

IA, biais et équité (en actuariat et en assurance)

Arthur Charpentier

avec Laurence Barry, Marie-Pier Côté, Olivier Côté,
Agathe Fernandes-Machado, Ewen Gallic, François Hu , Philipp Ratz
(et Ana Patrón Piñerez, Mulah Moriah, etc)

Marseille, Février 2025



AI, biases and fairness (in actuarial science and insurance)

Arthur Charpentier

with Laurence Barry, Marie-Pier Côté, Olivier Côté,
Agathe Fernandes-Machado, Ewen Gallic, François Hu, Philipp Ratz
(and Ana Patrón Piñerez, Mulah Moriah, etc)

Marseille, February 2025

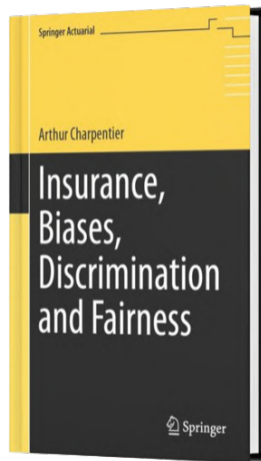


Preamble

- ▶ professor in Montréal, Canada (mathematics)
& Rennes, France (economics)
- ▶ talk based on a recent textbook
**Charpentier (2024) Insurance, Biases,
Discrimination and Fairness. Springer.**

⚠ disclaimer

This is part of ongoing work, and, despite my efforts, it might contain errors of any type. Concepts and results presented in those slides are probably either extremely vague, or wrong. All apologies.



What is a “discrimination”? (direct, or indirect)

- ▶ Discrimination is “*the act of treating different groups differently*,” Frees and Huang (2021)
- ▶ “*direct discrimination is intentional, whereas indirect discrimination is unintentional*,” Campbell and Smith (2023)
- ▶ “*Technology is neither good nor bad; nor is it neutral*,” Kranzberg (1986)

Definition 1: Discrimination, Merriam-Webster (2022)

Discrimination is the act, practice, or an instance of separating or distinguishing categorically rather than individually.

What is an “actuary”?

► “actuarial” ?

“To be an actuary is to be a specialist in generalization, and actuaries engage in a form of decision making that is sometimes called actuarial. Actuaries guide insurance companies in making decisions about large categories that have the effect of attributing to the entire category certain characteristics that are probabilistically indicated by membership in the category, but that still may not be possessed by a particular member of the category,” Schauer (2006).

PROFILES

PROBABILITIES

AND

STEREOTYPES

FREDERICK SCHAUER

The Belknap Press of Harvard University Press
Cambridge, Massachusetts
London, England

generalization is the stock in trade of the insurance industry. Indeed, the insurance industry has its own name for this kind of decisionmaking. To be an *actuary* is to be a specialist in generalization, and actuaries engage in a form of decisionmaking that is sometimes called *actuarial*. Actuaries guide insurance companies in making decisions about large categories (teenage males living in northern New Jersey) that have the effect of attributing to the entire category certain characteristics (carelessness in driving) that are probabilistically indicated by membership in the category, but that still may not be possessed by a particular member of the category (this *particular* teenage male living in northern New Jersey).

Occasionally the actuarial generalizations of the insurance industry become controversial. One example is the use of generalizations about the comparative safety of different neighborhoods as a basis for setting the rates for homeowners' insurance or determining the willing-

What is an “actuarial model” (as in most actuarial textbooks)?

- ▶ linear regression on categories - “**segmentation**”

$$\hat{y}(\text{man}) = \beta_0 + \beta_1 \mathbf{1}_{\text{urban}} + \beta_2 \mathbf{1}_{\text{young}} + \beta_3 \mathbf{1}_{\text{man}} = \hat{y}(\text{woman}) + \beta_3$$

$+ \beta_3$ ceteris paribus

- ▶ Poisson regression (frequency) on categories, or not

$$\hat{y}(\text{man}) = \exp [\beta_0 + \beta_1 \mathbf{1}_{\text{urban}} + \beta_2 \mathbf{1}_{\text{young}} + \beta_3 \mathbf{1}_{\text{man}}] = \hat{y}(\text{woman}) \cdot \exp[\beta_3]$$

$\times e^{\beta_3}$ ceteris paribus

$$\hat{y}(\text{man}) = \exp [\beta_0 + \beta_1 \mathbf{1}_{\text{urban}} + \beta_2 \text{age} + \beta_3 \mathbf{1}_{\text{man}}] = \hat{y}(\text{woman}) \cdot \exp[\beta_3]$$

If β_3 small, $e^{\beta_3} \approx 1 + \beta_3$, i.e. “ $\beta_3 = 0.2$ ” \longleftrightarrow “+20% for men”

Thus “**interpretation**” is simple (if we do not discuss what “ceteris paribus” means).

Why could there be a problem?

- ▶ **Econometrics** is dead, long live “**artificial intelligence**”
- ▶ “**Machine learning**” context, i.e. black boxes, with less intuitive interpretation
- ▶ “**Big data**” context, i.e. easy to get proxies for protected/sensitive variables

y	urban	age	race	y	urban	age	zip	lastname	model	credit
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

It is possible to predict the “**race**” based on non-protected variables, e.g. names and geolocation, see “**Bayesian Improved Surname Geocoding (BISG)**”, Elliott et al. (2009), Imai and Khanna (2016)

Machine learning in one slide

Given a dataset $\{(y_i, \mathbf{x}_i, \mathbf{1}_{\text{man},i})\}$, we want to solve

$$\hat{m} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i, \mathbf{1}_{\text{man},i})) + \lambda \operatorname{penalty}(m) \right\}$$

How does $\mathbf{1}_{\text{man}}$ influence \hat{m} ? \leftarrow “**interpretability**” and “**explainability**”

Classical answer, based on [Friedman \(2001\)](#)'s “**partial dependence plots**”

$$\operatorname{pdp}(\mathbf{1}_{\text{man}}) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\mathbf{x}_i, 1) \quad \text{v.s.} \quad \operatorname{pdp}(\mathbf{1}_{\text{woman}}) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\mathbf{x}_i, 0)$$

Explainability is a challenging problem when variables are correlated...

Why not consider techniques from **causal inference** ?

Where could there be a problem?

Ratemaking is an issue, but also **underwriting**,

“**Redlining**”, for loans, but also insurance, **Kerner (1968)**

*“use of a red line around the questionable areas on territorial maps centrally located in the Underwriting Division for ease of reference by all Underwriting personnel [...] mark off certain areas * * * to denote a lack of interest in business arising in these areas In New York these are called K.O. areas meaning knock-out areas; in Boston they are called redline districts. Same thing – don’t write the business.”*

to requests for information reveal clearly that business in certain geographic territories is restricted. For example, one underwriting guide states:

“An underwriter should be aware of the following situations in his territory:

1. The blighted areas.
2. The redevelopment operations.
3. Peculiar weather conditions which might make for a concentration of windstorm or hail losses.
4. The economic makeup of the area.
5. The nature of the industries in the area, etc.

“This knowledge can be gathered by drives through the area, by talking to and visiting agents, and by following local newspapers as to incidents of crimes and fires. A good way to keep this information available and up to date is by *the use of a red line* around the questionable areas on territorial maps centrally located in the Underwriting Division for ease of reference by all Underwriting personnel.” (Italics added.)

A New York City insurance agent at our hearings put it more pointedly:

“[M]ost companies mark off certain areas * * * to denote a lack of interest in business arising in these areas In New York these are called K.O. areas—meaning knock-out areas; in Boston they are called redline districts. Same thing—don’t write the business.”

What is a “actuarial fairness”?

► “Actuarial fairness” ?

... *“on an actuarially fair basis; that is, if the costs of medical care are a random variable with mean m , the company will charge a premium m , and agree to indemnify the individual for all medical costs,”* Arrow (1963).

“**actuarially fair premiums**” = “**expected losses**”
of the insured risk, see also Frezal and Barry (2020).

