

Bias, fairness and discrimination-free insurance pricing

Katrien Antonio

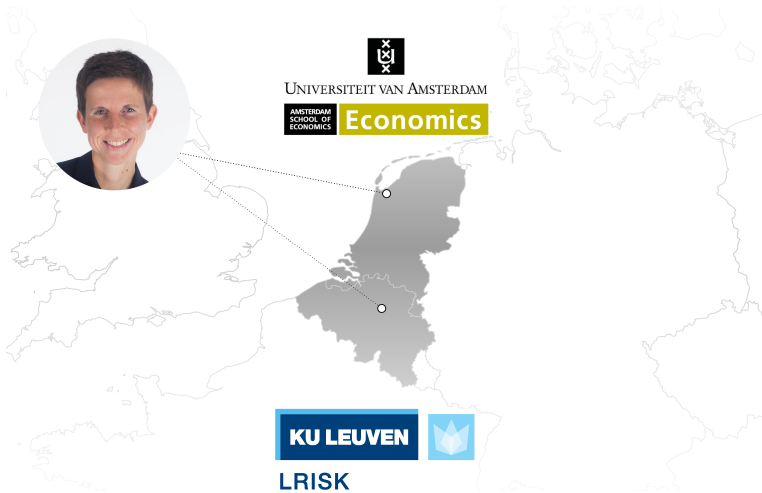
LRisk - KU Leuven and ASE - University of Amsterdam

September 22, 2022



AMSTERDAM
SCHOOL OF
ECONOMICS

Economics



My personal website: <https://katrienantonio.github.io>

This workshop's mission is threefold

3

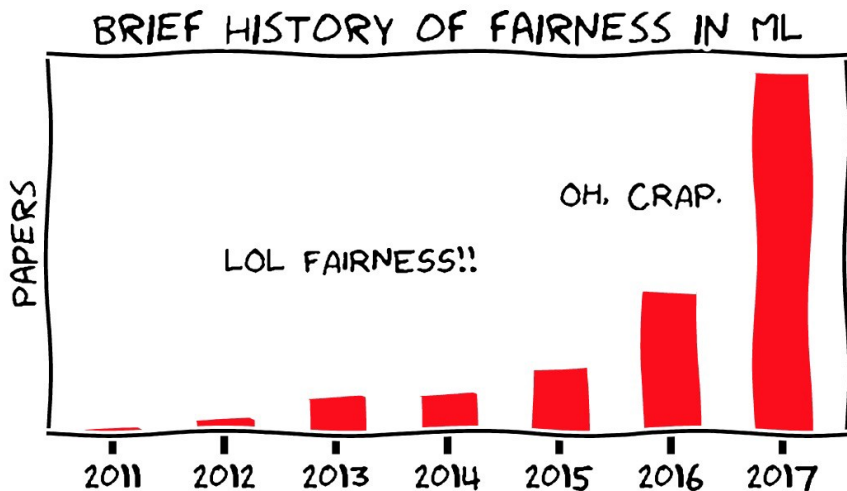
We aim to:

- (1) demystify the (exploding) literature on fair machine learning
- (2) explore in- and post-processing methods to establish fairness or to remove discrimination by proxy in insurance pricing
- (3) with GLM (\sim statistical learning) and GBM (\sim machine learning) based pricing methods.

A widely expanding literature

In machine learning ...

4



A widely expanding literature

In machine learning ...

Fairness Through Awareness

Cynthia Dwork* Moritz Hardt[†] Toniann Pitassi[‡] Omer Reingold[§]
 Richard Zemel[¶]

November 30, 2011

Algorithmic decision making and the cost of fairness

Sam Corbett-Davies Emma Pierson Avi Feller
 Stanford University Stanford University Univ. of California, Berkeley
 scorbett@stanford.edu emmap1@stanford.edu afeller@berkeley.edu

Sharad Goel Aziz Huq
 Stanford University University of Chicago
 sgoel@stanford.edu huq@uchicago.edu

Equality of Opportunity in Supervised Learning

Moritz Hardt Eric Price Nathan Srebro

October 11, 2016

**Fairness in Criminal Justice Risk Assessments:
The State of the Art**

Richard Berk^{a,b}, Hoda Heidari^c, Shahin Jabbari^c,
 Michael Kearns^c, Aaron Roth^c

**Fair prediction with disparate impact:
A study of bias in recidivism prediction instruments**

Alexandra Chouldechova *

**The Measure and Mismeasure of Fairness:
A Critical Review of Fair Machine Learning***

Sam Corbett-Davies Sharad Goel
 Stanford University Stanford University

August 14, 2018

HUMAN DECISIONS AND MACHINE PREDICTIONS*

JON KLEINBERG HIMABINDU LAKKARAJU
 JURE LESKOVEC JENS LUDWIG
 SENDIL MULLAINATHAN

The Frontiers of Fairness in Machine Learning

Alexandra Chouldechova* Aaron Roth[†]

October 23, 2018

A Survey on Bias and Fairness in Machine Learning

NINAREH MEHRABI, FRED MORSTATTER, NRIPSUTA SAXENA, KRISTINA LERMAN,
 and ARAM GALSTYAN, USC-ISI

FAIRNESS IN MACHINE LEARNING: A SURVEY

A PREPRINT

Simon Caton Christian Haas
 University College Dublin University of Nebraska at Omaha
 Dublin, Ireland Omaha, US
 simon.caton@ucd.ie christianhaas@unomaha.edu

DOI:10.1145/3270000

A group of industry, academic, and government experts convene in Philadelphia to explore the roots of algorithmic bias.

BY ALEXANDRA CHOULDECHOVA AND AARON ROTH

A Snapshot of the Frontiers of Fairness in Machine Learning

(Picture taken from A. Charpentier, presentation on Sept 9, see [here](#).)

