Olivier Côté, M.Sc., ACAS\*1, Marie-Pier Côté, Ph.D., FSA, ACIA1, and Arthur Charpentier, Ph.D.2

<sup>1</sup>École d'actuariat, Université Laval <sup>2</sup>Département de Mathématiques, Université du Québec à Montréal

October 20, 2025

<sup>\*</sup>Corresponding author: Olivier Côté, olivier.cote.12@ulaval.ca

## **Executive Summary**

According to actuarial standards of practice, insurance pricing relies on grouping policyholders by risk to set adequate premiums. Modern predictive models, especially machine learning, excel at detecting statistical associations to differentiate risks, but they can learn spurious or undesired correlations. This raises concerns when socioeconomic or demographic factors may (intentionally or inadvertently) affect the fairness of insurance pricing.

Fairness in insurance is difficult to operationalize due to its ambiguity. Fairness metrics from the machine learning literature lack the segment-specific relevance actuaries require and are expressed in abstract units that obscure real-world consequences. For actuaries to intervene, proxy effects and unfair biases must be quantified in insurance-relevant terms: dollars and people.

In this paper, we focus on fairness in actuarial pricing. We study the situation where insurance rates should be fair with respect to a categorical (or discretized) sensitive variable, such as race or economic status, and the latter is fully observed (despite the possible privacy challenges). Our main contributions are listed below.

- We argue that actuarial fairness, solidarity, and causality form the three core dimensions of fairness in insurance pricing:
  - Actuarial fairness aligns premiums with expected losses, mitigating crosssubsidies,
  - Solidarity aligns premiums across protected groups, mitigating disparities,
  - Causality ensures models capture only true risk factors, mitigating proxy effects.
- We translate these dimensions into a five-point spectrum of premiums:
  - The best-estimate premium is the most accurate predictor of losses using all available information, including the sensitive variable,
  - The unaware premium is the most accurate predictor of losses using all information except the sensitive variable,
  - The aware premium is the most accurate predictor of losses when controlling for the sensitive variable.
  - The corrective premium is the most accurate predictor that enforces similar premium distributions across levels of the sensitive variable,
  - The hyperaware premium is the most accurate approximation of the corrective premium that does not directly discriminate on the sensitive variable.
- We define actuarially relevant local metrics that quantify the potential monetary impact of unfairness at the policyholder level. Proxy vulnerability is the difference between unaware and aware premiums. It locally measures how much the allowed variables pick up the signal of a missing sensitive variable.

- We define post pricing local metrics to evaluate the fairness of any pricing structure relative to the estimated spectrum.
- We partition policyholders to expose the segments in which unfair discrimination is most severe.
- We integrate these components into a fairness assessment framework that partitions
  the policyholders, pinpoints segments most affected by unfairness, and evaluates
  local metrics to diagnose unfairness and guide intervention.
- We illustrate our approach with a large case study inspired by industry practice. The
  analysis relies a real dataset of approximately 768,000 vehicles insured in Québec
  (2016–2017), covering at-fault material damage claims. We examine the fairness
  of a pseudo commercial price with respect to discretized credit score: low (vulnerable group) vs high. This sensitive variable measures the policyholder's economic
  precariousness.
  - Proxy vulnerability is both material and skewed: while most policyholders may receive a modest rebate, a vulnerable minority of them could face 15–30% overpricing if the regulation only requires that the sensitive variable be omitted,
  - Our integrated framework (Fig. 14) illustrates that fairness in insurance pricing can be assessed efficiently, with minimal analyst effort. The framework provides simultaneous diagnostics from the three fairness dimensions, translates unfairness into dollar terms at the individual level, and highlights disparities across population segments.
- We provide additional information and the complete code illustrated on a comprehensive simulated data example in the online supplementary material.

Designed for routine portfolio monitoring, our toolbox delivers valuable insights whether or not the sensitive attribute is included in pricing, provided it is available for assessment. The toolbox's scalability, across large datasets and rich covariate sets, makes fairness operationalizable for actuaries: intuitive, practical, and encompassing the three fairness dimensions.

## **Contents**

1	Introduction	5
2	Scope, notation, and setup	6
3	,	<b>7</b> 7 8 9
4	The spectrum of fairness 4.1 The five families as five benchmarks	10 11 12 13
5	5.2 Post pricing local metrics	17 17 17 21 21 22 22 22 23
6	Exposing systematic disparities through partitioning  6.1 Pre-pricing policyholder partitioning by proxy vulnerability	<b>24</b> 24 26
7	Discussion	28
Α	Glossary	32

### **Section 1. Introduction**

Machine learning is now central to actuarial science for its predictive power (Embrechts and Wüthrich, 2022). It scales to large datasets, captures complex interactions, and excels at finding associations helpful to the predictive task (Frees et al., 2016). However, some associations are spurious; others reflect sensitive, though predictive, features. This raises ongoing debates over the legitimacy and fairness of relying on such associations in pricing.

Ratemaking datasets tend to subtly encode sensitive traits: geographic location may hint at ethnicity, and occupation often reflects gender (Bender et al., 2022a). Combined with inequalities, e.g., racial wealth gaps or credit access disparities (Bender et al., 2022b), non-sensitive variables associated with sensitive attributes can act as proxies, indirectly targeting protected subpopulations. Even without explicit use of protected features in ratemaking models, such proxy effects can perpetuate disparities.

The abundance of data amplifies the likelihood that combinations of input covariates inform on protected traits. This can be problematic when socioeconomic or demographic factors may (intentionally or inadvertently) affect insurance pricing.

Following ASOP No. 12 and 53 (ASB, 2005, 2017), insurance pricing requires grouping policyholders by risk to set adequate and financially sound premiums. Actuarial fairness ensures solvency; in contrast, other types of fairness may erode competitiveness, deterring their adoption without regulatory pressure.

Still, actuaries are expected to test for proxy effects in their models. A recent survey (Cavanaugh et al., 2024) indicates that U.S. regulators broadly agree that "insurers should test to ensure that their models do not use data and information that act as proxies for disallowed rating variables". The result is a tension: actuaries must balance risk-based pricing with ill-defined fairness notions. Even the Actuarial Standards Board acknowledges in ASOP No. 12 that there is "no general agreement on what constitutes an 'equitable' classification system or 'fair' discrimination" (ASB, 2005).

Fairness in insurance lacks an operationalizable definition and meaningful metrics. Its meaning is ambiguous, and debates hinge on speculation about variable behavior in black-box models. Standard group fairness metrics offer little segment-level insight and are difficult to translate in dollars or impacted policyholders. Fairness is discussed in theory (Lindholm et al., 2024b), but unfairness unfolds in practice. Empirical studies remain scarce, and Fahrenwaldt et al. (2024) call for high-quality datasets to move the field forward.

To address these challenges, we develop a framework that maps fairness debates onto actuarially meaningful benchmarks. We focus on fairness in actuarial pricing, acknowledging its unique challenges. Our contributions are both conceptual and applied: our methodology enables detection and quantification of unfairness using industry data. We illustrate on a large-scale industry dataset how potential unfairness manifest in practice.

The remainder of this paper is structured as follows. We first set the notation and scope in §2. In §3, we present the three dimensions of fairness in insurance pricing: actuarial fairness, solidarity, and causality. We then translate these dimensions as five ratemaking

benchmarks in §4 covering the *spectrum* of fairness viewpoints. From the estimated spectrum, we define in §5 actuarially relevant local metrics prior to pricing (risk spread, proxy vulnerability, fairness range, and parity cost) and post pricing (commercial loading, commercial burden, implied propensity and excess lift). In §6, we propose a method to detect systematic disparities by partitioning the portfolio to reveal vulnerable subpopulations. We then present two use cases: pre-pricing detection of proxy-vulnerable individuals (§6.1) and post-pricing monitoring of commercial loadings (§6.2). We illustrate the tools practicality with a case study focused on protecting individuals in precarious economic situation using a large-scale Canadian auto insurance dataset and provide the code illustrated on a simulated data example in the online supplementary material.

## Section 2. Scope, notation, and setup

In this article, we focus on a one-period fairness goal, independently of the notion of *intent*. Fairness is assessed from an output-based perspective. We assume the absence of unobserved confounders and selection bias, though these issues warrant discussion (see Côté et al., 2024, 2025).

The random variable Y represents the loss cost in a property and casualty insurance coverage. Variables available for pricing this coverage are denoted by the random vector  $\mathbf{X}$ , assumed measured without error.

Fairness is always relative to some pre-specified prohibited (or sensitive) variable D, which is here taken to be a single, categorical and fully observed random variable. Examples of sensitive attributes include gender in Europe, race in Texas, religion in California, or credit score in Ontario. Even though this might create privacy concerns, we assume that the insurer collects D, so that this variable is fully observed in our dataset. We further refer to "protected groups" as the subpopulations formed by the different levels of D and to "vulnerable groups" as those historically disadvantaged among the protected groups.

Suppose we have a portfolio of n policyholders  $(\mathbf{x}_i, d_i, e_i, Y_i)_{i=1,\dots,n}$ , where  $e_i$  is exposure to risk, measured in vehicle-years. Let  $\pi(\mathbf{x}, d)$  be the yearly commercial price for a policy with characteristics  $\mathbf{x}$  and d, including all loads, adjustments, and profit margins.

We introduce below the setup of the real data case study used for illustrating our method.

**Case study.** We study fairness regarding policyholder's economic precariousness in auto insurance premiums for material damage in at-fault accidents (Chapter B2) in the province of Québec, Canada. The data, obtained through a partnership with an insurer, includes over 768 000 vehicles insured from 2016 to 2017. The vector **X** comprises 16 explanatory variables<sup>1</sup>, including driver information, vehicle characteristics, and territorial information.

The response variable Y is the claim amount for at-fault accidents. It is highly zero-inflated, with approximately 97% of observations not filing any claim. The annual average claim is around  $\sum Y_i/\sum e_i \approx \$190$ , where exposure  $e_i \in (0,1]$  is measured in vehicle-years.

<sup>&</sup>lt;sup>1</sup>Strict anonymization and confidentiality measures were applied. A pre-selection, performed under senior actuarial oversight, reduced the size of the dataset from more than 30 candidate variables to just 16, balancing multicollinearity control with the inclusion of essential risk factors for the models.

We take credit score as D, given its link to economic precariousness (Bank of Canada, 2024). Furthermore, Prince and Schwarcz (2019) argue in favor of considering credit score as sensitive. They explain that insurers' reliance on credit data disproportionately impacts low-income and ethnic minority policyholders, characteristics generally protected against discrimination. We construct a binary variable D, where D=1 indicates high credit risk (economic precariousness) and represents about 40% of the sample.