“governments must recognise that there is a difference between unfair discrimination and insurers differentiating prices according to risk,”
Swiss Re (2015), cited in Meyers and Van Hoyweghen (2018)

THE AMERICAN ECONOMIC REVIEW

VOLUME LIII

DECEMBER 1963

NUMBER 5

UNCERTAINTY AND THE WELFARE ECONOMICS OF MEDICAL CARE

By KENNETH J. ARROW*

the latter. Suppose, therefore, an agency, a large insurance company plan, or the government, stands ready to offer insurance against medical costs on an actuarially fair basis; that is, if the costs of medical care are a random variable with mean m , the company will charge a premium m , and agree to indemnify the individual for all medical costs. Under these circumstances, the individual will certainly prefer to take out a policy and will have a welfare gain thereby.

Will this be a social gain? Obviously yes, if the insurance agent is suffering no social loss. Under the assumption that medical risks on different individuals are basically independent, the pooling of them reduces the risk involved to the insurer to relatively small proportions.

What is a “actuarial fairness”?

“During the first two decades of the century, many companies transformed these qualitative understandings into formal written categories, albeit incompletely and in an idiosyncratic fashion. Although such systems of classifying risk remained relatively primitive, they nonetheless helped to guide firms in making decisions about risks—especially in setting rates, which were quantitative expressions of insurers’ qualitative understandings of danger. Both rates and categories of risk were in a constant state of flux, as underwriters entered into a dialogue with prospective customers and the landscape,” Tebeau (2003)

Eating Smoke

Fire in Urban America, 1800–1950

Mark Tebeau

CHAPTER TWO

The Business of Safety

The American Fire Insurance Industry, 1800–1850

As they cultivated new business early in the nineteenth century, insurance firms began to focus on developing a better understanding of the problem of fire and on setting guidelines for everyday business practices. During the first three decades of the century, several activities became central to the fire insurance business: surveying a risk, corresponding with field representatives and customers about hazards and rates, and compiling records of surveys and transactions in ledgers, and classifying danger. By the 1810s, companies transformed such informal procedures into formal written guidelines and organizational structures. In particular, the industry diversified its risks, and underwriters established rudimentary distinctions between different sorts of property, manufacturing activities, and construction methods. Initially such divisions resided in the minds of company secretaries—in an expanding qualitative knowledge base about fire danger that they developed from their own experience. During the first two decades of the century, many companies transformed these qualitative understandings into formal written categories, albeit incompletely and in an idiosyncratic fashion. Although such systems of classifying risk remained relatively primitive, they nonetheless helped to guide firms in making decisions about risks—especially in setting rates, which were quantitative expressions of insurers’ qualitative understandings of danger. Both rates and categories of risk were in a constant state of flux, as underwriters entered into a dialogue with prospective customers and the landscape.

The Johns Hopkins University Press
Baltimore

What is a “actuarial fairness”?

"Indeed, the rationale that proscribing the use of certain rating variables is in the public interest because, under imperfect risk assessment systems, actuarial fairness is not achieved for some -- albeit unidentifiable - individuals is fundamentally contradictory. It promotes a remedy for unfairness to some that increases the unfairness overall (by the same actuarial yardstick) and redistributes it."

“Indeed, the rationale that proscribing the use of certain rating variables is in the public interest because, under imperfect risk assessment systems, actuarial fairness is not achieved for some – albeit unidentifiable - individuals is fundamentally contradictory. It promotes a remedy for unfairness to some that increases the unfairness overall (by the same actuarial yardstick) and redistributes it,” Casey et al. (1976), cited in Walters (1981)

So “actuarial fairness” has to do with “accuracy”?

Following [Arrow \(1963\)](#), “**actuarially fair premiums**” = “**expected losses**”

▶ but still, there is no “**law of one price**” in insurance, [Froot et al. \(1995\)](#)

→ with different models and different portfolio, we can have two different premiums

▶ estimating “**expected losses**” means maximizing “**accuracy**”

average losses / empirical losses

$$\bar{y} = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - \gamma)^2 \right\} \text{ or } \mathbb{E}[Y] = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \sum_y (y - \gamma)^2 \mathbb{P}[Y = y] \right\}$$

least squares

i.e. we want to minimize the error between observed losses y and predictions \hat{y} .

with binary observations $y \in \{0, 1\}$, hard to assess if $\hat{y} = 12.2486\%$ is accurate or not...

So “actuarial fairness” has to do with “accuracy”?

“If we are asked to find the probability holding for an individual future event, we must first incorporate the case in a suitable reference class,” Reichenbach (1971)

“When we speak of the ‘probability of death’, the exact meaning of this expression can be defined in the following way only. We must not think of an individual, but of a certain class as a whole, e.g., ‘all insured men forty-one years old living in a given country and not engaged in certain dangerous occupations’. A probability of death is attached to the class of men or to another class that can be defined in a similar way. The phrase ‘probability of death’, when it refers to a single person, has no meaning for us at all,” von Mises (1928, 1939)

THE THEORY OF PROBABILITY

*An Inquiry into the Logical and Mathematical
Foundations of the Calculus of Probability*

By HANS REICHENBACH

PROFESSOR OF PHILOSOPHY IN THE UNIVERSITY OF CALIFORNIA AT LOS ANGELES

UNIVERSITY OF CALIFORNIA PRESS
BERKELEY AND LOS ANGELES • 1949

§ 71. Attempts at a Single-Case Interpretation of Probability

After the discussion of the frequency meaning of probability, the investigation must turn to linguistic forms in which the concept of probability refers to an individual event. It is on this ground that the frequency interpretation has been questioned. Some logicians have argued that such usage is based on a different concept of probability, which is not reducible to frequencies. Is the existence of two disparate concepts of probability an inescapable consequence of the usage of language?

The first interpretation of the probability of single events is the *degree of expectation* with which an event is anticipated. The feeling of expectation certainly represents a psychological factor the existence of which is indispensable; it even shows degrees of intensity corresponding to the degrees of probability. Difficulty, however, arises from the fact that the degree of expectation varies from person to person and depends on more factors than the degree of the probability of the event to which the expectation refers. Apart from the probability of an event, emotional associations will anticipate it with too-certain expectations, whereas pessimistic persons will think of it in terms of too-uncertain expectations.

So “actuarial fairness” has to do with “accuracy”?

As explained in [Van Calster et al. \(2019\)](#), “*among patients with an estimated risk of 20%, we expect 20 in 100 to have or to develop the event,*”

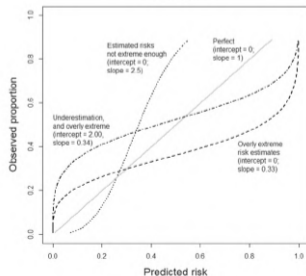
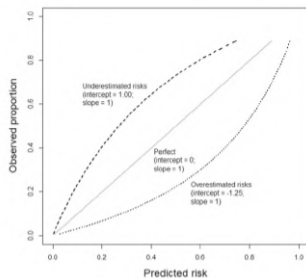
- If 40 out of 100 in this group are found to have the disease, the risk is **underestimated**
- If we observe that in this group, 10 out of 100 have the disease, we have **overestimated** the risk.

The prediction $\hat{m}(\mathbf{X})$ of Y is a well-calibrated prediction if

20 out of 100 (proportion $y = 1$)

$$\mathbb{E}[Y \mid \hat{Y} = \hat{y}] = \hat{y}, \forall \hat{y}$$

↑
estimate risk $\hat{y} = 20\%$



So “actuarial fairness” has to do with “accuracy”?

“Suppose the Met Office says that the probability of rain tomorrow in your region is 80%. They aren’t saying that it will rain in 80% of the land area of your region, and not rain in the other 20%. Nor are they saying it will rain for 80% of the time. What they are saying is there is an 80% chance of rain occurring at any one place in the region, such as in your garden. [...] A forecast of 80% chance of rain in your region should broadly mean that, on about 80% of days when the weather conditions are like tomorrow’s, you will experience rain where you are. [...] If it doesn’t rain in your garden tomorrow, then the 80% forecast wasn’t wrong, because it didn’t say rain was certain. But if you look at a long run of days, on which the Met Office said the probability of rain was 80%, you’d expect it to have rained on about 80% of them.” **McConway (2021)**



The nature of probability

Kevin McConway, Emeritus Professor of Applied Statistics at The Open University, helps to explain the nature of probability and how weather forecasting and horse racing are unlikely partners when it comes to beating the odds.