A widely expanding literature

6

... and in actuarial science!

“Insurance is particularly interesting because the **entire industry is based on discrimination**. Here, we use the word discrimination in an **entirely neutral way**, taking it to mean the act of treating different groups differently. (...) The modern-day insurance industry is founded on the ability **to differentiate, or discriminate, among risks**, known as risk classification.” (Frees & Huang, 2021, NAAJ)

In fact, insurers discriminate among customers via a diverse set of actions, e.g.:

- the decision to insure
- the coverage offered
- by charging different prices.

A widely expanding literature

7

... and in actuarial science!

- ▶ Frees & Huang (2021, NAAJ) argue how the meaning of **actuarial fairness** or **fair insurance systems**, may depend on:
 - historical context, with a shift of responsibility from individual to a pool
 - the nature of the pool, e.g. a mutual company vs a stock insurance company
 - whether the insurance product can be viewed as a social or type of public good, or not.
- ▶ Today, I will interpret actuarial fairness as a pricing system where each customer should pay a **premium proportional to their own risk**.

Sources I used to put this workshop together:

- Mehrabi et al. (2022), [A survey on bias and fairness in machine learning](#), ACM Computing Surveys
- Caton & Haas (2021), [Fairness in machine learning: a survey](#)
- E. Frees & F. Huang (2021), [The discriminating \(pricing\) actuary](#), North American Actuarial Journal
- M. Lindholm, R. Richman, A. Tsanakas & M. Wüthrich (2022), [Discrimination-free insurance pricing](#), ASTIN Bulletin
- A. Charpentier (2022), [Insurance: discrimination, biases and fairness](#), Institut Louis Bachelier

“In the context of decision-making **fairness** is the **absence of any prejudice or favoritism toward an individual or group based on their inherent or acquired characteristics**. Thus, an unfair algorithm is one whose decisions are skewed toward a particular group of people.” (Mehrabi et al., 2019)

Mehrabi et al. (2019) identify two potential sources of unfairness:

- from biases in the data
- from the learning algorithms.

“Bias can exist in many shapes and forms, some of which can lead to unfairness in different downstream learning tasks.” (Mehrabi et al., 2019)

Most important **sources of bias** discussed in the literature: from

- **data to algorithm**, e.g. *measurement bias* or *omitted variable bias*
- **algorithm to user**, e.g. *algorithmic bias* where bias is not present in the data but purely added by the algorithm
- **user to data**, e.g. when data are user-generated, any inherent biases in users might be reflected in the data they generate.

► **Explainable** discrimination

- when differences in treatment and outcomes amongst different groups can be justified and explained via some (acceptable, non-protected) attributes.

► **Direct** discrimination

- when protected attributes of individuals **explicitly** result in non-favorable outcomes toward them.

► **Indirect** discrimination

- individuals appear to be treated based on seemingly neutral and non-protected attributes
- however, protected groups, or individuals still get to be treated unjustly as a result of **implicit effects from their protected attributes**.

Indirect discrimination: by proxy and disparate impact

► Proxy discrimination:

- arises from correlation between protected and unprotected characteristics
- the implicit ability to infer protected characteristics from other (legitimately used) policyholder features
- e.g. geographic area serves as a substitute for a protected variable such as race, or proxy produced by an AI that summarizes the effects of many variables.

► Disparate impact:

- a systematic disadvantage resulting for a group that is protected by a nondiscrimination provision.

Protected features?

“Grouping, or classifying, insureds into homogeneous categories for the purposes of risk sharing is at the heart of the insurance function. Many variables that insurers use are seemingly innocuous (e.g., blindness for auto insurance), yet others can be viewed as *wrong* (e.g., religious affiliation), *unfair* (e.g., onset of cancer for health insurance), *sensitive* (e.g., marital status), or *mysterious* (e.g., Artificial Intelligence produced).” (Frees & Huang, 2021, NAAJ)

Protected variables:

- not permitted in risk classification
- their choice is a normative one, after societal debate with many actors.

Discrimination

Protected or sensitive features?

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

Table 2.1: A factor is considered “permitted” (•) when there are no laws or regulatory policies in the state or province that prohibit insurers from using that factor. Otherwise, it will be “prohibited” (X). In North Carolina, age is only allowed when giving a discount to drivers 55 years of age and older. In Pennsylvania, credit score can be used for new business and to reduce rates at renewal, but not to increase rates at renewal. In Alberta, credit score and driver’s license seniority cannot be used for mandatory coverage (but can be used on optional coverage). In Labrador, age cannot be used before 55, and beyond that, it must be a discount (as in North Carolina).

(Picture taken from A. Charpentier, report on *Insurance: discrimination, biases and fairness*, see [here](#).)

Sensitive features?

Structure to identify whether or not a variable contains sensitive information, from Frees & Huang (2021), inspired by Avraham (2018) and Prince and Schwarcz (2020):

Property	Explanation
Control	If a policyholder has control over a certain characteristic, e.g. Type of car, it is deemed appropriate to use for risk classification.
Mutability	A variable for which the value changes over time is considered to be fair, as individuals have the chance to be on the winning and losing side during their lifetime. An example is the variable Age.
Causality	When there is a proven causal connection between the variable and the insured event, it is fair to use the variable for risk classification. Note that a lot of research needs to be done to confirm causality.
Statistical discrimination	If a variable makes no significant contribution to the predictive accuracy it is best not to use it.
Limiting or reversing the effects of past discrimination	A variable used to preserve a historical negative stereotype is deemed unfair.
Inhibiting socially valuable behaviour	A variable that, when used to classify risks, will hinder socially desirable behaviour is considered as an unfair variable.

