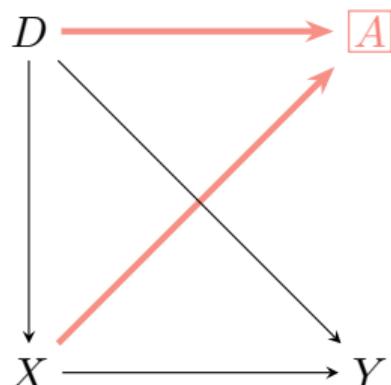
- The insurer's portfolio is conditional on  $A = 1$ , where  **$A$  is the insured's inclusion**;  
1 if included, 0 otherwise.

Selection bias in insurance : why  
portfolio-specific fairness

fails to extend market-wide, *available on SSRN*

# Portfolio selection bias

11/53



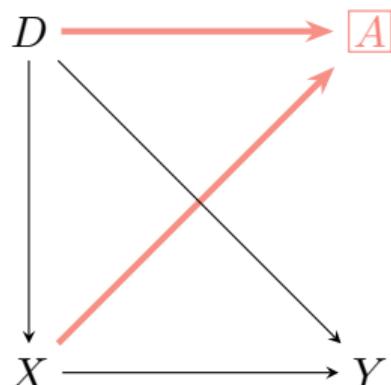
- The insurer's portfolio is conditional on  $A = 1$ , where  **$A$  is the insured's inclusion**; 1 if included, 0 otherwise.
- The path  $D \rightarrow A \leftarrow X$  represents portfolio-specific **selection bias**.

Selection bias in insurance : why  
portfolio-specific fairness

fails to extend market-wide, *available on SSRN*

# Portfolio selection bias

11 / 53



Selection bias in insurance : why  
portfolio-specific fairness

fails to extend market-wide, *available on SSRN*

- The insurer's portfolio is conditional on  $A = 1$ , where  **$A$  is the insured's inclusion**;  
1 if included, 0 otherwise.
- The path  $D \rightarrow A \leftarrow X$  represents portfolio-specific **selection bias**.
- Should insurers focus on **portfolio or market-wide** fairness ?

## Separation into components according to $D$

12/53

Fairness studies are always centered on  $D$ .

## Separation into components according to $D$

12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the part generated  
by  $D$  and the residual independent of  $D$

## Separation into components according to $D$

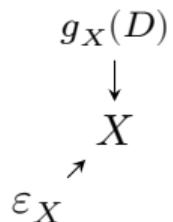
12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the part generated by  $D$  and the residual independent of  $D$ , i.e.,

$$X = g_X(D) + \varepsilon_X$$

where  $\varepsilon_X \perp\!\!\!\perp (D, g_X(D))$ ,



## Separation into components according to $D$

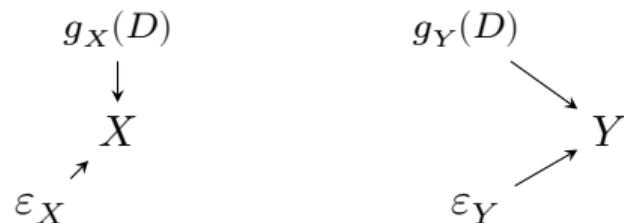
12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the part generated by  $D$  and the residual independent of  $D$ , i.e.,

$$X = g_X(D) + \varepsilon_X \quad \text{and} \quad Y = g_Y(D) + \varepsilon_Y,$$

where  $\varepsilon_X \perp\!\!\!\perp (D, g_X(D))$ ,  $\varepsilon_Y \perp\!\!\!\perp (D, g_Y(D))$



## Separation into components according to $D$

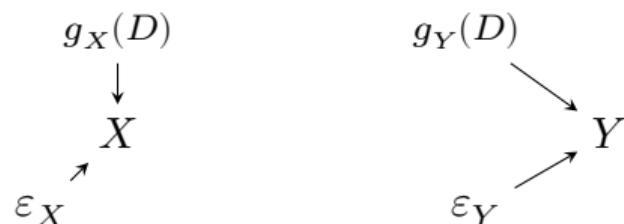
12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the part generated by  $D$  and the residual independent of  $D$ , i.e.,

$$X = g_X(D) + \varepsilon_X \quad \text{and} \quad Y = g_Y(D) + \varepsilon_Y,$$

where  $\varepsilon_X \perp\!\!\!\perp (D, g_X(D))$ ,  $\varepsilon_Y \perp\!\!\!\perp (D, g_Y(D))$ , and  $\varepsilon_X, \varepsilon_Y$  should be “large”.



## Separation into components according to $D$

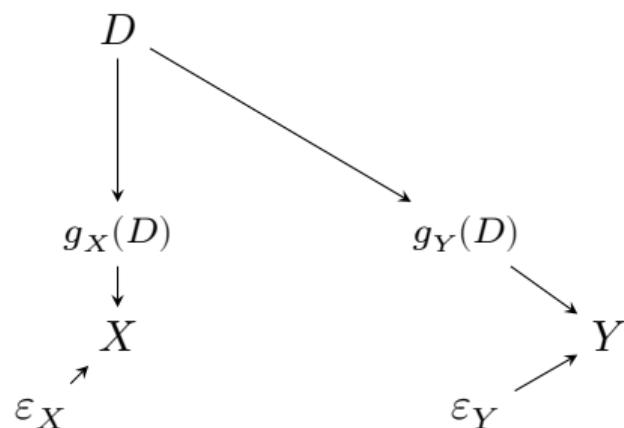
12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the **part generated by  $D$**  and the **residual independent of  $D$** , i.e.,

$$X = g_X(D) + \varepsilon_X \quad \text{and} \quad Y = g_Y(D) + \varepsilon_Y,$$

where  $\varepsilon_X \perp\!\!\!\perp (D, g_X(D))$ ,  $\varepsilon_Y \perp\!\!\!\perp (D, g_Y(D))$ , and  $\varepsilon_X, \varepsilon_Y$  should be “large”.



## Separation into components according to $D$

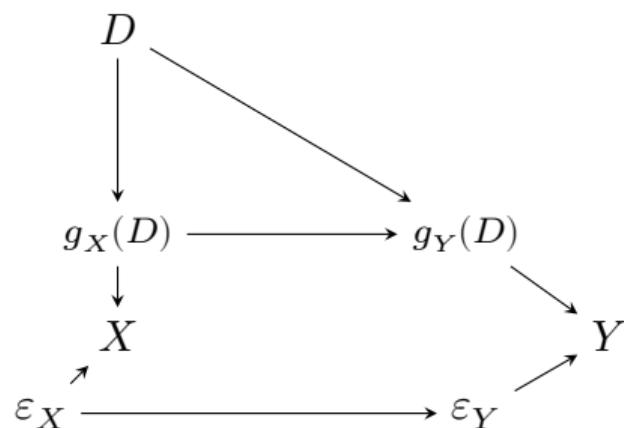
12/53

Fairness studies are always centered on  $D$ .

We decompose  $X$  and  $Y$  into : the **part generated by  $D$**  and the **residual independent of  $D$** , i.e.,

$$X = g_X(D) + \varepsilon_X \quad \text{and} \quad Y = g_Y(D) + \varepsilon_Y,$$

where  $\varepsilon_X \perp\!\!\!\perp (D, g_X(D))$ ,  $\varepsilon_Y \perp\!\!\!\perp (D, g_Y(D))$ , and  $\varepsilon_X, \varepsilon_Y$  should be “large”.



# What is a premium?

13/53

Discrimination is defined relative to a score, **premium**, or formula.

# What is a premium?

13/53

**Discrimination is defined relative to a score, premium, or formula.**

## Definition (premium)

Let a premium  $\mu$  be a deterministic function  $(x, d) \mapsto \mu(x, d)$

Examples of premiums include estimates of  $\mathbb{E}[Y|X, D]$  and [final market prices](#).

# What is a premium?

13/53

**Discrimination is defined relative to a score, premium, or formula.**

## Definition (premium)

Let a premium  $\mu$  be a deterministic function  $(x, d) \mapsto \mu(x, d)$

Examples of premiums include estimates of  $\mathbb{E}[Y|X, D]$  and final market prices.

A premium can be placed in a causal graph (Rosenbaum and Rubin, 1983) if it encapsulates a causal relationships.

$$A \longrightarrow B$$

**(a) Causal relationship**

# What is a premium?

13/53

**Discrimination is defined relative to a score, premium, or formula.**

## Definition (premium)

Let a premium  $\mu$  be a deterministic function  $(x, d) \mapsto \mu(x, d)$

Examples of premiums include estimates of  $\mathbb{E}[Y|X, D]$  and final market prices.

A premium can be placed in a causal graph (Rosenbaum and Rubin, 1983) if it encapsulates a causal relationships.

$$A \longrightarrow B$$

(a) Causal relationship

$$A \longrightarrow \mu(A) \longrightarrow B$$

(b) **Sufficient** premium  $\mu(A)$

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

According to Fredman (2022), recurring elements from legal definitions are :

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

According to Fredman (2022), recurring elements from legal definitions are :

- 1 “equal treatment”, i.e., no direct discrimination;

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

According to Fredman (2022), recurring elements from legal definitions are :

- 1 “equal treatment”, i.e., no direct discrimination;
- 2 “treatment must have disparate results”;

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

According to Fredman (2022), recurring elements from legal definitions are :

- 1 “equal treatment”, i.e., no direct discrimination;
- 2 “treatment must have disparate results”;
- 3 “the disparate impact can be justified if [...].”

# Indirect discrimination from the legal perspective

14/53

“A situation in which an **apparently neutral provision**, criterion or practice would put persons of a racial or ethnic origin at a **particular disadvantage compared with other persons**, unless that provision, criterion or practice is **objectively justified** by a legitimate aim and the means of achieving that aim are appropriate and necessary.”

—European Union (2000)

According to Fredman (2022), recurring elements from legal definitions are :

- 1 “**equal treatment**”, i.e., no direct discrimination;
- 2 “treatment must have **disparate results**”;
- 3 “the disparate impact can be justified if [...].”

We formalize for elements 1 and 2.

# Defining indirect discrimination

15/53

## Definition (indirect discrimination)

A premium  $\mu(X, D)$  **discriminates indirectly** on  $D$  if :

- 1 it does not discriminate directly on  $D$ ;
- 2 it exhibits at least one **source of disparate impact**.

# Defining indirect discrimination

15/53

## Definition (indirect discrimination)

A premium  $\mu(X, D)$  **discriminates indirectly** on  $D$  if :

- 1 it does not discriminate directly on  $D$ ;
- 2 it exhibits at least one **source of disparate impact**.

# Defining indirect discrimination

15/53

## Definition (indirect discrimination)

A premium  $\mu(X, D)$  **discriminates indirectly** on  $D$  if :

- 1 it does not discriminate directly on  $D$ ;
  - 2 it exhibits at least one **source of disparate impact**.
- 
- A **disparate impact** implies a decision  $\mu(X, D)$  affecting protected groups  $D$ .

# Defining indirect discrimination

15/53

## Definition (indirect discrimination)

A premium  $\mu(X, D)$  **discriminates indirectly** on  $D$  if :

- 1 it does not discriminate directly on  $D$ ;
  - 2 it exhibits at least one **source of disparate impact**.
- 
- A **disparate impact** implies a decision  $\mu(X, D)$  affecting protected groups  $D$ .
    - ▶ We formalise these **sources** as the paths between  $\mu(X, D)$  and  $D$ .

# Defining indirect discrimination

15/53

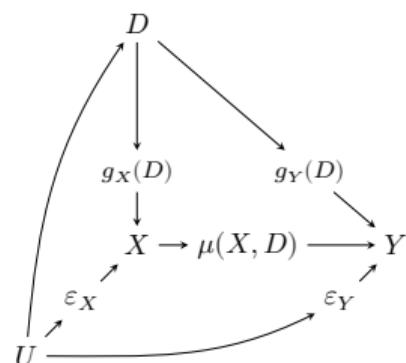
## Definition (indirect discrimination)

A premium  $\mu(X, D)$  **discriminates indirectly** on  $D$  if :

- 1 it does not discriminate directly on  $D$ ;
  - 2 it exhibits at least one **source of disparate impact**.
- 
- A **disparate impact** implies a decision  $\mu(X, D)$  affecting protected groups  $D$ .
    - ▶ We formalise these **sources** as the paths between  $\mu(X, D)$  and  $D$ .
  - We assume there is no direct discrimination.

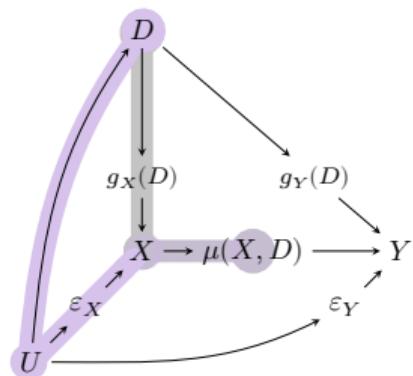
# Sources of disparate impact

16 / 53



# Sources of disparate impact

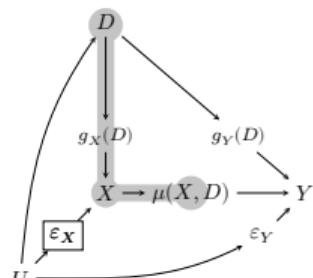
16 / 53



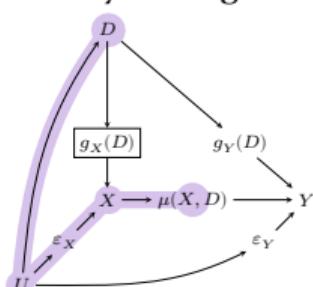
Disp. Impact through  $X$

# Sources of disparate impact

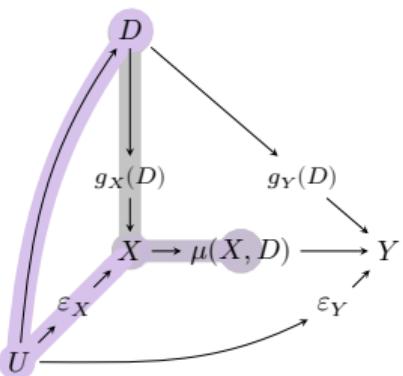
16 / 53



Purely through  $X$



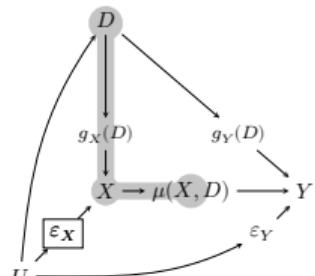
Through an conf.  
of  $D \rightarrow X$



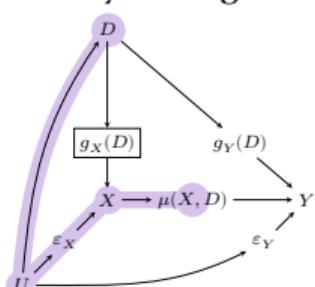
Disp. Impact through  $X$

# Sources of disparate impact

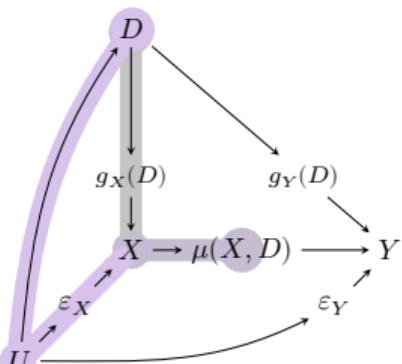
16 / 53



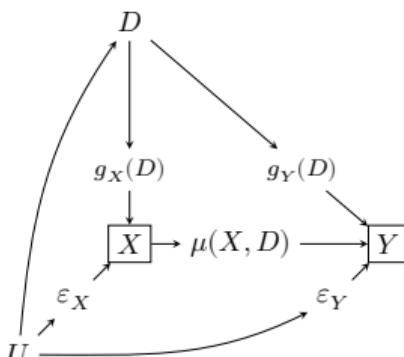
Purely through  $X$



Through an conf.  
of  $D \rightarrow X$

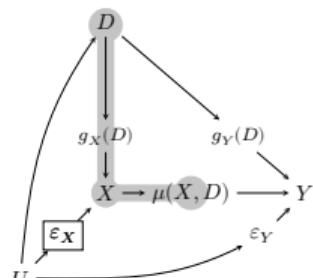
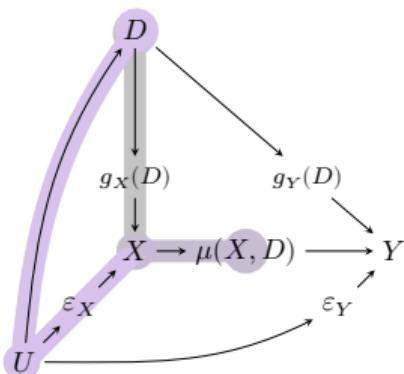
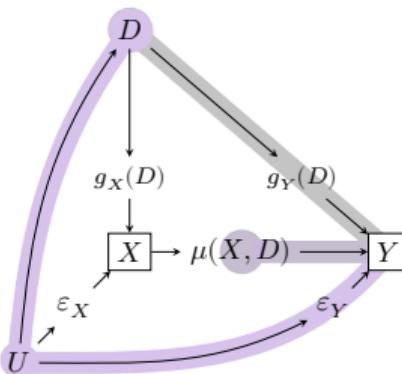


Disp. Impact through  $X$



# Sources of disparate impact

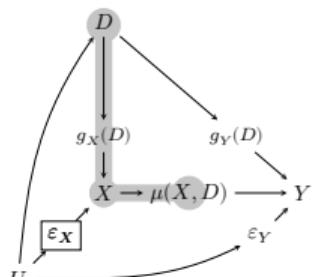
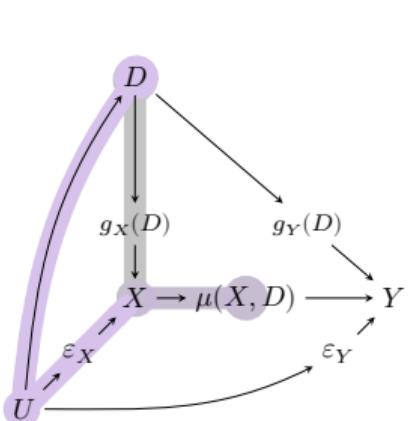
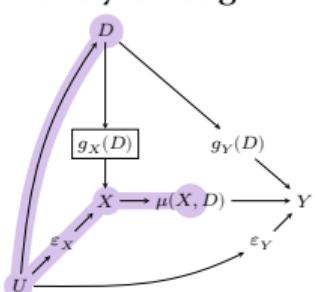
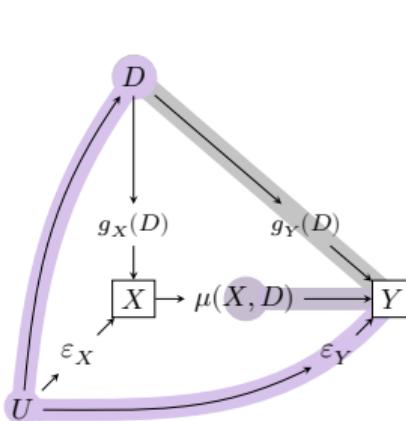
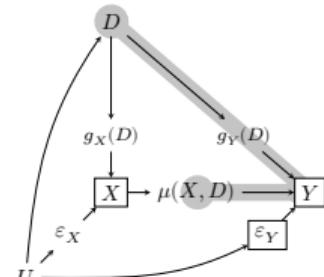
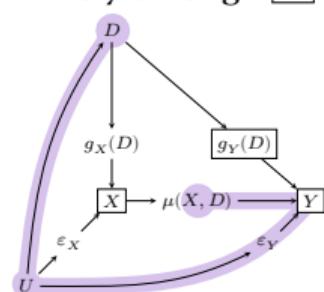
16 / 53

Purely through  $X$ Disp. Impact through  $X$ Disp. Impact through  $Y$ 

Through an conf.  
of  $D \rightarrow X$

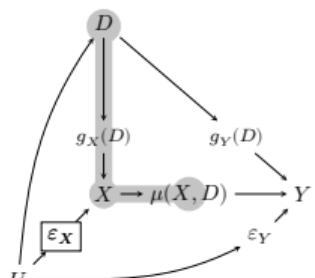
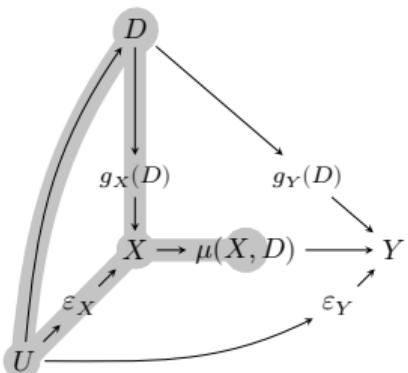
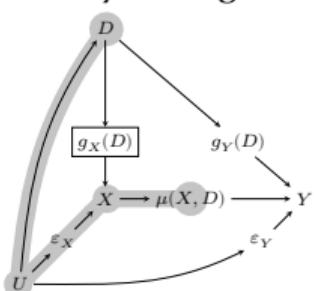
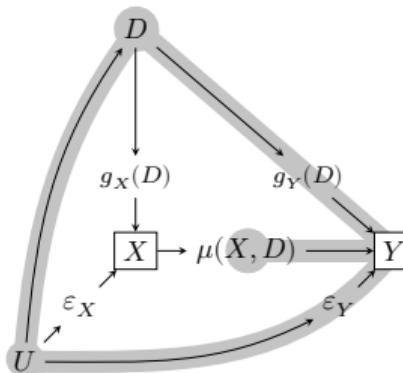
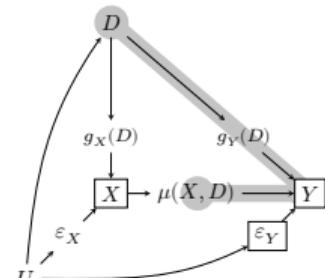
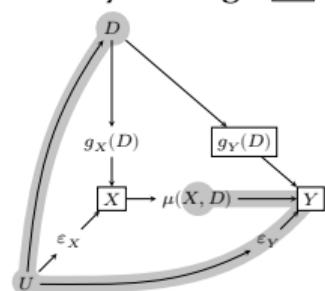
# Sources of disparate impact

16 / 53

Purely through  $X$ Disp. Impact through  $X$ Through an 🕵️ conf.  
of  $D \rightarrow X$ Disp. Impact through  $Y$ Purely through  $Y$ Through a 🕵️ conf.  
of  $D \rightarrow Y$

# Sources of disparate impact

16 / 53

Purely through  $X$ Disp. Impact through  $X$ Through an conf.  
of  $D \rightarrow X$ Mitigation leads to  
**Solidarity**  
(Levelled premiums)Mitigation leads to  
**Actuarial Fairness**  
(Levelled profits)Purely through  $\boxed{Y}$ Through a conf.  
of  $D \rightarrow Y$

# Discrimination is inevitable

17 / 53

**Proposition.** Unless  $D \perp\!\!\!\perp Y$  and  $D \perp\!\!\!\perp Y|X$ , all premiums disparately impact  $D$ .

## Discrimination is inevitable

17 / 53

**Proposition.** Unless  $D \perp\!\!\!\perp Y$  and  $D \perp\!\!\!\perp Y|X$ , all premiums disparately impact  $D$ .

### Corollary.

Unless  $D \perp\!\!\!\perp Y$  and  $D \perp\!\!\!\perp Y|X$ , all premiums **must discriminate** either directly or indirectly on  $D$ .