As one of the top five performing weather forecasting centres in the world, Met Office forecasts are highly valued. Continuing improvements in accuracy with, for example, four day forecasts today being as accurate as a one day forecast back in the 1960s, enable the public and society to take a wider range of weather related decisions with more confidence. The chaotic nature of weather does mean that there are unavoidable limitations to what we can predict, however, by calculating the confidence in a weather forecast we aim to give people a clear picture of any uncertainties.

Beating the odds

Weather forecasting and horse racing have more in common than you might think. Both involve predicting uncertain events. Will it rain on my wedding tomorrow? Will the horse win the next race? And there can be consequences of getting the prediction wrong – soaked guests, or lost money in bets. Nobody expects a racing tipster to make perfect predictions of all the winners – there’s too much uncertainty. Weather, with its chaotic nature and many variables, is undoubtedly even more complex, and that adds to the potential uncertainty. Many people are familiar with expressing the uncertainty in the outcome of a horse race in terms of odds, and we can do something very similar with weather forecasts using probability, which expresses the chance of particular weather occurring.

Probability is a way of expressing the uncertainty of an event in terms of a number on a scale. One very common way of doing this is on a scale going from 0% to 100%, where impossible events are given a probability of 0% and events that will certainly happen are given a probability of 100%.

Other events, that might or might not happen, are given intermediate values on the scale. So an event that is unlikely to happen is not given a probability halfway along the scale at 50%, an event that is pretty likely to happen, but could possibly not happen, might have a probability of 95%.

THE SCALE OF PROBABILITY

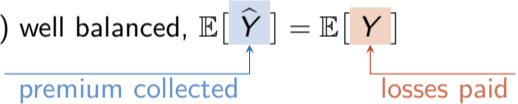


The long-run meaning of probability is all very well, but it doesn’t make so much sense in contexts where things cannot be repeated exactly. In horse racing, you can’t imagine the same horse running exactly the same race again and again and counting up how often it wins. And when the Met Office gives a probability of rain for your region tomorrow, they aren’t really talking about long-run exact repetitions of tomorrow, tomorrow’s only going to happen once.

So “actuarial fairness” has to do with “accuracy”?

This concept goes beyond the simple issue of personalization (discussed in [Barry and Charpentier \(2020\)](#))

There are usually classical assumptions for “model” \hat{y} ,

- ▶ (globally) well balanced, $\mathbb{E}[\hat{Y}] = \mathbb{E}[Y]$

- ▶ (locally) well balanced, $\mathbb{E}[\hat{Y} | \hat{Y} = \hat{y}] = \mathbb{E}[Y | \hat{Y} = \hat{y}] = \hat{y}, \forall \hat{y}$ (“calibration”)

Discrimination? Individual vs. Group Treatment

“**Discrimination** is the act, practice, or an instance of separating or distinguishing categorically rather than individually,” Merriam-Webster (2022).

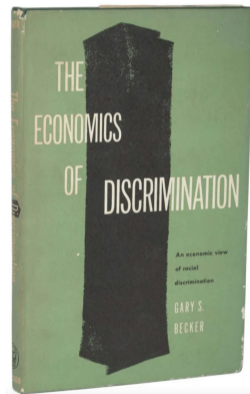
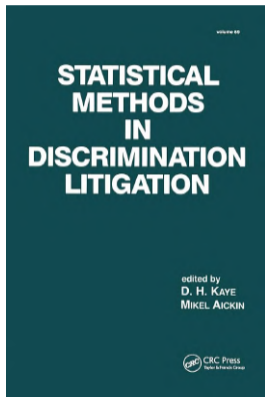
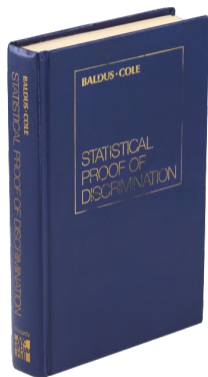
- ▶ “**Ten Oever**” judgement (*Gerardus Cornelis Ten Oever v Stichting Bedrijfspensioenfonds voor het Glazenwassers – en Schoonmaakbedrijf*, in April 1993), the Advocate General **Van Gerven (1993)** argued that “*the fact that women generally live longer than men has no significance at all for the life expectancy of a specific individual and it is not acceptable for an individual to be penalized on account of assumptions which are not certain to be true in his specific case,*” as mentioned in **De Baere and Goessens (2011)**.
 - ▶ **Schanze (2013)** used the term “**injustice by generalization,**” from **Britz (2008)** (“**Generalisierungsunrecht**”)
- Actuarial pricing is essentially discriminatory... and unfair.

“At the core of insurance business lies discrimination”.

- ▶ *”What is unique about insurance is that **even statistical discrimination** which by definition is absent of any malicious intentions, poses significant moral and legal challenges. Why? Because on the one hand, policy makers would like insurers to treat their insureds equally, without discriminating based on race, gender, age, or other characteristics, even if it makes statistical sense to discriminate (...) On the other hand, **at the core of insurance business lies discrimination** between risky and non-risky insureds. But riskiness often statistically correlates with the same characteristics policy makers would like to prohibit insurers from taking into account. ”*
Avraham (2017)
- ▶ *“Machine learning won’t give you anything like gender neutrality ‘for free’ that you didn’t explicitly ask for,”* Kearns and Roth (2019)

Quantifying discrimination, isn't it an old problem?

See [Becker \(1957\)](#) or [Baldus and Cole \(1980\)](#), among (many) others.



Several papers over the past 15 years revisited several notions and concepts.

Is there a (simple) way to quantify unfairness ?

- ▶ classical fairness concept are related to so called “**group fairness**”, where we have a statistical (overall perspective),
- ▶ in some problems, we focus on discrimination in “continuous outcomes”,
 - ▶ $\hat{m}(\mathbf{x}_i, s_i) \in [0, 1]$ (score) that could also be denoted \hat{y}_i
 - ▶ $\hat{m}(\mathbf{x}_i, s_i) \in \mathbb{R}_+$ (premium) that could also be denoted \hat{y}_i
 - classical in insurance modeling
- ▶ in some problems, we focus on discrimination in binary decisions $\hat{y}_i \in \{0, 1\}$, usually obtained as
 - ▶ $\hat{y}_i = \mathbf{1}(\hat{m}(\mathbf{x}_i, s_i) > \text{threshold}) \in \{0, 1\}$ (class) that could also be denoted
 - classical in computer science

Several definitions of “fairness” or “non-discriminatory”

demographic parity $\rightarrow \mathbb{E}[\hat{Y} | S = A] \stackrel{?}{=} \mathbb{E}[\hat{Y} | S = B]$

sensitive (green arrow pointing to $S = A$)

sensitive (yellow arrow pointing to $S = B$)

score \hat{y} (blue arrow pointing from $S = A$ to $S = B$)

equalized odds $\rightarrow \mathbb{E}[\hat{Y} | Y = y, S = A] \stackrel{?}{=} \mathbb{E}[\hat{Y} | Y = y, S = B], \forall y$

outcome y (orange arrow pointing from $Y = y$ in the first term to $Y = y$ in the second term)

score \hat{y} (blue arrow pointing from $S = A$ to $S = B$)

calibration $\rightarrow \mathbb{E}[Y | \hat{Y} = u, S = A] \stackrel{?}{=} \mathbb{E}[Y | \hat{Y} = u, S = B], \forall u$

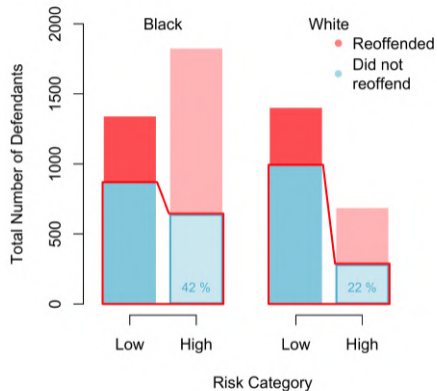
outcome y (orange arrow pointing from $\hat{Y} = u$ in the first term to $\hat{Y} = u$ in the second term)

score \hat{y} (blue arrow pointing from $S = A$ to $S = B$)

Isn't it a problem to have several definitions?