A multitude of fairness criteria and corresponding metrics has been proposed in the ML literature, almost exclusively linked to classification problems (with target Y being 0 or 1).

A selection: (with Y the target, \hat{Y} the prediction and D the protected feature)

- **demographic or statistical parity:**

a predictor \hat{Y} satisfies demographic parity if $\mathbb{P}(\hat{Y} \mid D = 0) = \mathbb{P}(\hat{Y} \mid D = 1)$

- **equal opportunity:**

a predictor \hat{Y} satisfies equal opportunity with respect to protected attribute D and outcome Y , if $\mathbb{P}(\hat{Y} = 1 \mid D = 0, Y = 1) = \mathbb{P}(\hat{Y} = 1 \mid D = 1, Y = 1)$. Thus, equal true positive rates for protected and unprotected groups.

See e.g. <http://research.google.com/bigpicture/attacking-discrimination-in-ml/> for visuals and <https://developers.google.com/machine-learning/glossary/fairness>

for detailed illustrations and definitions.

A multitude of fairness criteria and corresponding metrics has been proposed in the ML literature, almost exclusively linked to classification problems (with target Y being 0 or 1).

A selection:

- **fairness through awareness:**

an algorithm is fair if it gives similar predictions to similar individuals

- **fairness through unawareness:**

an algorithm is fair as long as any protected attributes D are not explicitly used in the decision-making process.

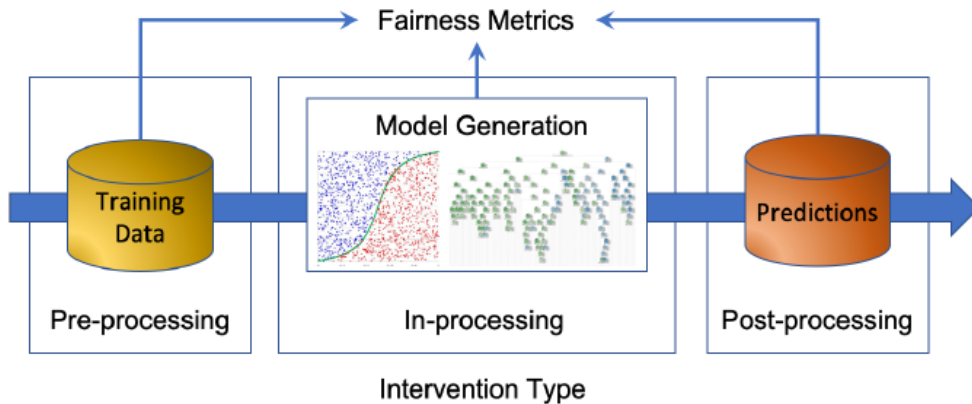
Some reflections:

- very much focused on classification problems!
- impossible to satisfy some of these fairness conditions simultaneously, except in some very special cases.
- users must decide on where to place emphasis, but be mindful of the trade off between any fairness measure and model accuracy.

Mitigation strategies

Methods for fair machine learning

19



(Picture taken from Caton & Haas, see [here](#).)

Mitigation strategies

Methods for fair machine learning

With **pre-processing**:

- transform the data with the aim to remove bias/discrimination from the training data
- then fit a learning model on the *repaired* data

With **in-processing**: (*~ Marchi, Antonio, et al., 2022, ongoing*)

- try to find a balance between multiple model objectives, e.g. accuracy and fairness

With **post-processing**: (*~ Lindholm et al., 2022, ASTIN*)

- apply transformations to model output to improve prediction fairness
- only needs access to the predictions and sensitive attribute information.

In ongoing research with Marchi (KU Leuven), Avanzi and Zhou (Uni of Melbourne), we aim to establish fair insurance pricing **via regularization** (\sim Lasso)

$$\min_{w_h} \left\{ \sum_{i=1}^n \mathcal{L}(h_{w_h}(x_i), y_i) + \lambda \cdot \Phi(h_{w_h}, \{x_i\}_{i=1}^n, \{y_i\}_{i=1}^n, \{d_i\}_{i=1}^n) \right\}. \quad (1)$$

where the data are $(x_i, y_i, d_i) \in \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{D}^n$, a generic loss function $\mathcal{L} : \mathbb{R}^2 \rightarrow \mathbb{R}$, a predictive model $h_{w_h} \in \mathcal{F}$ and a fairness criterion $\Phi : \mathcal{F} \times \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{D}^n \rightarrow \mathbb{R}$.

Hence, the goal is to balance accuracy (via \mathcal{L}) and fairness (via Φ).

Lindholm et al. (2022):

- propose a **post-processing** pricing adjustment formula that explicitly addresses discrimination by proxy
- with the goal **to remove indirect discrimination** – if it happens to exist – from insurance pricing models
- purely reason from an **actuarial** rather than a legal perspective
- assume **knowledge of a policyholder's discriminatory features**.

No connection with fairness, or a metric to express fairness. In fact, **open discussion** going on the connection between insurance pricing and notions of fairness.

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

$n_{i,j}$	Small	Large	Row total
Female	83	14	97
Male	47	60	107
Column total	130	74	204

Table: Fictional claim counts.

$e_{i,j}$	Small	Large	Row total
Female	567	83	650
Male	269	253	522
Column total	836	336	1172

Table: Fictional exposures.

Here:

- $i \in \{0, 1\}$ refers to Gender with $i = 1$ a *male* driver
- $j \in \{0, 1\}$ refers to Type of car with $j = 1$ a *large* car.

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

Let us focus on the following set of questions:

- how to estimate the overall expected claim frequency, say $\hat{\lambda}$, in this portfolio?
- what is $\hat{\lambda}_{ij}$ when risks are classified using both Gender and Type of car?
- what if prices can not discriminate based on Gender, hence: Gender is a protected feature?