We analyze a modified version of the insurer's pricing function. The modifications include:

- 1. Retrofitting the prices from available covariates  $(\mathbf{X}, D)$ ,
- 2. Integrating observations unavailable at the time of developing the prices,
- 3. Rescaling to match average loss to preserve premium scale confidentiality,
- 4. Making additional adjustments to protect segmentation confidentiality.

We refer to this adjusted tariff as the *pseudoprice*, denoted PseudoPrice( $\mathbf{x}, d$ ), which serves as the focal point for our fairness analysis throughout the paper<sup>2</sup>.

## Section 3. The dimensions of fairness in actuarial pricing

In our framework, three dimensions are needed to evaluate fairness relative to a sensitive D in actuarial pricing: actuarial fairness (alignment with expected loss), solidarity (redistribution across protected groups beyond pure risk), and causality (justifiable use of information). Together, the dimensions aim to cover all facets of fairness in actuarial pricing, each representing a distinct angle. Any fairness-related critique reduces to a breach of one dimension or to an explicit trade-off among them. The three dimensions of fairness are summarized in Fig. 1 and presented in §§3.1–3.3.

#### 3.1. Actuarial fairness

Actuarial fairness underpins viable insurance pricing. Originating in economic theory (Arrow, 1963), it requires policyholders to contribute to the insurance pool in proportion to their own risk. Specifically, 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).

A model is actuarially fair if it captures all risk differences and is locally balanced, keeping risk estimates unbiased and aligned with observed losses at all portfolio scales (Denuit et al., 2024). Each group, including protected ones, must have self-sustaining loss ratios. This avoids cross-subsidies and aims at a constant expected profit margin across policyholders. Actuarial fairness is about aligning premiums with expected losses.

Actuarial fairness reflects the predictive performance of the pure premium model and the lack of non-risk based commercial adjustments. The criteria of loss ratio parity (see, e.g., Bender et al., 2025) and sufficiency (see, e.g., Mosley and Wenman, 2022) align with this dimension.

<sup>&</sup>lt;sup>2</sup>While the pseudoprice is constructed to maintain realism within the case study, its non-equivalence to the actual pricing function precludes any conclusion regarding the fairness of the partner insurer's pricing.

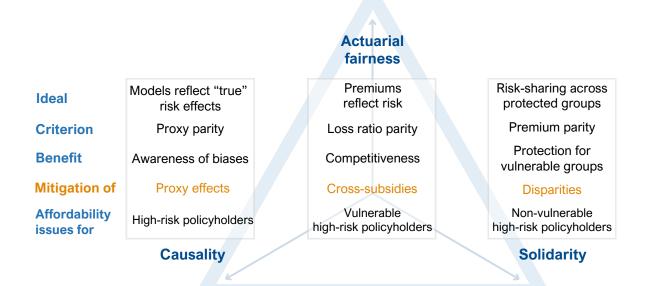


Figure 1. The three dimensions of fairness in actuarial pricing: their ideal, corresponding criterion and benefit. We clarify what they aim at mitigating, and which policyholder group might face affordability issues if this dimension is prioritized.

## 3.2. Solidarity

Solidarity is the foundation of insurance. In actuarial modeling with variables  $(\mathbf{X}, D)$ , the solidarity dimension of fairness pertains to prohibited variables D, allowing risk differentiation based on  $\mathbf{X}$  as long as premium distributions are similar across D. We intentionally create cross-subsidies if risk differs across groups of D, with the intent to promote societal welfare. Solidarity aligns with premium parity explained in Bender et al. (2025) and with demographic parity<sup>3</sup> of premiums: equal average premiums (weak parity) or identical premium distributions (strong parity) across protected groups. This is also referred to as the independence fairness criterion in Mosley and Wenman (2022).

## **Complement 1 – The shrinking homogeneous pool**



Initially, all policyholders were pooled in a collective fight against risk. By relying on data, insurers were able to create smaller, "homogeneous" pools, segmenting the original solidarity to better capture risk heterogeneity. With big data, the pools shrank further, and Barry (2020) discusses a shift toward fairness rooted in individualized pricing. In an unrealistic extreme, oracle insurers – capable of perfectly predicting both the amount and timing of individual claims – might charge each policyholder precisely their discounted future claim amount, questioning the very concept of insurance risk transfer. Increasingly granular risk factors widen the separation between actuarial fairness and solidarity.

<sup>&</sup>lt;sup>3</sup>See definition in Chapter 8 of Charpentier (2024) (Def. 8.5, Def. 8.6, Prop. 8.1).

#### 3.3. Causality

Causal inference is the art and science of reflecting the true impact of some variable on a given target, stripped of proxy effects. In a causal pricing model, each rating variable's influence corresponds to its causal effect on the claim risk Y, and does not include any indirect link through sensitive attributes D.

For causal pricing, only true and allowed risk factors, causally linked to Y, belong in the model, and their premium influence must match their actual risk contribution. While controllable risk factors are valuable for prevention, non-controllable ones, such as age, are also admissible when their causal link to risk is well supported.

Proxy and causality are about effects, not specific variables. A variable's use in a model – not the variable itself – determines its role as a proxy. Talking about proxy effects rather than "proxies" emphasizes this distinction. Even valid risk factors may induce proxy effects. Causal inference tools can isolate the "true risk component" of a rating factor.

Common causal inference strategies include using control variables to mitigate bias during training, ensuring the model is uncontaminated by D. The causality dimension of fairness aligns with the proxy-free fairness in Charpentier (2024), the proxy parity fairness criterion of Côté et al. (2024), and the proxy discrimination metric of Lindholm et al. (2024a).

## Complement 2 - Aligned tools: causal thinking and actuarial judgment



Statistical models capture associations—how claims vary with age or vehicle type. Causal models ask what happens when a variable changes. All causal models are statistical, but not all statistical models support causal claims.

A variable "causes" another if intervening on it shifts the distribution of what follows. Causation reflects a consequential distributional shift, not deterministic outcome.

Causal modeling depends on assumptions—some testable, some not—often made explicit in a causal graph. These diagrams clarify which variables influence others (see, e.g., Moodie and Stephens, 2022).

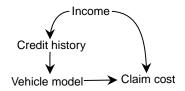


Figure 2. Illustrative causal graph for insurance pricing.

Figure 2 encodes one such structure: vehicle model directly affects claim cost; credit history correlates with claim cost through income, an unobserved common cause. Estimating the effect of vehicle model requires controlling for income or credit history, the latter being feasible via credit score data.

Causal graphs reveal valid signals and potential biases. **Actuaries know: rating factor choice is never solely about fit**. Causal thinking formalizes that intuition. Many causal assumptions mirror those already made, explicitly or not, in practice.

Causality matters for fairness. Proxy effects hide in variables that appear overly predictive but reflect something sensitive. Discrimination lies in data, not models (Charpentier, 2024). Causal reasoning helps separate valid from spurious signals, which is key to detecting unfairness along the causality dimension: **proxy effects**.

#### **Complement 3 – Revisiting the history of risk classification**

ΔΔ

Flat-rated pricing (first private auto in 1887, Insurance Information Institute, 2023) maximized solidarity but ignored risk heterogeneity and causality. Manual pricing tables (Kormes, 1935) added experience-based risk tables, improving actuarial fairness and likely causality. Multivariate analysis, such as minimum bias procedure of (Bailey and Simon, 1960), refined segmentation and risk factor use based on insurer-specific experience.

Usage-based insurance (first policy sold in 1998 by Progressive, Brobeck and Hunter, 2021) aligned premiums with behavior. Although it dilutes the predictive power of protected attributes (Boucher and Pigeon, 2024), behavioral data may still encode disparities, undermining solidarity as pictured in Fig. 3.

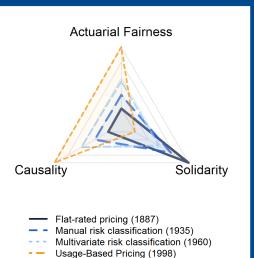


Figure 3. The dimensions of fairness with increasing segmentation capacity.

Risk differences across protected groups may exist due to historical and socioe-conomic factors. As data granularity increases, so does the potential for actuarial justification in perpetuating these disparities.

## **Section 4.** The spectrum of fairness

Côté et al. (2025) describe disparate impact as the association between premiums and sensitive attributes, lying between two extremes: **solidarity**, which aligns premiums across protected groups, and **actuarial fairness**, which aligns premiums with risk. Solidarity levels premiums; actuarial fairness levels profits. This tension is the core of fairness debates.

Nuances exist. Premium disparities relative to D are not uniformly problematic. They range from no association (**solidarity**), to justified association (**causality**), to association inflicted by proxy effects, and up to direct exploitation of D for maximal predictive accuracy (**actuarial fairness**). The challenge is in pinpointing where the disparate impact falls along this spectrum, and whether it crosses the blurry line between fair and unfair.

The five premium families of Côté et al. (2025) span this spectrum: best-estimate, unaware, aware, hyperaware, and corrective. These families represent how fairness considerations regarding D influence a pricing structure, and how permissible variables  $\mathbf{X}$  are leveraged in relation to D. Each family offers distinct trade-offs between actuarial fairness, causality, and solidarity. We use this spectrum to investigate potential unfairness.

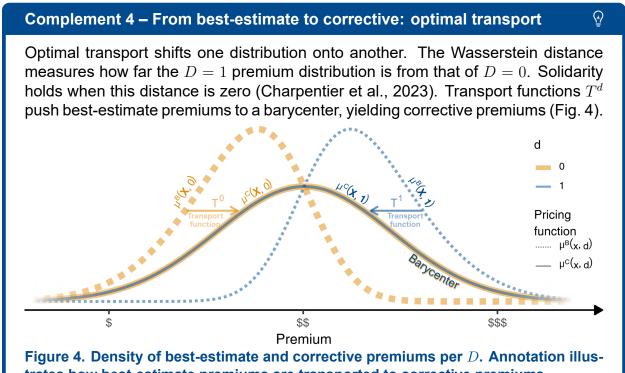
This section is structured as follows: we present the five benchmark premiums, one for each family, in §4.1 and we explain how we estimate them in §4.2. In §4.3, we reveal trade-offs between fairness dimensions in our Case study.

#### The five families as five benchmarks 4.1.

We present one model per family to obtain five benchmarks. We focus on the expectation. in line with the "expected value of all future costs" of Casualty Actuarial Society (1988).

First, the **best-estimate premium**  $\mu^B(\mathbf{x}, d)$  aligns with **actuarial fairness**, grouping risks following **X** and *D* to set premiums. One example is  $\mu^B(\mathbf{x}, d) = \mathsf{E}(Y|\mathbf{X} = \mathbf{x}, D = d)$ .

The corrective premium  $\mu^C(\mathbf{x}, d)$  leverages both **X** and D to satisfy solidarity, that is, equal premium distribution per protected group. One example is  $\mu^C(\mathbf{x},d) = T^d\{\mu^B(\mathbf{x},d)\}$ , where  $T^d$  is a transport function detailed in Complement 4 and the online material.



trates how best-estimate premiums are transported to corrective premiums.