From Feller et al. (2016),

- for White people, among those who did not re-offend (y), 22% were wrongly classified (\hat{y}),
- for Black people, among those who did not re-offend, 42% were wrongly classified,
- **Problem**, since $42\% \gg 22\%$

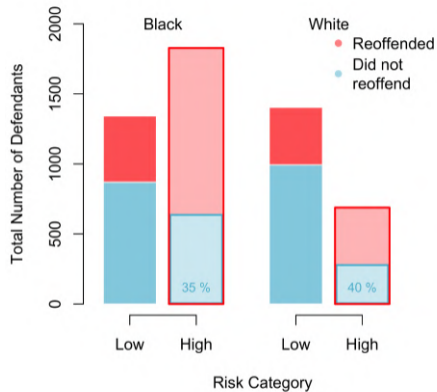


$$\mathbb{P}[\hat{Y} = \text{high} \mid Y = \text{no}, S = \text{black}] = 42\% \stackrel{?}{=} \mathbb{P}[\hat{Y} = \text{high} \mid Y = \text{no}, S = \text{white}] = 22\%$$

Isn't it a problem to have several definitions?

From Dieterich et al. (2016),

- for White people, among those who were classified as high risk (\hat{y}), 40% did not re-offend (y),
- for Black people, among those who were classified as high risk (\hat{y}), 35% did not re-offend (y),
- **No problem**, since $35 \approx 40\%$



$$\mathbb{P}[Y = \text{no} \mid \hat{Y} = \text{high}, S = \text{black}] = 35\% \stackrel{?}{=} \mathbb{P}[Y = \text{no} \mid \hat{Y} = \text{high}, S = \text{white}] = 40\%.$$

Is it always possible to have a sensitive-free model (with respect to ...)?

For **decisions** ($\hat{y} \in \{0, 1\}$, e.g., “obtain a loan”), decision \hat{y}

$$\text{demographic parity} \rightarrow \mathbb{P}[\hat{Y} = 1 \mid S = A] \stackrel{?}{=} \mathbb{P}[\hat{Y} = 1 \mid S = B]$$

those decisions are usually based on **scores**, and **thresholds**

$$\text{demographic parity} \rightarrow \mathbb{E}[\hat{m}(\mathbf{X}, S) > t \mid S = A] \stackrel{?}{=} \mathbb{E}[\hat{m}(\mathbf{X}, S) > t \mid S = B]$$

One can achieve **demographic parity**, simply selecting **different thresholds**

$$\text{demographic parity} \rightarrow \mathbb{E}[\hat{m}(\mathbf{X}, S) > t_A \mid S = A] \stackrel{?}{=} \mathbb{E}[\hat{m}(\mathbf{X}, S) > t_B \mid S = B]$$

(with that strategy, usually impossible to achieve **equalized odds**)

Is it always possible to have a sensitive-free model (with respect to ...)?

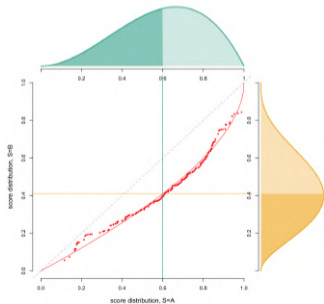
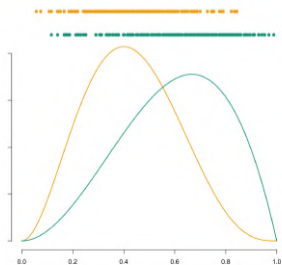
For **decisions** ($\hat{y} \in \{0, 1\}$, e.g., “obtain a loan”), we considered

$$\text{demographic parity} \rightarrow \mathbb{E}[\hat{Y} | S = A] \stackrel{?}{=} \mathbb{E}[\hat{Y} | S = B]$$

and we can consider the analogous for **scores** (possibly used to assess premiums),

$$\text{demographic parity} \rightarrow \mathbb{E}[\hat{m}(\mathbf{X}, S) | S = A] \stackrel{?}{=} \mathbb{E}[\hat{m}(\mathbf{X}, S) | S = B]$$

↑ score \hat{y} ↑



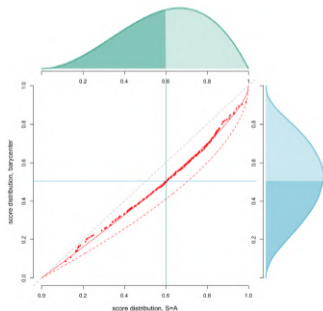
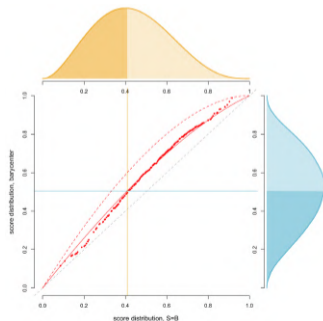
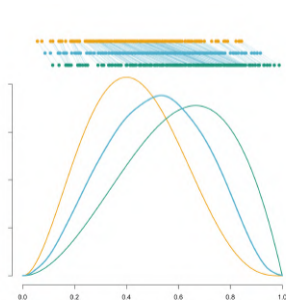
- ▶ individual in group **A** with a score $\hat{y}(A) = 60\%$ corresponding to quantile α (here 0.5)
- ▶ in group **B**, the same quantile α corresponds to $\hat{y}(B) = 40\%$

Is it always possible to have a sensitive-free model (with respect to ...)?

- ▶ To get a fair model (**neutral with respect to s**), consider an average between the two models,

score in group A with quantile α score in group B with quantile α

$$\hat{y}^* = \mathbb{P}[S = A] \cdot \hat{y}(A) + \mathbb{P}[S = B] \cdot \hat{y}(B)$$



“In order to treat some persons equally, we must treat them differently”

- ▶ Supreme Court Justice Harry Blackmun stated, in 1978,

“In order to get beyond racism, we must first take account of race. There is no other way. And in order to treat some persons equally, we must treat them differently,” Knowlton (1978), cited in Lippert-Rasmussen (2020)

- ▶ In 2007, John G. Roberts of the U.S. Supreme Court submits

“The way to stop discrimination on the basis of race is to stop discriminating on the basis of race,” Sabbagh (2007) and Turner (2015)

See philosophical discussions about **affirmative action**, e.g., Rubinfeld (1997); Pojman (1998); Anderson (2004)

“In order to treat some persons equally, we must treat them differently”

- ▶ In 2021, in Europe, first sentence of chapter 7, “*State of Democracy, Human Rights and the Rule of Law,*”

“The strategic goal of the Council of Europe in the field of anti-discrimination, diversity and inclusion is to ensure genuine equality and full access to rights and opportunities for all members of society,”
Council of Europe (2021)

- ▶ In 2025, in the U.S. “*Ending Radical And Wasteful Government DEI Programs And Preferencing,*” Executive order,

“Nearly every Federal agency and entity submitted “Equity Action Plans” to detail the ways that they have furthered DEIs infiltration of the Federal Government. The public release of these plans demonstrated immense public waste and shameful discrimination. That ends today. Americans deserve a government committed to serving every person with equal dignity and respect,” The White House (2025)

“Neutral with respect to some sensitive attribute?”

What does “**neutral with respect to s** ” really means ?

We have seen that accuracy was assessed with respect to data in the portfolio,

$$\bar{y} = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - \gamma)^2 \right\} \text{ or } \mathbb{E}[Y] = \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \sum_y (y - \gamma)^2 \mathbb{P}[Y = y] \right\}$$

based on observations from the insurer’s portfolio. Technically, should we consider

- ▶ expected values / probabilities / independence properties based on \mathbb{P} (portfolio)
- ▶ expected values / probabilities / independence properties based on \mathbb{Q} (market)

(ongoing work *Why portfolio-specific fairness should fail to extend market-wide: Selection bias in insurance* with M.P. Côté & O. Côté)

Should we ask for neutrality “in the portfolio” or for some “targeted population” ?