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

Overall expected claim frequency:

$$\hat{\lambda} = \frac{\sum_{i,j} n_{i,j}}{\sum_{i,j} e_{i,j}} = \frac{n_{\bullet}}{e_{\bullet}} = \frac{204}{1172} = 0.174.$$

Price discrimination, or risk classification, using both Gender and Type of car leads to **best-estimate prices**:

$\hat{\lambda}_{i,j}$	Small	Large
Female	0.146	0.169
Male	0.175	0.237

where, for instance, $\hat{\lambda}_{0,0} = \frac{83}{567} = 0.146$. Try to summarize some findings here!

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

Some initial observations:

- claim frequencies higher for men than for women
- male drivers with a large car are high risk
- claim frequencies higher for large car drivers compared to small cars.

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

Now we treat Gender as a protected feature, and proceed by **simply ignoring** the variable in the pricing:

$$\hat{\lambda}_{\bullet,j} = \frac{n_{\bullet,j}}{e_{\bullet,j}} = \frac{n_{0,j} + n_{1,j}}{e_{0,j} + e_{1,j}} \text{ with } j \in \{0, 1\}$$

which results in

$$\begin{aligned} \hat{\lambda}_{\bullet,0} &= \frac{130}{836} = 0.156 \text{ for a } \textit{Small} \text{ car} \\ \hat{\lambda}_{\bullet,1} &= \frac{74}{336} = 0.220 \text{ for a } \textit{Large} \text{ car.} \end{aligned}$$

These prices adhere to the **unawareness** principle and are called **unawareness prices**.

Are these prices truly not discriminating based on Gender?

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

Type of car is very informative for Gender in our example.

Indeed,

$$\begin{aligned}\hat{\mathbb{P}}(\text{Male} \mid \text{Large}) &= \frac{e_{1,1}}{e_{0,1} + e_{1,1}} = \frac{253}{336} = 0.75 \\ \hat{\mathbb{P}}(\text{Male} \mid \text{Small}) &= \frac{e_{1,0}}{e_{0,0} + e_{1,0}} = \frac{269}{836} = 0.32.\end{aligned}$$

Thus, higher propensity of drivers with a *Large* car to be *Male* drivers.

At portfolio level, $\hat{\mathbb{P}}(\text{Male}) = \frac{522}{1172} = 0.45$ and $\hat{\mathbb{P}}(\text{Female}) = \frac{650}{1172} = 0.55$.

Indirect discrimination

Motivating example, from Baeten (2021), inspired by Lindholm et al. (2022)

In fact, the unawareness price $\hat{\lambda}_{\bullet,1}$ for a driver of a *Large* car can be rewritten as:

$$\begin{aligned}
 \hat{\lambda}_{\bullet,1} &= \frac{n_{0,1} + n_{1,1}}{e_{0,1} + e_{1,1}} = \frac{n_{0,1}}{e_{0,1}} \cdot \frac{e_{0,1}}{e_{0,1} + e_{1,1}} + \frac{n_{1,1}}{e_{1,1}} \cdot \frac{e_{1,1}}{e_{0,1} + e_{1,1}} \\
 &= \hat{\lambda}_{0,1} \cdot \frac{e_{0,1}}{e_{0,1} + e_{1,1}} + \hat{\lambda}_{1,1} \cdot \frac{e_{1,1}}{e_{0,1} + e_{1,1}} \\
 &= \hat{\lambda}_{0,1} \cdot \hat{\mathbb{P}}(\text{Female} \mid \text{Large}) + \hat{\lambda}_{1,1} \cdot \hat{\mathbb{P}}(\text{Male} \mid \text{Large}) \\
 &= 0.169 \cdot 0.25 + 0.237 \cdot 0.75 \\
 &= 0.22.
 \end{aligned}$$

Not only information about the influence of *Type* of car on producing a claim is used, but also about the propensity of drivers with a specific *Type* to be *male* or *female*.

The correlation between *Type* of car and *Gender* is exploited.

Motivating example

31

From Baeten (2021), inspired by Lindholm et al. (2022)

Similarly, the unawareness price $\hat{\lambda}_{\bullet,0}$ for a driver of a *Small* car becomes:

$$\begin{aligned}\hat{\lambda}_{\bullet,0} &= \frac{n_{0,0} + n_{1,0}}{e_{0,0} + e_{1,0}} = \frac{n_{0,0}}{e_{0,0}} \cdot \frac{e_{0,0}}{e_{0,0} + e_{1,0}} + \frac{n_{1,0}}{e_{1,0}} \cdot \frac{e_{1,0}}{e_{0,0} + e_{1,0}} \\ &= \hat{\lambda}_{0,0} \cdot \frac{e_{0,0}}{e_{0,0} + e_{1,0}} + \hat{\lambda}_{1,0} \cdot \frac{e_{1,0}}{e_{0,0} + e_{1,0}} \\ &= \hat{\lambda}_{0,0} \cdot \hat{\mathbb{P}}(\text{Female} \mid \text{Small}) + \hat{\lambda}_{1,0} \cdot \hat{\mathbb{P}}(\text{Male} \mid \text{Small}) \\ &= 0.146 \cdot 0.678 + 0.175 \cdot 0.322 \\ &= 0.156.\end{aligned}$$

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

We see the **potential for indirect discrimination** reflected in the **unawareness price**.

Indeed, in our example:

- men cause on average more claims than women
- in the sub-population of *Large* car drivers men are more prevalent, compared to the Gender composition at portfolio level
- hence, the unawareness price for *Large* is leveraged, and vice versa for *Small* car drivers.