The **unaware premium**  $\mu^U(\mathbf{x})$  is the best non-directly discriminatory approximation of the best-estimate premium, an example of unaware premium being  $\mu^U(\mathbf{x}) = \mathsf{E}(Y|\mathbf{X}=\mathbf{x})$ .

The **hyperaware premium**  $\mu^H(\mathbf{x})$  is the best non-directly discriminatory approximation of the corrective premium. One example is  $\mu^H(\mathbf{x}) = \mathsf{E}\{\mu^C(\mathbf{X},D)|\mathbf{X}=\mathbf{x}\}.$ 

Finally, the **aware premium**  $\mu^A(\mathbf{x})$  captures the effect of **X** on Y when controlling for D. An example is  $\mu^A(\mathbf{x}) = \mathsf{E}_D\{\mu^B(\mathbf{x},D)\}$ , a discrimination-free price of Lindholm et al. (2022).

The best-estimate premium sets the baseline with both  $\mathbf{X}$  and D in a "paying for your own risk" approach. The unaware, aware, and hyperaware premiums omit D but handle proxies differently: the first reflects proxy effects, the second resists them, and the third uses them toward premium parity. The hyperaware and corrective explicitly pursue solidarity. The best-estimate and unaware are solely risk-focused. The aware permits parity shifts when causally justified by **X**. Together, these benchmarks reveal the trade-offs in fair pricing regarding D across fairness dimensions. We summarize them in Table 1.

Premium	Best-estimate	Unaware	Aware	Hyperaware	Corrective	
Notation $\mu^B(\mathbf{x},d)$		$\mu^U(\mathbf{x}) \qquad  \mu^A(\mathbf{x})$		$\mu^H(\mathbf{X})$	$\mu^C(\mathbf{x},d)$	
Formula	$E(Y \mathbf{X}=\mathbf{x},D=d)$	$E(Y \mathbf{X}=\mathbf{x})$	$E_D\{\mu^B(\mathbf{x},D)\}$	$E\{\mu^C(\mathbf{x},D) \mathbf{X}=\mathbf{x}\}$	$T^d\{\mu^B(\mathbf{x},d)\}$	
Direct discrimination	<b>✓</b>	×	×	×	~	
Proxy discrimination	_	<b>✓</b>	×	✓	_	
Demographic disparities	<b>✓</b>	<b>✓</b>	<b>✓</b>	×	×	
Dimension prioritized	Actuarial fairness	Actuarial fairness	Causality	Solidarity	Solidarity	

Table 1. Properties of the five fair premiums from §4.1.

#### 4.2. Estimation of the five premiums

We give an example of procedure to estimate the spectrum of fairness. For additional details and examples, see the online supplement.

1. **Best-estimate:** Fit a lightgbm (Ke et al., 2017) that includes  $(\mathbf{X}, D)$  to predict Y using an appropriate distribution (e.g., Tweedie):

$$\widehat{\mu}^B(\mathbf{x}, d) = \widehat{\mathsf{E}}(Y \mid \mathbf{X} = \mathbf{x}, D = d).$$

Alternatively, one can rely on a (directly discriminating) technical price.

## **Complement 5 – Technical or data-driven best-estimate premium?**

Insurers typically start pricing by estimating **indicated rates**, the actuary's best estimate of risk-based prices. As discussed in §4.1, the last four fairness families are derived from  $\widehat{\mu}^B$ . The choice of anchor for  $\widehat{\mu}^B$  shapes the fairness assessment:

- a) Indicated rates as  $\widehat{\mu}^B$ : Indicated rates can be used as  $\widehat{\mu}^B$  if the sensitive variable D is included as a rating variable. This approach benefits from actuarial oversight, but ties fairness benchmarking to actuarial choices. Consequently, biases within technical pricing may remain undetected.
- **b)** Data-driven  $\widehat{\mu}^B$ : To guard against institutional or analyst-induced bias, actuaries may estimate  $\widehat{\mu}^B$  directly from data using flexible algorithms like lightgbm (Ke et al., 2017), detached from technical or commercial pricing.
- 2. Unaware: Train a second lightgbm to approximate  $\widehat{\mu}^B(\mathbf{X}, D)$  using only **X**:

$$\widehat{\mu}^U(\mathbf{x}) = \widehat{\mathsf{E}}\{\widehat{\mu}^B(\mathbf{X}, D) \mid \mathbf{X} = \mathbf{x}\}.$$

3. **Aware:** Compute the empirical proportions of D in the training set. For each  $\mathbf{x}$ , average group-specific best-estimate premiums weighted by empirical frequencies:

$$\widehat{\mu}^{A}(\mathbf{x}) = \sum_{d} \widehat{\mu}^{B}(\mathbf{x}, d) \, \widehat{\mathsf{Pr}}(D = d).$$

4. **Corrective:** Train the optimal transport function  $\widehat{T}^d$  of best-estimate premiums using Equipy (Fernandes Machado et al., 2025). The corrective premium is then:

$$\widehat{\mu}^C(\mathbf{x},d) = \widehat{T}^d \{ \widehat{\mu}^B(\mathbf{x},d) \}.$$

5. **Hyperaware:** Train a last lightGBM model to regress  $\widehat{\mu}^C$  on **X** only:

$$\widehat{\mu}^{H}(\mathbf{X}) = \widehat{\mathsf{E}}\{\widehat{\mu}^{C}(\mathbf{X}, D) \mid \mathbf{X} = \mathbf{X}\},\$$

removing any direct discrimination on *D* while partly preserving parity corrections.

#### **Complement 6 – Implementation tips**

**\$**0

- If *D* was discretized, store bin definitions for future use.
- Scale each of the five premiums by a constant to align with revenue targets.

All models are estimated on the same data, with the same features and overall target profit. Aside from natural estimation variability, differences reflect only the fairness goal.

## 4.3. Deviations within the spectrum

Deviations of a commercial price from a fairness benchmark suggest either misalignment with its intent or potential for predictive gain without fairness sacrifice. The meaning of the spectrum emerges only through the lens of the three fairness dimensions (§3), explored next.

Case study (Cont'd). We estimate the spectrum of fairness following the methodology recommended in §4.2. We obtain our best-estimate premium  $\widehat{\mu}^B$  purely from data.

id	GenderMainDriver	DrivExp (year)	DriverAge (year)	YearlyMileage (km)	Location (from ZipCode)	ОссТуре	hasPropertyIns	Economic precariousness (D)
1	Female	19	35	10 000	Island of Montreal	Full time	Yes	1
2	Male	25	42	10 000	Capitale-Nationale	None	Yes	0
3	Female	56	80	5 000	Laurentians	None	Yes	0
4	Male	8	24	20 000	Island of Montreal	Full time	Yes	1
5	Male	0	42	15 000	Island of Montreal	Full time	No	1
6	Male	3	19	15 000	Centre-du-Québec	Full time	No	0

Table 2. Partial description of six profiles for the analysis in the Case study.

In the left panel of Fig. 5, we display the lightgbm estimated propensity  $\widehat{\Pr}(D=1|\mathbf{X}=\mathbf{x})$  with dashed line indicating an observed high credit risk D=1 for individuals 1, 4, and 5. In the right panel of Fig. 5, we show the estimated spectrum and the pseudoprice for six individuals, partly described in Tab. 2. These graphs offer initial intuition on fairness:

- For individual 5, the pseudoprice lies outside the fairness spectrum. While methodology and commercial strategy (e.g., marketing or customer experience) may explain this, any deviation from the spectrum warrants close attention.
- For individuals with D=1, the best-estimate premium (red) is the highest of the spectrum. The corrective premium  $\widehat{\mu}^C(\mathbf{x},d)$  (blue) exceeds  $\widehat{\mu}^B(\mathbf{x},d)$  only if D=0.
- Fair premium ranges vary. For individual 2, all premiums closely align, suggesting fairness adjustments have little matter for some policyholders. The intuition that higher risk appears linked to wider premium range will be exemplified in Fig. 14.
- Individual i=5 shows higher  $\widehat{\Pr}(D=1|\mathbf{X}=\mathbf{x}_i)$  than individual i=4, pulling the unaware premium (orange cross) closer to the best-estimate (red triangle), illustrating proxy discrimination. For outlier cases like i=6, with D=0 despite high propensity  $\widehat{\Pr}(D=1|\mathbf{X}=\mathbf{x}_6)$ , the unaware premium is far from the best-estimate.
- Higher propensity for D=1 in individuals  $\{4,5,6\}$  versus  $\{1,2,3\}$ , aligns with higher risk estimates, reinforcing the core motivation for fairness. Vulnerable individuals (high credit risk D=1) tend to present allowed covariates **X** associated with higher claim propensity, driving up premiums via both protected and allowed variables<sup>4</sup>.

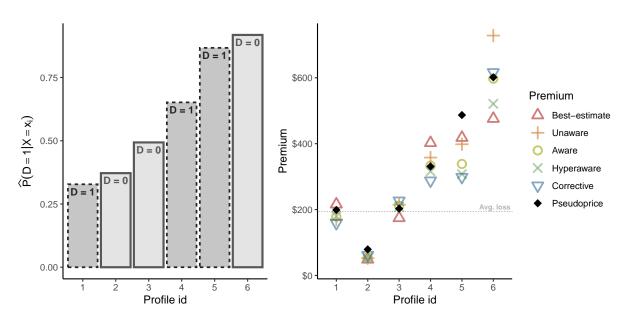


Figure 5. Propensity to observe D=1 (left) and estimated spectrum of premiums along with the pseudoprice (right) for six individuals in the Case study.

<sup>&</sup>lt;sup>4</sup>In practice, biased covariates (e.g., uneven law enforcement of traffic violations) may further inflate premiums for vulnerable groups (Leong et al., 2024).

To formalize these observations, we assess each premium's alignment with the fairness dimensions introduced in §3. With a binary sensitive variable, Wasserstein distance measures distributional differences between the two protected groups (D = 0 and D = 1).

- 1. For **actuarial fairness**, we assess loss ratio parity via the distance between its distributions across groups. Because every value of *D* yields a single loss ratio, we partition the data into 100 random subsamples to allow estimation of distributional differences. This sampling is only used for this dimension.
- For causality, we assess proxy parity. Following Lindholm et al. (2024a); Côté et al. (2024), we treat deviations from the aware premium as proxy effects and compare their distributions across groups.
- 3. For **solidarity**, we compare premium distributions between protected groups.

These metrics are illustrative, they can be adapted or refined depending on the context.

Table 3 presents the Wasserstein distances for all premiums and the pseudoprice for atfault material damage (Chapter B2) coverage. A Wassertein distance of zero implies the corresponding fairness criterion is satisfied. This provides context for the pseudoprice's position within the fairness spectrum.