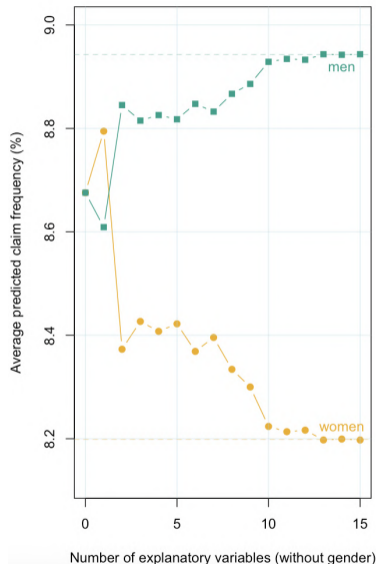
Discrimination in the data, or in the model?

On a French motor dataset, average claim frequencies are **8.94%** (men) and **8.20%** (women).

Consider some logistic regression to estimate annual claim frequency, on k explanatory variables **excluding gender**.

	men	women
$k = 0$	8.68%	8.68%
$k = 2$	8.85%	8.37%
$k = 8$	8.87%	8.33%
$k = 15$	8.94%	8.20%
empirical	8.94%	8.20%

Models simply tend to reproduce what was observed in the data (see “**is-ought**” problem, in [Hume \(1739\)](#)).



Discrimination in the data, or in the model?

David Hume's "**is-ought**" problem, in [Hume \(1739\)](#)



what **is** observed, what is **statistically normal**

$\pi(\mathbf{x}) = \mathbb{E}_{\mathbb{P}}[Y|\mathbf{X} = \mathbf{x}]$ where \mathbb{P} is the historical probability

\neq what **should be**, what we expect from an **ethical norm**

$\pi(\mathbf{x}) = \mathbb{E}_{\mathbb{P}^*}[Y|\mathbf{X} = \mathbf{x}]$ where \mathbb{P}^* is some "fair" probability

"keep in mind that machine learning can only be used to memorize patterns that are present in your training data. You can only recognize what you've seen before. Using machine learning trained on past data to predict the future is making the assumption that the future will behave like the past," [Chollet \(2021\)](#)

Classical **clausula rebus sic stantibus** ("with things thus standing") in predictive modeling (statistics and machine learning)

Discrimination in the data, or in the model?

- ▶ change the training data to de-bias (through weights) : **pre-processing**
if we can draw i.i.d. copies of a random variable X_i 's, under probability \mathbb{P} , then

$$\frac{1}{n} \sum_{i=1}^n h(x_i) \rightarrow \mathbb{E}_{\mathbb{P}}[h(X)], \text{ as } n \rightarrow \infty \text{ "law of large numbers"}$$

but if we want to reach $\mathbb{E}_{\mathbb{Q}}[h(X)]$, consider

$$\frac{1}{n} \sum_{i=1}^n \underbrace{\frac{d\mathbb{Q}(x_i)}{d\mathbb{P}(x_i)}}_{\text{weight } \omega_j} h(x_i) \rightarrow \mathbb{E}_{\mathbb{Q}}[h(X)], \text{ as } n \rightarrow \infty.$$

- ▶ keep the biases data, but distort the outcome : **post-processing**
- ▶ add a fairness constraint (penalty) in the optimization problem : **in-processing**
as classical adversarial techniques, [Grari et al. \(2021\)](#)

Discrimination, with different perspectives

- ▶ Regulatory perspective, “**group fairness**” (discussed previously)
- ▶ Policyholders perspective, “**individual fairness**”

A decision satisfies individual fairness if “*had the protected attributes (e.g., race) of the individual been different, other things being equal, the decision would have remained the same.*”

- ▶ also named “**counterfactual fairness**” in [Kusner et al. \(2017\)](#), and should be related to classical causal inference problem, (conditional) average treatment effect (the “treatment” being the sensitive attribute),

“*other things being equal*” ? **ceteris paribus** ? See “revolving variable” in [Kilbertus et al. \(2017\)](#). Consider a man ($s = A$) with height $x = 6'3$ (or 190 cm). If that person had been a woman ($s = B$) would she have height $x = 6'3$?

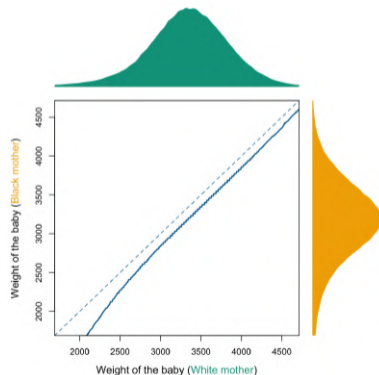
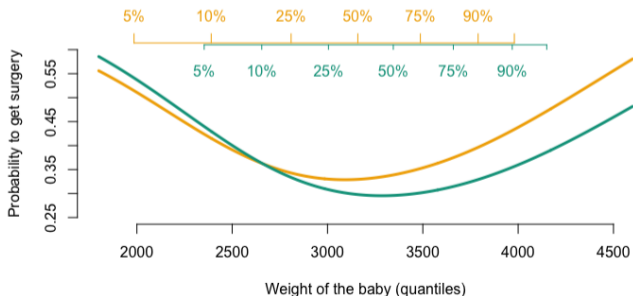
(hint: no, consider similar quantiles, as discussed previously, see [Charpentier et al. \(2023a\)](#))

Discrimination, with different perspectives

$$\mathbb{P}[Y = \text{yes} \mid S = \text{white}, \mathbf{X} = \mathbf{x}] \stackrel{?}{=} \mathbb{P}[Y = \text{yes} \mid S = \text{black}, \mathbf{X} = \mathbf{x}], \forall \mathbf{x}$$

Diagram illustrating the equation above. The word "sensitive" is written in green above "white" and "black". A red bracket labeled "having surgery" is positioned below the "Y = yes" parts of both sides of the equation. A yellow arrow labeled "sensitive" points from the word "sensitive" to the "black" variable.

e.g., x is the “weight of the baby” (at birth), in kg.



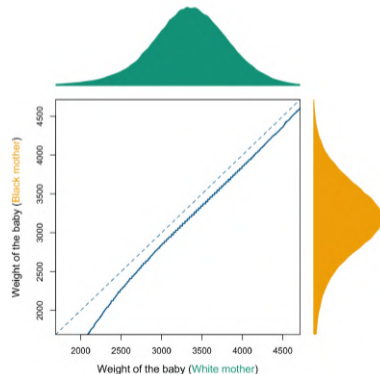
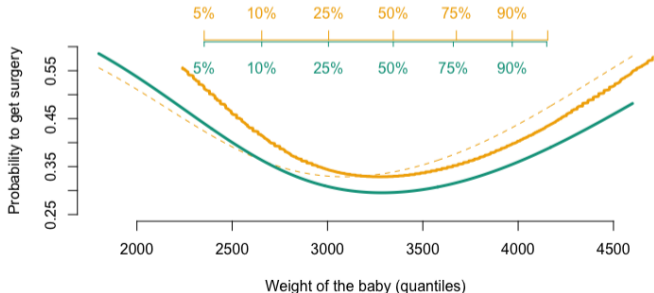
Discrimination, with different perspectives

$$\mathbb{P}[Y = \text{yes} \mid S = \text{white}, X = x] \stackrel{?}{=} \mathbb{P}[Y = \text{yes} \mid S = \text{black}, X = x^*], \forall x$$

sensitive (pointing to 'white') sensitive (pointing to 'black')

having surgery (with arrows pointing to 'yes' in both sides)

e.g., x is the “weight of the baby” (at birth), in kg.



What if we neither observe nor collect sensitive personal information (s) ?

September 27, 2023, the Colorado Division of Insurance exposed a new proposed regulation entitled **Concerning Quantitative Testing of External Consumer Data and Information Sources, Algorithms, and Predictive Models Used for Life Insurance Underwriting for Unfairly Discriminatory Outcomes**. Use of **BIFSG** (Bayesian Improved First Name Surname and Geocoding), from **Elliott et al. (2009)**. Consider 12 people living near Atlanta, GA (Fulton & Gwinnett counties),

	last	first	county	city	zipcode	whi	bla	his	asi
1									
2	RADLEY	OLIVIA	Fulton	Fairburn	30213	14	83	1	0
3	BOORSE	KEISHA	Fulton	Atlanta	30331	97	0	3	0
4	MAZ	SAVANNAH	Gwinnett	Norcross	30093	5	6	76	13
5	GAULE	NATASHIA	Gwinnett	Snellville	30078	67	19	14	0
6	MCMELLEN	ISMAEL	Gwinnett	Lilburn	30047	73	15	6	3
7	WASHINGTON	BRYN	Gwinnett	Norcross	30093	0	95	3	0

(ongoing *Predicting Unobserved Multi-Class sensitive Attributes : Enhancing Calibration with Nested Dichotomies for Fairness* with A.M. Patrón Piñerez, A. Fernandes Machado, & E. Gallic)

Can we use aggregate data related to sensitive information (\bar{y})?