## Discrimination is inevitable

17 / 53

**Proposition.** Unless  $D \perp\!\!\!\perp Y$  and  $D \perp\!\!\!\perp Y|X$ , all premiums disparately impact  $D$ .

### Corollary.

Unless  $D \perp\!\!\!\perp Y$  and  $D \perp\!\!\!\perp Y|X$ , all premiums **must discriminate** either directly or indirectly on  $D$ .

The **impossibility** theorem by Kleinberg et al. (2016) **extends to discrimination**.

## Examples from the literature

**18**/53

Premium of Lindholm et al. (2022)

Premium of Frees and Huang (2023)

## Examples from the literature

18 / 53

### Premium of Lindholm et al. (2022)

The premium is

$$\pi^L(X) = \mathbb{E}_D\{\mathbb{E}(Y|X, D)\}$$

### Premium of Frees and Huang (2023)

For  $X^* = (\mathbf{I} - H_D)X$ , the premium is

$$\pi^F(X, D) = X^\top(\mathbf{I} - H_D)(X^{*\top}X^*)^{-1}X^{*\top}Y.$$

## Examples from the literature

18 / 53

### Premium of Lindholm et al. (2022)

The premium is

$$\pi^L(X) = \mathbb{E}_D\{\mathbb{E}(Y|X, D)\}$$

It is a **causal inference** technique.

### Premium of Frees and Huang (2023)

For  $X^* = (\mathbf{I} - H_D)X$ , the premium is

$$\pi^F(X, D) = X^\top(\mathbf{I} - H_D)(X^{*\top}X^*)^{-1}X^{*\top}Y.$$

$X^*$  is an **estimator of  $\varepsilon_X$** .

## Examples from the literature

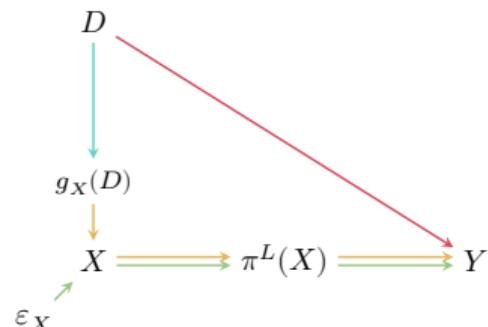
**18**/<sub>53</sub>

### Premium of Lindholm et al. (2022)

The premium is

$$\pi^L(X) = \mathbb{E}_D\{\mathbb{E}(Y|X, D)\}$$

It is a **causal inference** technique.

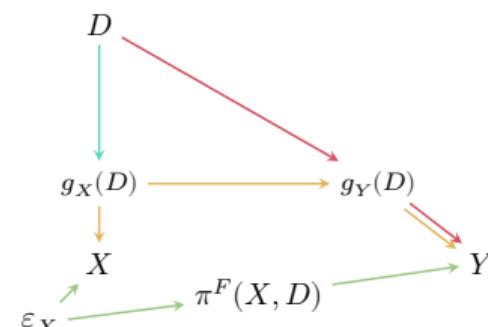


### Premium of Frees and Huang (2023)

For  $X^* = (\mathbf{I} - H_D)X$ , the premium is

$$\pi^F(X, D) = X^\top(\mathbf{I} - H_D)(X^{*\top}X^*)^{-1}X^{*\top}Y.$$

$X^*$  is an **estimator of  $\varepsilon_X$** .



# Discrimination for the premium $\pi^L(X)$

19 / 53

Data generating mechanism

	$D$ $X \longrightarrow Y$	$D$ $\downarrow$ $X \longrightarrow Y$	$D$ $\searrow$ $X \longrightarrow Y$	$D$ $\downarrow$ $X \longrightarrow Y$	$D$ $\nearrow$ $\downarrow$ $\searrow$ $U$ $X \longrightarrow Y$
Discr.	∅	∅	∅	∅	∅
	Direct	Indirect	Indirect	Indirect	Indirect
	∅	∅	∅	∅	∅

# Discrimination for the premium $\pi^L(X)$

19 / 53

Data generating mechanism

		D	D	D	D	D
		$X \longrightarrow Y$	$X \downarrow D \longrightarrow Y$	$X \xrightarrow{D} Y$	$X \downarrow D \xrightarrow{D} Y$	$X \xrightarrow{D} Y$
Discr.		∅	∅	∅	∅	∅
	Indirect	∅	∅	∅	∅	∅
X	Purely Conf.	∅	∅	∅	∅	∅
Y	Purely Conf.	∅	∅	∅	∅	∅

# Discrimination for the premium $\pi^F(X, D)$

20/53

Data generating mechanism

		D	D	D	D	D
		$X \longrightarrow Y$	$\downarrow X \longrightarrow Y$	$X \longrightarrow Y$	$\downarrow X \longrightarrow Y$	$\downarrow X \longrightarrow Y$
Discr.	Direct	∅	∅	∅	∅	∅
	Indirect	∅	∅	∅	∅	∅
X	Purely	∅	∅	∅	∅	∅
	👻 Conf.	∅	∅	∅	∅	∅
Y	Purely	∅	∅	∅	∅	∅
	👻 Conf.	∅	∅	∅	∅	∅

The table illustrates five data generating mechanisms for the premium  $\pi^F(X, D)$ . The columns represent different causal structures between the discriminant variable  $D$  and the outcome  $Y$ , and the rows represent different types of discrimination and their impact on the premium calculation.

- Discr. (Discrimination) Rows:**
  - Direct:** Shows no discrimination in all cases.
  - Indirect:** Shows no discrimination in all cases.
- X (Explanatory Variable) Columns:**
  - Purely:** Shows no discrimination in all cases.
  - 👻 Conf. (Confounding):** Shows discrimination in all cases.
- Y (Outcome) Columns:**
  - Purely:** Shows no discrimination in all cases.
  - 👻 Conf.:** Shows discrimination in all cases.

The causal structures are as follows:

- Column 1:  $D \perp\!\!\!\perp X \longrightarrow Y$
- Column 2:  $D \perp\!\!\!\perp X \downarrow X \longrightarrow Y$
- Column 3:  $D \not\perp\!\!\!\perp X \longrightarrow Y$
- Column 4:  $D \perp\!\!\!\perp X \downarrow X \longrightarrow Y$
- Column 5:  $D \not\perp\!\!\!\perp X \not\perp\!\!\!\perp Y$  (with a latent variable  $U$  influencing both  $D$  and  $Y$ )

# Best-estimate

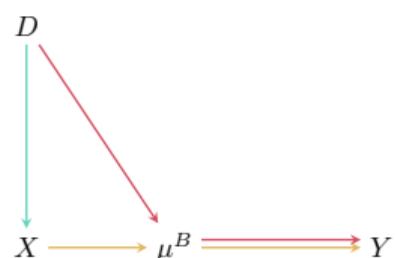
# premiums

21/53

## The best-estimate

premium  $(X, D) \mapsto$

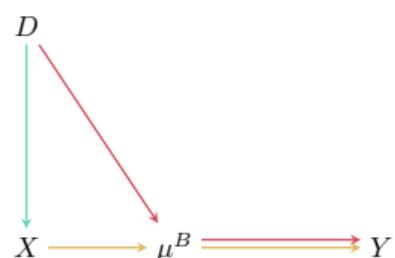
$\mu^B(X, D)$  uses  $(X, D)$  to predict  $Y$ .



$$\mu^B(x, d) = E(Y|X = x, D = d)$$

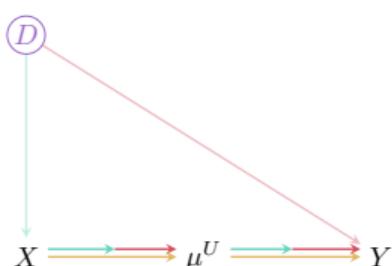
**Best-estimate, unaware****premiums****21**/53

The **best-estimate**  
 premium  $(X, D) \mapsto$   
 $\mu^B(X, D)$  uses  $(X, D)$  to  
 predict  $Y$ .



$$\mu^B(x, d) = E(Y|X = x, D = d)$$

The **unaware** premium  
 $X \mapsto \mu^U(X)$  uses  $X$ ,  
 ignoring  $D$ , to predict  $Y$ .

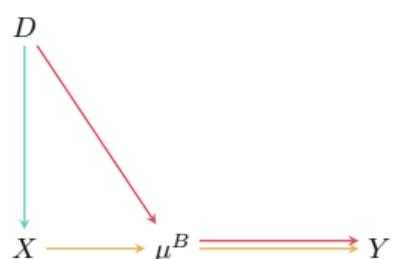


$$\mu^U(x) = E(Y|X = x)$$

# Best-estimate, unaware, and aware premiums

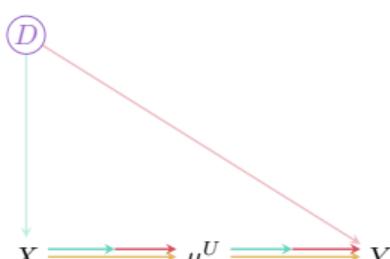
21/53

The **best-estimate**  
premium  $(X, D) \mapsto$   
 $\mu^B(X, D)$  uses  $(X, D)$  to  
predict  $Y$ .



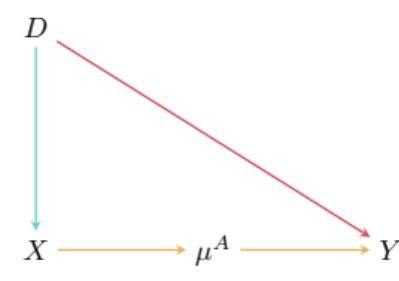
$$\mu^B(x, d) = E(Y|X = x, D = d)$$

The **unaware** premium  
 $X \mapsto \mu^U(X)$  uses  $X$ ,  
ignoring  $D$ , to predict  $Y$ .



$$\mu^U(x) = E(Y|X = x)$$

The **aware** premium  
 $X \mapsto \mu^A(X)$  uses  $X$  to  
predict  $Y$  while controlling  
for  $D$ .

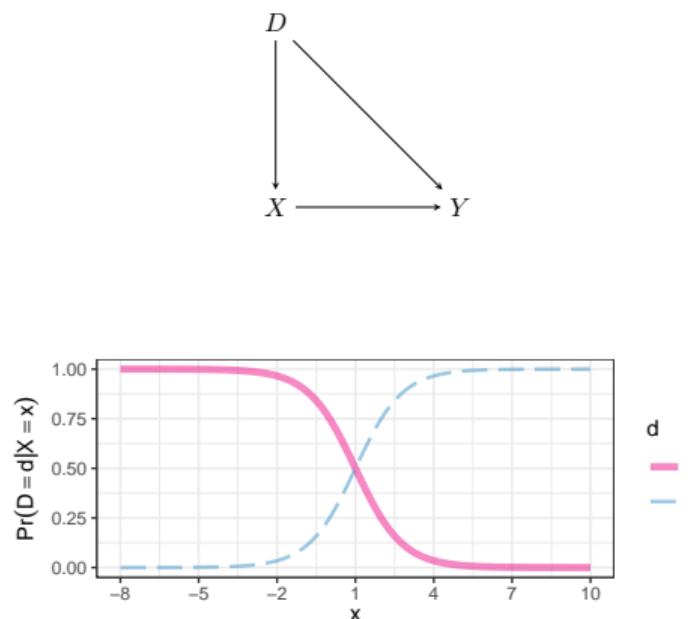


$$\mu^A(x) = E_D\{E(Y|X = x, D)\}$$

## ▶ Configuration of the example

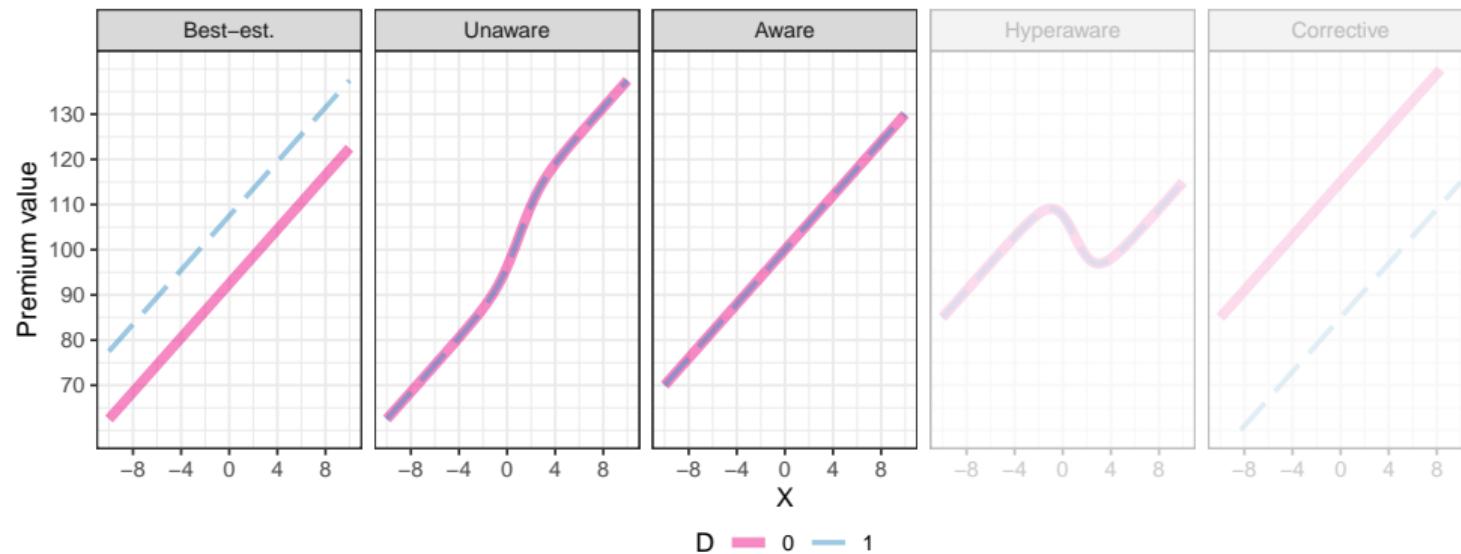
22/53

- Let  $D \in \{0, 1\}$  be Bernoulli with  $\Pr(D = 1) = 0.5$ .
- The variables  $X$  and  $Y$  are Gaussian and the DAG is satisfied.



# Premiums in terms of $x$ and $d$

23/53

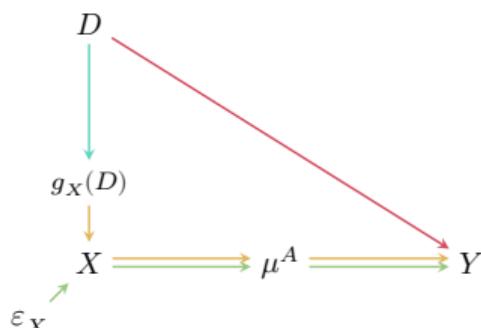


# Hyperaware and corrective premiums

24/53

## The **aware** premium

$X \mapsto \mu^A(X)$  uses all of  $X$  to predict  $Y$  while correcting for the omitted variable bias associated with  $D$ .

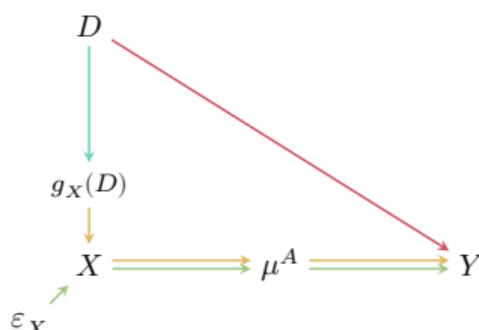


# Hyperaware and corrective premiums

24/53

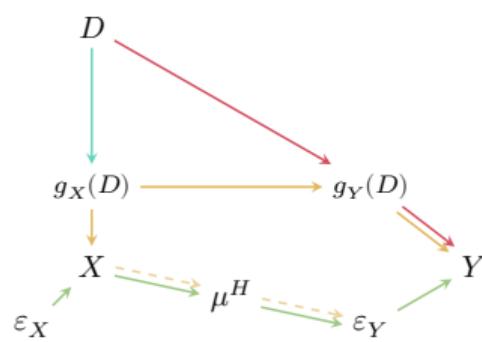
## The aware premium

$X \mapsto \mu^A(X)$  uses all of  $X$  to predict  $Y$  while correcting for the omitted variable bias associated with  $D$ .



## A hyperaware premium

$X \mapsto \mu^H(X)$  use all of  $X$  to predict  $\varepsilon_Y$ .

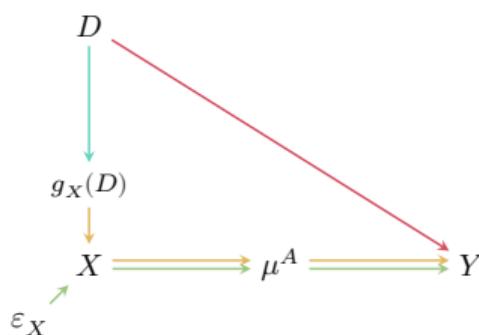


# Hyperaware and corrective premiums

24/53

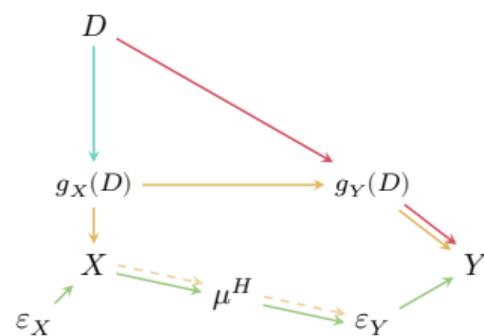
## The aware premium

$X \mapsto \mu^A(X)$  uses all of  $X$  to predict  $Y$  while correcting for the omitted variable bias associated with  $D$ .



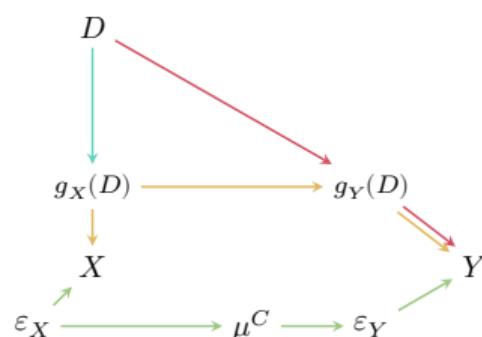
## A hyperaware premium

$X \mapsto \mu^H(X)$  use all of  $X$  to predict  $\varepsilon_Y$ .



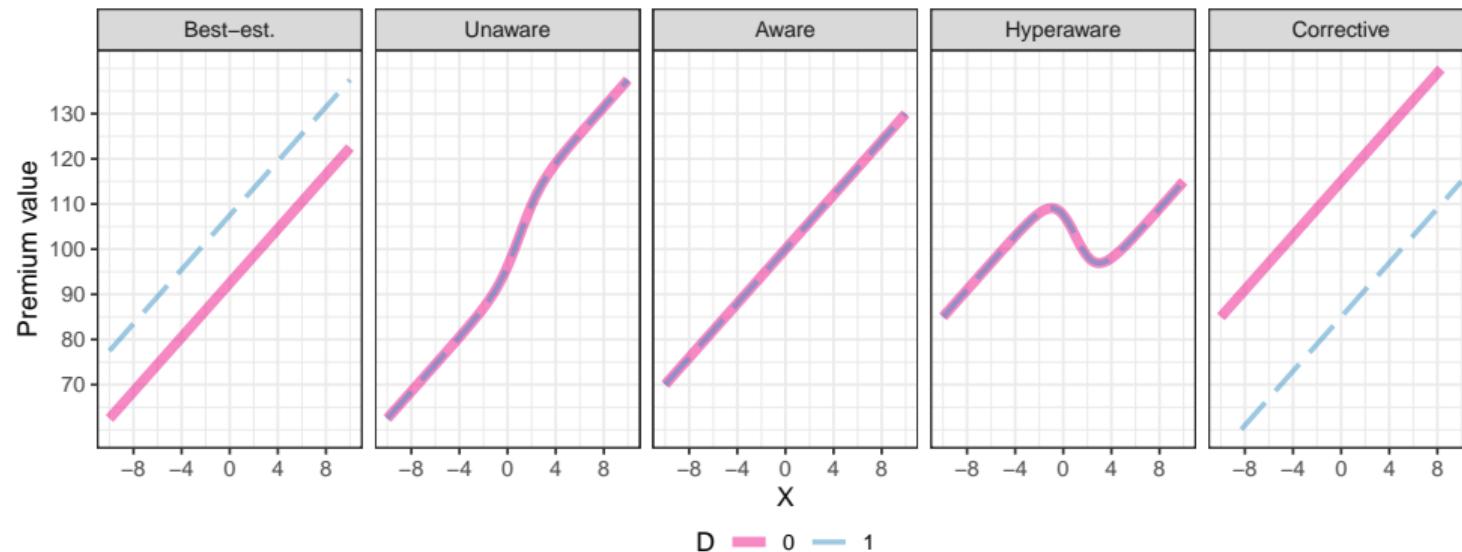
## The corrective premium

$\varepsilon_X \mapsto \mu^C(\varepsilon_X)$  models the residual information independent from  $D$ .



# Premiums in terms of $x$ and $d$

25/53



# Our classification : a few examples

26/53

		Type of intervention		
		<i>Pre-processing</i>	<i>In-processing</i>	<i>Post-processing</i>
<b>Aware</b>	Plečko et al. (2021)		Lindholm et al. (2024)	Pope and Sydnor (2011)
	Charpentier et al. (2023a)		Gabric et al. (2024)	Aseervatham et al. (2016) <a href="#">Lindholm et al. (2022)</a>
<b>Hyperaware</b>	Zemel et al. (2013)		Grari et al. (2022)	
<b>Corrective</b>	Calmon et al. (2017)			Charpentier et al. (2023b)
	Komiyama and Shimao (2017)	Fermanian and Guegan (2021)		Hu et al. (2023)
		<a href="#">Frees and Huang (2023)</a>		Hu et al. (2024)

# Our classification : a few examples

26/53

		Type of intervention ( <i>when</i> )		
		<i>Pre-processing</i>	<i>In-processing</i>	<i>Post-processing</i>
Expected premium ( <i>what</i> )	Aware	Plečko et al. (2021) Charpentier et al. (2023a)	Lindholm et al. (2024) Gabric et al. (2024)	Pope and Sydnor (2011) Aseervatham et al. (2016) <b>Lindholm et al. (2022)</b>
	Hyperaware	Zemel et al. (2013)	Grari et al. (2022)	
	Corrective	Calmon et al. (2017) Komiyama and Shimao (2017) <b>Frees and Huang (2023)</b>	Fermanian and Guegan (2021)	Charpentier et al. (2023b) Hu et al. (2023) Hu et al. (2024)

# So far

27 / 53

- We detailed the causal graph used for fairness in insurance to adequately place the different fair premiums in the graph.

## So far

27 / 53

- We detailed the causal graph used for fairness in insurance to adequately place the different fair premiums in the graph.
- We defined direct and indirect discrimination for a premium.