Fig. 6 displays a radar plot ranking premiums by their alignment with each fairness dimension. The closer a premium is to a triangle's vertex, the stronger its adherence to that principle. As expected, the best-estimate, aware, and corrective premiums strongly adhere to actuarial fairness, causality, and solidarity, respectively. The three non-directly discriminatory premiums (unaware, aware, hyperaware) cluster together, with the unaware leaning toward actuarial fairness and the hyperaware tilting toward solidarity, as expected.

Table 3. Wasserstein distance between distributions for D=0 and D=1 of loss ratios (actuarial fairness), deviations from the aware family (causality), and premiums (solidarity) in the Case study.

	Actuarial Fairness	Causality	Solidarity
Best-estimate	0.036	67.054	116.938
Unaware	0.282	12.800	62.684
Aware	0.343	0	49.884
Hyperaware	0.392	10.071	49.840
Corrective	0.604	50.167	0.877
PseudoPrice	0.175	35.773	85.654

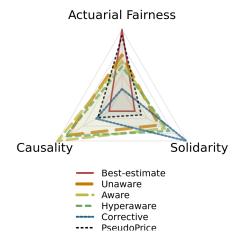


Figure 6. Alignment of premiums with fairness dimensions in the Case study.

The pseudoprice aligns with a best-estimate strategy, which is unsurprising given the sensitive variable's inclusion in industry-wide rates during the study period. When a variable is not deemed sensitive, insurers legitimately prioritize actuarial fairness.

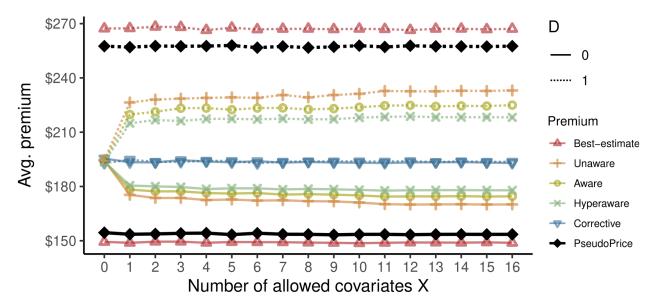


Figure 7. Evolution of the mean premium for D=0 (solid) or 1 (dashed) by family as a function of how many covariates are allowed in the vector  ${\bf X}$  in the Case study. The best-estimate, corrective, and pseudoprice directly discriminate on D.

We illustrate in Fig. 7 how average premiums per protected group evolve when covariates are sequentially introduced. The covariates in  $\mathbf{X}$  are added in order of variable importance in the best-estimate  $\mathtt{lightgbm}$ ; the first is the main driver's experience in years  $\mathtt{DrivExp}$  which is strongly correlated with credit score. Both the best-estimate and corrective premiums directly discriminate on D but with opposing intents: the best-estimate reflects all risk differences in D, while the corrective offsets group disparities through corrective direct discrimination. Even with abundant proxy effects, the unaware premium's ability to capture risk differences from the prohibited variable D is limited, as evidenced by the persistent gap between its average (orange) and that from the best-estimate premium (red). The ordering of premium families aligns with intuition.

Fig. 6 reminds us that fairness dimensions cannot be satisfied simultaneously. Classical fairness criteria (e.g., independence, proxy parity, and sufficiency), each tied to a distinct fairness dimension in §3, are fundamentally incompatible: extensive results, for example, by Kleinberg et al. (2016); Charpentier (2024); Lindholm et al. (2024b) demonstrate that satisfying one typically violates another. In insurance setups, Bender et al. (2025) illustrate this impossibility with actuarial examples.

Universal fairness breaks on one truth: disparities *do exist* to begin with as seen, for example, in Fig. 7. No premium can be deemed universally fair; fairness is – and will remain – an elusive ideal, requiring ongoing governance of trade-offs across the three dimensions. All dimensions align only in the trivial case where the sensitive variable is entirely unrelated to the rest of the dataset (see, for example, Côté et al., 2025).

## Section 5. Actuarially relevant local fairness metrics

The Case study suggests that deviations from the fair premium spectrum provide meaningful insights at the individual level. In this section, we first interpret key deviations within the spectrum (pre-pricing) and then examine deviations from a given tariff  $\pi(\mathbf{x}, d)$  to the benchmarks (post-pricing).

#### 5.1. Pre-pricing local metrics

First, the **risk spread** measures the range of best-estimates for different sensitive attribute values. Building on proxy effects of Lindholm et al. (2024a), we interpret the **proxy vulnerability** as the deviation between the unaware and aware benchmark premiums. Next, the **fairness range** reflects the overall range of the spectrum. Lastly, we define the **parity cost**, the overcharge experienced when going from loss ratio parity to premium parity.

#### 5.1.1. Risk spread

For a segment  $\mathbf{x}$ , the **risk spread**, denoted  $\Delta_{\mathsf{risk}}(\mathbf{x})$ , measures the range of data-driven risk estimates across different values of the sensitive attribute D:

$$\Delta_{\mathrm{risk}}(\mathbf{x}) = \max\left\{\mu^B(\mathbf{x},1), \mu^B(\mathbf{x},0)\right\} - \min\left\{\mu^B(\mathbf{x},1), \mu^B(\mathbf{x},0)\right\} = \left|\mu^B(\mathbf{x},1) - \mu^B(\mathbf{x},0)\right|,$$

the within-segment premium gap between D=1 and D=0. It captures how much the best-estimate premium attributes risk differences to D within a given segment  $\mathbf{x}$ .

The risk spread represents the model's incentive to capture risk differences driven by D for the segment  $\mathbf{x}$ . The risk spread is positive; a larger value indicates a greater potential for disparate treatment across protected groups should pricing differentiate on D.

#### 5.1.2. Proxy vulnerability

For a segment  $\mathbf{x}$ ,  $\mathbf{proxy}$   $\mathbf{vulnerability}$  quantifies the unintended price shift between a model that ignores D and one that controls for it. We define it as:

$$\Delta_{\mathsf{proxy}}(\mathbf{x}) = \mu^U(\mathbf{x}) - \mu^A(\mathbf{x}). \tag{1}$$

A large proxy vulnerability indicates that a segment  $\mathbf{x}$  is prone to proxy effects, where seemingly neutral variables in  $\mathbf{x}$  serve as proxies for D. This occurs when a significant risk spread exists and  $\mathbf{X}$  informs on D. A positive value means the unawareness model overcharges the segment, while a negative value results from underpricing due to proxy. Studying proxy vulnerability highlight segments that are most exposed to potential proxy discrimination, providing a best guess regarding its monetary magnitude and direction.

Because the proxy phenomenon captures risk differences across groups of D indirectly, it is bounded by the distance between the aware premium and surrounding best-estimates:

$$\min\{\mu^B(\mathbf{x},1),\mu^B(\mathbf{x},0)\} - \mu^A(\mathbf{x}) \leq \Delta_{\text{proxy}}(\mathbf{x}) \leq \max\{\mu^B(\mathbf{x},1),\mu^B(\mathbf{x},0)\} - \mu^A(\mathbf{x}).$$

Proxy vulnerability arises from the interplay between risk spread (potential direct discrimination on D) and propensity (ability to exploit it when using only  $\mathbf{x}$ ).

Case study (Cont'd). The top left panel of Fig. 8 plots proxy vulnerability  $\widehat{\Delta}_{\text{proxy}}(\mathbf{x}_i)$  against the estimated propensity  $\widehat{\Pr}(D=1|\mathbf{X}=\mathbf{x}_i)$  for individuals i in the test set. Black triangles relate to observed low credit score  $(d_i=1)$  and purple circles to high credit score  $(d_i=0)$ . The upward spline trend indicates that proxy vulnerability increases with the likelihood of belonging to the D=1 group. Colors suggest that the model for D|X (horizontal axis) is good, with more black triangles (true d=1) on the right.

The top right panel of Fig. 8 shows boxplots of the proxy vulnerability, expressed as a percentage of the aware premium. Proxy vulnerability is higher for the D=1 subpopulation, and often exceeds 10%. Low credit score individuals ( $d_i=1$ ) are riskier, and the unaware premium  $\mu^U$  captures this even when D is unobserved or excluded. This outcome is unavoidable when predictive covariates correlate with protected traits.

The middle row of Fig. 8 plots the relationship between propensity and aware premium, using both original (left) and normalized rank (right) scales. Though the aware premium is constructed to be invariant to the propensity to observe D=1 (§4.1), a dependence persists: individuals with higher propensity for D=1 tend to have higher aware premiums. In this case study, vulnerable individuals (D=1) are riskier and more likely to exhibit values of allowed covariates  ${\bf x}$  linked to higher claim risk. The normalized rank plot confirms this: the upper tail is mainly populated by D=1 individuals.

The bottom line of Fig. 8 depicts proxy vulnerability (color scale, low in purple, high in black) as a function of risk spread (x-axis) and propensity for D=1 (y-axis). Patterns reveal that proxy vulnerability arise when risk differs by protected group (large risk spread) and when  $\bf x$  allows indirect inference of D (propensity near 0 or 1).

We depict in Fig. 9 a map of Québec aggregated by forward sortation area, our chosen geographic unit. Each area is colored by the Tail Value-at-Risk  $TVaR_{0.95}$  of proxy vulnerability  $\widehat{\Delta}_{\text{proxy}}(\mathbf{x})$ , computed as the average of the top 5% of values within that unit. This highlights regions where the most vulnerable individuals are concentrated, with zooms around Québec City and Montréal. We see that some regions, such as Alma, Montréal-Nord or St-Georges, exhibit high proxy vulnerability (darker purple). Representing proxy vulnerability on the map and supporting the analysis with census data may help to pinpoint sensitive demographics in these specific regions, for example the large proportion (42%) of immigrants in Montréal-Nord<sup>5</sup>.

This is a compelling illustration of the materiality of proxy effects and why they warrant scrutiny. Even without explicit use of a prohibited attribute, its statistical imprint propagates through associated covariates, sustaining disparities in ways that evade direct detection. Proxy effects are material and their potential impact is not evenly distributed.

<sup>&</sup>lt;sup>5</sup>Montréal en statistiques (2018) Profil sociodémographique, Recensement 2016: Arrondissement de Montréal-Nord.