Motivating example

33

From Baeten (2021), inspired by Lindholm et al. (2022)

Lindholm et al. (2022) replace the conditional probabilities by **unconditional probabilities** to obtain a **discrimination-free** price:

$$\hat{\lambda}_{\bullet,j}^{\text{DF}} = \hat{\lambda}_{0,j} \cdot \hat{\mathbb{P}}(\text{Female}) + \hat{\lambda}_{1,j} \cdot \hat{\mathbb{P}}(\text{Male}).$$

The discrimination-free expected claim frequencies for *Small* and *Large* car drives then becomes:

$$\hat{\lambda}_{\bullet,0}^{\text{DF}} = 0.146 \cdot 0.55 + 0.175 \cdot 0.45 = 0.168 > \hat{\lambda}_{\bullet,0} = 0.156$$

$$\hat{\lambda}_{\bullet,1}^{\text{DF}} = 0.169 \cdot 0.55 + 0.237 \cdot 0.45 = 0.199 < \hat{\lambda}_{\bullet,1} = 0.22.$$

Note: the resulting price list still discriminates on the basis of Type of car (as it should do).

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

The **best estimate** and **unawareness** prices in this example are **unbiased**.

Indeed, assume claim cost is 1 EUR, then at portfolio level:

$$\sum_{i,j} \hat{\lambda}_{i,j} \cdot e_{i,j} = 567 \cdot 0.146 + 83 \cdot 0.169 + 269 \cdot 0.175 + 253 \cdot 0.237 = 204.$$

and

$$\sum_j \hat{\lambda}_{\bullet,j} \cdot e_{\bullet,j} = 836 \cdot 0.156 + 336 \cdot 0.22 = 204.$$

Thus, at portfolio level, the total premium volume (here: expected claim counts) equals the observed total loss (here: observed claim counts).

Motivating example

From Baeten (2021), inspired by Lindholm et al. (2022)

However, the **discrimination-free** prices **imply a bias** at portfolio level. Indeed,

$$\hat{\lambda}_{\bullet,0}^{\text{DF}} \cdot (e_{0,0} + e_{1,0}) + \hat{\lambda}_{\bullet,1}^{\text{DF}} \cdot (e_{0,1} + e_{1,1}) = 0.168 \cdot 836 + 0.199 \cdot 336 = 207.3 > 204.$$

To de-bias, one option is to adjust $\hat{\mathbb{P}}(\text{Female})$ and $\hat{\mathbb{P}}(\text{Male})$ so that portfolio bias is removed.

In our example, choose $\hat{\mathbb{P}}^*(\text{Female}) = 0.65$ and $\hat{\mathbb{P}}^*(\text{Male}) = 0.35$, then

$$\begin{aligned}\hat{\lambda}_{\bullet,0}^{\text{DF}^*} &= 0.162 \cdot 0.65 + 0.175 \cdot 0.35 = 0.167 \\ \hat{\lambda}_{\bullet,1}^{\text{DF}^*} &= 0.169 \cdot 0.65 + 0.237 \cdot 0.35 = 0.193,\end{aligned}$$

and at portfolio level

$$\hat{\lambda}_{\bullet,0}^{\text{DF}^*} \cdot (e_{0,0} + e_{1,0}) + \hat{\lambda}_{\bullet,1}^{\text{DF}^*} \cdot (e_{0,1} + e_{1,1}) = 0.167 \cdot 836 + 0.193 \approx 204.$$

Let \mathbf{D} correspond to some protected (or: sensitive) features and let the non-protected features be denoted by \mathbf{X} .

Let Y denote the target variable of interest, e.g. number of claims or claim sizes.

Lindholm et al. (2022) then consider three pricing formulas:

- the **best-estimate** price, where both \mathbf{D} and \mathbf{X} are used
- the **unawareness** price, simply ignoring \mathbf{D}
- the **discrimination-free** price, where best-estimate prices are averaged over discriminatory covariates.

The **best-estimate** price for Y wrt (\mathbf{X}, \mathbf{D}) is:

$$\mu(\mathbf{X}, \mathbf{D}) := \mathbb{E}[Y \mid \mathbf{X}, \mathbf{D}].$$

Hereby, the price list is

- in general not discrimination-free, unless $\mu(\mathbf{X}, \mathbf{D})$ reduces to $\mu(\mathbf{X})$ because of independence between \mathbf{X} and \mathbf{D}
- obtained from some predictive model that **can use both** \mathbf{X} and \mathbf{D} .

The **unawareness** price for Y wrt (\mathbf{X}) is:

$$\mu(\mathbf{X}) := \mathbb{E}[Y \mid \mathbf{X}].$$

Hereby, the price list

- attempts at avoiding discrimination by simply ignoring the protected features in \mathbf{D}
- may produce indirect discrimination via (\mathbf{D} proxied by \mathbf{X})

$$\mu(\mathbf{X}) = \int_{\mathbf{d}} \mu(\mathbf{X}, \mathbf{d}) d\mathbb{P}(\mathbf{D} = \mathbf{d} \mid \mathbf{X}),$$

where the conditional probability enables inference of protected features \mathbf{D} via unprotected features \mathbf{X} .

A discrimination-free price for Y wrt (\mathbf{X}) is:

$$h^*(\mathbf{X}) := \int_{\mathbf{d}} \mu(\mathbf{X}, \mathbf{d}) d\mathbb{P}^*(\mathbf{D} = \mathbf{d}).$$

Hereby, the price list

- averages best-estimate prices over discriminatory covariates $\sim \mathbb{P}^*(\mathbf{D} = \mathbf{d})$
- is free of direct discrimination, since $h^*(\mathbf{X})$ does not explicitly use \mathbf{D}
- is free of indirect discrimination, since the possible explanatory power that \mathbf{X} may have for \mathbf{D} is removed
- is in general not unbiased.