## So far

27 / 53

- We detailed the causal graph used for fairness in insurance to adequately place the different fair premiums in the graph.
- We defined direct and indirect discrimination for a premium.
- We define five families of premiums according to their fairness properties, therefore giving a **new categorization** of fairness methodologies.

## So far

27 / 53

- We detailed the causal graph used for fairness in insurance to adequately place the different fair premiums in the graph.
- We defined direct and indirect discrimination for a premium.
- We define five families of premiums according to their fairness properties, therefore giving a **new categorization** of fairness methodologies.

How does this framework make a difference **in practice** ?

# A scalable toolbox for exposing indirect discrimination

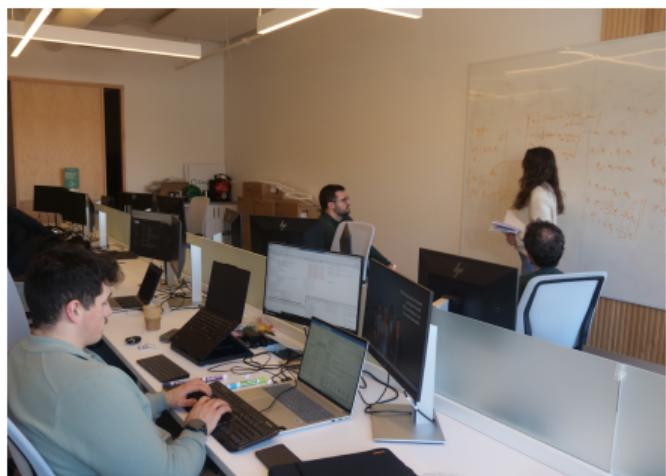
---

- 1 Exploiting causal graphs for fairness
- 2 A scalable toolbox for exposing indirect discrimination
  - What does it mean to be fair ?
  - Can one benchmark for fairness ?
  - Is unfairness truly material ?
  - Where indirect discrimination concentrates ?

## Context around this project

## Context around this project

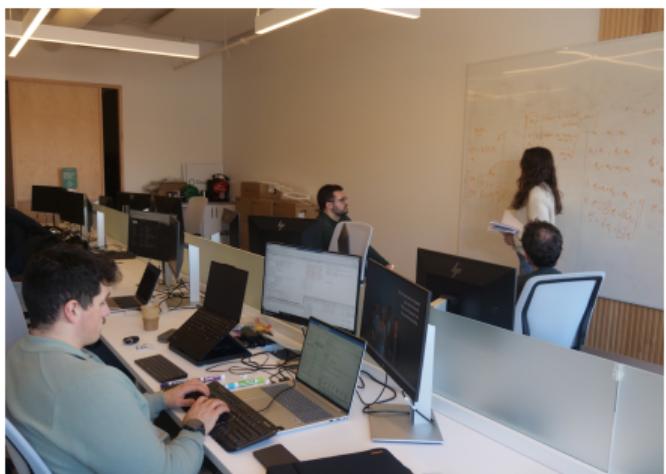
28/53



**Figure 2** – Journée de collaboration  
au local **Desjardins** du MILA, 20  
février 2025, Montréal

# Context around this project

28 / 53



**Figure 2 - Journée de collaboration au local **Desjardins** du MILA, 20 février 2025, Montréal**

 Outlook

---

Congratulations on Achieving Your ACAS! CRM:0100683

---

De CAS Portal - NO REPLY <cobaltams@casact.org>

Date Ven 2025-06-06 12:26

À Olivier Côté <olivierc\_actuarial@outlook.com>

Olivier Côté, ACAS  
870 av Marguerite-Bourgeoys  
Lévis QC G1S3X2  
Canada

Dear Olivier,

On behalf of the Casualty Actuarial Society, it is our pleasure to notify you that **you are now an Associate of the CAS**. You may begin using the designation "Associate, Casualty Actuarial Society" or "ACAS" immediately. CAS credentials are unmatched for their rigor, integrity, and relevance, making it the gold standard among property & casualty actuarial credentials. Congratulations on attaining this milestone in your actuarial career!

Sincerely,



David Cummings, FCAS  
President



Expertise. Insight.  
Solutions.®



Victor Carter-Bey, DM  
Chief Executive Officer

# Motivation

29/53

- There is growing interest in ensuring fairness in insurance pricing.
  - ▶ E.g. Financial Services Regulatory Authority of Ontario (2024)

# Motivation

29/53

- There is growing interest in ensuring fairness in insurance pricing.
  - ▶ E.g. Financial Services Regulatory Authority of Ontario (2024)
- Fairness in insurance is difficult to operationalize due to its ambiguity.
  - ▶ « there is no general agreement on what constitutes an “equitable” classification system or “fair” discrimination. »
    - ASOP No. 12, ASB (2005)

# Motivation

29/53

- There is growing interest in ensuring fairness in insurance pricing.
  - ▶ E.g. Financial Services Regulatory Authority of Ontario (2024)
- Fairness in insurance is difficult to operationalize due to its ambiguity.
  - ▶ « there is no general agreement on what constitutes an “equitable” classification system or “fair” discrimination. »
    - ASOP No. 12, ASB (2005)
- Current fairness metrics are not suited for actuarial needs.
  - ▶ They don't identify problematic segments or variables.
  - ▶ Their units are abstract.

# Motivation

29/53

- There is growing interest in ensuring fairness in insurance pricing.
  - ▶ E.g. Financial Services Regulatory Authority of Ontario (2024)
- Fairness in insurance is difficult to operationalize due to its ambiguity.
  - ▶ « there is no general agreement on what constitutes an “equitable” classification system or “fair” discrimination. »
    - ASOP No. 12, ASB (2005)
- Current fairness metrics are not suited for actuarial needs.
  - ▶ They don't identify problematic segments or variables.
  - ▶ Their units are abstract.
- Fairness is debated in theory, but **unfairness unfolds in practice**.
  - ▶ Few Canadian insurance case studies exist.

# Vulnerability to proxies in Québec auto insurance

30/53

**Objective :** Evaluate a given **pseudo price** for fairness regarding financial vulnerability in material damage premiums for at-fault accidents (Chapter B2) in Quebec.

# Vulnerability to proxies in Québec auto insurance

30/53

**Objective :** Evaluate a given **pseudo price** for fairness regarding financial vulnerability in material damage premiums for at-fault accidents (Chapter B2) in Quebec.

**Data :**  $\approx$  768,000 insured vehicles in Quebec, from 2016-2017. Data obtained via insurer partnership.

Note : Personal data anonymized; strict confidentiality measures applied.

# Data overview

31/53

Notation	Concept	Domain	Notes
$Y$	Claim amount (\$)	$\mathbb{R}^+$	$\bar{Y} \approx 200$ , with 97% at 0

# Data overview

31/53

Notation	Concept	Domain	Notes
$Y$	Claim amount (\$)	$\mathbb{R}^+$	$\bar{Y} \approx 200$ , with 97% at 0
$D$	Low credit indicator	$\{0, 1\}$	1 indicates low credit, with $\bar{D} \approx 0.40$

# Data overview

31/53

Notation	Concept	Domain	Notes
$Y$	Claim amount (\$)	$\mathbb{R}^+$	$\bar{Y} \approx 200$ , with 97% at 0
$D$	Low credit indicator	$\{0, 1\}$	1 indicates low credit, with $\bar{D} \approx 0.40$
$X$	Policyholder info	Dim. 16	E.g., gender, driving experience, mileage, education, occupation
	Geographic info	Dim. 4	E.g., FSA and territorial risk score
	Vehicle info	Dim. 4	E.g., vehicle age, new purchase, vehicle risk score
	Policy info	Dim. 3	E.g., home insurance, endorsements

# Data overview

31/53

Notation	Concept	Domain	Notes
$Y$	Claim amount (\$)	$\mathbb{R}^+$	$\bar{Y} \approx 200$ , with 97% at 0
$D$	Low credit indicator	$\{0, 1\}$	1 indicates low credit, with $\bar{D} \approx 0.40$
$X$	Policyholder info	Dim. 16	E.g., gender, driving experience, mileage, education, occupation
	Geographic info	Dim. 4	E.g., FSA and territorial risk score
	Vehicle info	Dim. 4	E.g., vehicle age, new purchase, vehicle risk score
	Policy info	Dim. 3	E.g., home insurance, endorsements

We aim to analyze **PseudoPrice(x,d)**.

# Data overview

31/53

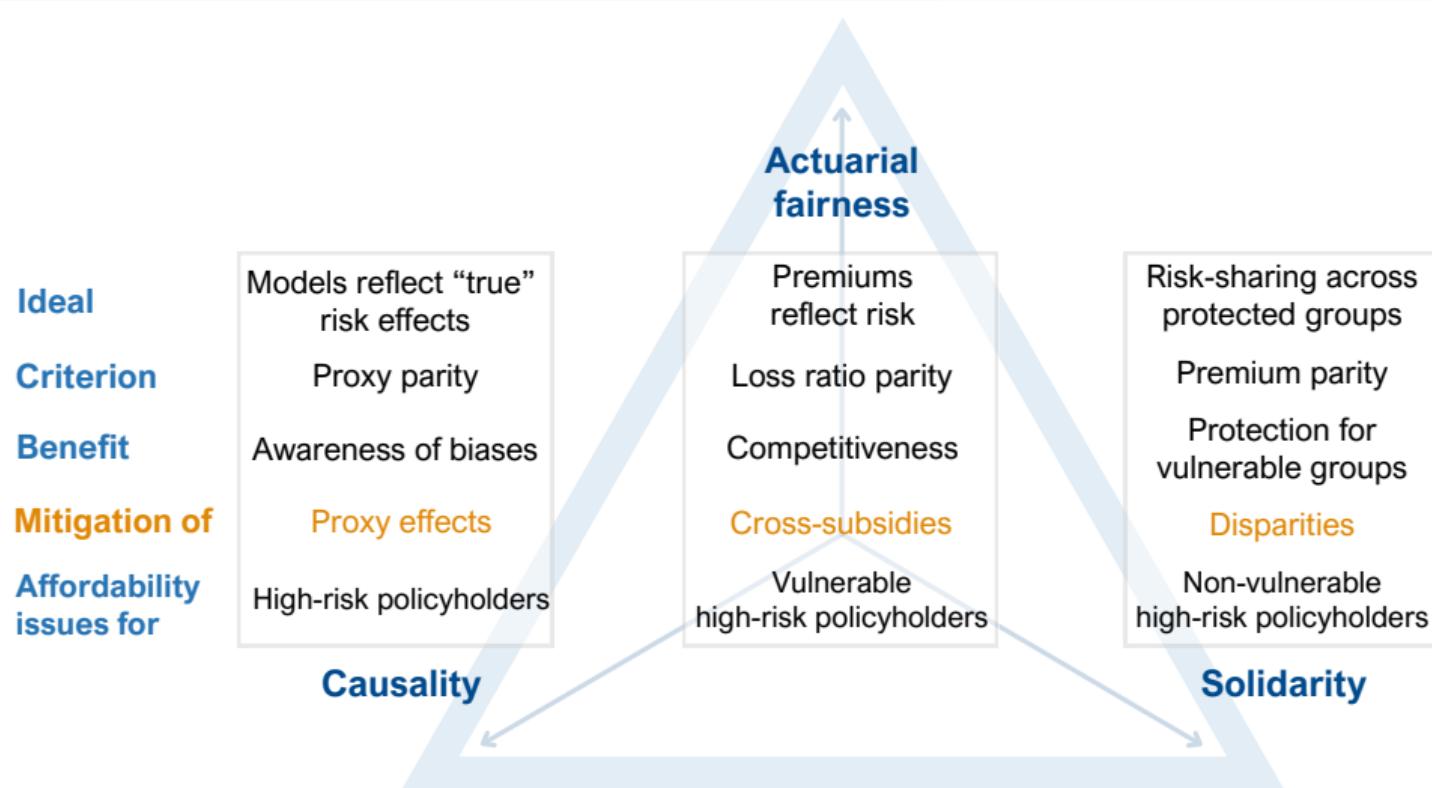
Notation	Concept	Domain	Notes
$Y$	Claim amount (\$)	$\mathbb{R}^+$	$\bar{Y} \approx 200$ , with 97% at 0
$D$	Low credit indicator	$\{0, 1\}$	1 indicates low credit, with $\bar{D} \approx 0.40$
$X$	Policyholder info	Dim. 16	E.g., gender, driving experience, mileage, education, occupation
	Geographic info	Dim. 4	E.g., FSA and territorial risk score
	Vehicle info	Dim. 4	E.g., vehicle age, new purchase, vehicle risk score
	Policy info	Dim. 3	E.g., home insurance, endorsements

We aim to analyze **PseudoPrice(x,d)**.

The **non-equivalence of the pseudoprice to the actual pricing function** precludes any conclusion regarding the fairness of the partner insurer's pricing.

# The dimensions of fairness

32/53



# Actuarial Fairness

33/53

A premium is actuarially fair if “it represents an unbiased estimate of the expected value of all future costs associated with the risk transfer” (Casualty Actuarial Society, 1988).

## Actuarial Fairness

33/53

A premium is actuarially fair if “it represents an unbiased estimate of the expected value of all future costs associated with the risk transfer” (Casualty Actuarial Society, 1988).

- Premium matches expected cost.

## Actuarial Fairness

33/53

A premium is actuarially fair if “it represents an unbiased estimate of the expected value of all future costs associated with the risk transfer” (Casualty Actuarial Society, 1988).

- Premium matches expected cost.
- Self-sustaining loss ratios (no cross-subsidies).

# Actuarial Fairness

33/53

A premium is actuarially fair if “it represents an unbiased estimate of the expected value of all future costs associated with the risk transfer” (Casualty Actuarial Society, 1988).

- Premium matches expected cost.
- Self-sustaining loss ratios (no cross-subsidies).
- Avoids **non-risk-based adjustments**.

# Solidarity

34/53

Solidarity is the foundation of insurance : pooling and sharing risks among **all** policyholders.

# Solidarity

34/53

Solidarity is the foundation of insurance : pooling and sharing risks among **all** policyholders.

In the presence of admissible variables  $X$  and sensitive ones  $D$ ,

- we can leverage the effect of  $X$  on the response variable  $Y$
- while aiming for solidarity on  $D$  : equal premiums (in expectation or distribution) across all protected groups.

This is referred to as **demographic parity of premiums**.

## Causality and proxy effects

35/53

Avoiding proxy effects requires two actions :

- Exclude factors that do not determine risk,
- Limit effect of risk factors to their “**true**” risk relevance.

## Causality and proxy effects

35/53

Avoiding proxy effects requires two actions :

- Exclude factors that do not determine risk,
- Limit effect of risk factors to their “**true**” risk relevance.

Even valid risk factors can suffer from proxy effects. (e.g., territory and ethnicity)

A variable's use – not the variable itself – determines its role as a proxy.

## Causality and proxy effects

35/53

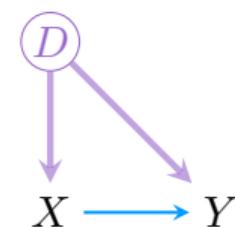
Avoiding proxy effects requires two actions :

- Exclude factors that do not determine risk,
- Limit effect of risk factors to their “**true**” risk relevance.

Even valid risk factors can suffer from proxy effects. (e.g., territory and ethnicity)

A variable's use – not the variable itself – determines its role as a proxy.

In fairness analysis with respect to  $D$ , causality seeks to identify the **effect of  $X$  on  $Y$**  without **proxy effects from  $D$** .



# Five fair benchmarks

36/53

Premium	Best-estimate	Unaware	Aware	Hyperaware	Corrective
Notation	$\mu^B(\mathbf{x}, d)$	$\mu^U(\mathbf{x})$	$\mu^A(\mathbf{x})$	$\mu^H(\mathbf{x})$	$\mu^C(\mathbf{x}, d)$
Formula	$\mathbb{E}(Y \mathbf{X} = \mathbf{x}, D = d)$	$\mathbb{E}(Y \mathbf{X} = \mathbf{x})$	$\mathbb{E}_D\{\mu^B(\mathbf{x}, D)\}$	$\mathbb{E}\{\mu^C(\mathbf{x}, D) \mathbf{X} = \mathbf{x}\}$	$\mathcal{T}^{d \rightarrow *}\{\mu^B(\mathbf{x}, d)\}$
Direct discrimination	✓	✗	✗	✗	✓
Proxy discrimination	-	✓	✗	✓	-
Demographic disparities	✓	✓	✓	✗	✗
Pillar	AF	AF	C	S	S

# Five fair benchmarks

36/53

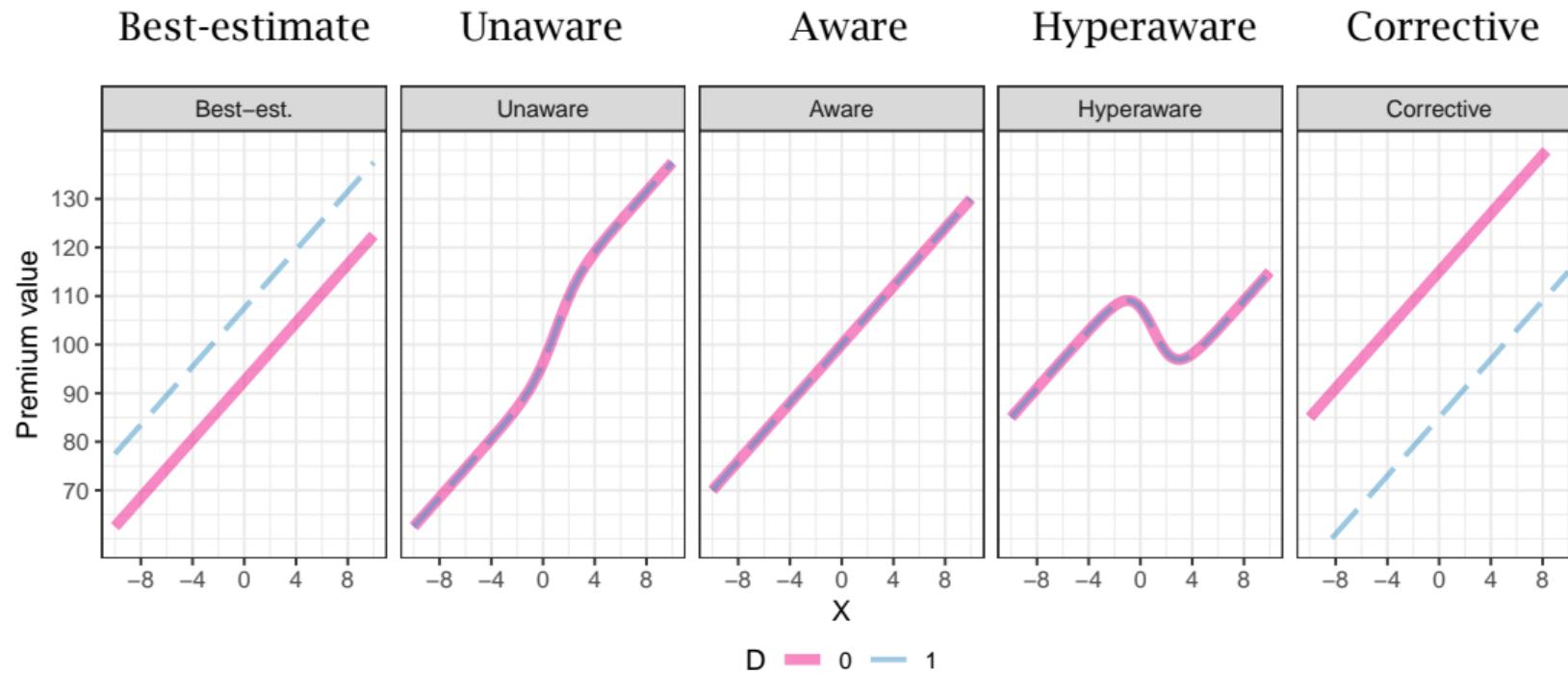
Premium	Best-estimate	Unaware	Aware	Hyperaware	Corrective
Notation	$\mu^B(\mathbf{x}, d)$	$\mu^U(\mathbf{x})$	$\mu^A(\mathbf{x})$	$\mu^H(\mathbf{x})$	$\mu^C(\mathbf{x}, d)$
Formula	$\mathbb{E}(Y \mathbf{X} = \mathbf{x}, D = d)$	$\mathbb{E}(Y \mathbf{X} = \mathbf{x})$	$\mathbb{E}_D\{\mu^B(\mathbf{x}, D)\}$	$\mathbb{E}\{\mu^C(\mathbf{x}, D) \mathbf{X} = \mathbf{x}\}$	$\mathcal{T}^{d \rightarrow *}\{\mu^B(\mathbf{x}, d)\}$
Direct discrimination	✓	✗	✗	✗	✓
Proxy discrimination	-	✓	✗	✓	-
Demographic disparities	✓	✓	✓	✗	✗
Pillar	AF	AF	C	S	S

We estimate the premium spectrum using :

- **lightGBM** (Ke et al., 2017) to learn conditional expectations,
- **empirical marginals** of  $D$  for population-level integration, and
- **optimal transport mappings** via Equipy (Fernandes Machado et al., 2025).