Sex Bias in Graduate Admissions: Data from Berkeley

Measuring bias is harder than is usually assumed, and the evidence is sometimes contrary to expectation.

F. J. Becke, E. A. Hannel, J. W. O'Connell

Examining whether discrimination because of sex or ethnic identity is being practiced against persons seeking passage from one social status or locus to another is an intractable problem in our society today. It is legally important and morally imperative. It is also often quite difficult. This article is an exploration of some of the issues of measurement and assessment involved in one example of the general problem, by means of which we hope to shed some light on the difficulties. We will proceed in a straightforward and logical way, even though we have been misled by an unorthodox, case-study approach to the problem that we do not think is quite likely that other persons interested in questions of bias might proceed in just the same way, and several examples of the mistakes in our own procedure may be instructive.

Data and Assumptions

The particular body of data chosen for examination here consists of applications for admission to graduate study at the University of California, Berkeley, for the fall 1973 quarter. In the admission cycle for that quarter, the Graduate Division at Berkeley received approximately 15,000 applications, most of which were taken with care or transferred to a different program or department. The applications were made for fall 1973 (see Table 1) and are sufficiently complete to permit a

decision to admit or to deny admission. The question we wish to pursue is whether or not the decision to admit or to deny was influenced by the sex of the applicant. We cannot begin with any certainty the influence on the admissions in the Graduate Admissions Office, or on the faculty reviewing committees, or on any other administrative personnel participating in the chain of actions that led to a decision on an individual application. We can, however, say that if the admissions decision and the sex of the applicant are statistically associated and the results of a series of applications, we may judge that they have occurred, and we may then seek to find whether discrimination existed. By "we" we mean here a pattern of covariation between a particular decision and a particular sex of applicant, of the relation in our own procedure may be instructive.

The simplest approach (which we shall call approach A) is to examine the aggregate data for the sexes. This approach would surely be the most direct, for it would require no assumptions about the relative importance of the sex of the applicant in the admissions process. It would require only the sex of the applicant and the number of applications received for each sex. This approach would be the most direct, for it would require no assumptions about the relative importance of the sex of the applicant in the admissions process. It would require only the sex of the applicant and the number of applications received for each sex. This approach would be the most direct, for it would require no assumptions about the relative importance of the sex of the applicant in the admissions process. It would require only the sex of the applicant and the number of applications received for each sex.

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

that bias existed in the fall 1973 admissions. On that occasion, we are sure of the fact for the responsible parties to see whether they give evidence of discrimination. Now, the question is whether an application for admission to graduate study is influenced, in a statistically significant way, by the faculty of the department to which the prospective student applies. Let us now examine such data of the department for indications of bias. Among the 101 departments we find 16 that either had no women applicants or that they admitted to no applicants of either sex. For other departments, there are cases where either sex was admitted, but based on the remaining 85. For a more full and detailed study of the 85 with bias sufficiently large to occur by chance (see the next section) is a hard task. There are 85 departments. The deficit in the number of women admitted to these four (over the assumption for calculating expected frequency) is given above) is 26. Looking further, we find six departments bias in the opposite direction, at the same probability level; these account for a deficit of six.

We present this investigation by computing the expected frequencies of male and female applicants admitted and rejected. Table 1, on the assumption that men and women applicants have equal chances of admission to the program (1 and 2). This computation also gives in Table 1 the number of men and women who were admitted to the 85 departments. This is a large number, and it is unlikely that so large a number of male applicants would occur by chance alone. The chi-square value for this table is 1018, and the probability of observing such a large or larger value is approximately 10^{-20} . We should on this evidence judge that

from Becke et al. (1975), discussed as an illustration of

"Simpson's paradox"

from Becke et al. (1975), discussed as an illustration of "Simpson's paradox"

from Becke et al. (1975), discussed as an illustration of "Simpson's paradox"

from Becke et al. (1975), discussed as an illustration of "Simpson's paradox"

from Becke et al. (1975), discussed as an illustration of "Simpson's paradox"

from Becke et al. (1975), discussed as an illustration of "Simpson's paradox"

Table 1. Decision on application to Graduate Division for fall 1973, by sex of applicant, department, and admission status. (From Becke et al., 1975, p. 1018. $\chi^2 = 1018$, $df = 1$, $P < 0.01$.)

Applicants	Outcome				Difference
	Observed	Expected	Observed	Expected	
Men	204	204	271	271	0
Women	436	307	173	273	273

Table 2. Admission data by sex of applicant to two hypothetical departments. For sex, $\chi^2 = 571$, $df = 1$, $P < 0.01$ (repeated).

Applicants	Outcome				Difference	
	Admitted	Denied	Admitted	Denied	Admitted	Denied
Men	100	200	100	200	0	0
Women	100	100	100	100	0	0
Men	100	100	100	100	0	0
Women	100	200	100	100	0	0
Men	200	100	100	200	200	200
Women	200	400	276	273	276	273

of the total population of applicants) we obtain $\bar{y} = 0.6$, while the remaining 68 departments have a corresponding $\bar{y} = 0.9$. The significance of \bar{y} under the hypothesis of an association can be calculated. All three values obtained are highly significant. The effect may be defined by means of an analogy. Picture a cabinet with two different sized bins. A school of fish, all of identical size (assumption 1), swim toward the net and seek to pass. The female fish all try to get through the small net, while the male fish all try to get through the large net. On the other side of the net all the fish are male. Assumption 2 is that the sex of the fish had no relation to the size of the net they tried to get through. It is false. To take another

Table 3. Proportion of applicants that are women (relative proportion of applicants) in the 83 departments. Six of the last indicates relative number of applicants to the department.

Applicants	Observed		Expected		Difference	
	Admitted	Denied	Admitted	Denied	Admitted	Denied
Men	100	200	100	200	0	0
Women	100	100	100	100	0	0
Men	100	100	100	100	0	0
Women	100	200	100	100	0	0
Men	200	100	100	200	200	200
Women	200	400	276	273	276	273

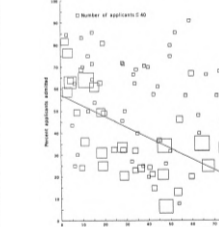


Fig. 1. Proportion of applicants that are women (relative proportion of applicants) in the 83 departments. Six of the last indicates relative number of applicants to the department.

Can we use aggregate data related to sensitive information (\bar{S}) ?

	Total	Men	Women	Proportions
Total	5233/12763 \sim 41%	3714/8442 \sim 44%	1512/4321 \sim 35%	66%-34%
Top 6	1745/4526 \sim 39%	1198/2691 \sim 45%	557/1835 \sim 30%	59%-41%
A	597/933 \sim 64%	512/825 \sim 62%	89/108 \sim 82%	88%-12%
B	369/585 \sim 63%	353/560 \sim 63%	17/ 25 \sim 68%	96%- 4%
C	321/918 \sim 35%	120/325 \sim 37%	202/593 \sim 34%	35%-65%
D	269/792 \sim 34%	138/417 \sim 33%	131/375 \sim 35%	53%-47%
E	146/584 \sim 25%	53/191 \sim 28%	94/393 \sim 24%	33%-67%
F	43/714 \sim 6%	22/373 \sim 6%	24/341 \sim 7%	52%-48%

Data from [Bickel et al. \(1975\)](#). Formalized as follows: S is the (binary) genre, \hat{Y} the admission decision, and X the program (category),

Can we use aggregate data related to sensitive information (\bar{S}) ?

$$\begin{aligned} \mathbb{P}[\hat{Y} = \text{yes} \mid S = \text{men}] &\geq \mathbb{P}[\hat{Y} = \text{yes} \mid S = \text{women}] \\ \mathbb{P}[\hat{Y} = \text{yes} \mid X = x, S = \text{men}] &\leq \mathbb{P}[\hat{Y} = \text{yes} \mid X = x, S = \text{women}], \forall x. \end{aligned}$$

The diagram includes several annotations: a green arrow labeled "sensitive" points to the $S = \text{men}$ term in the first equation; a yellow arrow labeled "sensitive" points to the $S = \text{women}$ term in the first equation; a red bracket labeled "overall admission" spans the inequality in the first equation; and a blue bracket labeled "conditional on program" spans the inequality in the second equation.