In fact, the problem that arises with the unawareness prices is referred to in econometrics as **omitted variable bias**.

When features \mathbf{X} included in a linear regression model are correlated with omitted features \mathbf{Z} , the \mathbf{X} will (partially) proxy for the \mathbf{Z} and the estimated regression parameters β will be biased.

For more details:

- Pope & Sydnor (2011) on *Implementing Anti-Discrimination Policies in Statistical Profiling Models* in the American Economic Journal: Economic Policy
- [our Colab](#).

Simulated example

We now demonstrate technical insurance pricing with:

- Generalized Linear Models (GLMs)
- Gradient Boosting Machines (GBMs).

We use the set-up from Lindholm et al. (2022) and consider the best-estimate, the unawareness and discrimination-free pricing formulas.

Set up

Consider three covariates, X_1 , X_2 (unprotected) and D (protected), with

- $D \in \{\text{female}, \text{male}\}$
- $X_1 \in \{15, \dots, 80\}$, the age of the policyholder
- $X_2 \in \{\text{non-smoker}, \text{smoker}\}$.

For the \mathbb{P} distribution of (\mathbf{X}, D) we assume:

$$\mathbb{P}(D = \text{female}) = 0.45$$

$$\mathbb{P}(X_2 = \text{smoker}) = 0.3$$

$$\mathbb{P}(D = \text{female} \mid X_2 = \text{smoker}) = 0.8.$$

Set up

The example assumes **three types of health related claims**:

(1) type 1 (related to giving birth)

$$\lambda_1(\mathbf{X}, D) := \exp(\alpha_0 + \alpha_1 \mathbf{1}_{\{X_1 \in [20, 40]\}} \mathbf{1}_{\{D = \text{female}\}})$$

(2) type 2 (cancer related)

$$\lambda_2(\mathbf{X}, D) := \exp(\beta_0 + \beta_1 X_1 + \beta_2 \mathbf{1}_{\{X_2 = \text{smoker}\}} + \beta_3 \mathbf{1}_{\{D = \text{female}\}})$$

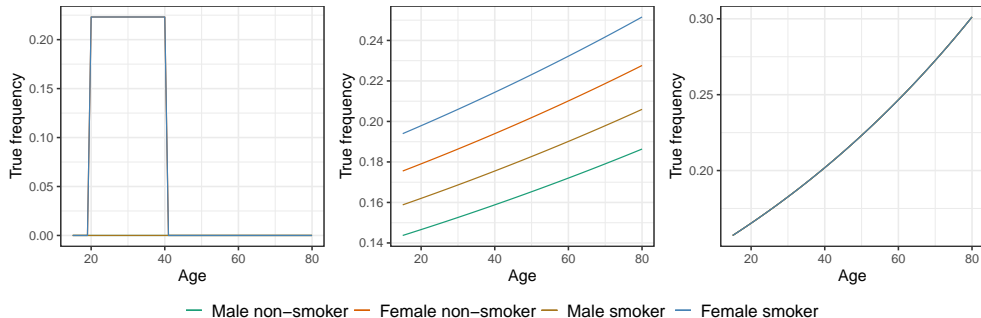
(3) type 3 (other diseases)

$$\lambda_3(\mathbf{X}, D) := \exp(\gamma_0 + \gamma_1 X_1).$$

Stylized example

Set up

44



Assumptions on previous sheet with $(\alpha_0, \alpha_1) = (-40, 38.5)$,
 $(\beta_0, \beta_1, \beta_2, \beta_3) = (-2, 0.004, 0.1, 0.2)$ and $(\gamma_0, \gamma_1) = (-2, 0.01)$.

Let us assume **deterministic claim amounts**, namely $(c_1, c_2, c_3) = (0.5, 0.9, 0.1)$.

Using these assumptions, the pricing formulas become:

- the best-estimate price

$$\mu(\mathbf{X}, D) = c_1 \lambda_1(\mathbf{X}, D) + c_2 \lambda_2(\mathbf{X}, D) + c_3 \lambda_3(\mathbf{X}, D)$$

- the unawareness price

$$\mu(\mathbf{X}) = \sum_{d \in \{\text{female}, \text{male}\}} (c_1 \lambda_1(\mathbf{X}, d) + c_2 \lambda_2(\mathbf{X}, d) + c_3 \lambda_3(\mathbf{X}, d)) \mathbb{P}(D = d \mid \mathbf{X})$$

- a discrimination-free price

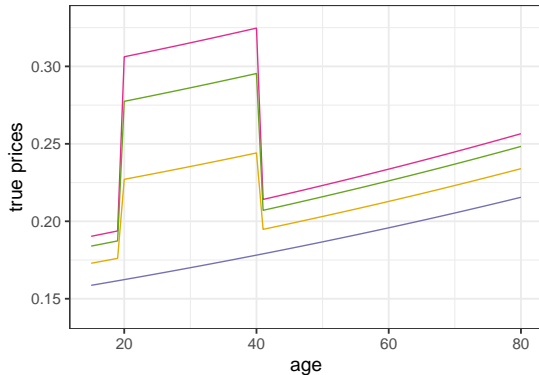
$$h^*(\mathbf{X}) = \sum_{d \in \{\text{female}, \text{male}\}} (c_1 \lambda_1(\mathbf{X}, d) + c_2 \lambda_2(\mathbf{X}, d) + c_3 \lambda_3(\mathbf{X}, d)) \mathbb{P}(D = d).$$

To calculate the unawareness prices, we use:

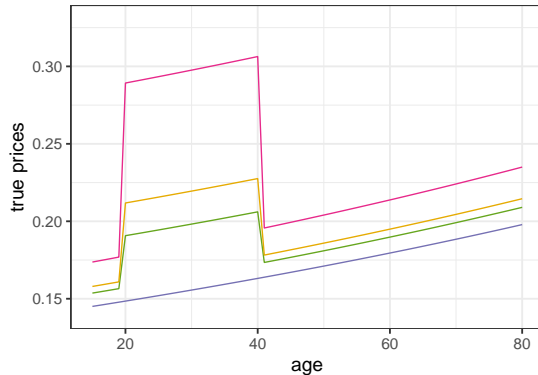
$$\begin{aligned}\mathbb{P}(D = d \mid \mathbf{X}) &= \frac{\mathbb{P}(D = d, \mathbf{X} = \mathbf{x})}{\mathbb{P}(\mathbf{X} = \mathbf{x})} \\ &= \frac{\mathbb{P}(D = d, X_2 = x_2)}{\mathbb{P}(X_2 = x_2)},\end{aligned}$$

assuming X_1 is independent of (D, X_2) .

Smokers



Non-smokers



— Best-estimate price (male) — Best-estimate price (female) — Unawareness price — Discrimination-free price

Some reflections:

- smokers are charged a higher price than non-smokers, due to a higher expected number of type 2 claims
- for smokers: unawareness price $>$ discrimination-free price, and vice versa for non-smokers
- if a policyholder indicates to be a smoker, the unawareness price implicitly incorporates the higher price for women (on type 1 and type 2 claims) because - in this portfolio - smokers are more likely women
- explore the impact on the unawareness prices when - say - $\mathbb{P}(D = \text{female} \mid X_2 = \text{smoker}) = 0.2$ instead of the assumed 0.8.

We obtain the following Gender-Smoking habits probabilities

$$\begin{aligned}\mathbb{P}(D = \text{female} \cap X_2 = \text{smoker}) &= \mathbb{P}(D = \text{female} \mid X_2 = \text{smoker}) \cdot \mathbb{P}(X_2 = \text{smoker}) \\ &= 0.8 \cdot 0.3 = 0.24\end{aligned}$$

$$\mathbb{P}(D = \text{female} \cap X_2 = \text{non-smoker}) = 0.21$$

$$\mathbb{P}(D = \text{male} \cap X_2 = \text{smoker}) = 0.06$$

$$\mathbb{P}(D = \text{male} \cap X_2 = \text{non-smoker}) = 0.49,$$

using Bayes rule.

We generate 100 000 policyholder records, and distribute their Gender-Smoking habits risk profiles along this probability distribution.

Age of the policyholder is simulated from a specified probability distribution.

Number of claims per type is sampled from a POI distribution with mean λ_j (for $j = 1, 2, 3$).

Excerpt from this data set:

#	Gender	Smoker	Age	Type_1	Type_2	Type_3
1	0	0	61	0	0	0
2	1	0	46	0	0	0
3	0	0	68	0	0	0
4	1	0	46	0	0	0
5	0	1	39	0	1	0
6	0	0	70	0	0	1

For the **best-estimate frequencies**, we calibrate three POI GLMs: (for $j = 1, 2, 3$)

$$\begin{aligned}N_j &\sim \text{POI}(\lambda_j) \\ \lambda_j &= \exp\left(\beta_0^j + \beta_1^j X_1 + \beta_2^j X_2 + \beta_3^j D\right).\end{aligned}$$

The **best-estimate prices** then follow from:

$$\hat{\mu}(\mathbf{X}, D) = c_1 \hat{\lambda}_1(\mathbf{X}, D) + c_2 \hat{\lambda}_2(\mathbf{X}, D) + c_3 \hat{\lambda}_3(\mathbf{X}, D),$$

with the c 's denoting the (fixed, type-specific) claim amounts.

For the **unawareness frequencies**, we drop the D variable from the linear predictor:

$$\begin{aligned}N_j &\sim \text{POI}(\lambda_j) \\ \lambda_j &= \exp\left(\beta_0^j + \beta_1^j X_1 + \beta_2^j X_2\right).\end{aligned}$$

For the **discrimination-free prices**, we average the best-estimate prices over the values of D , using the empirical probabilities:

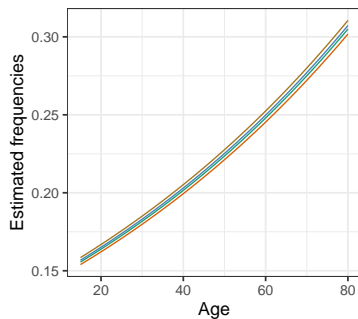
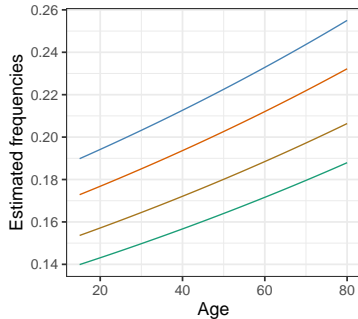
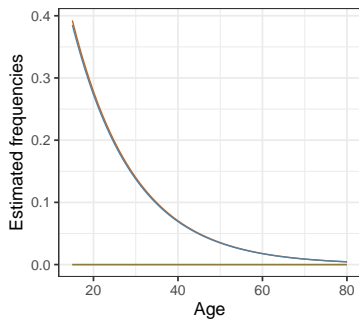
$$\hat{h}^*(\mathbf{x}) = \sum_d \hat{\mu}(\mathbf{x}, d) \frac{n_d}{n},$$

with $n = 100\,000$ the sample size.