## Example : Premiums in terms of $x$ and $d$

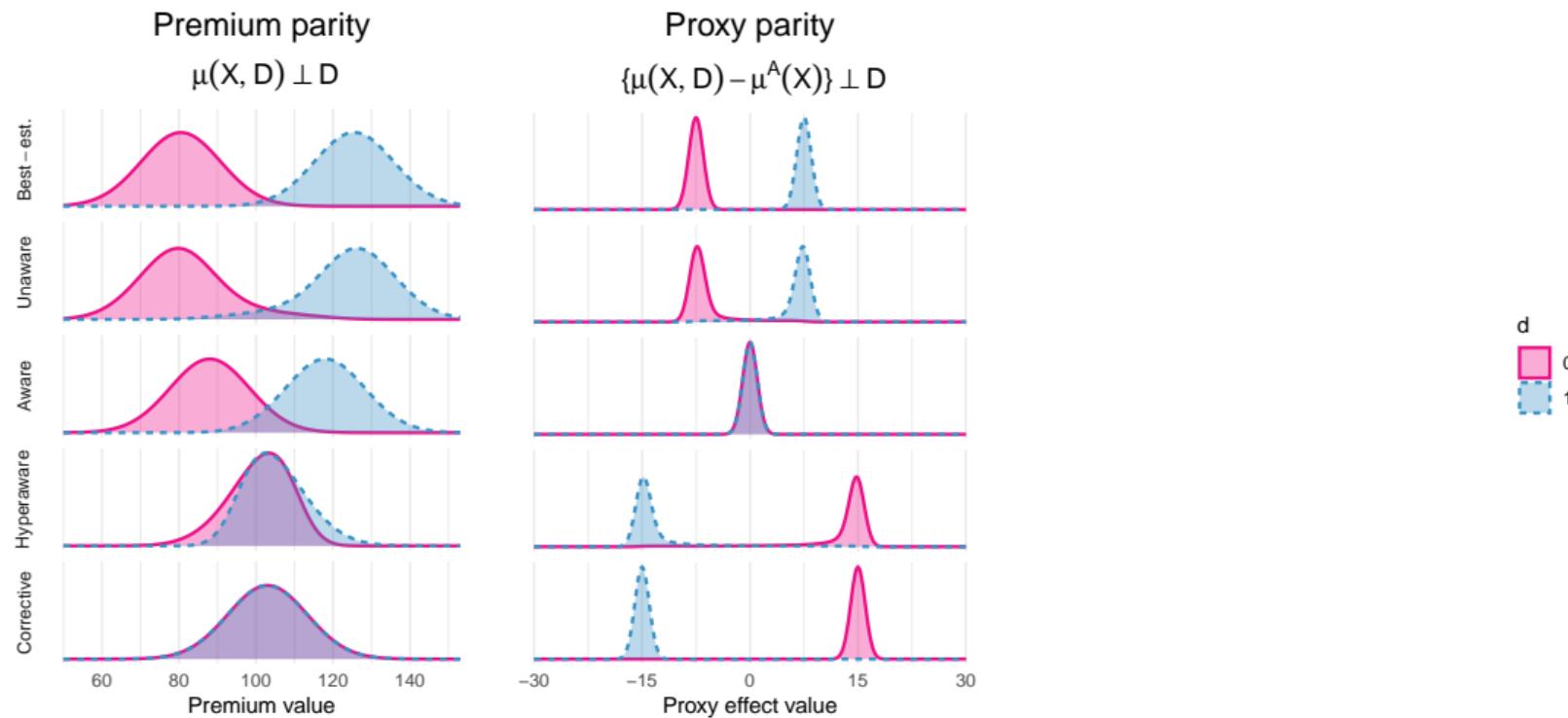
37 / 53

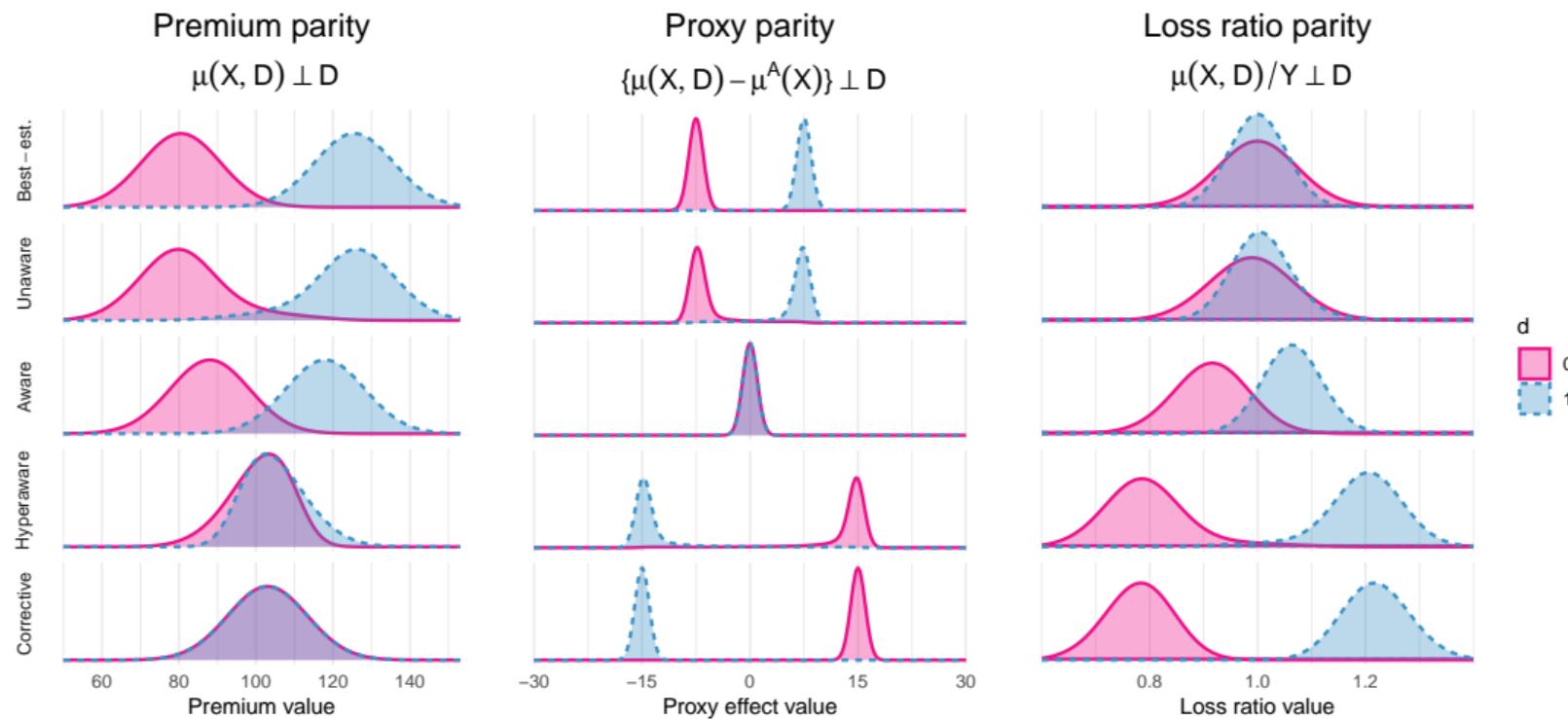


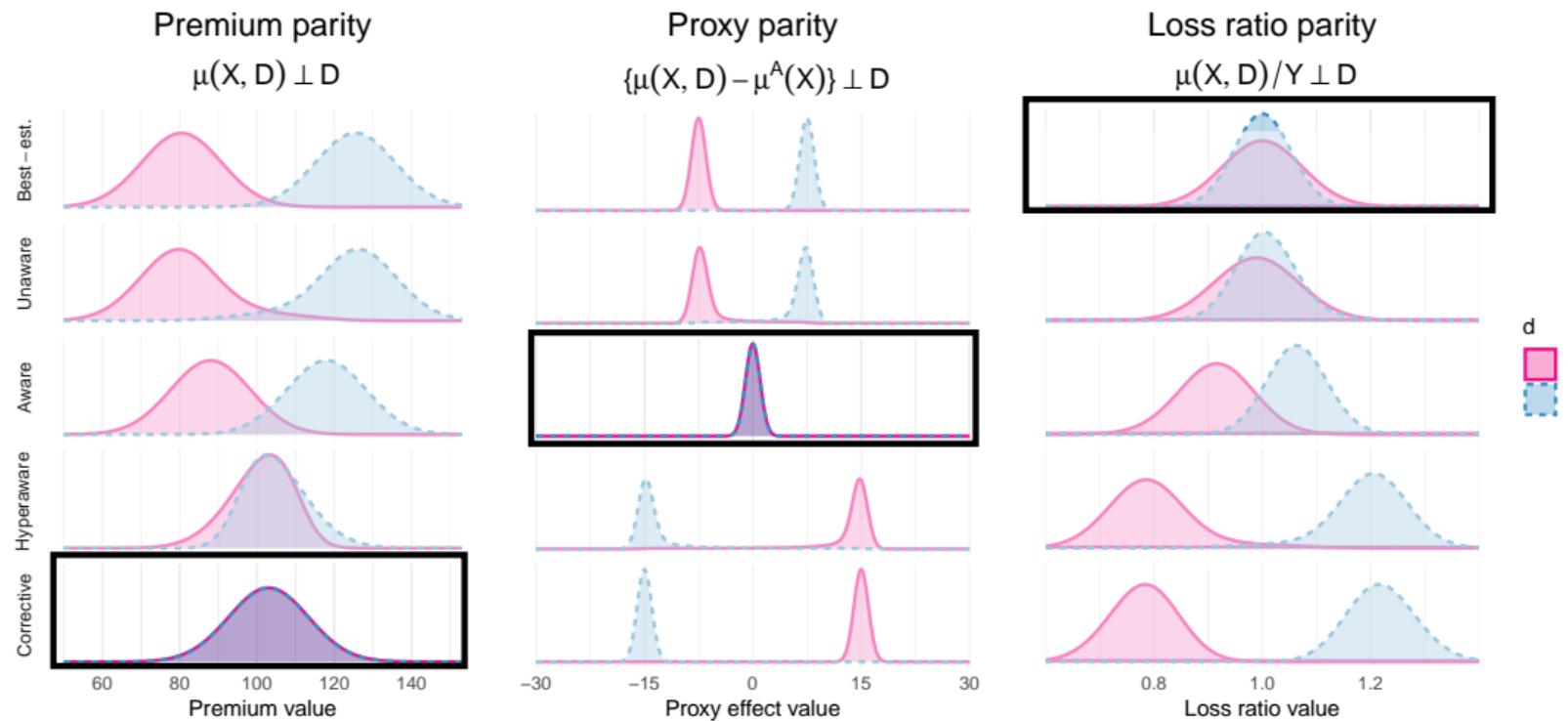
## Premium parity

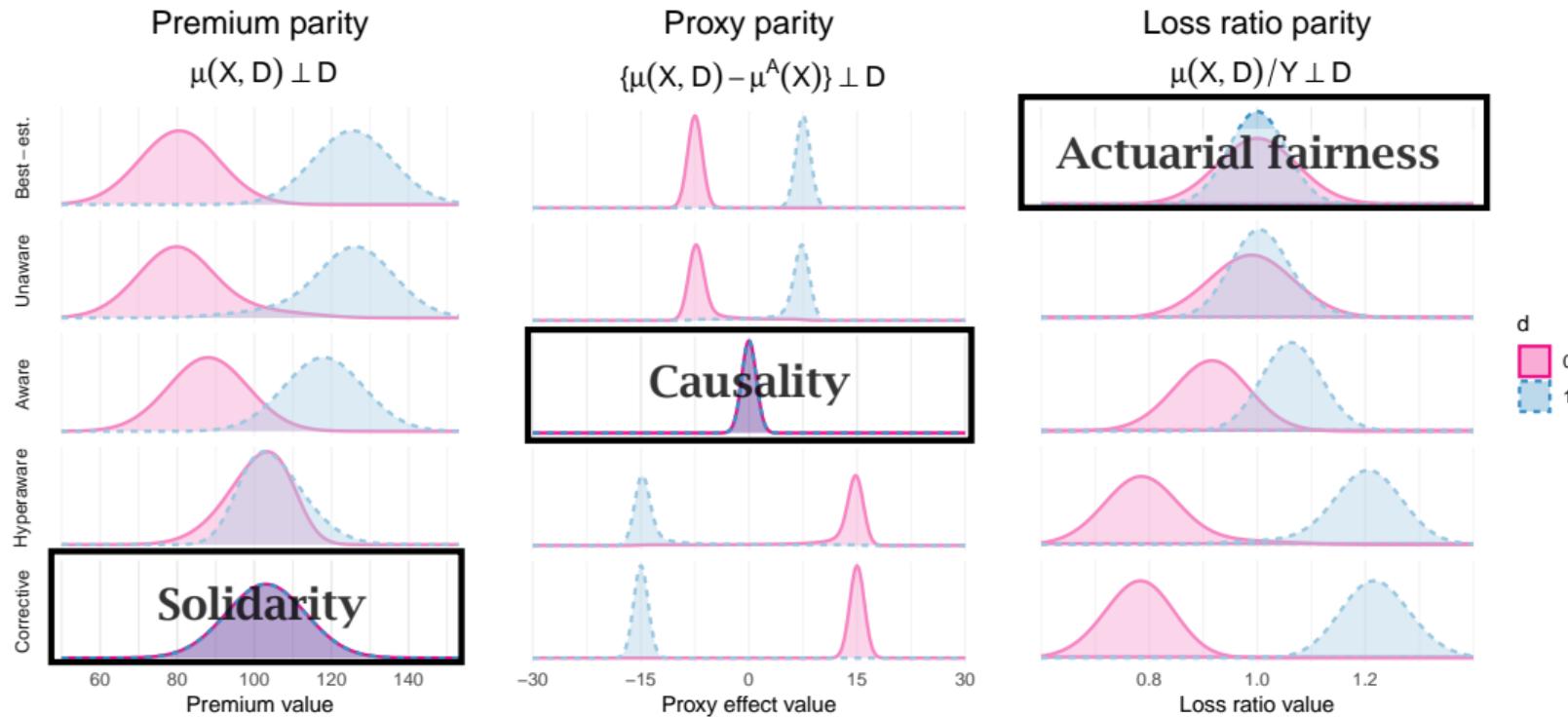
$$\mu(X, D) \perp D$$

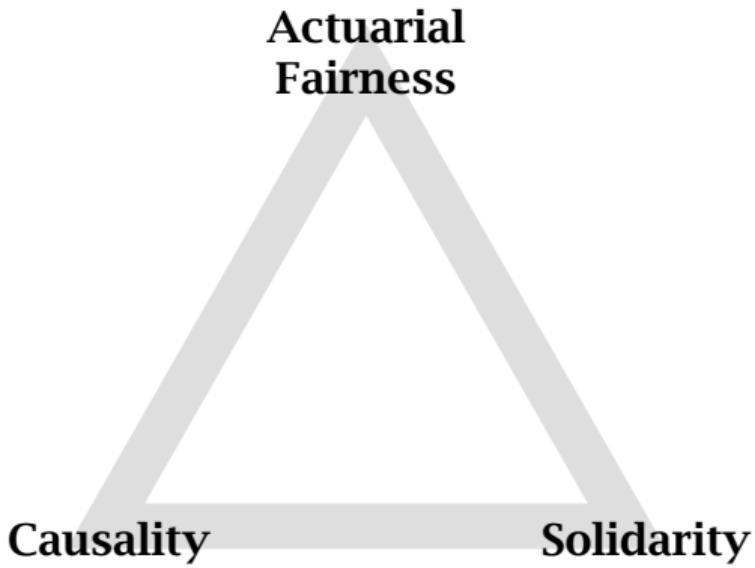












Actuarial  
Fairness

Causality

Solidarity

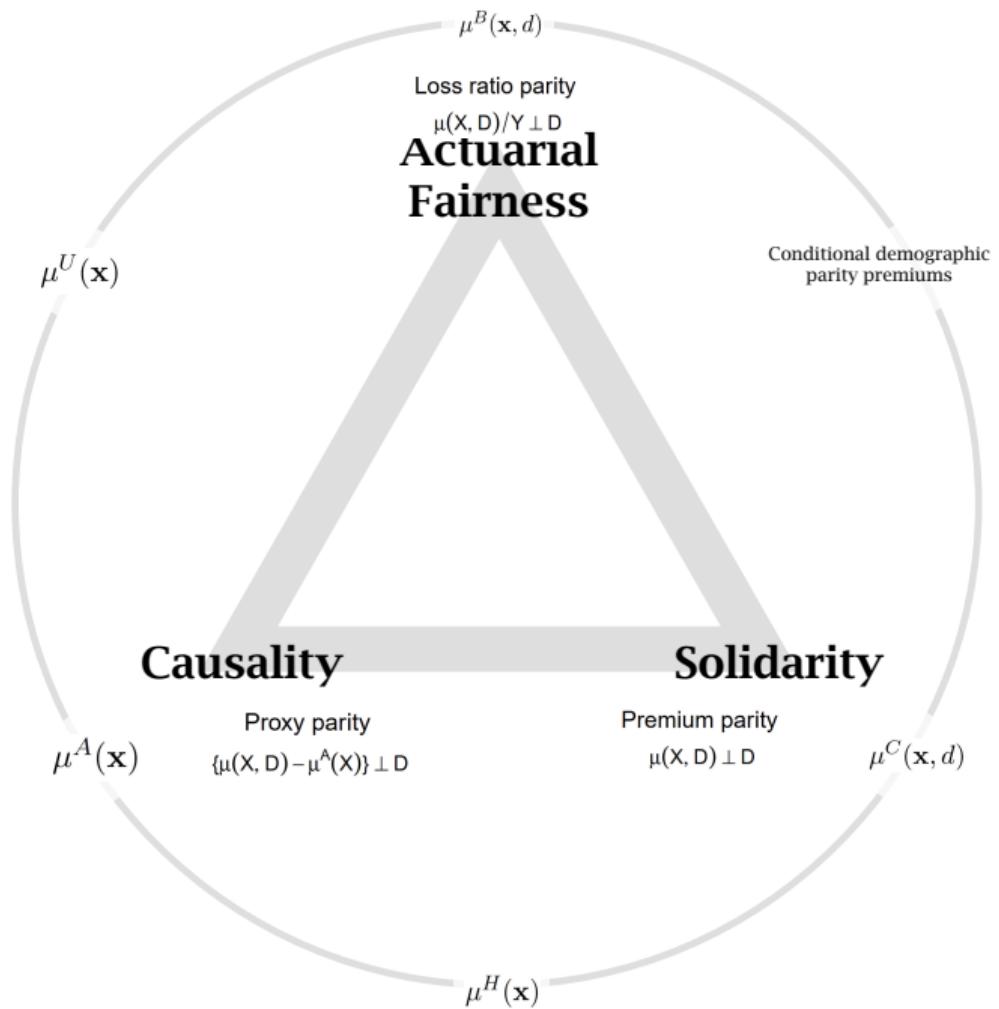
**Actuarial  
Fairness**

**Causality**

Proxy parity  
 $\{\mu(X, D) - \mu^A(X)\} \perp D$

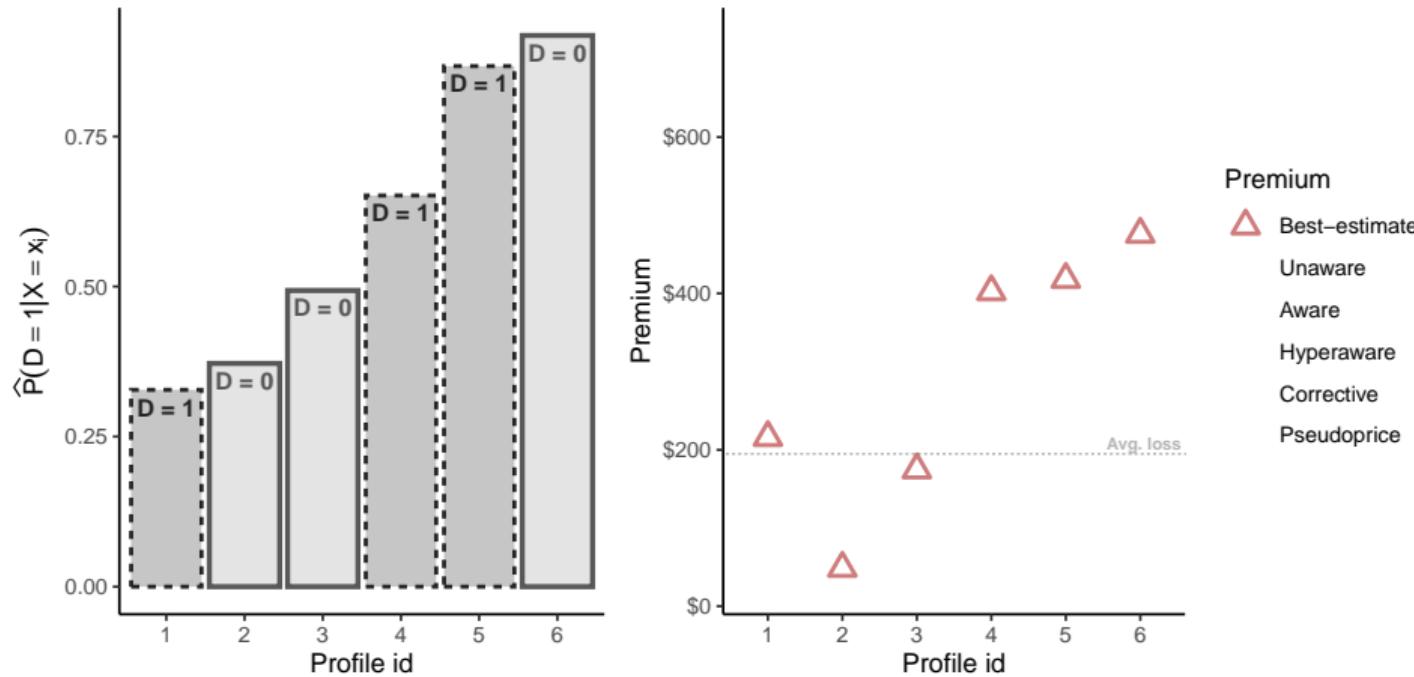
**Solidarity**

Premium parity  
 $\mu(X, D) \perp D$



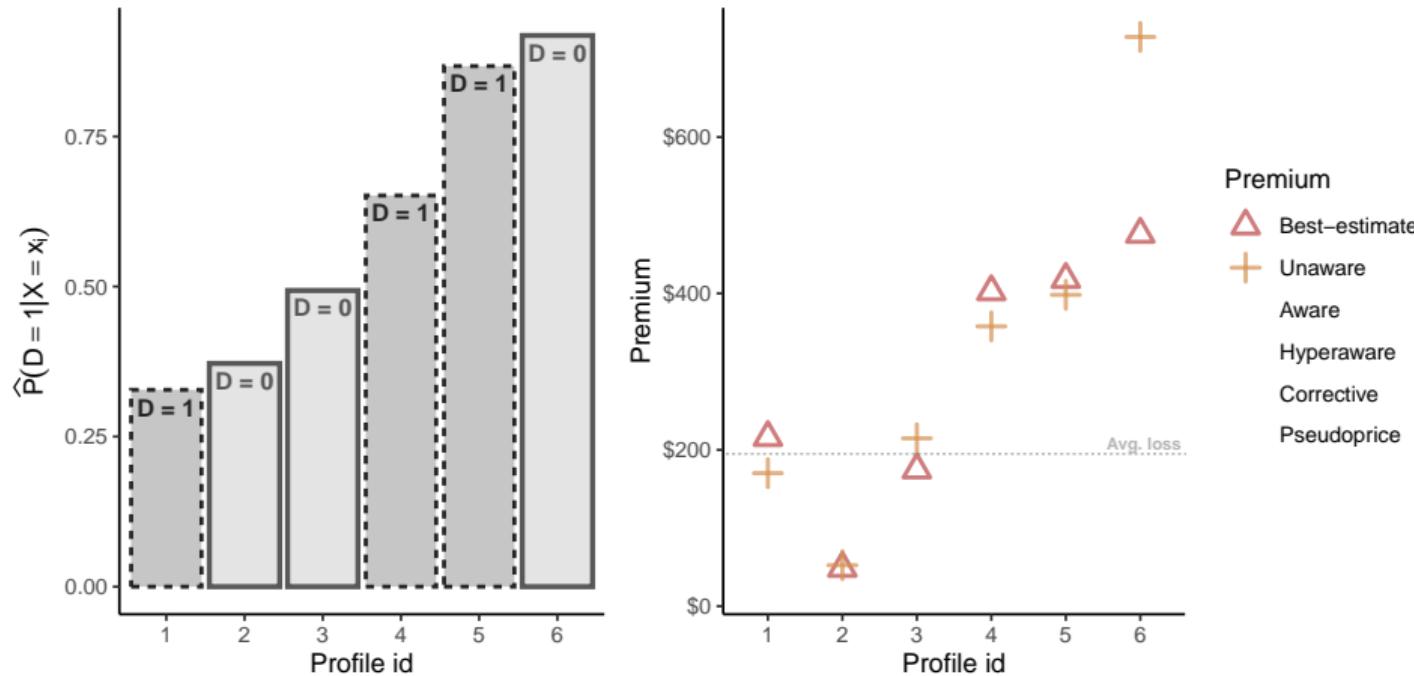
## Fairness range for six profiles

40/53



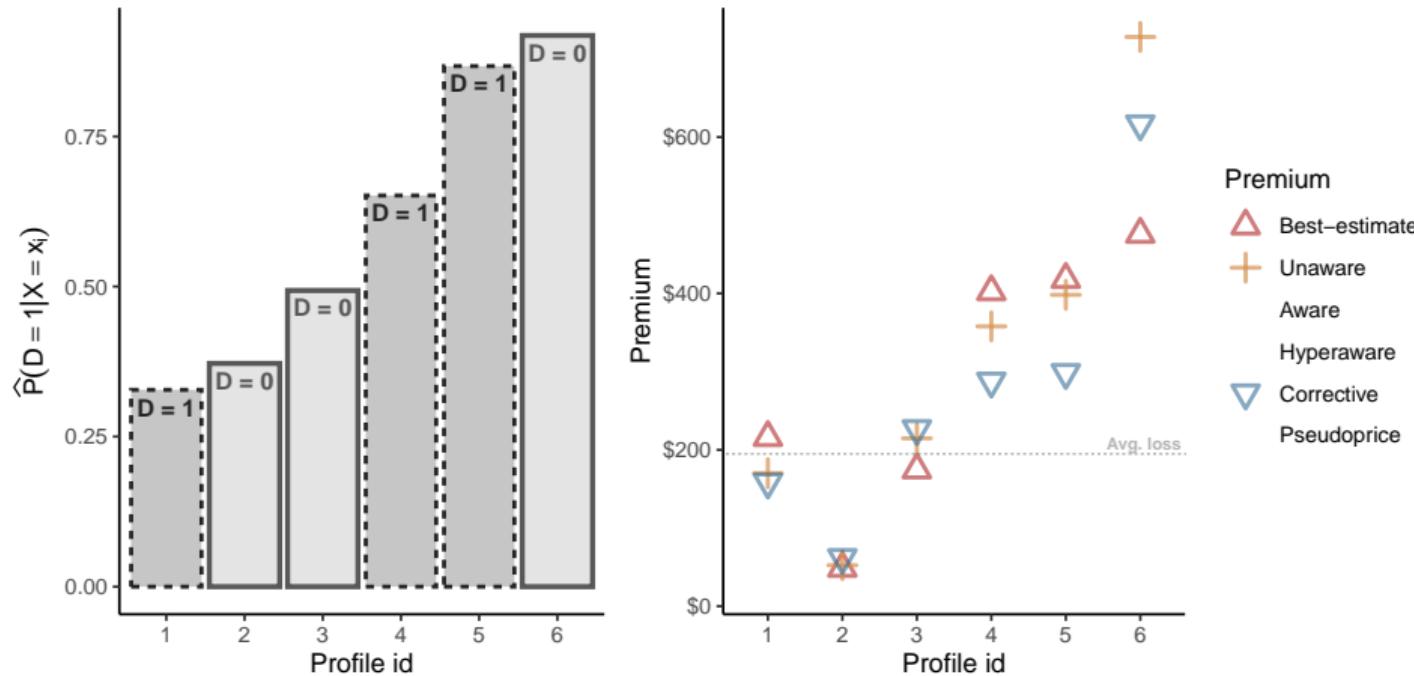
# Fairness range for six profiles

40/53