“the bias in the aggregated data stems not from any pattern of discrimination on the part of admissions committees, which seems quite fair on the whole, but apparently from prior screening at earlier levels of the educational system. Women are shunted by their socialization and education toward fields of graduate study that are generally more crowded, less productive of completed degrees, and less well funded, and that frequently offer poorer professional employment prospects,” Bickel et al. (1975)

What if we collect s but we miss an important predictor (x) ?

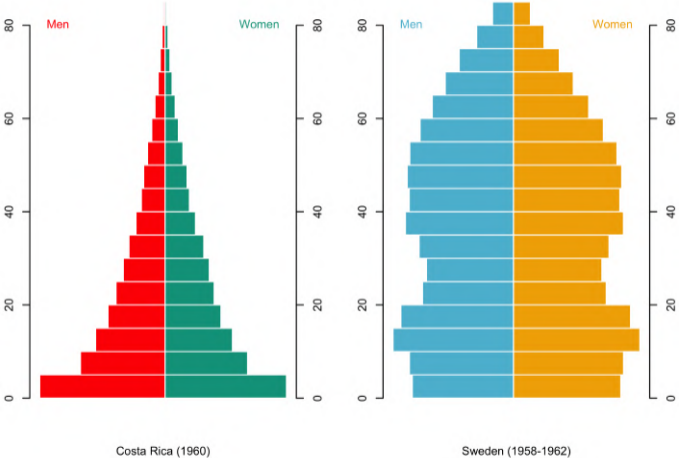
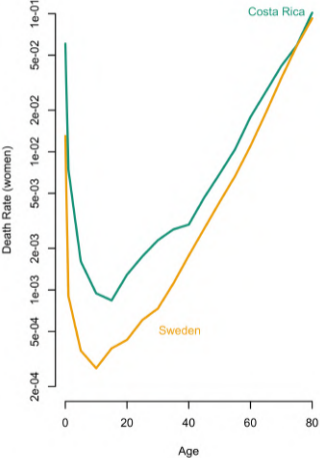
Simpson's paradox can also be seen as an **omitted variable bias** problem,

$$\begin{cases} y_i = \beta_0 + \mathbf{x}_1^\top \beta_1 + \mathbf{x}_2^\top \beta_2 + \varepsilon_i & \text{true model} \\ y_i = b_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \eta_i & \text{estimated models} \end{cases}$$

$$\begin{aligned} \hat{\mathbf{b}}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top [\mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 + \varepsilon] \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_1 \beta_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \beta_2 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon \\ &= \beta_1 + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \beta_2}_{\beta_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon}_{\nu_i}, \end{aligned}$$

so that $\mathbb{E}[\hat{\mathbf{b}}_1] = \beta_1 + \beta_{12} \neq \beta_1$.

What if we collect s but we miss an important predictor (x) ?



Overall mortality rate for women, **8.12‰** in Costa Rica, against **9.29‰** in Sweden.

Disentangling correlations

BBC

Some diverse areas of England face car insurance 'ethnicity penalty'

By Maryam Ahmed
BBC Verify

Quote A



Teacher
Aged 30
Male

Car: Ford Fiesta

Address: Princes End area of Sandwell, near Birmingham

Black, Asian & minority ethnic population: 11%

Average quote: £1,975

Quote B



Teacher
Aged 30
Male

Car: Ford Fiesta

Address: Great Bridge area of Sandwell, near Birmingham

Black, Asian & minority ethnic population: 44%

Average quote: £2,796

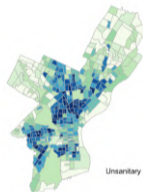


BBC

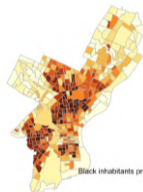
See some diverse areas of England face car insurance 'ethnicity penalty' (remove from the BBC website since)



Red area (too risky)



Unsanitary index (0-100)



Black inhabitants proportion (%)

y , x and s can easily be correlated variables

spurious correlations problem ?

Need to use causal models to avoid indirect discrimination

Multiple sensitive attributes, “robbing Peter to pay Paul”?

$$\mathbb{E}[\hat{m}(\mathbf{X}, S_1, S_2) \mid S_1 = A] \neq \mathbb{E}[\hat{m}(\mathbf{X}, S_1, S_2) \mid S_1 = B]$$

sensitive attribute 1

$$\mathbb{E}[\hat{m}(\mathbf{X}, S_1, S_2) \mid S_2 = C] \approx \mathbb{E}[\hat{m}(\mathbf{X}, S_1, S_2) \mid S_2 = D]$$

sensitive attribute 2

Distort model \hat{m} to achieve fairness with respect to $S_1 \rightarrow$ model \tilde{m}

$$\mathbb{E}[\tilde{m}(\mathbf{X}, S_1, S_2) \mid S_1 = A] = \mathbb{E}[\tilde{m}(\mathbf{X}, S_1, S_2) \mid S_1 = B]$$

sensitive attribute 1

$$\mathbb{E}[\tilde{m}(\mathbf{X}, S_1, S_2) \mid S_2 = C] \neq \mathbb{E}[\tilde{m}(\mathbf{X}, S_1, S_2) \mid S_2 = D]$$

sensitive attribute 2

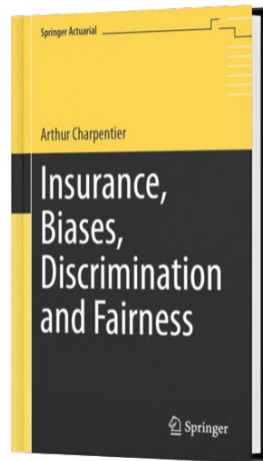
“The myth of the actuary” (objectivity vs. subjectivity)

- ▶ The rhetoric of insurance exclusion – numbers, objectivity and statistics – forms what Brian Glenn calls “*the myth of the actuary*,” “*a powerful rhetorical situation in which decisions appear to be based on objectively determined criteria when they are also largely based on subjective ones*” or “**the subjective nature of a seemingly objective process.**” “*Virtually every aspect of the insurance industry is predicated on stories first and then numbers,*” Glenn (2000, 2003)
- ▶ Importance of **interpretation** and **explainability** of models

Conclusion (?)

- ▶ dealing with discrimination in insurance is tricky since actuarial pricing is deeply related to the idea of focusing on groups, and not individuals
- ▶ if we do not address properly those questions, there is no way we can get fair models
- ▶ not collecting and not using protected attributes is clearly not a good strategy
- ▶ there are still important questions that should be addressed by regulators, that should provide guidelines

To go further, **Charpentier (2024) Insurance, Biases, Discrimination and Fairness. Springer.**





Laurence
Barry



Marie-Pier
Côté



Olivier
Côté



Agathe
Fernandes



Ewen
Gallic



François
Hu



Philipp
Ratz



Ana
Patrón



Mulah
Moriah



Arthur
Charpentier

✉ charpentier.arthur@uqam.ca

✉ francois.hu@milliman.com

References

- Anderson, T. H. (2004). *The pursuit of fairness: A history of affirmative action*. Oxford University Press.
- Arrow, K. J. (1963). Uncertainty and the welfare economics of medical care. *The American Economic Review*, 53(5):941–973.
- Avraham, R. (2017). Discrimination and insurance. In Lippert-Rasmussen, K., editor, *Handbook of the Ethics of Discrimination*, pages 335–347. Routledge.
- Baldus, D. C. and Cole, J. W. (1980). *Statistical proof of discrimination*. McGraw-Hill.
- Barry, L. and Charpentier, A. (2020). Personalization as a promise: Can big data change the practice of insurance? *Big Data & Society*, 7(1):2053951720935143.
- Becker, G. S. (1957). *The economics of discrimination*. University of Chicago press.
- Bickel, P. J., Hammel, E. A., and O'Connell, J. W. (1975). Sex bias in graduate admissions: Data from Berkeley. *Science*, 187(4175):398–404.
- Britz, G. (2008). *Einzelfallgerechtigkeit versus Generalisierung: verfassungsrechtliche Grenzen statistischer Diskriminierung*. Mohr Siebeck.
- Campbell, C. and Smith, D. (2023). Distinguishing between direct and indirect discrimination. *The Modern Law Review*, 86(2):307–330.