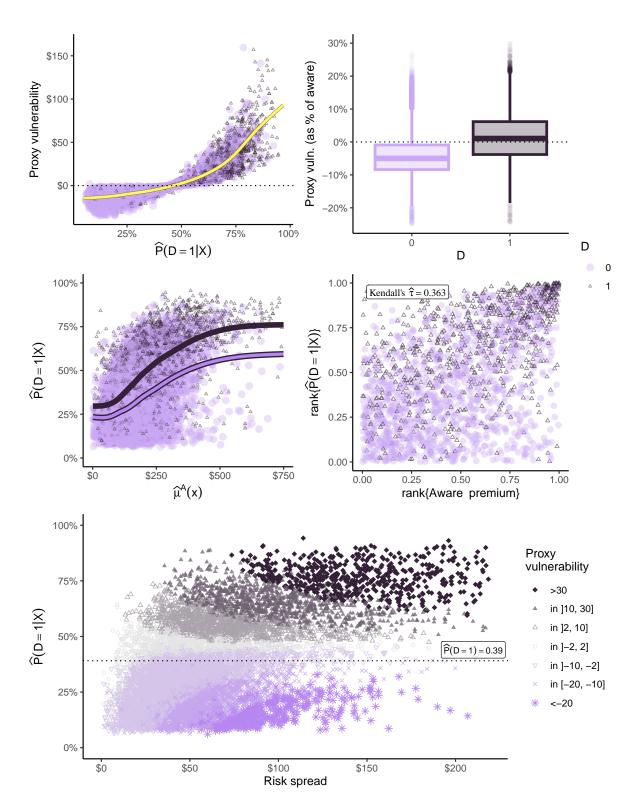


Figure 8. Dashboard of proxy vulnerability for the Case study: proxy vulnerability (in CAN\$) as a function of the propensity to observe D=1 (top left); proxy vulnerability in percentage of aware premium per protected group (top right); propensity as a function of aware premiums on their original scale (middle left) and on the scale of their normalized ranks (middle right); and estimated propensity to observe D=1 in terms of risk spread and colored by proxy vulnerability interval (bottom).



Figure 9. Geographic distribution of the empirical 95% TVaR of proxy vulnerability, assessed at the forward sortation area level, in the Case study. The top-left panel shows the entire province, the right panel zooms in on its central region, and the bottom panel provides a detailed view of the island of Montréal.

#### **5.1.3.** Fairness range

For a specific segment  $(\mathbf{x}, d)$ , the **fairness range**, denoted  $\Delta_{\text{fair}}(\mathbf{x}, d)$ , is defined as

$$\begin{split} \Delta_{\text{fair}}(\mathbf{x},d) &= \max\{\mu^B(\mathbf{x},d), \mu^U(\mathbf{x}), \mu^A(\mathbf{x}), \mu^H(\mathbf{x}), \mu^C(\mathbf{x},d)\} - \\ &\qquad \qquad \min\{\mu^B(\mathbf{x},d), \mu^U(\mathbf{x}), \mu^A(\mathbf{x}), \mu^H(\mathbf{x}), \mu^C(\mathbf{x},d)\}. \end{split}$$

The fairness range measures how much prices vary across fairness methods. A large value indicates that pricing is sensitive to fairness considerations for the segment.

#### 5.1.4. Parity cost

The **parity cost** is the (monetary) cost for a policyholder of enforcing demographic parity compared to a "pay for your own risk" approach. For a segment  $(\mathbf{x}, d)$ , it is defined as:

$$\Delta_{\mathrm{parity}}(\mathbf{X},d) = \mu^C(\mathbf{X},d) - \mu^B(\mathbf{X},d).$$

A higher parity cost signals that larger adjustments are needed to enforce demographic parity. It quantifies how solidarity objectives conflict with actuarial fairness for an individual.

Case study (Cont'd). Fig. 10 shows fairness range and parity cost. On the left, fairness range tracks risk spread, revealing that sensitivity to fairness adjustments follows potential for disparate treatment. Conditional on risk spread, the D=1 group shows greater sensitivity to fairness adjustments. On the right, parity cost reflects discounts for highrisk (D=1) and surcharges for low-risk (D=0) individuals. Point size reflects sample density. The distribution for D=1 centers on large discounts; for D=0, on small surcharges. Thus, achieving demographic parity in involves imposing modest levies on the non-vulnerable group (D=0) to fund substantial rebates for the vulnerable group (D=1).

This captures the redistributive effect of the corrective premium relative to its risk-based counterpart. Critics of demographic parity often highlight cross-subsidies as problematic; the parity cost explicitly quantifies them.

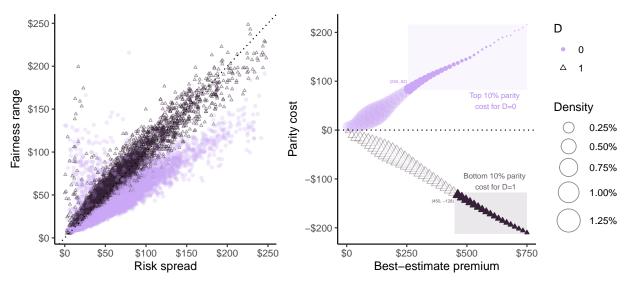


Figure 10. Fairness range as a function of risk spread (left) and parity cost as a function of best-estimate premiums (right) for the Case study. Colors denote protected group.

#### **Complement 7 – Subsidizing fairness without cross-subsidies**



A public scheme could fund fairness by reimbursing insurers the **parity cost**: the gap between corrective and best-estimate premiums. Vulnerable individuals (D=1) pay the reduced corrective rate; others (D=0) retain their actuarially fair price. This lowers prices for protected groups without burdening others, enabling fairness without cross-subsidies or insurer loss.

#### 5.2. Post pricing local metrics

Commercial pricing includes adjustments unrelated to risk or fairness. Bender et al. (2025) advise assessing final prices by their "actual impact on policyholders". We introduce local fairness metrics to evaluate a given commercial price  $\pi(\mathbf{X}, D)$  relative to the spectrum.

#### 5.2.1. Commercial loading and commercial burden

For a segments  $(\mathbf{x}, d)$ , the **commercial loading** is defined as

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

where  $\mu$  denotes a reference premium intended to best represent indicated rates. If direct use of D is allowed, set  $\mu = \mu^B$  to assess actuarial fairness, aligning, e.g., with the "deviation from indicated rates" metric of the Financial Services Regulatory Authority of Ontario (2024). If not, or to monitor proxy effects, use  $\mu = \mu^A$  to examine causality.

The **commercial burden**, denoted  $\rho_{\text{burden}}(\mathbf{x}, d; \pi, \mu)$ , is the commercial loading as a percentage of  $\mu$ . High burdens may raise affordability concerns for low-income policyholders:

$$\rho_{\mathrm{burden}}(\mathbf{x},d;\pi,\mu) = \frac{\pi(\mathbf{x},d)}{\mu(\mathbf{x},d)} - 1 = \frac{\Delta_{\mathrm{load}}(\mathbf{x},d;\pi;\mu)}{\mu(\mathbf{x},d)}.$$

If the reference premium is the best-estimate  $\mu^B$ , the commercial burden equals the markup over the claim costs, i.e., the inverse expected loss ratio.

#### 5.2.2. Implied propensity

The **implied propensity**, denoted  $\tilde{P}_D(\mathbf{x};\pi)$ , is the implicit weight of D=1 for segment  $\mathbf{x}$  when expressing a (non-directly discriminatory)  $\pi$  as a linear combination of best-estimate premiums across values of D:

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

Because  $\pi$  is unconstrained,  $\tilde{P}_D(\mathbf{x};\pi)$  may lie outside [0,1]. Solving for  $\tilde{P}_D$  yields:

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

It is well-defined when  $\mu^B(\mathbf{x},1) \neq \mu^B(\mathbf{x},0)$ . Values outside [0,1] reveal targeting of a protected group. Naturally,  $\tilde{P}_D(\mathbf{x};\mu^U) = \Pr(D=1|\mathbf{X}=\mathbf{x})$ . An implied propensity aligned with  $\Pr(D=1|\mathbf{x})$  or  $1-\Pr(D=1|\mathbf{x})$  reflects proxy effects or solidarity, respectively.

#### 5.2.3. Excess lift

For directly discriminatory pricing function, we define the excess lift for segment **x** as

$$\Delta_{\text{excess}}(\mathbf{x};\pi) = \left|\pi(\mathbf{x},1) - \pi(\mathbf{x},0)\right| - \left|\mu^B(\mathbf{x},1) - \mu^B(\mathbf{x},0)\right| = \left|\pi(\mathbf{x},1) - \pi(\mathbf{x},0)\right| - \Delta_{\text{risk}}(\mathbf{x}).$$

The excess lift quantifies how "excessively" a pricing function  $\pi$  differentiate on D for a segment relative to the "pure risk" best-estimate premium. By construction,  $\Delta_{\sf excess}(\mathbf{x}; \mu^B) = 0$  for every segment  $\mathbf{x}$ . Strictly positive values signal over-differentiation on D; negative values signal under-differentiation between D=0 and D=1 (e.g., from smoothing, regulatory caps, or solidarity efforts) and imply cross-subsidization within that segment.

Case study (Cont'd). Fig. 11 shows commercial burden  $\widehat{\rho}(\mathbf{x},d; \mathsf{PseudoPrice}, \widehat{\mu}^A)$  (top) and excess lift  $\widehat{\Delta}_{\mathsf{excess}}(\mathbf{x}; \mathsf{PseudoPrice})$  (bottom), plotted against the aware premium. Though  $\mu^B$  reflects industry norms during the study period, we use  $\mu^A$  as the reference premium to assess targeting of D, whether directly or via proxy effects.

In the top panel, both groups show high variability. Commercial burden is on average higher for vulnerable population than for others. In the bottom panel, excess lift declines with  $\hat{\mu}^A$ . Excess lift on D is positive for low  $\hat{\mu}^A$ , but negative for large premiums, possibly due to fixed costs, solidarity efforts, or rigid model (e.g., models without interactions).

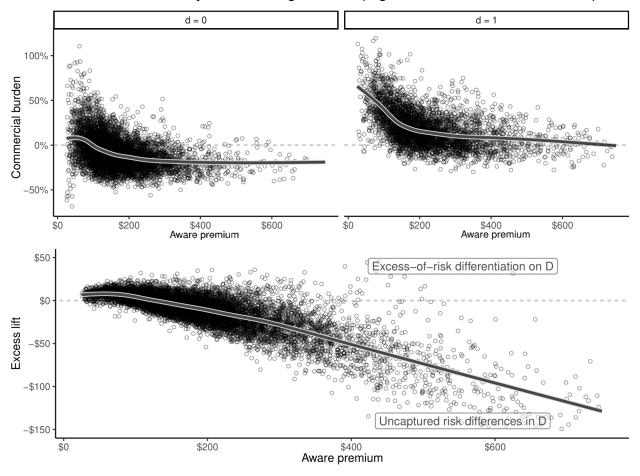


Figure 11. Commercial burden  $\widehat{\rho}(\mathbf{x},d; \mathsf{PseudoPrice}, \widehat{\mu}^A)$  by protected group (top) and excess lift  $\widehat{\Delta}_{\mathsf{excess}}(\mathbf{x}; \mathsf{PseudoPrice})$  (bottom) plotted against aware premiums for the Case study.