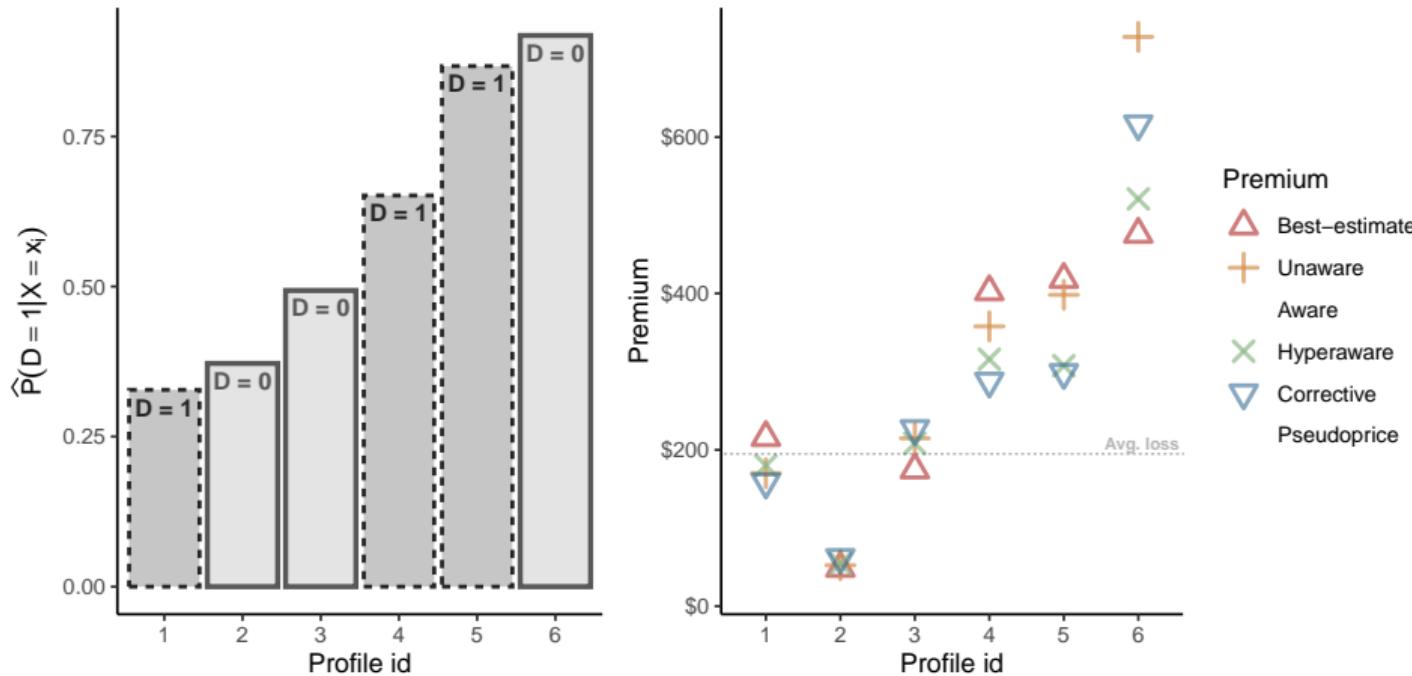
# Fairness range for six profiles

40/53



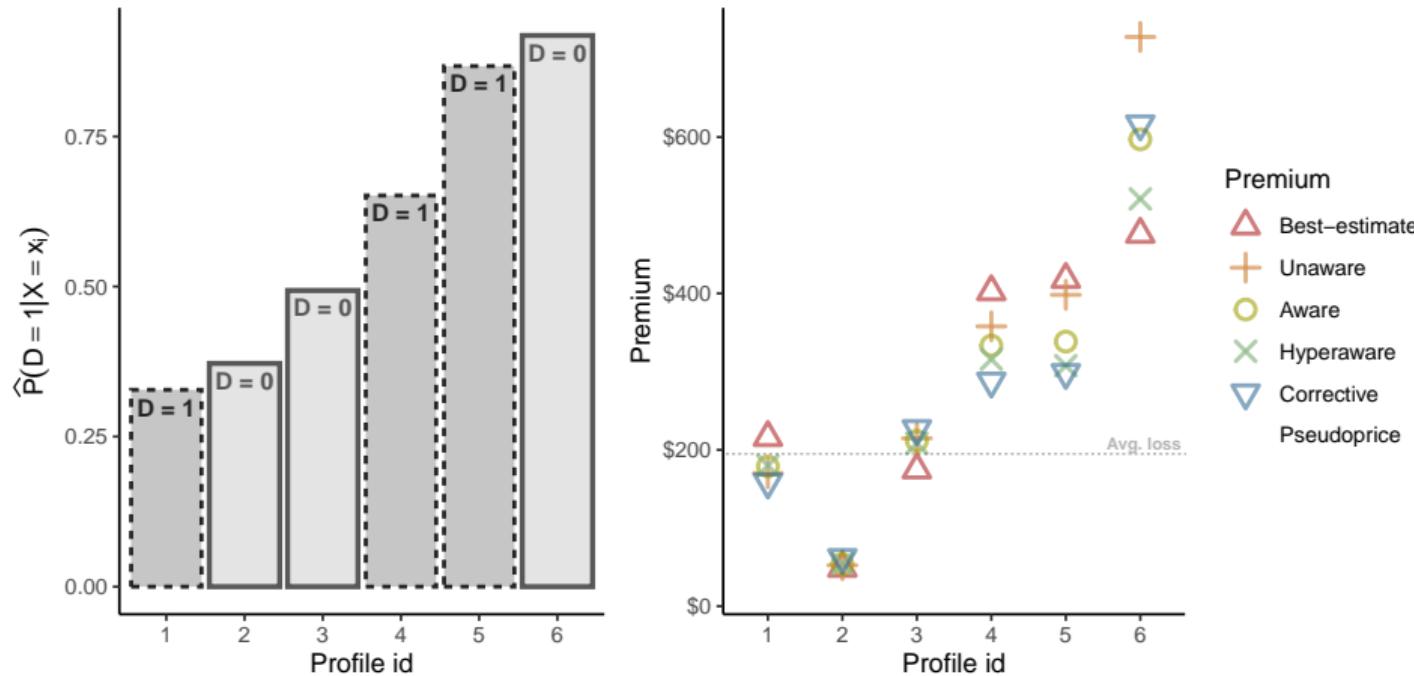
# Fairness range for six profiles

40/53



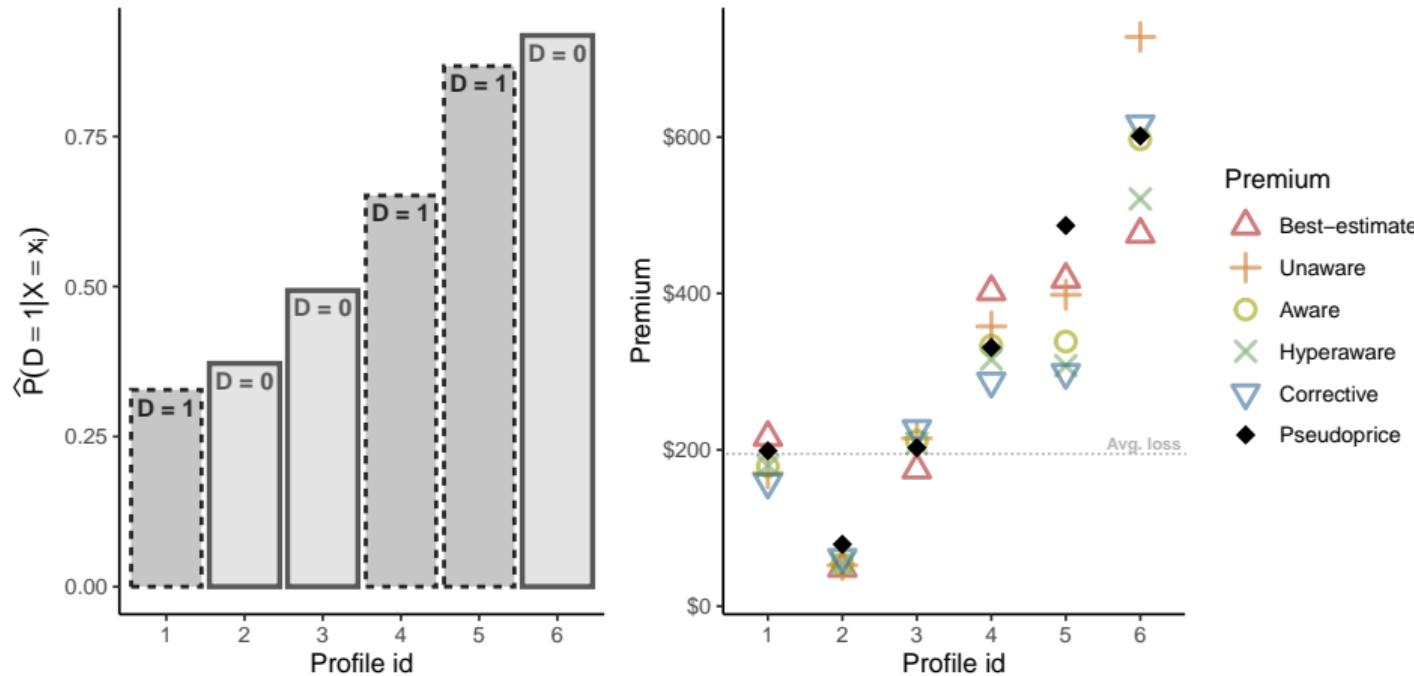
## Fairness range for six profiles

40/53



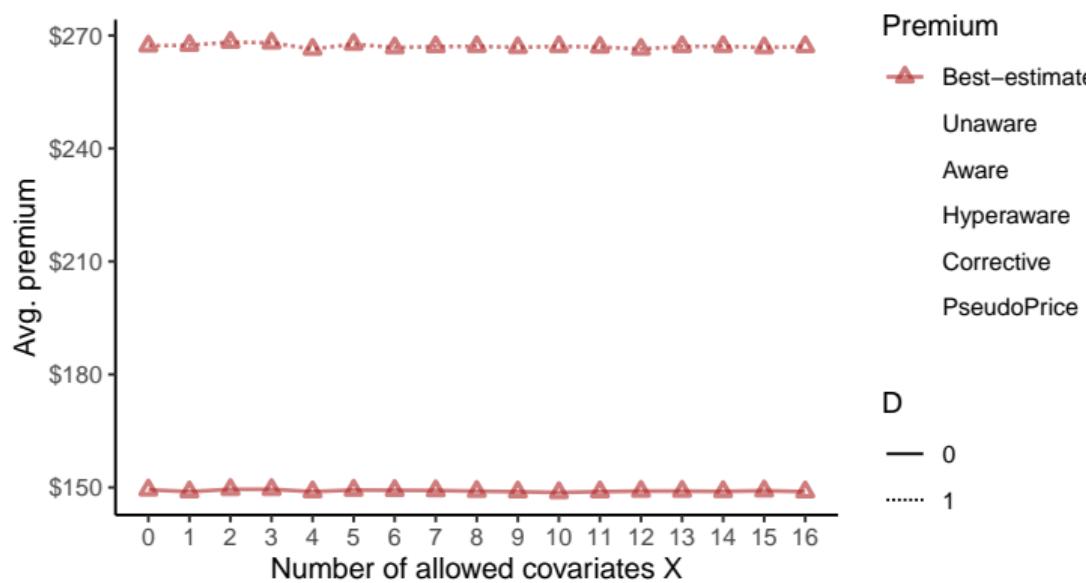
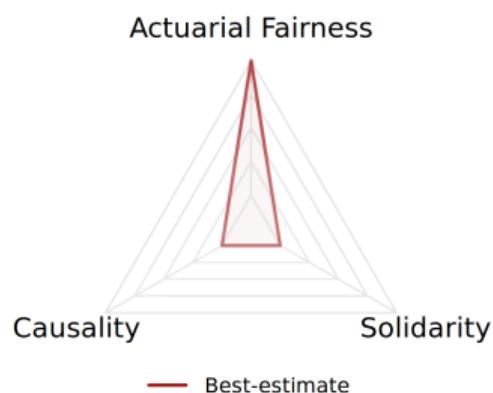
# Fairness range for six profiles

40/53



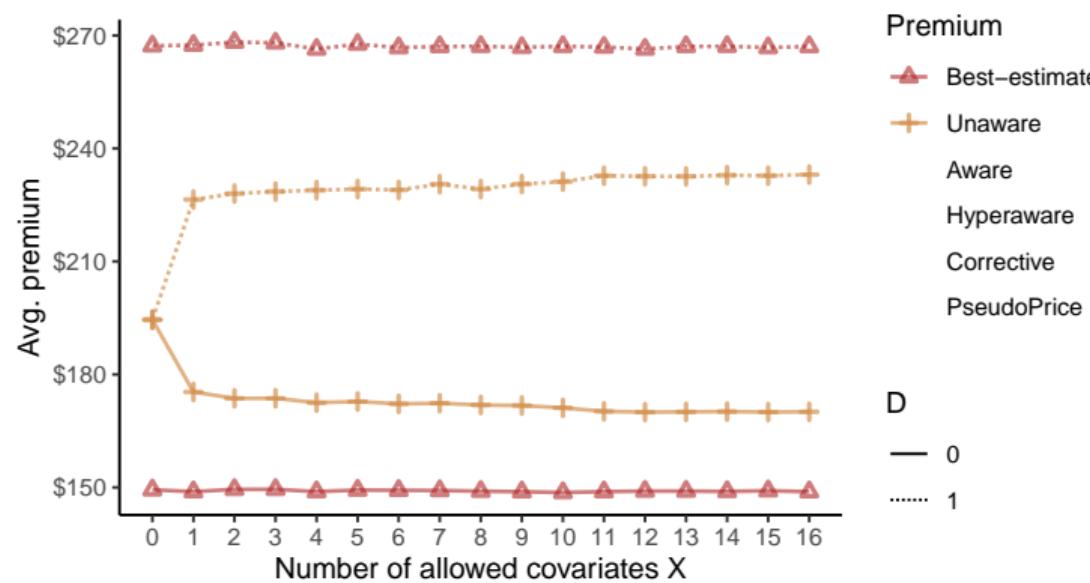
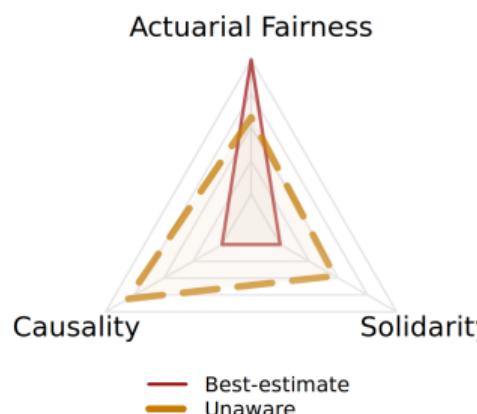
# Grasping behavior of families

41/53



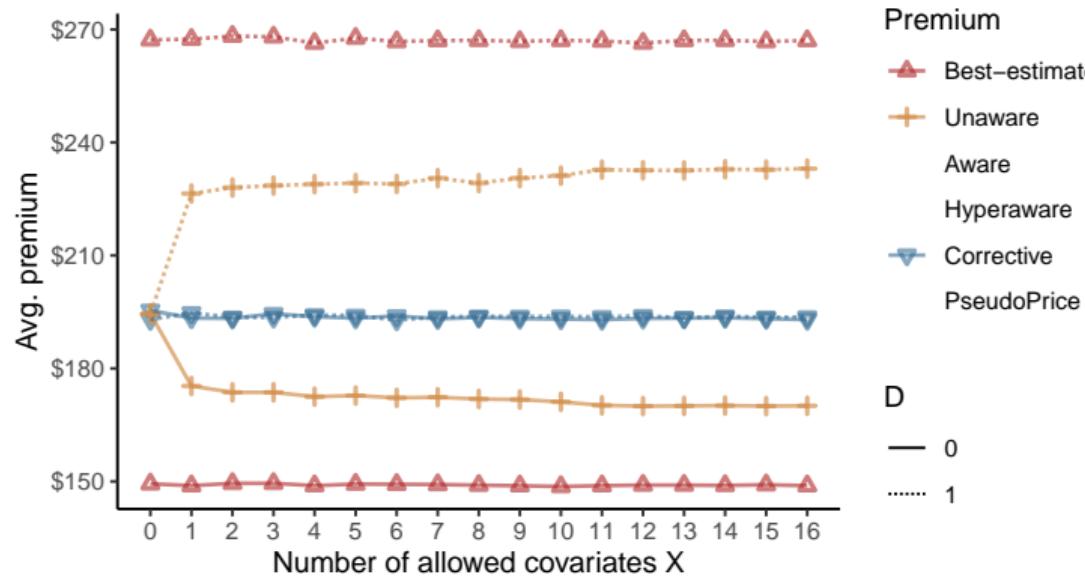
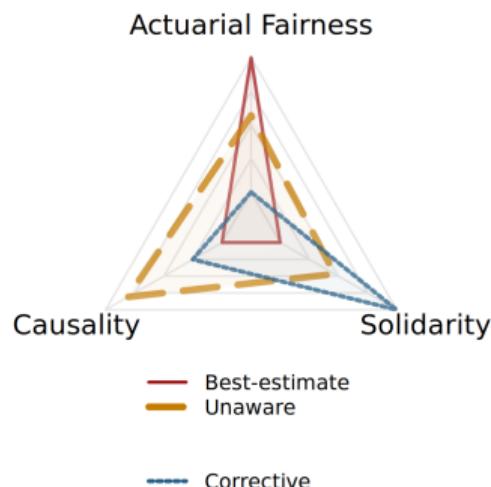
# Grasping behavior of families

41/53



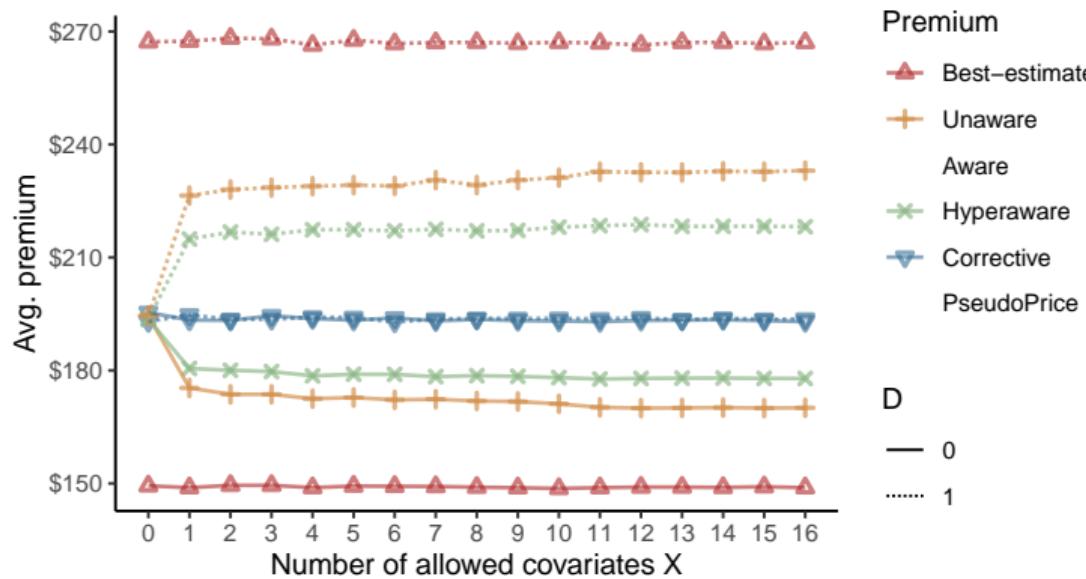
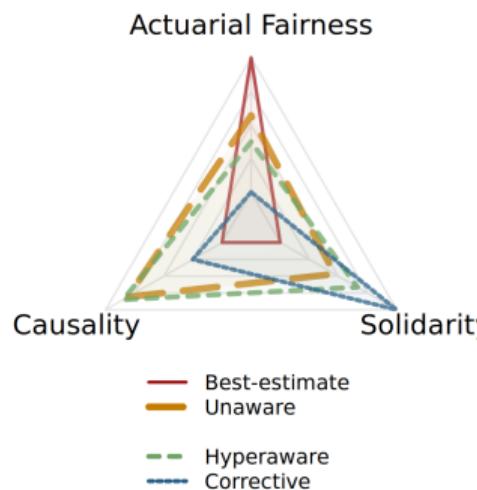
# Grasping behavior of families