Estimated prices

GLMs

53

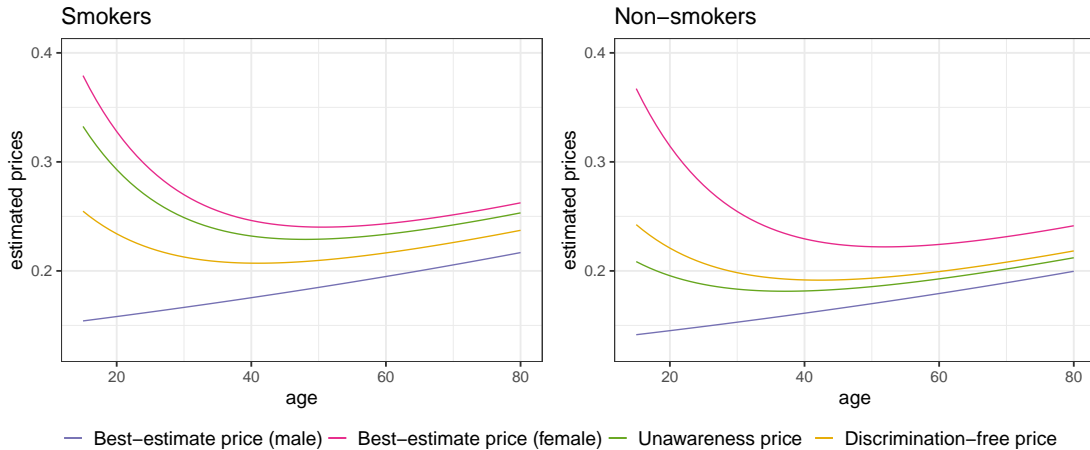


— Male non-smoker — Female non-smoker — Male smoker — Female smoker

Estimated prices

GLMs

54



Lindholm et al. (2022) analyze the simulated data set with neural network regression.

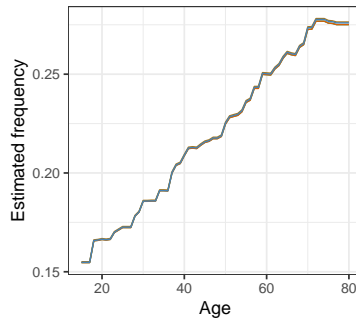
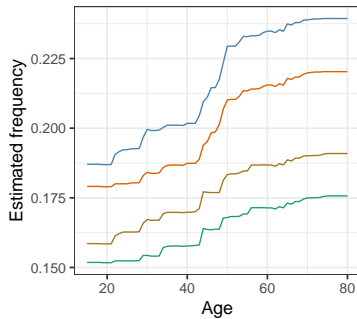
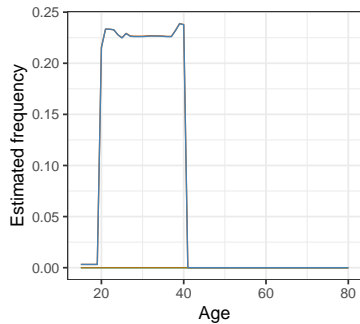
We opt for a **Gradient Boosting Machine (GBM)** (see Henckaerts et al., 2021):

- iteratively combines weak learners into a powerful predictor
- at each iterative step a new tree is fit using information from previously grown trees
- tuning parameters: shrinkage, interaction depth, number of trees, minimum number of observations in a node and the bag fraction.

Estimated prices

GBMs - estimated frequencies

56

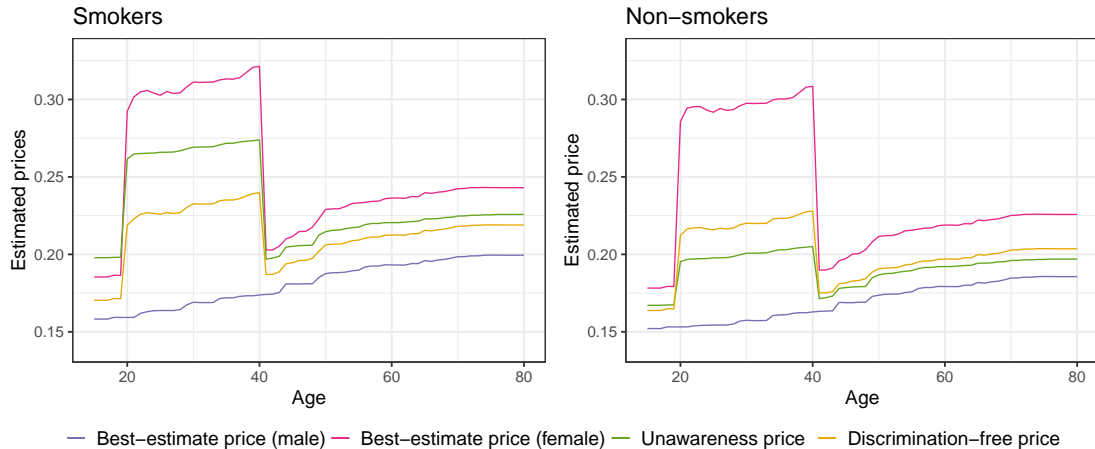


— Male non-smoker — Female non-smoker — Male smoker — Female smoker

Estimated prices

GBMs - estimated prices

57



Comparing the assumed technical prices and the calibrated ones:

- GLM prices are a poor approximation to the true prices, because of difficulties capturing the highly nonlinear birthing-related effects
- GBMs perform better
- unawareness price discriminates indirectly by learning the gender D from smoking habits X_2 .

Throughout the workshop we

- took a quick tour of the literature on bias and fairness in machine learning
- explored a recently proposed strategy to adjust insurance prices for discrimination by proxy, when information on protected features is available.