

Algorithmic Bias

from discrimination discovery
to fairness-aware data mining

Sara Hajian

Francesco Bonchi @francescobonchi

Carlos Castillo @chatox

Part I: Introduction and context

Part II: Discrimination discovery

Part III: Fairness-aware data mining

Part IV: Challenges and directions for future research

Discussion and further questions

 **Part I: Introduction and context**

Part II: Discrimination discovery

Part III: Fairness-aware data mining

Part IV: Challenges and directions for future research

Discussion and further questions

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

A dangerous reasoning

To discriminate is to treat someone differently

(Unfair) discrimination is based on group membership, not individual merit

People's decisions include objective and subjective elements

Hence, they can be discriminate

Algorithmic inputs include only objective elements

Hence, they cannot discriminate?

A stable but unequal marriage

National Resident Match Program (NRMP): US program to match interns to hospitals

V1: 1952 → 1997: algorithm favors hospitals and allows strategic manipulation

V2: 1998 → present: new algorithm is resistant to strategic manipulation



On the web: race and gender stereotypes reinforced

- Results for "CEO" in Google Images: 11% female, US 27% female CEOs
 - Also in Google Images, "doctors" are mostly male, "nurses" are mostly female
- Google search results for professional vs. unprofessional hairstyles for work

Image results:
"Unprofessional
hair for work"



Image results:
"Professional
hair for work"

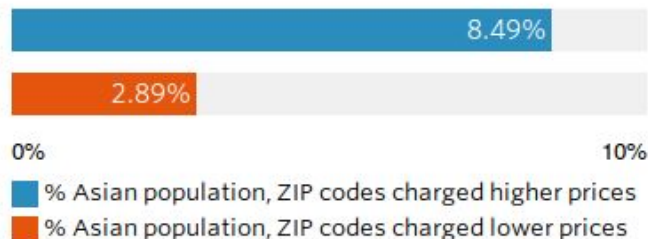
More on web
discrimination
later ...

Geography and race: the "Tiger Mom Tax"

Pricing of SAT tutoring by The Princeton Review in the US doubles for Asians, due to geographical price discrimination

Asians More Likely To Be Among Those Charged Higher Prices By The Princeton Review

Asians make up 4.9 percent of the U.S. population overall. But they account for more than 8 percent of the population in areas where The Princeton Review charges higher prices for its SAT prep packages.



Judiciary use of COMPAS scores



COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) is a 137-questions questionnaire and a predictive model for "risk of recidivism" and "risk of violent recidivism." The model is a proprietary secret of Northpointe, Inc.

Prediction accuracy of recidivism for blacks and whites is about the same (63% and 59%), but errs by being too lenient with whites and too harsh with blacks:

- **Blacks that did not reoffend**
were classified as **high risk** twice as much as **whites that did not reoffend**
- **Whites who did reoffend**
were classified as **low risk** twice as much as **blacks who did reoffend**

Supreme Court of Wisconsin (US) — July 13th, 2016

The court ruled that judges are allowed to use COMPAS scores, but they:

- **must** receive them accompanied by disclaimers and criticisms
- **can** use them as a factor to give non-prison alternatives to prison
- **can** use them as a factor to impose terms and conditions in parole
- **cannot** use them to make sentences longer or shorter