41/53



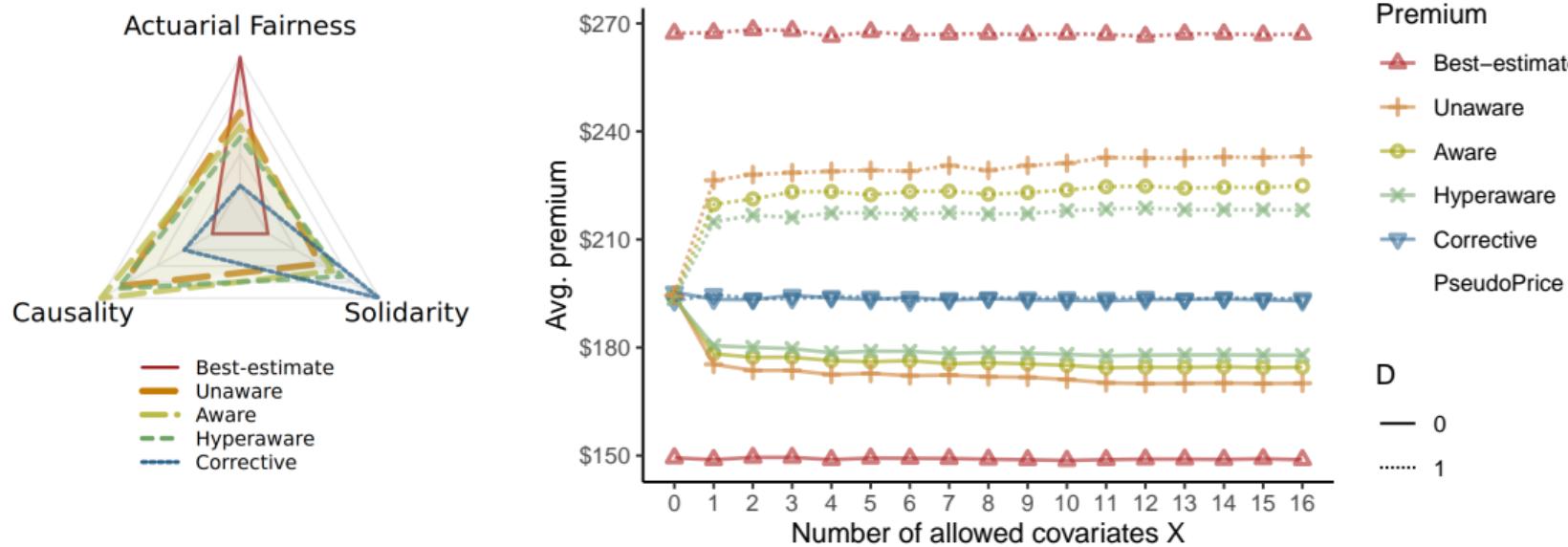
# Grasping behavior of families

41/53



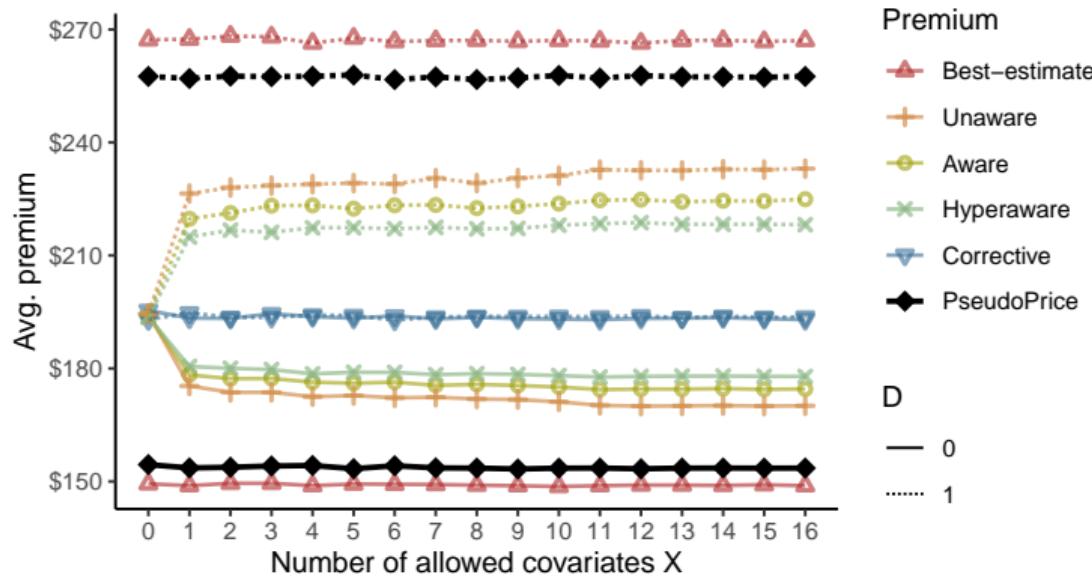
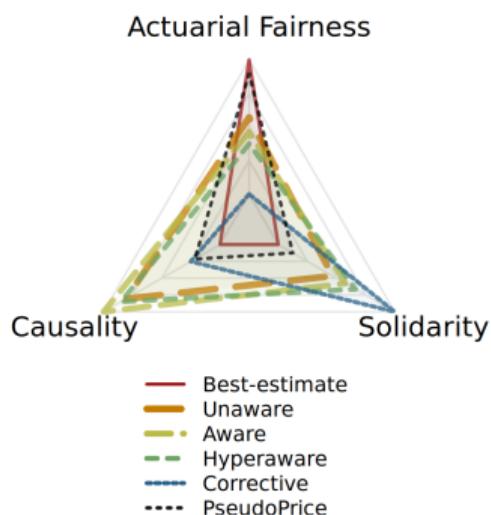
# Grasping behavior of families

41/53



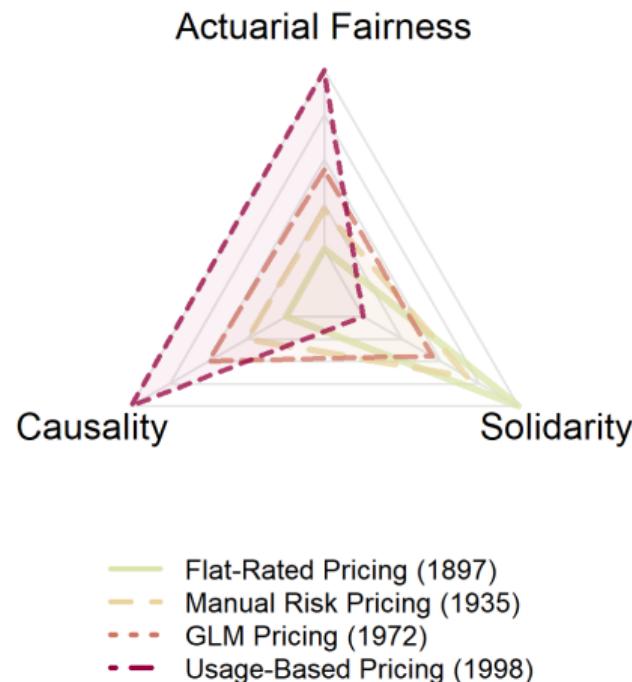
# Grasping behavior of families

41/53



# The history of risk classification revisited

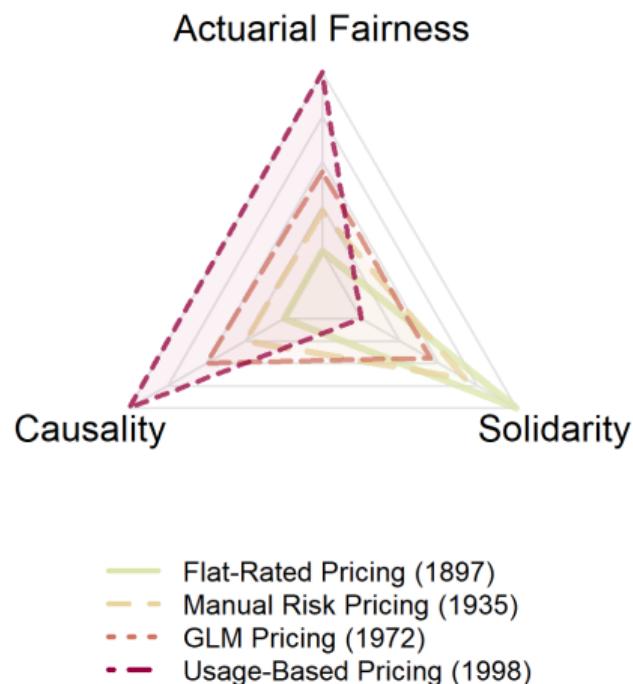
42/53



# The history of risk classification revisited

42/53

As data granularity increases, so does the potential for actuarial justification in perpetuating disparities.



## Prior-pricing local fairness metrics

43/53

The potential for disparate treatment on  $D$  is the **risk spread** :

$$\Delta_{\text{risk}}(\mathbf{x}) = \sup_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d) - \inf_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d).$$

## Prior-pricing local fairness metrics

43/53

The potential for disparate treatment on  $D$  is the **risk spread** :

$$\Delta_{\text{risk}}(\mathbf{x}) = \sup_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d) - \inf_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d).$$

The vulnerability of a segment  $\mathbf{x}$  to proxy effect is the **proxy vulnerability** :

$$\Delta_{\text{proxy}}(\mathbf{x}) = \mu^U(\mathbf{x}) - \mu^A(\mathbf{x}).$$

## Prior-pricing local fairness metrics

43/53

The potential for disparate treatment on  $D$  is the **risk spread** :

$$\Delta_{\text{risk}}(\mathbf{x}) = \sup_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d) - \inf_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d).$$

The vulnerability of a segment  $\mathbf{x}$  to proxy effect is the **proxy vulnerability** :

$$\Delta_{\text{proxy}}(\mathbf{x}) = \mu^U(\mathbf{x}) - \mu^A(\mathbf{x}).$$

The **fairness range** shows how much prices vary across fairness methods

$$\Delta_{\text{fair}}(\mathbf{x}, d) = \sup_{a \in \{B, U, A, H, C\}} \mu^a(\mathbf{x}, d) - \inf_{b \in \{B, U, A, H, C\}} \mu^b(\mathbf{x}, d)$$

The **parity cost** is the (monetary) cost of enforcing premium parity :

$$\Delta_{\text{cost}}(\mathbf{x}, d) = \mu^C(\mathbf{x}, d) - \mu^B(\mathbf{x}, d).$$

## Prior-pricing local fairness metrics

43/53

The potential for disparate treatment on  $D$  is the **risk spread** :

$$\Delta_{\text{risk}}(\mathbf{x}) = \sup_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d) - \inf_{d \in \mathcal{D}} \mu^B(\mathbf{x}, d).$$

The vulnerability of a segment  $\mathbf{x}$  to proxy effect is the **proxy vulnerability** :

$$\Delta_{\text{proxy}}(\mathbf{x}) = \mu^U(\mathbf{x}) - \mu^A(\mathbf{x}).$$

The **fairness range** shows how much prices vary across fairness methods

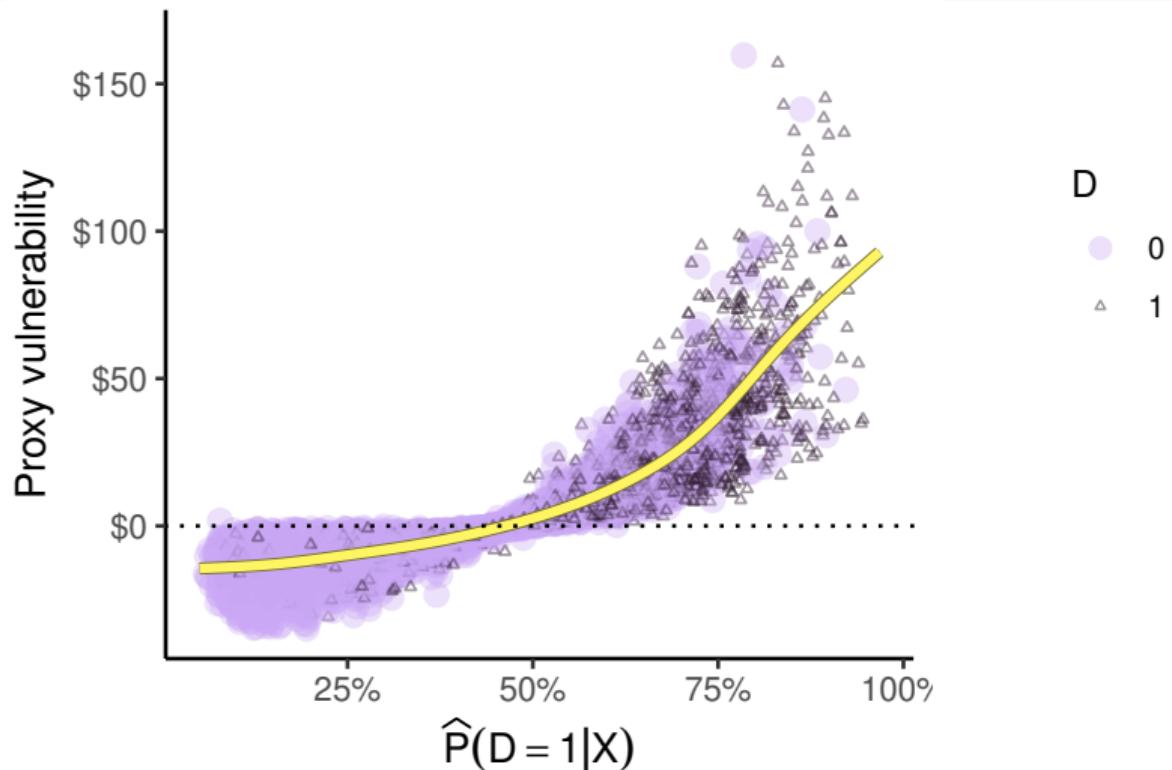
$$\Delta_{\text{fair}}(\mathbf{x}, d) = \sup_{a \in \{B, U, A, H, C\}} \mu^a(\mathbf{x}, d) - \inf_{b \in \{B, U, A, H, C\}} \mu^b(\mathbf{x}, d)$$

The **parity cost** is the (monetary) cost of enforcing premium parity :

$$\Delta_{\text{cost}}(\mathbf{x}, d) = \mu^C(\mathbf{x}, d) - \mu^B(\mathbf{x}, d).$$

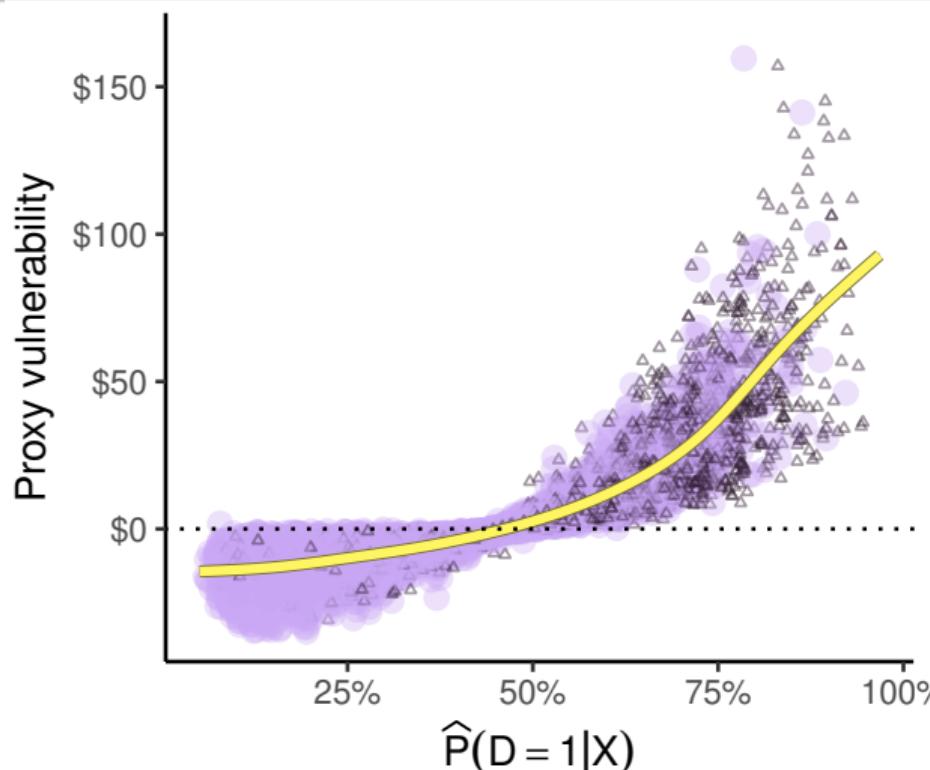
# Visualizing the proxy vulnerability

44/53



# Visualizing the proxy vulnerability

44/53

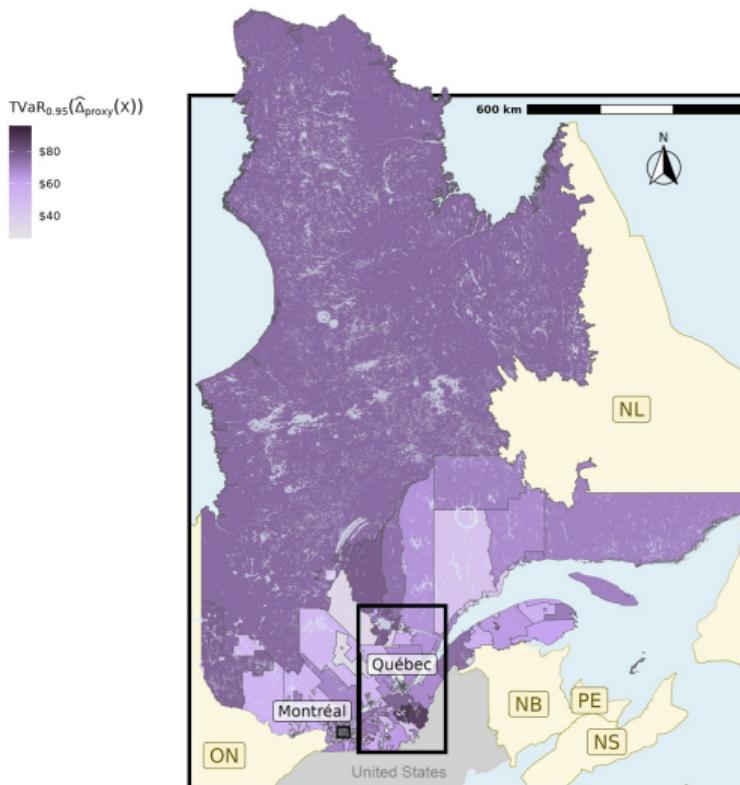


D  
● 0  
△ 1

Proxy vulnerability is both material and skewed

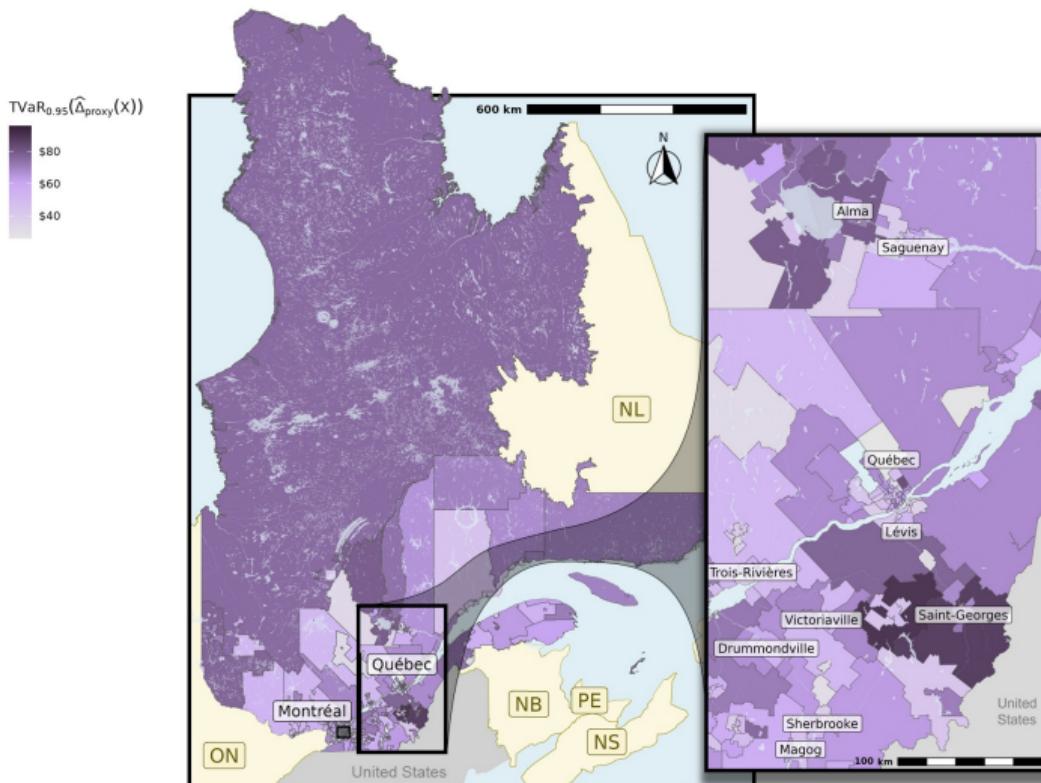
# Geographic distribution of the 95% TVaR of proxy vulnerability