## Section 6. Exposing systematic disparities through partitioning

Fairness metrics on broad groups can mislead. A model may seem fair with respect to gender while masking disparities for smaller vulnerable groups, like single mothers. Tailored fairness assessment for specific subpopulations helps to target corrections. To detect systematic fairness disparities, we propose in this section a simple methodology to partition policyholders following a relevant fairness metric.

To create these supervised partitions, we use decision trees (see, e.g., Loh, 2014) for simplicity and interpretability. Five essential components guide this process: **data**, **feature space**, **response variable**, **loss function**, and **algorithm**, further detailed in the online supplement. The complexity of partitioning should align with the analyst's need for fairness granularity rather than be dictated by loss minimization alone.

We leverage the partitioning algorithm with relevant local fairness quantities to differentiate segments depending on a given fairness quantity: we identify segments with high proxy vulnerability upfront in §6.1, and we detect commercial loading in rates in §6.2.

### 6.1. Pre-pricing policyholder partitioning by proxy vulnerability

We identify segments most exposed to potential proxy discrimination by partitioning policyholders based on proxy vulnerability (defined in §5.1.2).

Case study (Cont'd). We apply a regression tree to predicted proxy vulnerability  $\widehat{\Delta}_{proxy}(\mathbf{x})$ , using all allowed variables  $\mathbf{X}$  for partitioning. The resulting regularized tree<sup>6</sup>, depicted in Fig. 12, reveals the following:

- The highest proxy vulnerability leaves (numbered 1–5) are split by hasPropertyIns (property insurance indicator), DrivExp (driving experience), NbTraffViolation (count of traffic violations), and OccType (type of occupation).
- The indicator hasPropertyIns likely captures property ownership (clearly associated with credit score), DrivExp correlates with age or financial constraints, and a high value of NbTraffViolation may reflect low risk aversion.
- Inexperienced drivers, such as DrivExp < 5 and hasPropertyIns == 'N' for nodes 1 and 2, stand out as a group that is vulnerable to high proxy effects. Their true (causal) risk may be inflated by implicit inference of their credit score.
- Proxy vulnerability averages to zero but hides wide disparities. In leaf node 1, the median overcharge is \$70 on \$540 losses (13%). Elsewhere, the proxy rebate is at most \$23. Proxy effects are both material and asymmetric.

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

<sup>&</sup>lt;sup>6</sup>Our partitioning with evtree uses a subsample of 50,000 observations and evaluates a population of 150 candidate evolutionary trees; the final model is the single best tree (no ensemble). See Grubinger et al. (2014) for more details.

25

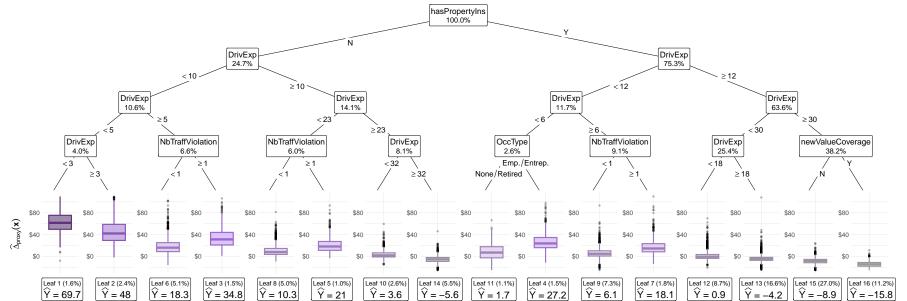


Figure 12. Optimal partitioning of  $\widehat{\Delta}_{proxy}(x)$  for the Case study.

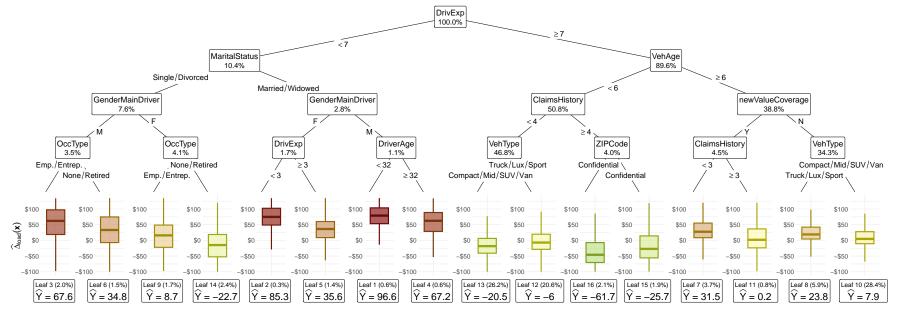


Figure 13. Optimal partitioning of  $\widehat{\Delta}_{\text{load}}(\mathbf{x})$  for the Case study.

#### 6.2. Post-pricing policyholder partitioning by commercial loading

Constructing a regression tree on commercial loading (§5.2.1) reveals rating patterns.

Case study (Cont'd). We fit a regression tree to  $\widehat{\Delta}_{load}(\mathbf{x}; \widehat{\mu}^A)$  and depict it in Fig. 13. The leaves with highest commercial loading (numbered 1–6) are split by  $\mathtt{DrivExp}$  (driving experience of the main driver),  $\mathtt{MaritalStatus}$  (marital status),  $\mathtt{GenderMainDriver}$  (gender of the main driver),  $\mathtt{DriverAge}$  (age of the main driver), and  $\mathtt{OccType}$  (type of occupation). Highest loadings occur among inexperienced married drivers (leaves 1, 2, 4, and 5) and single males with little driving experience (leaves 3 and 6).

Unlike proxy vulnerability, commercial loading stems from more than just proxy effects: differences in technical pricing methodology, lag in historical loss modeling for Y, commercial strategies, new business discounts, capping, and, most importantly, direct discrimination on D. Pricing decisions, when compounded, may produce unintended disparities, disadvantaging groups beyond the insurer's intent and/or awareness.

Combined with the pre-pricing partition on proxy vulnerabilty (§6.1), the two partitions may help track model behavior across flagged subgroups. We integrate these components into a structured fairness assessment framework depicted in Fig. 14, combining the partitions illustrated in Figs. 12 and 13 with the actuarial metrics defined in Section 5. Our toolbox guides actuaries in pinpointing which segments (partitioning columns) may warrant premium adjustment and which fairness indicators (rows) should be evaluated.

The tool assembles key actuarial components for each subgroup:

- Demographic summaries on selected covariates;
- Classical pricing diagnostics: expected losses, premiums, predictive performance;
- Basic fairness assessment: group disparities in premiums and losses;
- Pre- and post-pricing local fairness indicators;
- The partitioning rule identifying the subgroup at the bottom;
- Optional enrichment from external data (e.g., census).

Analysis reveals a high proxy vulnerability among groups with elevated commercial loading. The alignment between pre-pricing proxy vulnerability and post-pricing commercial loading suggests that proxy vulnerability is captured by the pseudoprice – a predictable outcome given the sensitive attribute's acceptability during the study period. Both partitions point to inexperienced drivers (low <code>DrivExp</code>) as a critical group, offering a clear path to intervene in the pseudoprice ratemaking algorithm or to apply post-processing adjustments for mitigating unfairness with respect to credit score.

Integrating the three dimensions of fairness in model assessment may form part of future actuarial standards. See the online supplement for code applied to simulated data.

Bender et al. (2025) groups biases as systemic, statistical, and human. Assuming D is not a true risk driver, risk spread flags segment-level systemic bias (measurement, sampling, label), and parity cost is the dollar cost to undo it. Proxy vulnerability estimates statistical bias from omitting D. With a data-driven spectrum<sup>7</sup>, pre-pricing metrics capture systemic and (potential) statistical bias; prior-pricing metrics can reflect all three to guide mitigation.

<sup>&</sup>lt;sup>7</sup>A spectrum anchored on a data-driven best-estimate premium (Option 1 of Complement 5).