Self-perpetuating algorithmic biases

Credit scoring algorithm suggests Joe has high risk of defaulting

Hence, Joe needs to take a loan at a higher interest rate

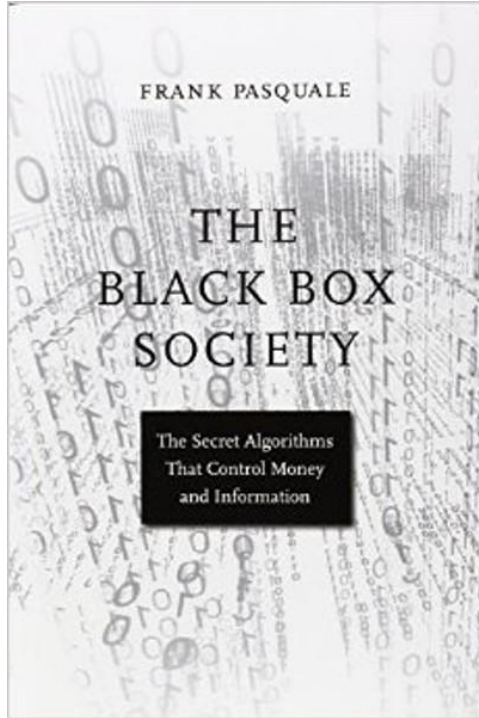
Hence, Joe has to make payments that are more onerous

Hence, Joe's risk of defaulting has increased

The same happens with stop-and-frisk of minorities
further increasing incarceration rates



To make things worse ...



Algorithms are "black boxes" protected by

Industrial secrecy

Legal protections

Intentional obfuscation

Discrimination becomes invisible

Mitigation becomes impossible

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

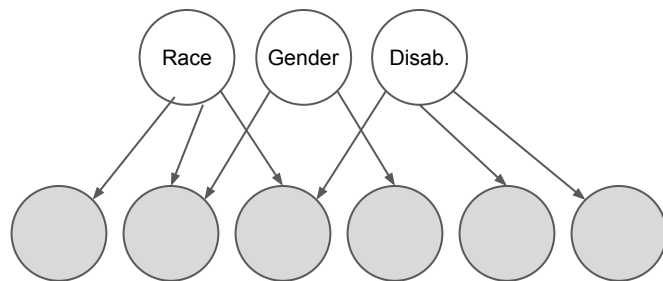
Some sources of algorithmic bias

Data as a social mirror

Protected attributes redundantly encoded in observables

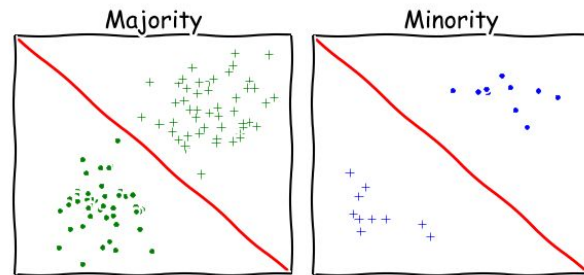
Correctness and completeness

Garbage in, garbage out (GIGO)



Sample size disparity: learn on majority

Errors concentrated in the minority class



Data mining assumptions might not hold

Data mining assumptions are not always observed in reality

- Variables might not be independently identically distributed

- Samples might be biased

- Labels might be incorrect

Errors might be concentrated in a particular class

Sometimes, we might be seeking more simplicity than what is possible

Two areas of concern: data and algorithms

Data inputs:

- Poorly selected (e.g., observe only car trips, not bicycle trips)
- Incomplete, incorrect, or outdated
- Selected with bias (e.g., smartphone users)
- Perpetuating and promoting historical biases (e.g., hiring people that "fit the culture")



Algorithmic processing:

- Poorly designed matching systems
- Personalization and recommendation services that narrow instead of expand user options
- Decision making systems that assume correlation implies causation
- Algorithms that do not compensate for datasets that disproportionately represent populations
- Output models that are hard to understand or explain hinder detection and mitigation of bias

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Data analysis techniques for discrimination discovery: multidisciplinary approaches in the economic, legal, and statistical domains

Parallels with privacy (conceptually)

Legal concepts

Anti-discrimination legislation typically seeks **equal access** to employment, working conditions, education, social protection, goods, and services

Anti-discrimination legislation is very diverse and includes **many legal concepts**

Genuine occupational requirement (male actor to portray male character)

Disparate impact and **disparate treatment**

Burden of proof and **situation testing**

Group under-representation principle

Discrimination: treatment vs impact

Modern legal frameworks offer various levels of **protection** for being discriminated by belonging to a particular class of: gender, age, ethnicity, nationality, disability, religious beliefs, and/or sexual orientation