References

- Casey, B., Pezier, J., and Spetzler, C. (1976). *The Role of Risk Classification in Property and Casualty Insurance: A Study of the Risk Assessment Process : Final Report*. Stanford Research Institute.
- Charpentier, A. (2024). *Insurance: biases, discrimination and fairness*. Springer Verlag.
- Charpentier, A., Flachaire, E., and Gallic, E. (2023a). Causal inference with optimal transport. In Thach, N. N., Kreinovich, V., Ha, D. T., and Trung, N. D., editors, *Optimal Transport Statistics for Economics and Related Topics*. Springer Verlag.
- Charpentier, A., Hu, F., and Ratz, P. (2023b). Mitigating discrimination in insurance with wasserstein barycenters. *BIAS, 3rd Workshop on Bias and Fairness in AI, International Workshop of ECML PKDD*.
- Chollet, F. (2021). *Deep learning with Python*. Simon and Schuster.
- Côté, O., Côté, M.-P., and Charpentier, A. (2024). A fair price to pay: exploiting causal graphs for fairness in insurance. *forthcoming*.
- Council of Europe (2021). *State of Democracy, Human Rights and the Rule of Law: A Democratic Renewal for Europe*. Council of Europe.
- De Baere, G. and Goessens, E. (2011). Gender differentiation in insurance contracts after the judgment in case c-236/09, *Association Belge des Consommateurs Test-Achats asbl v. conseil des ministres*. *Colum. J. Eur. L.*, 18:339.

References

- Dieterich, W., Mendoza, C., and Brennan, T. (2016). Compas risk scales: Demonstrating accuracy equity and predictive parity. *Northpointe Inc*, 7(7.4):1.
- Elliott, M. N., Morrison, P. A., Fremont, A., McCaffrey, D. F., Pantoja, P., and Lurie, N. (2009). Using the census bureau's surname list to improve estimates of race/ethnicity and associated disparities. *Health Services and Outcomes Research Methodology*, 9:69–83.
- Feller, A., Pierson, E., Corbett-Davies, S., and Goel, S. (2016). A computer program used for bail and sentencing decisions was labeled biased against blacks. it's actually not that clear. *The Washington Post*, October 17.
- Fernandes Machado, A., Charpentier, A., Flachaire, E., Gallic, E., and Hu, F. (2024a). From uncertainty to precision: Enhancing binary classifier performance through calibration. *arXiv preprint arXiv:2402.07790*.
- Fernandes Machado, A., Charpentier, A., Flachaire, E., Gallic, E., and Hu, F. (2024b). Post-calibration techniques: Balancing calibration and score distribution alignment. *Thirty-Eighth Annual Conference on Neural Information Processing Systems*.
- Fernandes Machado, A., Charpentier, A., Flachaire, E., Gallic, E., and Hu, F. (2024c). Probabilistic scores of classifiers, calibration is not enough. *arXiv preprint arXiv:2408.03421*.

References

- Frees, E. W. and Huang, F. (2021). The discriminating (pricing) actuary. *North American Actuarial Journal*, pages 1–23.
- Frezal, S. and Barry, L. (2020). Fairness in uncertainty: Some limits and misinterpretations of actuarial fairness. *Journal of Business Ethics*, 167:127–136.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232.
- Froot, K. A., Kim, M., and Rogoff, K. S. (1995). The law of one price over 700 years. *National Bureau of Economic Research Cambridge*.
- Glenn, B. J. (2000). The shifting rhetoric of insurance denial. *Law and Society Review*, pages 779–808.
- Glenn, B. J. (2003). Postmodernism: the basis of insurance. *Risk Management and Insurance Review*, 6(2):131–143.
- Grari, V., Lamprier, S., and Detyniecki, M. (2021). Fairness without the sensitive attribute via causal variational autoencoder.
- Hu, F., Ratz, P., and Charpentier, A. (2023). Fairness in multi-task learning via wasserstein barycenters. *Joint European Conference on Machine Learning and Knowledge Discovery in Databases – ECML PKDD*.

References

- Hu, F., Ratz, P., and Charpentier, A. (2024). A sequentially fair mechanism for multiple sensitive attributes. *Annual AAAI Conference on Artificial Intelligence*.
- Hume, D. (1739). *A Treatise of Human Nature*. Cambridge University Press Archive.
- Imai, K. and Khanna, K. (2016). Improving ecological inference by predicting individual ethnicity from voter registration records. *Political Analysis*, 24(2):263–272.
- Kearns, M. and Roth, A. (2019). *The ethical algorithm: The science of socially aware algorithm design*. Oxford University Press.
- Kerner, O. (1968). *Report of The National Advisory Commission on Civil Disorder*. Bantam Books.
- Kilbertus, N., Rojas Carulla, M., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. (2017). Avoiding discrimination through causal reasoning. *Advances in neural information processing systems*, 30.
- Knowlton, R. E. (1978). Regents of the university of california v. bakke. *Arkansas Law Review*, 32:499.
- Kranzberg, M. (1986). Technology and history: "kranzberg's laws". *Technology and culture*, 27(3):544–560.
- Kusner, M. J., Loftus, J., Russell, C., and Silva, R. (2017). Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pages 4066–4076.

References

- Lippert-Rasmussen, K. (2020). *Making sense of affirmative action*. Oxford University Press.
- Merriam-Webster (2022). *Dictionary*.
- Meyers, G. and Van Hoyweghen, I. (2018). Enacting actuarial fairness in insurance: From fair discrimination to behaviour-based fairness. *Science as Culture*, 27(4):413–438.
- Pojman, L. P. (1998). The case against affirmative action. *International Journal of Applied Philosophy*, 12(1):97–115.
- Reichenbach, H. (1971). *The theory of probability*. University of California Press.
- Rubinfeld, J. (1997). Affirmative action. *Yale Law Journal*, 107:427.
- Sabbagh, D. (2007). *Equality and transparency: A strategic perspective on affirmative action in American law*. Springer.
- Schanze, E. (2013). Injustice by generalization: notes on the Test-Achats decision of the european court of justice. *German Law Journal*, 14(2):423–433.
- Schauer, F. (2006). *Profiles, probabilities, and stereotypes*. Harvard University Press.
- Swiss Re (2015). Life insurance risk selection: Required differentiation or unfair discrimination? *Sigma*.
- Tebeau, M. (2003). *Eating smoke: Fire in urban America, 1800–1950*. John Hopkins University Press.
- The White House (2025). Ending radical and wasteful government dei programs and preferencing.

References

- Turner, R. (2015). The way to stop discrimination on the basis of race. *Stanford Journal of Civil Rights & Civil Liberties*, 11:45.
- Van Calster, B., McLernon, D. J., Van Smeden, M., Wynants, L., and Steyerberg, E. W. (2019). Calibration: the achilles heel of predictive analytics. *BMC medicine*, 17(1):1–7.
- Van Gerven, G. (1993). Case c-109/91, Gerardus Cornelis Ten Oever v. Stichting bedrijfspensioenfonds voor het glazenwassers-en schoonmaakbedrijf. *EUR-Lex*, 61991CC0109.
- von Mises, R. (1928). *Wahrscheinlichkeit Statistik und Wahrheit*. Springer.
- von Mises, R. (1939). *Probability, statistics and truth*. Macmillan.
- Walters, M. A. (1981). Risk classification standards. In *Proceedings of the Casualty Actuarial Society*, volume 68, pages 1–18.