45/53



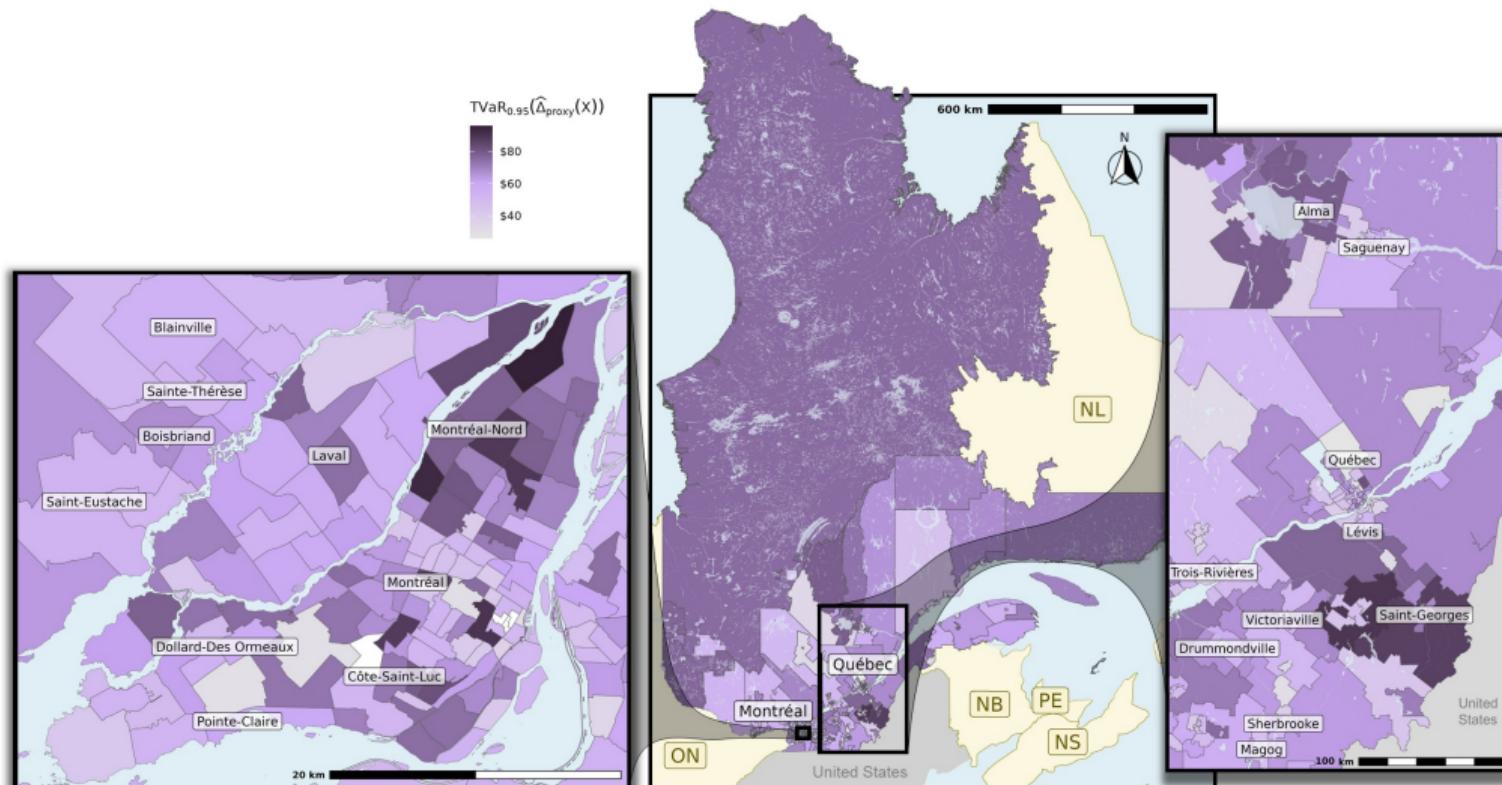
# Geographic distribution of the 95% TVaR of proxy vulnerability

45/53



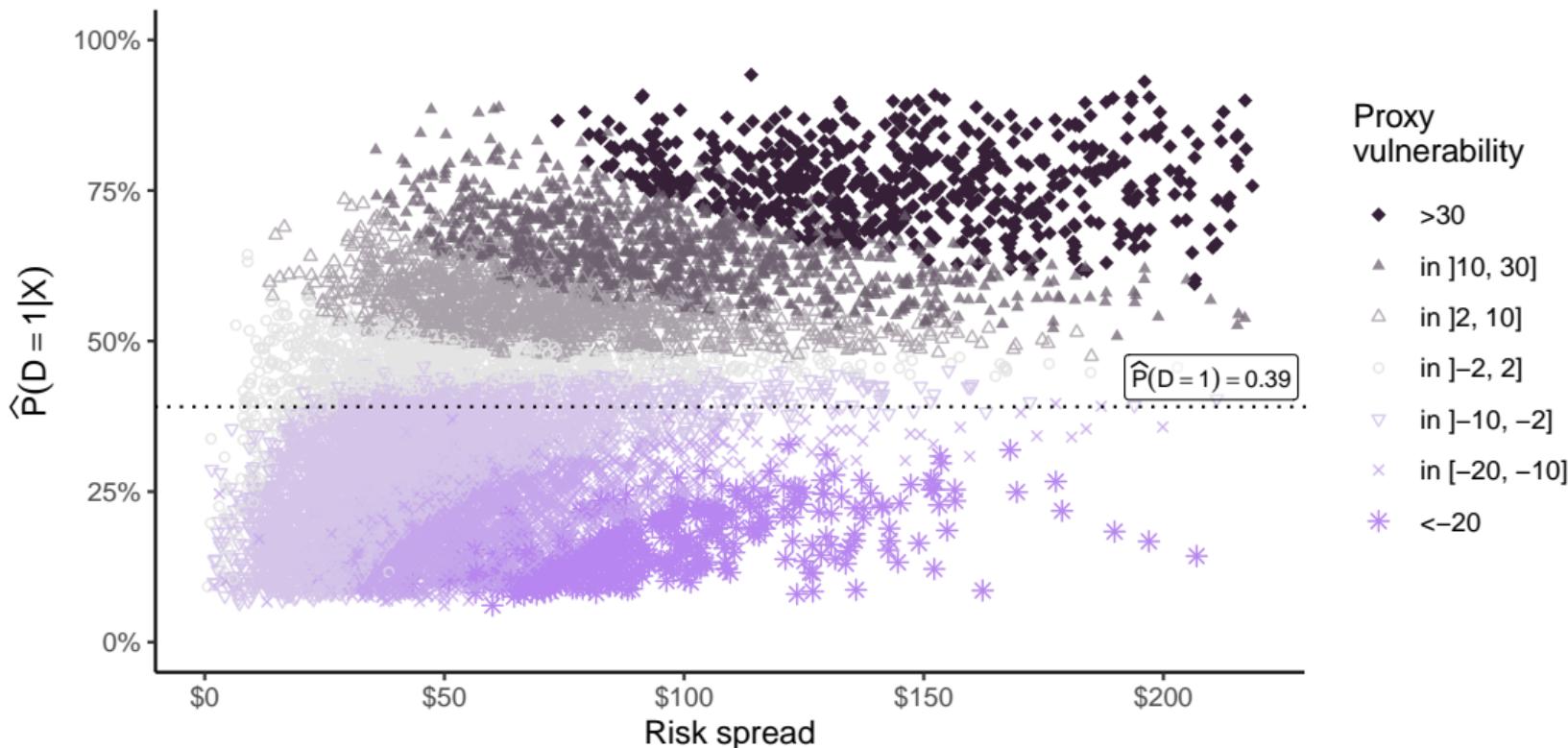
# Geographic distribution of the 95% TVaR of proxy vulnerability

45/53



# Decomposing proxy vulnerability

46/53



## Local metrics between commercial premium and fair benchmarks

47/53

The **commercial loading** relative to the reference premium  $\mu^A$  is :

$$\Delta_{\text{load}}(\mathbf{x}, d; \pi) = \pi(\mathbf{x}, d) - \mu^A(\mathbf{x}),$$

The **implied propensity** is the implicit\* weight of  $D = 1$  in  $\pi(\mathbf{x})$  :

$$\tilde{P}_D(\mathbf{x}; \pi) = \frac{\pi(\mathbf{x}) - \mu^B(\mathbf{x}, 0)}{\mu^B(\mathbf{x}, 1) - \mu^B(\mathbf{x}, 0)}.$$

The **excess lift** measures excessive price differentiation on  $D$  :

$$\Delta_{\text{excess}}(\mathbf{x}; \pi) = \left\{ \sup_d \pi(\mathbf{x}, d) - \inf_d \pi(\mathbf{x}, d) \right\} - \Delta_{\text{risk}}(\mathbf{x}).$$

## Revealing disparities via partitioning

48/53

**Problem :** Fairness metrics over large groups can hide disparities within vulnerable subgroups.

## Revealing disparities via partitioning

48/53

**Problem :** Fairness metrics over large groups can hide disparities within vulnerable subgroups.

**Solution :** Partition policyholders based on fairness-relevant metrics.

## Revealing disparities via partitioning

48/53

**Problem :** Fairness metrics over large groups can hide disparities within vulnerable subgroups.

**Solution :** Partition policyholders based on fairness-relevant metrics.

**Method :**

- Use (optimal) decision trees for simplicity and interpretability.
- Partition based on **proxy vulnerability**  $\Delta_{\text{proxy}}(\mathbf{X})$ .
- Partition based on **commercial loading**  $\Delta_{\text{load}}(\mathbf{X})$ .

## Revealing disparities via partitioning

48/53

**Problem :** Fairness metrics over large groups can hide disparities within vulnerable subgroups.

**Solution :** Partition policyholders based on fairness-relevant metrics.

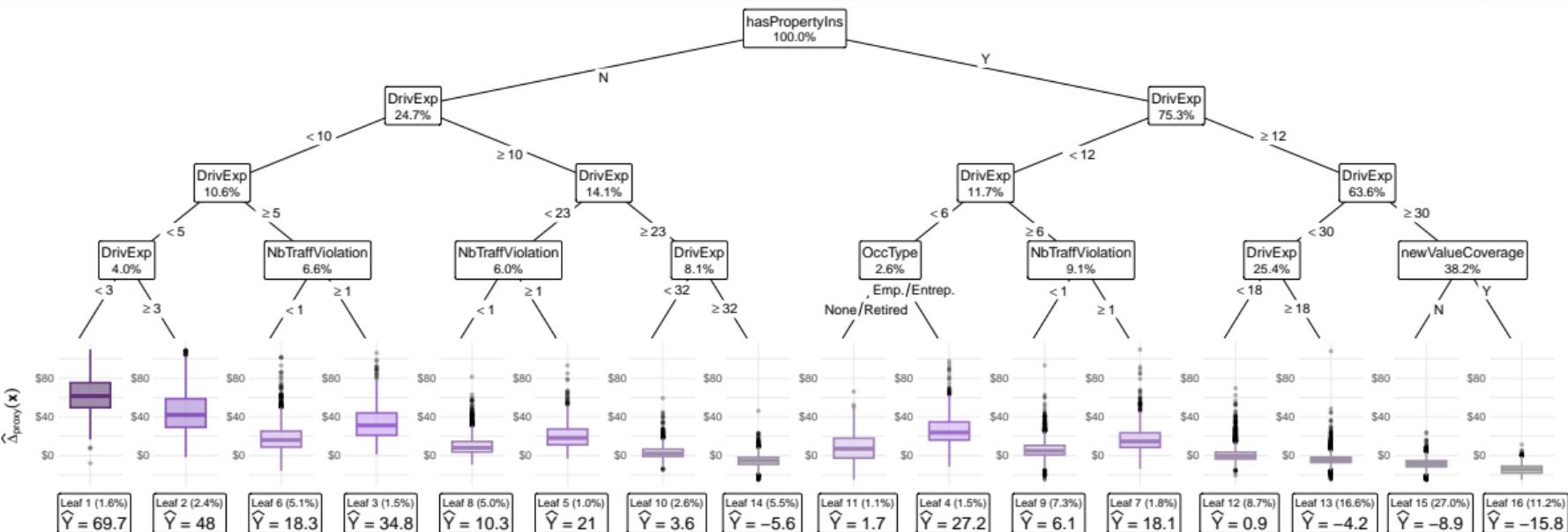
**Method :**

- Use (optimal) decision trees for simplicity and interpretability.
- Partition based on **proxy vulnerability**  $\Delta_{\text{proxy}}(\mathbf{X})$ .
- Partition based on **commercial loading**  $\Delta_{\text{load}}(\mathbf{X})$ .

**Goal :** Identify groups that are systematically **vulnerable** and **overloaded**.

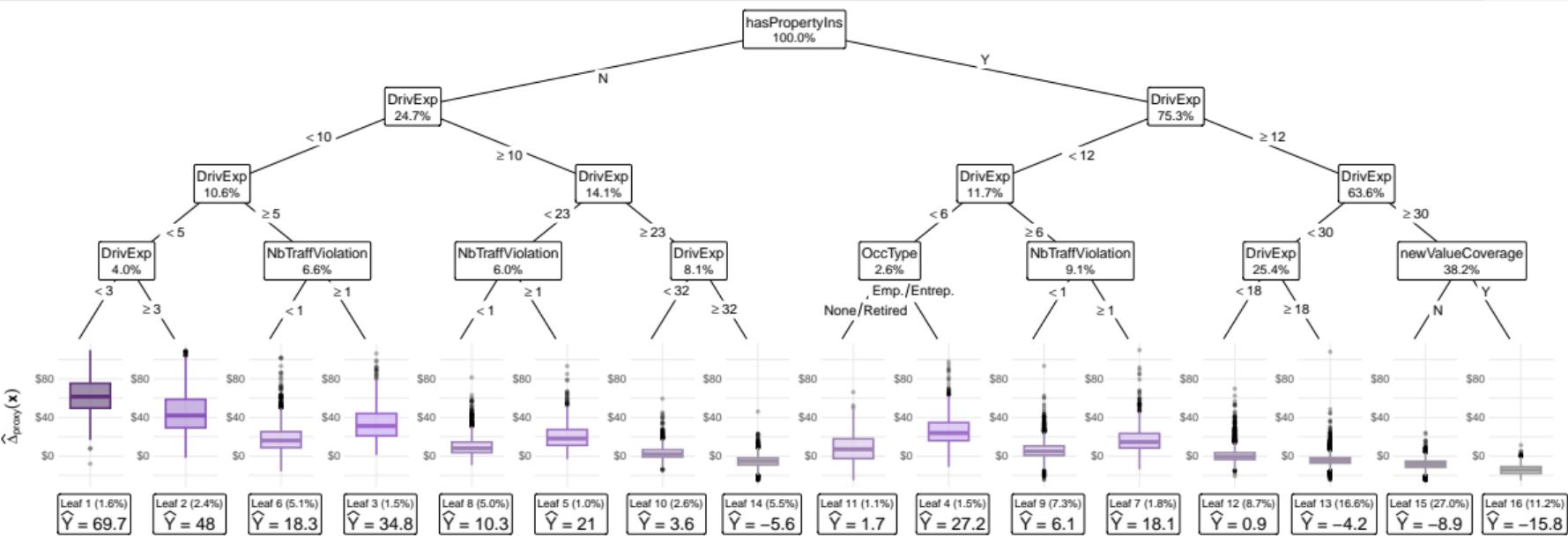
# Partitioning policyholders following proxy vulnerability

49/53



## Partitioning policyholders following proxy vulnerability

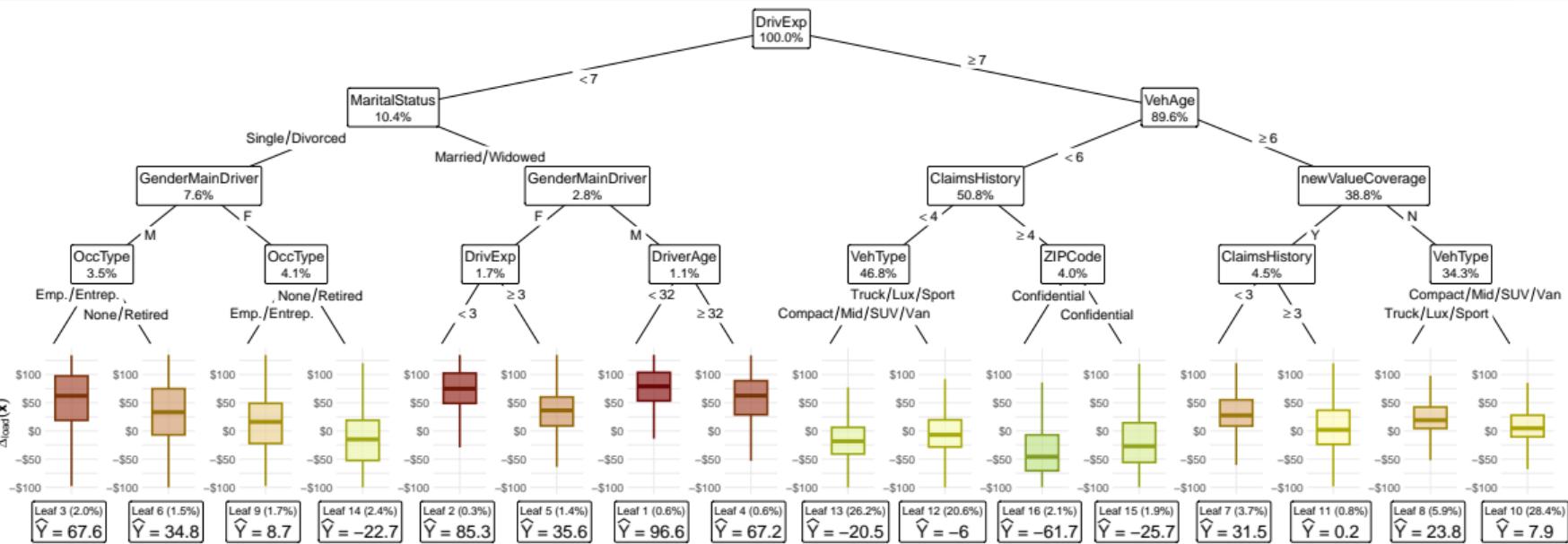
49/53



Leaves with high predictions signal **segments where proxy effects are most likely** to be exploited by a model.

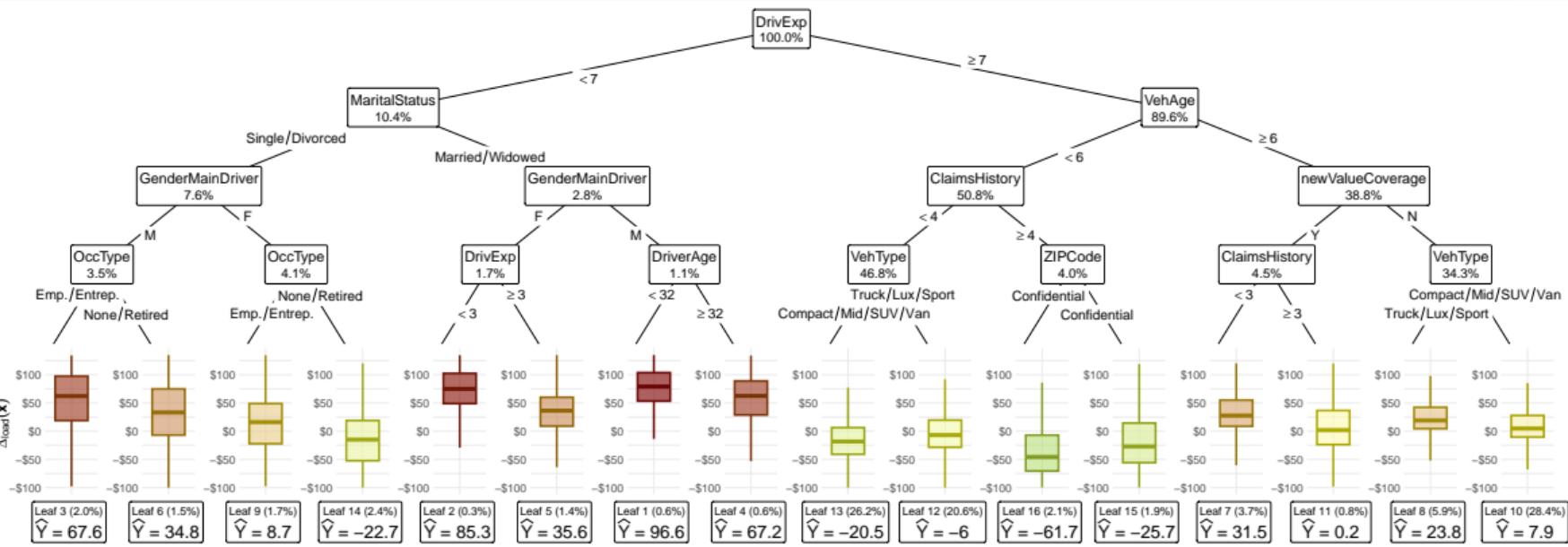
# Partitioning policyholders following commercial loading

50/53



# Partitioning policyholders following commercial loading

50/53



Monitoring commercial loadings means **monitoring the cumulative impact of pricing decisions** – likely beyond the insurer's intent or awareness.

# An integrated fairness assessment framework

51/53

Is the pseudo-price fair?

Section	Variable	Statistic	Subpo	All Data
Demographics	Exposure	Sum	All	164,064
	Credrisk (lvl %)	Level %	Credrisk = 1	37.5%
	DrivExp	Mean	All	#####
	VehAge	Mean	All	#####
	DrivAge	Mean	All	#####
	Zip Code (first character only)	Level %	J	#####
		Level %	G	#####
		Level %	H	#####
		Level %	Employed	60.6%
	OccType	Level %	Retired	18.7%
HasPropertyIns		Level %	None	10.3%
		Level %	Other	10.5%
		Level %	Y	#####
		Level %	lvl % : N	#####
	Loss	% of 0	All	97.2%
		Mean	All	189.25
	Severity	Mean	All	3743
		TVaR 0.95	All	19093
	Pseudoprice	VaR 0.05	All	68.40
		Mean	All	191.25
		VaR 0.95	All	401.59
		MAE	All	272.15
Loss and pseudo price metrics	Performance	Avg. Dev. Twee.	All	137.20
	Price	Mean	D=0	152.74
		Mean	D=1	255.31
		VaR 0.95	D=0	288.64
		VaR 0.95	D=1	497.85
	Loss Ratio	-	D=0	87.4%
		-	D=1	102.3%
	Est. P(D=1 X)	Mean	All	0.49
	Wass. Dist.	Prem(d = 0) vs Prem(d = 1)		103.99

# An integrated fairness assessment framework

51/53

Is the pseudo-price fair?

Section	Variable	Statistic	Subpop.	All Rule
Demographic	Exposure	Sum	All	94464
	CreditRisk (M+N)	Level %	CreditRisk = 1	37.5%
	DriveType	Mean	All	*****
	Voltage	Mean	All	*****
	Divide	Mean	All	*****
	Zip Code (first character only)	Level %	I	*****
		Level %	G	*****
		Level %	H	*****
		Level %	F	*****
		Level %	E	*****
Occupation	Profession	Mean	All	19.4%
	OccType	Level %	Retired	18.7%
		Level %	New	18.3%
		Level %	Other	18.5%
	HasProperty	Level %	N	*****
Loss and pseudo price	Severity	% of 0	All / N	0.0000
		Mean	All	97.2%
		Mean	All	189.35
		Mean	All	3743
		Mean	All	189.05
		Tall E/05	All	69.44
	Pseudoprice	Mean	All	191.25
		Tall E/05	All	481.58
	Performance	Mean	All	292.13
	Metrics	Sum / Avg / All	All	137.28
Loss and pseudo price group	Price	Mean	D=0	152.74
		Mean	D=1	255.31
		Tall E/05	D=0	288.69
		Tall E/05	D=1	449.68
	Loss Ratio	-	D=0	87.89
	Est. P[D=1 X]	Mean	All	0.48
Wass. Dist.	Mean(E[X=0] vs P[X=1])	All	185.98	

# An integrated fairness assessment framework

51/53

## Is the pseudo-price fair?

- Partitions define relevant subgroups (columns).