Disparate **treatment**:

Treatment depends on class membership

Disparate **impact**:

Outcome depends on class membership

Even if (apparently?) people are treated the same way

Disparate treatment: crossing legal red lines

In the US, "steering", i.e. offering more expensive financial services to disadvantaged individuals is forbidden [[Steel and Angwin 2010](#)]

In Europe, insurers cannot discriminate by gender, even if they have accurate, well-grounded, and transparent statistical models indicating costs are higher for women [[European Court of Justice \(2011\)](#)]

Disparate impact

Doctrine solidified in the US after [[Griggs v. Duke Power Co. 1971](#)] where a high school diploma was required for unskilled work, excluding black applicants

In algorithms, disparate impact can be caused by, e.g.:

- (i) how data are collected
- (ii) how data are labeled
- (iii) which features are used
- (iv) usage of proxies for a protected attribute

...

Proving discrimination is hard

In very rare cases there might be a confession that someone discriminated someone else because s/he was women, black, homosexual, etc.

In rare cases the evidence is almost comically overwhelming (e.g., Chinese applicants not interviewed for a position requiring fluency in Mandarin [UK, 1991](#))

In most cases it is hard to prove intentions. The burden of proof can be shared (as in the European Court of Justice) for discrimination cases:

the **accuser** must produce evidence of the **consequences**,

the **defendant** must produce evidence that the **process** was fair

Alternative: situational testing

Do-it-yourself situational testing (We are not encouraging you to try!)

"Gay" vs "Hetero" Job Interview (2015)

Hidden camera video about two job interviews for the same position.

Olle has a truck driver license, experience, and is enthusiastic about the job. "But" he mentions he is saving money to buy an apartment with his boyfriend.

Konrad says he has no experience, spends the whole day playing video games, and shows no enthusiasm for the job.

Who do you think got the job offer?



Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

Principles for quantifying discrimination

In the context of an algorithm generating a prediction:

Predictions for people with similar non-protected attributes should be similar

Differences should be mostly explainable by non-protected attributes

Two basic frameworks for measuring discrimination:

Discrimination at the **individual level**: consistency or individual fairness

Discrimination at the **group level**: statistical parity

Consistency or individual fairness

Consistency score

$$C = 1 - \sum_i \sum_{y_j \in \text{knn}(y_i)} |y_i - y_j|$$

Where $\text{knn}(y_i) = k$ nearest neighbors of y_i

A consistent or individually fair algorithm is one in which similar people experience similar outcomes

Note that similar people to an individual in a protected group may also belong to that protected group, and perhaps they are all treated equally badly

Statistical parity focuses on proportions

Example:

"Protected group" ~ "people with disabilities"

"Benefit granted" ~ "getting a scholarship"

group	benefit		
	denied	granted	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

Intuitively, if

a/n_1 , the fraction of people with disabilities that **does not** get a scholarship is much **larger** than

c/n_1 , the fraction of people without disabilities that **does not** get a scholarship,

then people with disabilities could claim they are being discriminated.

Simple discrimination measures

These measures compare the protected group against the unprotected group:

- Risk difference = $RD = p_1 - p_2$
- Risk ratio or relative risk = $RR = p_1 / p_2$
- Relative chance = $RC = (1-p_1) / (1-p_2)$
- Odds ratio = RR/RC

group	benefit		
	denied	granted	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

$$p_1 = a/n_1 \quad p_2 = c/n_2 \quad p = m_1/n$$

Simple discrimination measures

These measures compare the protected group against the unprotected group:

group	benefit		
	denied	granted	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

- **Risk difference** = $RD = p_1 - p_2$ ← Mentioned in UK law
- **Risk ratio** or relative risk = $RR = p_1 / p_2$ ← Mentioned by EU Court of Justice
- Relative chance = $RC = (1-p_1) / (1-p_2)$
- Odds ratio = RR/RC

US courts focus on selection rates:
 $(1-p_1)$ and $(1-p_2)$

Extended discrimination measures

These measures compare the **protected** group against the **entire population**:

- Extended risk difference = $p_1 - p$
- Extended risk ratio or extended lift = p_1 / p
- Extended chance = $(1-p_1) / (1-p)$

group	benefit		
	denied	granted	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

$$p_1 = a/n_1 \quad p_2 = c/n_2 \quad p = m_1/n$$

More on the *extended lift* later ...