					Proxy vulnerable groups				Commercially loaded groups									
Section	Variable	Statistic	Subpop	All Data	1	2	3	4-13	14	15	16	1	2	3	4-13	14	15	16
	Exposure	Sum	All	164,064	2,405	3,882	2,295	79,797	9,694	46,763	19,227	852	419	2,964	149,637	3,877	3,114	3,201
	Credrisk (lvl %)	Level %	Credrisk = 1	37.5%	81.6%	74.8%	70.1%	48.1%	33.6%	21.3%	18.3%	76.7%	73.3%	76.1%	35.9%	59.1%	36.2%	38.0%
	DrivExp VehAge	Mean Mean	All All	#######	1.44 5.22	3.53 4.93	6.82 4.31	#####	42.46 5.99	42.29 7.57	42.24 1.79	4.40 4.31	1.31 4.71	4.23 4.62	#####	3.70 5.52	30.87 2.33	31.18 2.33
	DrivAge	Mean	All	#######	21.89	22.30	24.36	#####	62.33	62.65	62.16	24.23	30.69	22.57	#####	20.52	50.08	51.42
Demographics	Zip Code (first	Level %	J	#######	####						####	####						####
grap	character only)	Level %	G	#######	####						####	####						####
gom	, , , , , , , , , , , , , , , , , , , ,	Level %	H	#######	####	####	####		####	####	####	####	####	####		####	####	####
De		Level % Level %	Employed Retired	60.6% 18.7%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
	ОссТуре	Level %	None	10.3%	58.2%	50.7%	26.1%	11.8%	5.3%	4.9%	3.4%	14.9%	25.0%	0.0%	8.4%	100.0%	2.8%	4.2%
		Level %	Other	10.5%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
	HasPropertyIns	Level %	Υ	#######	1000/	4000/	4000/		4000/	100%	100%	59%	56%	22%		24%	83%	83%
		Level % % of 0	lvl % : N All	97.2%	100% 94.1%	100% 95.1%	100% 96.2%	97.1%	100% 97.5%	97.8%	96.8%	41% 96.1%	44% 94.8%	78% 94.9%	97.4%	76% 95.8%	17% 96.1%	94.2%
ice	Loss	Mean	All	189.25	539.7	414.4	322.3	202.52	134	108.2	187.6	198.6	315.9	74.84	177.73	256.4	404	441.8
br pr	Severity	Mean	All	3743	4686	4544	4351	3948	3289	3046	3593	4126	4233	5441	3689	3809	2973	3672
pna	Severity	TVaR 0.95	All	19093	#####						#####	#####						#####
l pse	Ddi	VaR 0.05	All	68.40	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
and pseudo price	Pseudoprice	Mean VaR 0.95	All All	191.25 401.59	509.87	416.36	355.28	210.82	151.21	125.10	186.18	362.93	399.29	470.30	176.60	347.86	257.60	290.33
Loss	Performance	MAE	All	272.15	712.19	586.57	483.54		214.04	175.85	267.70	466.07	545.53		251.25	469.10		433.41
	metrics	Avg. Dev. Twee	e. All	137.20	226.26	200.34	192.22	150.27	117.06	106.12	138.29	221.77	233.75	233.01	132.23	174.53	157.82	180.45
Ļ.		Mean	D=0	152.74			290.66	171.48		114.67		301.75	326.45	369.51	145.13	274.23		252.51
Loss and price per group	Price	Mean VaR 0.95	D=1 D=0	255.31 288.64	530.92	447.06	382.84	253.27	182.59	163.65	234.15	381.53	425.76	502.00	232.75	398.71	321.53	351.98
orice ap		VaR 0.95	D=0 D=1	497.85	#####						#####	#####						#####
nd pric group	Logo Datio	-	D=0	87.4%	86.1%	95.4%	102.4%	88.6%	84.7%	80.5%	93.9%	76.3%	113.6%	84.9%	86.5%	102.1%	74.6%	111.6%
ss a	Loss Ratio	-	D=1	102.3%		100.5%		101.5%		101.9%					102.3%	77.3%		137.8%
Po	Est. P(D=1 X)	Mean	All	0.49	0.82	0.76	0.73	0.54	0.38	0.26	0.21	0.77	0.80	0.77	0.47	0.65	0.45 101.78	0.47
	Wass. Dist. Best-est.	Mean	vs Prem(d = 1) All	103.99 191.16	113.22 537.97	122.78 437.34	94.74 399.48	82.14 213.28	48.01 147.83	49.88 116.51	59.27 184.78	82.84 300.44	99.28 347.84	136.92 454.87	88.42 174.57	392.08	277.52	345.69
ss	Unaware	Mean	All	191.16	535.91	435.05	401.29	213.52		116.52		300.39	350.52		174.58		277.18	
Fairness Spectrum	Aware	Mean	All	191.16	468.35	389.01	366.31	209.20	153.09	125.43	200.63	271.32	317.39	400.77	176.35	368.23	283.84	351.65
Fa Spe	Hyperaware	Mean	All	191.16	429.22	360.74	339.62	204.61		132.10	216.09	247.98	288.18		177.49	350.26		357.14
	Corrective	Mean	All	191.16 -0.08	426.82 67.56	359.46 46.03	341.01 34.98	204.73 4.32	157.32 -5.51	131.99 -8.91	-16.07	247.65	290.23 33.13	370.30 53.38	177.40 -1.77	349.98 24.23	290.57 -6.66	357.76
Pre-pricing local metrics		Mean	D=1	7.85	68.59	48.37	36.34	7.66	-3.54	-8.02	-14.82	30.13	35.67	56.35	4.59	32.08	0.84	-5.59 4.04
met	Proxy	Mean	D=0	-4.84	63.01	39.13	31.81	1.23	-6.51	-9.15	-16.35	25.59	26.14	43.94	-5.33	12.86	-10.92	-11.51
cal	vulnerability	TVaR 0.95	All	57.97	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
lg lc		TVaR 0.95	D=1	76.71	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
riciı	Risk spread	TVaR 0.95 Mean	D=0 All	34.31 65.90	178.05	150.62	130.62	71.93	##### 51.53	42.80	65.42	##### 87.67	101.95	##### 157.52	60.30	138.71	91.90	118.66
e-p	Fairness spread	Mean	All	55.38			118.66	61.22	42.49	33.43	54.97	86.73	101.53		50.23	116.34	83.26	104.59
Pı	Parity cost	Mean	All	0.12	-111.15	-77.88	-58.47	-8.55	9.48	15.48	30.73	-52.79	-57.62	-84.57	2.83	-42.10	13.05	12.07
		Mean	All	0.09	-28.10	-20.99	-44.20	-2.46	3.37	8.59	1.39	62.50	51.45	15.43	2.03	-44.22	-19.92	-55.36
t-pricing local metrics	Commor=:-1	Mean	D=1 D=0	-8.66	-42.25	-34.20	-56.53	-9.62	3.02	10.49	-1.60	57.83	43.62	1.80	-5.88	-68.30	-23.91	-73.51
ing l	Commercial loading	Mean TVaR 0.95	D=0 All	5.35 72.32	34.74	18.12	-15.30	4.17	3.55	8.08	2.06	77.86	72.99	58.74	6.47	-9.37	-17.66	-44.23
pricing		TVaR 0.95	D=1	77.57	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
		TVaR 0.95	D=0	68.71	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
Pos	Comm. burden	Mean	All	10.7%	0.0%	-1.4%	-7.8%	5.3%	16.3%	24.0%	4.1%	24.5%	17.6%	9.3%	12.0%	-9.0%	-5.9%	-14.6%
H	Excess lift	Mean Prop > 0%	All	-11.53 56.5%	-76.78 36.5%	-57.63 37.0%	-43.98 24.4%	-14.43 50.8%	-4.43 61.1%	-0.16 68.8%	-9.38 58.6%	-19.51 93.8%	-26.73 86.0%	-56.78 64.6%	-8.98 58.6%	-51.46 19.8%	-14.26 29.6%	-33.75 8.7%
Extra summary statistics		Prop > 0%	D=1	48.3%	29.2%	28.5%	18.6%	44.8%	60.8%	69.9%	54.1%	92.0%	81.4%	55.3%	50.3%	8.4%	29.5%	5.9%
tati		Prop > 0%	D=0	61.5%	69.0%	62.1%	37.9%	56.3%		68.5%		99.6%	98.4%		63.3%	36.4%		10.5%
ıry s	Commercial	Prop > 9.2%	All	38.2%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
nma	burden	Prop > 9.2% Prop > 9.2%	D=1 D=0	30.3% 43.0%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
san		Prop > 9.2% Prop > 16.9%		25.9%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
xtra		Prop > 16.9%		19.2%	#####						#####	#####						#####
畄		Prop > 16.9%	D=0	29.9%	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####	#####
					AND < 3	3 <=	5 <= ) AND >= 1		ID 32 <= DrivExp	0 <= rage is N	30 <= /erage is Y	7 AND sus is	7 AND sus is	7 AND sus is		AND s is AND	< 6 ZIP ial]	< 6 ZIP ial]
					hasPropertyIns is N AND DrivExp < 3				AND 32 Driv	s is Y AND 30 <= newValueCoverage is N	asPropertyIns is Y AND 30 <= DrivExp AND newValueCoverage is Y	rat owe	DriverExp < 7 AND MaritalStatus is Married/Widowed AND	DriverExp < 7 AND MaritalStatus is Single/Divorced AND		DriverExp < 7 AND MaritalStatus is Single/Divorced AND	xp >= 7 AND VehAge < 6 msHistory >= 4 AND ZIF Code is [Confidential]	DriverExp >= 7 AND VehAge < a AND ClaimSHistory >= 4 AND ZI Code is [Confidential
			Raw dec		ns i Driv	N AND Drivex	ertyIns is N AND DrivExp < 10 NbTraffViolation		N AN	Y AN	Y AND	Exp alst Wido	DriverExp < MaritalStat rried/Widow	Exp alst ivor		Exp alst ivor	Veh = 4 nfid	Vel = 4 nfid
			rule fro		rtyI	i S	Ins is N DrivExp affViolat		-H	is 'ewVa	is 'ewVa	DriverExp MaritalSi rried/Wide	iver arit ied/h	DriverExp MaritalS ngle/Divo		DriverExp MaritalS ngle/Divo	AND ry v	AND Cy V
			p e	.0	obei	yIn:	yIn Di raft		Ins	Ins ID n	Ins ID ne	Dr.	Dr	Dr. Maingl		Dr Maing]	sto:	sto
					asPr	pert	pert		erty	ertyIr p AND	erty p AN	Σ	Σ	S		S	xp > nsHi	xp >= msHis Code
					Ë	hasPropertyIns	hasPropertyIns Dr NbTraff		hasPropertyIns	hasPropertyIns DrivExp AND ne	hasPropertyIns DrivExp AND ne						DriverExp >= 7 AND ND ClaimsHistory >= Code is [Con	/erE
						has	has		hasF	hasF	hasF						Driv AND C	Driv
																	⋖	Ø

Figure 14. Fairness monitoring table summarizing demographic, pricing, and fairness metrics for partitioned subpopulations in the Case study. Masking symbols ("#") preserve the partner insurer's confidentiality.

#### Section 7. Discussion

Our Case study grounds fairness in a large-scale, realistic setting. The fairness spectrum translates dimensions into pricing benchmarks. It provides context to judge any given (commercial) price. It also supports intuitive metrics, like proxy vulnerability by group or the share of policyholders facing high commercial burden, expressed in dollars and policyholders. By making fairness practical, we hope actuaries engage. The collaboration which made this Case study possible reflects insurer's interest in understanding fairness: its meaning, materiality, and actuarial relevance.

In our Case study, pre-pricing analysis shows proxy vulnerability is material and skewed: while many receive small rebates, some face 15–30% overpricing. The pseudoprice burdens the vulnerable group D=1 more than the D=0 group, but in a lesser extent than risk alone would justify, suggesting efforts toward solidarity.

Partitioning before and after pricing (§6) extends fairness analysis beyond large protected groups. Combined with local metrics (§5), it supports tools like Fig. 14 to surgically adjust for fairness within specific subgroups. In our Case study, the pseudoprice appears to capture proxy vulnerability, because commercially loaded groups also exhibit high proxy vulnerability. This surfaced fairness concerns in specific policyholder segments, like inexperienced drivers, that were not initially flagged.

In this article, we progress from fairness principles to detection, under assumptions that warrant scrutiny. In the current state of research, the three dimensions of fairness presented in §3 are necessary, but their exhaustivity remains an open question. Also, fairness dimensions are general – multiple premiums may reflect the same dimension. Both a corrective and a flat-rate price satisfy solidarity, suggesting at a broader range of models. We also ignored uncertainty in estimating benchmarks. How should we account for it?

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

- 1. Fairness often assumes **access to protected attributes**, which may be unavailable. Can we assess fairness without them? Predicting D (for example, with BIFSG as in Voicu, 2018) and Census data help, but are no substitute for direct access.
- 2. **Market dynamics** are ignored; portfolio fairness may conflict with market fairness (Côté et al., 2024). Can insurers contribute to market fairness using their own data?
- 3. Fairness is typically studied as a one-year objective, but its **long-term welfare** effects remain unclear (Shimao et al., 2022). Which fairness approach perpetuate, mitigate, or reverse disparities over time?
- 4. Seemingly neutral variables can mediate the link between protected traits and losses. Behavioral data may attenuate the impact of protected attributes on premiums by enriching X and detailing the causal risk chain (Boucher and Pigeon, 2024). This offers actuarial justification for disparities, but does it resolve proxy issues?
- 5. Insurance operates between law and statistics: one demands fairness case by case; the other defends differentiation at scale. Applying anti-discrimination **regulations** meets resistance where actuarial justification holds authority. Can regulations reconcile these perspectives to fairly serve insurers, regulators, and policyholders?