Variable	Modele	Subpo	AllData	Proxy vulnerable groups				Commercially insured groups				
				1	2	3	4	5	6	7	8	
Exposure	Sum	All	104,694	2,405	3,082	2,235	7,797	6,994	44,793	19,227	652	419
CreditRisk (M+N)	Level %	CreditRisk = 1	37.5%	81.0%	74.0%	70.3%	48.1%	34.9%	23.9%	18.3%	56.7%	73.3%
DriveTyp	Mean	All	1.44	3.53	6.82	3.78	4.246	4.229	42.24	4.49	1.31	4.23
Voltage	Mean	All	5.22	4.93	4.31	5.96	7.57	1.79	4.31	4.71	4.62	5.52
DivDige	Mean	All	21.89	22.33	24.36	22.04	22.05	22.16	22.16	24.23	20.49	22.57
Zip Code (first character only)	Level %	I										
	Level %	G										
	Level %	H										
	Level %	J										
	Level %	P										
	Level %	R										
	Level %	Refund										
GenType	Level %	New	18.9%	58.2%	59.7%	26.1%	11.0%	5.7%	4.9%	3.4%	14.9%	25.6%
	Level %	Other	18.5%									
	Level %	Old										
	Level %	Y										
	Level %	N										
	Level %	id %: N										
			100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Loss and accident price	%of0	All	97.2%	94.1%	95.1%	94.2%	93.1%	97.5%	97.9%	94.0%	96.1%	94.8%
	Mean	All	189.25	593.73	414.37	322.26	193.16	134.01	189.18	187.61	796.37	315.87
Severity	Mean	All	374.3	466.8	454.4	431.9	394.6	329.9	394.6	359.8	812.6	423.3
	Total Loss	All	60.44									
	Mean	All	191.25	593.67	416.38	355.20	230.03	151.21	125.18	186.18	362.83	399.29
	StdDev	All	481.59									
	Total Loss	All	272.39	742.39	586.57	403.58	286.85	234.09	195.88	247.79	606.87	545.53
	Mean	All	177.39	236.39	230.34	252.27	180.97	117.06	161.17	192.39	221.77	232.25
Performance metrics	Mean	All	148.1	212.1	174.2	152.1	135.2	118.6	175.6	161.51	301.75	326.45
	Total Loss	All	177.39	236.39	230.34	252.27	180.97	117.06	161.17	192.39	221.77	232.25
	Mean	All	102.9%	104.6%	103.5%	105.9%	103.0%	94.9%	105.9%	127.9%	87.6%	110.6%
	StdDev	All	102.9%	104.6%	103.5%	105.9%	103.0%	94.9%	105.9%	127.9%	87.6%	110.6%
	Loss Ratio	All	87.8%	88.1%	85.4%	88.4%	86.8%	84.2%	88.3%	93.9%	26.3%	31.4%
	Mean	All	0.48	0.82	0.76	0.73	0.54	0.38	0.28	0.21	0.77	0.80
	Est. P[D=1 X]	Mean	All	0.48	0.82	0.76	0.73	0.54	0.38	0.28	0.21	0.77
	Wass.Dist.	Mean	D[0:1] vs D[1:2]	185.99	115.22	122.79	94.74	67.14	48.01	49.66	92.84	99.20

## An integrated fairness assessment framework

51/53

## Is the pseudo-price fair?

- Partitions define relevant subgroups (columns).
  - Local metrics quantify dissection unfairness (rows).

# An integrated fairness assessment framework

51/53

## Is the pseudo-price fair?

- Partitions define relevant subgroups (columns).
- Local metrics quantify dissect unfairness (rows).

The toolbox guide **fairness assessment**.

Variable	Model	Subpop.	All Model	Proxy vulnerable groups				Commercially loaded groups				
				1	2	3	4	5	6	7	8	
Exposure	Sum	All	104,694	2,405	3,082	2,235	7,979	3,694	44,793	19,227	652	419
CreditRisk	Level %	CreditRisk = 1	37.9%	81.0%	74.0%	70.3%	48.3%	33.6%	23.9%	18.3%	56.7%	73.3%
DriveType	Mean	All	1.44	3.53	5.62	4.82	4.24	4.24	4.24	4.24	4.49	1.31
Voltage	Mean	All	5.22	4.93	4.31	5.96	7.57	1.79	4.31	4.71	4.62	5.52
Divide	Mean	All	21.89	22.33	24.36	22.05	22.05	22.05	22.05	22.05	24.23	20.49
Zip Code (first character only)	Level %	I	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Zip Code (first character only)	Level %	G	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Zip Code (first character only)	Level %	H	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Zip Code (first character only)	Level %	P	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Zip Code (first character only)	Level %	R	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
OnType	Refund	All	18.7%	9.9%	9.9%	9.9%	9.9%	9.9%	9.9%	9.9%	9.9%	9.9%
OnType	New	All	18.9%	58.2%	59.7%	26.1%	1.0%	5.7%	4.9%	3.4%	14.9%	25.6%
HasProperty	Yes	All	18.5%	58.2%	59.7%	26.1%	1.0%	5.7%	4.9%	3.4%	8.4%	18.0%
HasProperty	No	All	18.5%	58.2%	59.7%	26.1%	1.0%	5.7%	4.9%	3.4%	8.4%	18.0%
Ind % / N	All	All	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
Loss	Mean	All	97.29	94.1%	95.1%	94.2%	93.1%	95.9%	94.8%	94.9%	96.1%	94.8%
Severity	Mean	All	189.25	593.73	414.37	322.26	303.12	134.01	188.01	187.87	796.87	315.87
Severity	SD	All	374.3	466.8	454.4	431.3	394.0	309.8	334.6	329.5	812.6	433.2
PenPricing	Mean	All	191.25	593.87	416.38	355.20	393.08	151.21	225.18	186.18	362.83	399.29
Performance	Mean	All	481.59	1,020.25	742.25	742.25	742.25	742.25	742.25	742.25	742.25	742.25
Performance	SD	All	177.28	726.39	230.34	232.27	180.77	117.06	161.17	192.39	221.77	232.29
Loss and stand price	Mean	All	152.74	410.44	225.45	200.60	171.40	135.2	118.63	175.62	301.75	326.45
Loss and stand price	Price	All	255.23	530.93	447.86	322.84	302.87	182.89	182.89	182.89	360.87	545.53
Loss and stand price	Tariff	All	288.66	489.22	380.78	336.82	336.82	236.82	236.82	236.82	381.53	425.26
Loss Ratio	-	All	87.89	88.1%	95.4%	88.4%	88.0%	87.89	89.9%	89.9%	88.3%	88.3%
Loss Ratio	Mean	All	102.99	105.6%	103.5%	106.9%	105.9%	94.9%	105.9%	127.9%	77.6%	104.9%
Loss Ratio	SD	All	102.99	105.6%	103.5%	106.9%	105.9%	94.9%	105.9%	127.9%	77.6%	104.9%
Exp. P[D=0 D=1]	Mean	All	157.74	417.34	224.80	200.60	171.40	135.2	118.63	175.62	301.75	326.45
Exp. P[D=0 D=1]	SD	All	255.23	530.93	447.86	322.84	302.87	182.89	182.89	182.89	360.87	545.53
Exp. P[D=0 D=1]	Mean	All	288.66	489.22	380.78	336.82	336.82	236.82	236.82	236.82	381.53	425.26
Exp. P[D=0 D=1]	SD	All	288.66	489.22	380.78	336.82	336.82	236.82	236.82	236.82	381.53	425.26
Forwards	Indep. Spurious	All	191.16	327.87	477.34	224.80	200.60	171.40	135.2	118.63	301.75	347.85
Forwards	Indep. Spurious	Unknown	191.16	327.87	477.34	224.80	200.60	171.40	135.2	118.63	301.75	347.85
Forwards	Indep. Spurious	Aware	191.16	404.35	389.01	366.31	370.93	159.05	128.54	200.63	271.72	317.39
Forwards	Indep. Spurious	Hyperaware	191.16	429.22	380.78	336.82	284.81	177.02	223.18	216.09	247.88	268.38
Forwards	Indep. Spurious	Corrective	191.16	327.87	477.34	224.80	200.60	171.40	135.2	118.63	301.75	347.85
Proxy vulnerability	Mean	All	93.09	43.54	43.54	43.54	43.54	43.54	43.54	43.54	43.54	43.54
Proxy vulnerability	SD	All	7.88	40.89	61.17	36.38	7.88	-3.85	-4.62	-14.82	0.31	0.80
Proxy vulnerability	Mean	All	88.08	35.03	35.13	31.03	12.8	-6.81	-9.18	-16.18	-23.81	-43.83
Proxy vulnerability	SD	All	87.97	99.04	99.15	99.15	99.15	99.15	99.15	99.15	99.15	99.15
Risk spread	Mean	All	65.98	170.05	350.62	330.62	218.5	81.53	42.88	45.62	304.85	324.58
Risk spread	Mean	All	55.38	161.05	271.86	118.66	41.62	42.49	34.43	54.87	66.73	101.53
Risk spread	SD	All	72.27	94.38	214.22	114.22	114.22	114.22	114.22	114.22	114.22	114.22
Party cost	Mean	All	0.09	-20.18	-20.89	-44.23	-2.66	3.07	0.39	-1.29	42.83	31.65
Party cost	SD	All	-0.68	-62.28	-34.20	-56.88	-9.46	3.02	0.84	-1.49	78.71	63.82
Party cost	Mean	All	5.38	34.74	18.12	-15.93	4.37	3.55	0.08	2.04	77.80	72.89
Party cost	SD	All	69.77	72.27	49.85	49.85	49.85	49.85	49.85	49.85	49.85	49.85
Commercial loading	Mean	All	18.7%	8.0%	-1.4%	-7.8%	5.3%	16.3%	24.0%	4.1%	24.5%	17.6%
Commercial loading	SD	All	-11.7	-76.79	-14.63	-14.63	-4.41	-1.16	-9.38	-19.81	-26.77	-24.74
Commercial loading	Mean	All	54.76	27.04	27.04	27.04	27.04	27.04	27.04	27.04	27.04	27.04
Commercial loading	SD	All	48.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06
Commercial loading	Mean	All	63.59	69.0%	62.5%	37.8%	38.7%	61.2%	68.5%	59.7%	92.0%	81.4%
Commercial loading	SD	All	48.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06	23.06
Commercial loading	Mean	All	38.2%	38.2%	38.2%	38.2%	38.2%	38.2%	38.2%	38.2%	38.2%	38.2%
Commercial loading	SD	All	38.49	38.49	38.49	38.49	38.49	38.49	38.49	38.49	38.49	38.49
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	SD	All	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%	23.9%
Commercial loading	Mean	All										

## Question

Is the given pricing  
function  $\pi(\mathbf{x}, d)$  fair?

## Inputs

A pricing function:  $\pi(\mathbf{x}, d)$   
Data:  $\{\mathbf{X}_i, D_i, Y_i\}_{i=1}^n$

## Question

Is the given pricing  
function  $\pi(\mathbf{x}, d)$  fair?

## Inputs

A pricing function:  $\pi(\mathbf{x}, d)$   
Data:  $\{\mathbf{X}_i, D_i, Y_i\}_{i=1}^n$

What does it mean for pricing to be fair? → Can one benchmark fairness? → Is unfairness truly material? → Where does unfairness concentrate?

## Question

Is the given pricing function  $\pi(x, d)$  fair?

## Inputs

A pricing function:  $\pi(x, d)$   
Data:  $\{X_i, D_i, Y_i\}_{i=1}^n$

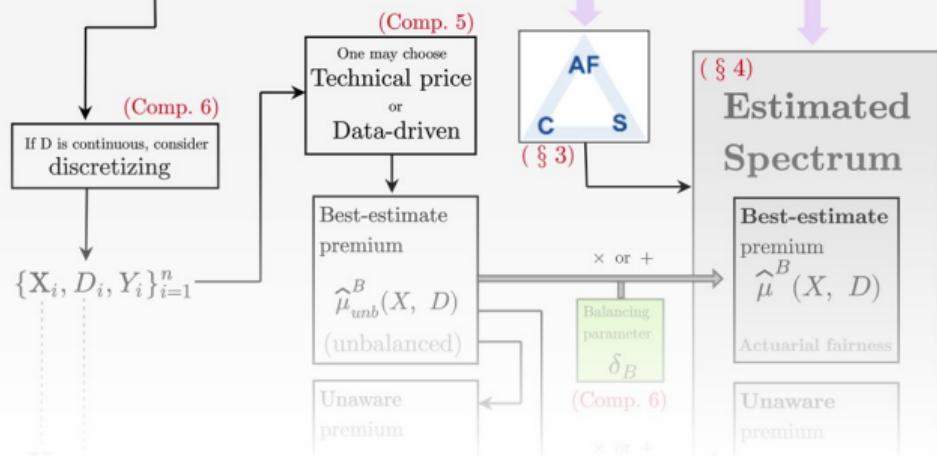
What does it mean for pricing to be fair?

Can one benchmark fairness?

Is unfairness truly material?

Where does unfairness concentrate?

## A Scalable toolbox for exposing indirect discrimination in insurance rates



## discrimination in insurance rates

(§ 5)  
Actuarially relevant local fairness metrics

(§ 6)  
Policyholder partitioning

Pre-pricing  
• Risk spread  
• Proxy vulnerability  
 $\Delta_{risk}(X)$

Pre-pricing following proxy  
 $\Delta_{proxy}(X)$

**Question**

Is the given pricing function  $\pi(x, d)$  fair?

**Inputs**

A pricing function:  $\pi(x, d)$   
Data:  $\{X_i, D_i, Y_i\}_{i=1}^n$

What does it mean for pricing to be fair?

Can one benchmark fairness?

Is unfairness truly material?

Where does unfairness concentrate?

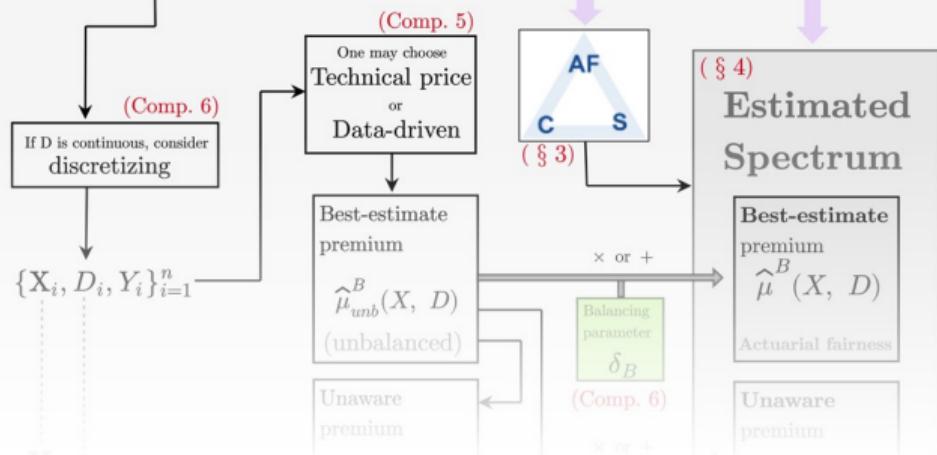
The toolbox guides assessment.



(Fig. 14)



# A Scalable toolbox for exposing indirect discrimination in insurance rates



# discrimination in insurance rates

(§ 5)  
Actuarially relevant local fairness metrics

## Pre-pricing

- Risk spread
- Proxy vulnerability

$$\Delta_{risk}(X)$$

$$\Delta_{proxy}(X)$$

(§ 6)  
Policyholder partitioning

Pre-pricing following proxy

## Moving forward

53/53

This study had a specific scope. Advancing fairness requires expanding it :

- 1 Fairness often assumes **access to protected attributes**, which may unavailable.

## Moving forward

53/53

This study had a specific scope. Advancing fairness requires expanding it :

- 1 Fairness often assumes **access to protected attributes**, which may unavailable.
- 2 Applying anti-discrimination **regulations** meets resistance where actuarial justification holds authority.

## Moving forward

53/53

This study had a specific scope. Advancing fairness requires expanding it :

- 1 Fairness often assumes **access to protected attributes**, which may unavailable.
- 2 Applying anti-discrimination **regulations** meets resistance where actuarial justification holds authority.
- 3 Fairness is typically studied as a one-year objective, but its **long-term welfare** effects remain unclear (Shimao et al., 2022).

## Moving forward

53/53

This study had a specific scope. Advancing fairness requires expanding it :

- 1 Fairness often assumes **access to protected attributes**, which may unavailable.
- 2 Applying anti-discrimination **regulations** meets resistance where actuarial justification holds authority.
- 3 Fairness is typically studied as a one-year objective, but its **long-term welfare** effects remain unclear (Shimao et al., 2022).
- 4 Portfolio fairness may conflict with market fairness (Côté et al., 2024).

# Thank you



## References i

53/53

Araiza Iturria, C. A., Hardy, M., and Marriott, P. (2024). A discrimination-free premium under a causal framework. *North American Actuarial Journal*, pages 1–21.

ASB (2005). Actuarial Standard of Practice No. 12 : Risk Classification (for All Practice Areas). [actuarialstandardsboard.org/asops/risk-classification-practice-areas/](http://actuarialstandardsboard.org/asops/risk-classification-practice-areas/). Accessed : March 31, 2025.

Aseervatham, V., Lex, C., and Spindler, M. (2016). How do unisex rating regulations affect gender differences in insurance premiums ? *The Geneva Papers on Risk and Insurance-Issues and Practice*, 41(1) :128–160.

Boucher, J.-P. and Pigeon, M. (2024). Balancing risk assessment and social fairness : an auto telematics case study. *Casualty Actuarial Society Research Paper Series on Race and Insurance Pricing*. Retrieved from [casact.org/publications-research/research/research-paper-series-race-and-insurance-pricing](http://casact.org/publications-research/research/research-paper-series-race-and-insurance-pricing). Accessed : March 31, 2025.

## References ii

53/53

Calmon, F., Wei, D., Vinzamuri, B., Natesan Ramamurthy, K., and Varshney, K. R. (2017). Optimized pre-processing for discrimination prevention. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30.

Casualty Actuarial Society (1988). Statement of Principles Regarding Property and Casualty Insurance Ratemaking.

[casact.org/statement-principles-regarding-property-and-casualty-insurance-ratemaking](https://casact.org/statement-principles-regarding-property-and-casualty-insurance-ratemaking). Originally published in 1988, rescinded in 2020, and reinstated in 2021 for reference. Accessed : March 31, 2025.

Charpentier, A., Flachaire, E., and Gallic, E. (2023a). Optimal transport for counterfactual estimation : A method for causal inference. In *Optimal Transport Statistics for Economics and Related Topics*, pages 45–89. Springer.

Charpentier, A., Hu, F., and Ratz, P. (2023b). Mitigating discrimination in insurance with wasserstein barycenters. *BIAS, 3rd Workshop on Bias and Fairness in AI, International Workshop of ECML PKDD*.

## References iii

53/53

- Côté, M.-P., Côté, O., and Charpentier, A. (2024). Selection bias in insurance : why portfolio-specific fairness fails to extend market-wide. *Available at SSRN* : 10.2139/ssrn.5018749.
- Côté, O., Côté, M.-P., and Charpentier, A. (2025). A fair price to pay : Exploiting causal graphs for fairness in insurance. *Journal of Risk and Insurance*, 92 :33-75.
- Embrechts, P. and Wüthrich, M. V. (2022). Recent challenges in actuarial science. *Annual Review of Statistics and Its Application*, 9(1) :119-140.
- European Union (2000). Council Directive 2000/43/EC of 29 June 2000 implementing the principle of equal treatment between persons irrespective of racial or ethnic origin.  
<https://eur-lex.europa.eu/eli/dir/2000/43/oj>. Accessed : December 4, 2024.
- Fermanian, J.-D. and Guegan, D. (2021). Fair learning with bagging. *Documents de travail du Centre d'Économie de la Sorbonne*.
- Fernandes Machado, A., Charpentier, A., and Gallic, E. (2024). Sequential conditional transport on probabilistic graphs for interpretable counterfactual fairness.

## References iv

53/53

- Fernandes Machado, A., Grondin, S., Ratz, P., Charpentier, A., and Hu, F. (2025). Equipy : Sequential fairness using optimal transport in python. *arXiv preprint*, arXiv :2503.09866.
- Financial Services Regulatory Authority of Ontario (2024). Proposed Guidance : Automobile Insurance Rating and Underwriting Supervision.  
[fsrao.ca/industry/auto-insurance/regulatory-framework/guidance-auto-insurance/proposed-guidance-automobile-insurance-rating-and-underwriting-supervision-guidance](http://fsrao.ca/industry/auto-insurance/regulatory-framework/guidance-auto-insurance/proposed-guidance-automobile-insurance-rating-and-underwriting-supervision-guidance). Accessed : March 31, 2025.
- Fredman, S. (2022). *Discrimination Law*. Oxford University Press, Oxford, 3rd edition.
- Frees, E. W. and Huang, F. (2023). The discriminating (pricing) actuary. *North American Actuarial Journal*, 27(1) :2-24.
- Gabric, L. J., Zhou, S., and Zhou, K. Q. (2024). A Bayesian approach to discrimination-free insurance pricing. *Available at SSRN 4785927*.

## References v

53/53

- Grari, V., Charpentier, A., Lamprier, S., and Detyniecki, M. (2022). A fair pricing model via adversarial learning. *arXiv preprint arXiv:2202.12008*.
- Hu, F., Ratz, P., and Charpentier, A. (2023). Fairness in multi-task learning via Wasserstein barycenters. In Koutra, D., Plant, C., Gomez Rodriguez, M., Baralis, E., and Bonchi, F., editors, *Machine Learning and Knowledge Discovery in Databases : Research Track*, pages 295–312, Cham. Springer Nature Switzerland.
- Hu, F., Ratz, P., and Charpentier, A. (2024). A sequentially fair mechanism for multiple sensitive attributes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(11) :12502–12510.
- Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y. (2017). LightGBM : A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30 :3146–3154.
- Kleinberg, J., Mullainathan, S., and Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *arXiv preprint, arXiv:1609.05807*.

## References vi

53/53

- Komiyama, J. and Shimao, H. (2017). Two-stage algorithm for fairness-aware machine learning. *arXiv preprint arXiv:1710.04924*.
- Lindholm, M., Richman, R., Tsanakas, A., and Wüthrich, M. V. (2022). Discrimination-free insurance pricing. *ASTIN Bulletin*, 52(1) :55–89.
- Lindholm, M., Richman, R., Tsanakas, A., and Wüthrich, M. V. (2024). What is fair? Proxy discrimination vs. demographic disparities in insurance pricing. *Scandinavian Actuarial Journal*, 2024(9) :935–970.
- Makhlof, K., Zhioua, S., and Palamidessi, C. (2024). When causality meets fairness : A survey. *Journal of Logical and Algebraic Methods in Programming*, 141 :101000.
- Plečko, D., Bennett, N., and Meinshausen, N. (2021). fairadapt : Causal reasoning for fair data pre-processing. *arXiv preprint arXiv:2110.10200*.
- Pope, D. G. and Sydnor, J. R. (2011). Implementing anti-discrimination policies in statistical profiling models. *American Economic Journal : Economic Policy*, 3(3) :206–31.

## References vii

53/53

- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1) :41–55.
- Shimao, H., Huang, F., and Khern-am nuai, W. (2022). Welfare implications of fairness regulations in insurance cost modeling : A multi-method study. Available at SSRN : 10.2139/ssrn.5112616.
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. (2013). Learning fair representations. In Dasgupta, S. and McAllester, D., editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 325–333.