Other measures of discrimination

Differences of mean

Difference of regression coefficients

Rank tests

Mutual information (between outcome and protected attribute)

Unexplained difference (residuals of predictions built with non-protected attributes)

Consistency (comparison of prediction with nearest neighbors)

Example: regression coefficients

$$P = \beta_0 + \beta_1 R + \beta_2 Y + \beta_3 E + \beta_4 W + \varepsilon$$



P: price of a product on eBay, determined by ...

R: buyer and seller have same race. Y: seller has a good rating. E, W: other.

β_1 positive means in-group discrimination, higher for low competition markets

Mean price paid: \$6.37 if seller and buyer are same race, \$5.93 (-7%) otherwise

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

Labor economic perspective

Who can discriminate

employers, other employees, customers

Issues

differences in wage, (un)employment rate, segregation into industries

Typical tool

regression of wage as a function of qualifications or seniority

Racial profiling

Typical tool

difference in proportions

Investigative tools

situational testing; natural experiments
(e.g. observe other motorists in a stop zone to see if police stops blacks more than whites)



Credit and consumer markets

A traditional area for discrimination studies

Example of indirect discrimination is **redlining**: the practice of denying credit to all inhabitants of a certain area (when area is correlated with a protected attribute)

Typical tool for mortgages: comparison of mortgage rejection rates, mortgage price distributions, mortgage default rates

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

Privacy and data protection are related

German *Ausländerzentralregister* (AZR) registers all foreigners living in Germany

Austrian national complained in court that AZR goes against free and non discriminatory movement of European citizens inside Europe

European Court of Justice decision in 2008:

For determining right of residence, AZR is legal

For demographics and statistics, only anonymized AZR data may be used

For criminal investigations, no AZR data may be used

Privacy and data protection legislation differ

Privacy legislation cares about **one action** (storage of personal data)
independently of the consequences

Discrimination legislation cares about **one consequence** (unfair treatment)
independently of the mechanism

A connection between privacy and discrimination

Finding if people having attribute X were discriminated is like inferring attribute X from a database in which:

- the attribute X was removed

- a new attribute (the decision), which is based on X, was added

This is similar to trying to reconstruct a column from a privacy-scrubbed dataset

More on this relationship later ...

Introduction and context

Motivation and examples of algorithmic bias

Sources of algorithmic bias

Legal definitions and principles of discrimination

Measures of discrimination

Discrimination in specific contexts

Discrimination and privacy

Resources

Some resources

- Presentations/keynotes/book
 - Suresh Venkatasubramanian: Keynote at ICWSM 2016
 - Ricardo Baeza: Keynote at WebSci 2016
 - Toon Calders: Keynote at EGC 2016
 - Discrimination and Privacy in the Information Society by Custers et al. 2013
- Groups/workshops/communities
 - Fairness, Accountability, and Transparency in Machine Learning (FATML) workshop and resources
 - Data Transparency Lab - <http://dtlconferences.org/>

Tools

DCUBE (databases) <http://kdd.di.unipi.it/dcube>

S. Ruggieri, D. Pedreschi, F. Turini. Data mining for discrimination discovery. TKDD 4(2), May 2010

S. Ruggieri, D. Pedreschi, F. Turini. DCUBE: Discrimination Discovery in Databases. In SIGMOD 2010 Demos.

Adfisher (web ads)