## Statement on the use of generative AI

We used generative AI tools to refine wording and syntax, draft or refactor small code snippets, and accelerate literature discovery (query formulation and complementing other bibliographic search tools). All AI outputs (text, code, and references) were independently reviewed, verified, and edited before inclusion; no analysis, modeling choices, or conclusions were delegated to these tools. We take full responsibility for the accuracy and integrity of all content in this publication.

## **Acknowledgements**

We are grateful for the support and helpful comments of the CAS Task Force members: Elizabeth Bellefleur-MacCaul, Mallika Bender, Denise Cheung, Jingwen Li, Shayan Sen, and Craig Sloss. We thank Marouane II Idrissi for discussions on projections, Fei Huang for discussions on encompassing frameworks for commercial prices, Agathe Fernandes Machados for assistance with Equipy, and Ewen Gallic for guidance on economic aspects.

The first author is thankful for the financial support from the Society of Actuaries' Hickman Scholars Program, Desjardins General Insurance Group, from MITACS, Fonds de recherche du Québec, the Chaire d'actuariat of Université Laval, Leadership Université Laval, and the Viger Family. The second author acknowledges the National Sciences and Engineering Research Council (NSERC) for funding (RGPIN-2019-04190). AC acknowledges the SCOR Foundation for Science and NSERC for funding (RGPIN-2019-07077).

## References

- Arrow, K. J. (1963). Uncertainty and the welfare economics of medical care. *The American Economic Review*, 53(5):941–973.
- ASB (2005). Actuarial Standard of Practice No. 12: Risk Classification (for All Practice Areas). actuarialstandardsboard.org/asops/risk-classification-practice-areas/. Accessed: October 16, 2025.
- ASB (2017). Actuarial Standard of Practice No. 53: Estimating Future Costs for Prospective Property/-Casualty Risk Transfer and Risk Retention. actuarialstandardsboard.org/asops/estimating-future-costs-prospective-propertycasualty-risk-transfer-risk-retention/. Accessed: October 16, 2025.
- Bailey, R. A. and Simon, L. J. (1960). Two studies in automobile insurance ratemaking. *ASTIN Bulletin*, 1(4):192–217.
- Bank of Canada (2024). Indicators of financial vulnerabilities. bankofcanada.ca/rates/indicators/indicators-of-financial-vulnerabilities/. Accessed: October 16, 2025.
- Barry, L. (2020). Insurance, big data and changing conceptions of fairness. *European Journal of Sociology*, 61(2):159–184.
- Bender, M., Dill, C., Hurlbert, M., Lindberg, C., and Mott, S. (2022a). Understanding potential influences of racial bias on the P&C insurance: four rating factors explored. *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. Accessed: October 16, 2025.
- Bender, M., Dillon, D. L., Harbage, R. A., Hurta, K., and Mullen, B. J. (2022b). Approaches to address racial bias in financial services: Lessons for the insurance industry. *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. Accessed: October 16, 2025.

- Bender, M., Margaret, B., Krafcheck, E., Sloss, C., Wang, G., and Woods, M. (2025). Practical application of bias measurement and mitigation techniques in insurance pricing. *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. Accessed: October 16, 2025.
- 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. Accessed: October 16, 2025.
- Brobeck, S. and Hunter, J. R. (2021). Watch Where You're Going: New Research on Telematics and Auto Insurance. *Consumer Federation of America*. Retrieved from consumerfed.org/reports/watch-where-youre-going/. Accessed: October 16, 2025.
- Casualty Actuarial Society (1988). Statement of Principles Regarding Property and Casualty Insurance Ratemaking. casact.org/statement-principles-regarding-property-and-casualty-insurance-ratemaking. Originally published in 1988, rescinded in 2020, and reinstated in 2021 for reference. Accessed: October 16, 2025.
- Cavanaugh, L., Merkord, S., Davis, T., and Heppen, D. (2024). Regulatory perspectives on algorithmic bias and unfair discrimination. *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. Accessed: October 16, 2025.
- Charpentier, A. (2024). Insurance, Biases, Discrimination and Fairness. Springer, New York.
- Charpentier, A., Hu, F., and Ratz, P. (2023). Mitigating discrimination in insurance with Wasserstein barycenters. *BIAS 2023, 3rd Workshop on Bias and Fairness in AI, International Workshop of ECML-PKDD*.
- 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.
- Denuit, M., Huyghe, J., Trufin, J., and Verdebout, T. (2024). Testing for auto-calibration with Lorenz and Concentration curves. *Insurance: Mathematics and Economics*, 117:130–139.
- Embrechts, P. and Wüthrich, M. V. (2022). Recent challenges in actuarial science. *Annual Review of Statistics and Its Application*, 9(1):119–140.
- Fahrenwaldt, M., Furrer, C., Hiabu, M. E., Huang, F., Jørgensen, F. H., Lindholm, M., Loftus, J., Steffensen, M., and Tsanakas, A. (2024). Fairness: plurality, causality, and insurability. *European Actuarial Journal*.
- 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. Accessed: October 16, 2025.
- Frees, E., Meyers, G., and Derrig, R. (2016). *Predictive Modeling Applications in Actuarial Science: Volume 2, Case studies in insurance*. Cambridge University Press.
- Grubinger, T., Zeileis, A., and Pfeiffer, K.-P. (2014). evtree: Evolutionary learning of globally optimal classification and regression trees in r. *Journal of Statistical Software*, 61(1):1–29.
- Insurance Information Institute (2023). Insurance 101: The history of insurance. iii.org/article/insurance-101. Accessed: October 16, 2025.
- 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.
- Kormes, M. (1935). The experience rating plan as applied to workmen's compensation risks part II. *Proceedings of the Casualty Actuarial Society (PCAS)*, 22:81–108.
- Leong, J., Moncher, R., and Jordan, K. (2024). A practical guide to navigating fairness in insurance pricing. 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. Accessed: October 16, 2025.
- 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. (2024a). Sensitivity-based measures of discrimination in insurance pricing. *Available at SSRN:* ssrn.com/abstract=4897265.
- Lindholm, M., Richman, R., Tsanakas, A., and Wüthrich, M. V. (2024b). What is fair? Proxy discrimination vs. demographic disparities in insurance pricing. *Scandinavian Actuarial Journal*, 2024(9):935–970.
- Loh, W.-Y. (2014). Fifty years of classification and regression trees. *International Statistical Review*, 82(3):329–348.
- Moodie, E. E. M. and Stephens, D. A. (2022). Causal inference: Critical developments, past and future. *Canadian Journal of Statistics*, 50(4):1299–1320.
- Mosley, R. and Wenman, R. (2022). Methods for quantifying discriminatory effects on protected classes in insurance. *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. Accessed: October 16, 2025.
- Prince, A. E. and Schwarcz, D. (2019). Proxy discrimination in the age of artificial intelligence and big data. *Iowa Law Review*, 105:1257.
- 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.
- Voicu, I. (2018). Using first name information to improve race and ethnicity classification. *Statistics and Public Policy*, 5(1):1–13.

## Appendix A. Glossary

Let Y denote the claim cost for a property-casualty coverage. Let **X** denote the vector of pricing covariates. Fairness is defined relative to a pre-specified (sensitive) variable D, taken here to be a single, binary, fully observed random variable. See §2 for details.

**Protected groups** — Levels of the sensitive attribute D.

**Vulnerable groups** — Subset of protected groups showing historic disadvantage.

**Fairness** — An elusive, contested ideal; a three-dimensional trade-off to manage (§3):

**Actuarial fairness** aligns premiums with expected losses, mitigating cross-subsidies, **Solidarity** aligns premiums across protected groups, mitigating disparities,

Causality ensures models capture only true risk factors, mitigating proxy effects.

**Spectrum of fair premiums** — The five fair premiums of the spectrum are (§4):

Premium	Best-estimate	Unaware	Aware	Hyperaware	Corrective
Notation	$\mu^B(\mathbf{x},d)$	$\mu^U(\mathbf{x})$	$\mu^A(\mathbf{x})$	$\mu^H(\mathbf{x})$	$\mu^C(\mathbf{x},d)$
Direct discrimination	✓	×	×	×	<b>✓</b>
Dimension prioritized	Actuarial fairness	Actuarial fairness	Causality	Solidarity	Solidarity

**Pre-pricing local metrics** reveal potential unfairness in the dataset (§5.1).

**Risk spread** — The spread of best-estimates across values of *D*:

$$\Delta_{\mathsf{risk}}(\mathbf{x}) = \left| \mu^B(\mathbf{x}, 1) - \mu^B(\mathbf{x}, 0) \right|.$$

**Proxy vulnerability** — The difference between ignoring D and controlling for it:

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

**Fairness range** — The range of estimates across the spectrum of fair premiums:

$$\begin{split} \Delta_{\text{fair}}(\mathbf{x},d) &= \max\{\mu^B(\mathbf{x},d), \mu^U(\mathbf{x}), \mu^A(\mathbf{x}), \mu^H(\mathbf{x}), \mu^C(\mathbf{x},d)\} - \\ &\quad \min\{\mu^B(\mathbf{x},d), \mu^U(\mathbf{x}), \mu^A(\mathbf{x}), \mu^H(\mathbf{x}), \mu^C(\mathbf{x},d)\}. \end{split}$$

Parity cost — The (monetary) cost of shifting from actuarial fairness to solidarity:

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

**Post-pricing local metrics** relate a given price  $\pi(\mathbf{X}, D)$  to the spectrum (§5.2):

**Commercial loading** — The difference between  $\pi$  and a chosen reference premium  $\mu$  from the spectrum:

$$\Delta_{\text{load}}(\mathbf{X}, d; \pi; \mu) = \pi(\mathbf{X}, d) - \mu(\mathbf{X}, d).$$

**Commercial burden** — The commercial loading as a percentage of  $\mu$ :

$$\rho_{\mathrm{burden}}(\mathbf{x},d;\pi,\mu) = \frac{\pi(\mathbf{x},d)}{\mu(\mathbf{x},d)} - 1 = \frac{\Delta_{\mathrm{load}}(\mathbf{x},d;\pi;\mu)}{\mu(\mathbf{x},d)}.$$

**Implied propensity** (Non-directly discr.  $\pi$ ) — The implicit weight on D=1:

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

**Excess lift** (Directly discr.  $\pi$ ) — The excess differentiation on D compared to  $\mu^B$ :

$$\Delta_{\mathrm{excess}}(\mathbf{x};\pi) = \left|\pi(\mathbf{x},1) - \pi(\mathbf{x},0)\right| - \left|\mu^B(\mathbf{x},1) - \mu^B(\mathbf{x},0)\right|.$$