<https://github.com/tadatitam/info-flow-experiments>

A. Datta, M.C. Tschantz, and A. Datta. Automated experiments on Ad privacy settings. In PETS, pp.92-112, 2015.

The screenshot shows the DCUBE web interface. The top part displays a SQL query for selecting discriminatory rules. The query filters for rules with a selection lift measure greater than 2.5 and a maximum size of B less than or equal to 2, ordered by selection lift in descending order.

```
SELECT d3pddecode(aset) AS aset, -- PD itemset A
       d3pnddecode(bset) AS bset, -- PND itemset B
       'class=bad' AS c, -- negative decision C
       TRUNC(r.ct.slift(),2) AS measure -- selection lift measure
FROM pdrule r JOIN pnditemsets pnd ON r.bset = pnd.id
WHERE r.c = d3encode('class=bad')
AND r.ct.slift() > 2.5 -- minimum measure value
AND pnd.len <= 2 -- maximum size of B
ORDER BY r.ct.slift() DESC; -- descending order
```

Below the query, the 'Query Result' section shows a table with 6 rows of results. The columns are ASET, BSET, C, and M... (Measure). The results show discriminatory rules for the 'class=bad' decision.

	ASET	BSET	C	M...
1	personal_status=female...	purpose=new_car employment=lt_1	class=bad	3,67
2	personal_status=female...	employment=from_1_lt_4 property_magnitude=real_estate	class=bad	2,9
3	personal_status=female...	purpose=furniture_or_equipment employment=from_1_lt_4	class=bad	2,75
4	personal_status=female...	savings_status=no_know	=gt_2d8 class=bad	2,61
5	personal_status=female...	installment_commitmen	class=bad	2,55
6	personal_status=female...	checking_status=lt_0 pr	class=bad	2,54

At the bottom right, there is a 'Snippets' menu with several options for analyzing discriminatory rules, such as 'List a-directly discriminatory rules per context' and 'List a-directly discriminatory rules and decode'.

Part I: Introduction and context

 **Part II: Discrimination discovery**

Part III: Fairness-aware data mining

Part IV: Challenges and directions for future research

Discussion and further questions

Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

The discrimination discovery task at a glance

Given a large database of historical decision records,
find discriminatory situations and practices.

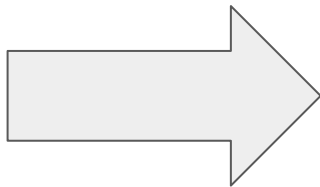
Discrimination discovery scenario

INPUT

Database of historical
decision records

A *criterion* of (unlawful)
discrimination

A set of *potentially*
discriminated groups



OUTPUT

A subset of *decision records*
and *potentially discriminated* people
for which the *criterion* holds

Why is discrimination discovery hard?

Many different concepts regarding discrimination

Including all the ones we mention in Part I of this tutorial

High dimensionality

There are a huge number of possible contexts that may, or may not, be theater for discrimination.

Hidden indirect discrimination

The features that may be the object of discrimination may not be directly recorded in the data

The original data may have been pre-processed due to privacy constraints

The German credit score

A small dataset used in many papers about discrimination
(like Zachary's karate club for networks people)

N = 1,000 records of bank account holders

Class label: good/bad creditor (grant or deny a loan)

Attributes: *numeric/interval-scaled*: duration of loan, amount requested, number of installments, age of requester, existing credits, number of dependents; *nominal*: result of past credits, purpose of credit, personal status, other parties, residence since, property magnitude, housing, job, other payment plans, own telephone, foreign worker; *ordinal*: checking status, saving status, employment

Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Data mining provides a powerful tool for discovering discrimination in historical decision records

Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

Group discr.

Individual discr.

Individual discr.

Ind./Group discr.

Group discr.

Group discr.

Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

D. Pedreschi, S. Ruggieri and F. Turini (2008). Discrimination-aware data mining. In KDD'08.

D. Pedreschi, S. Ruggieri, and F. Turini (2009). Measuring discrimination in socially-sensitive decision records. In SDM'09.

S. Ruggieri, D. Pedreschi, and F. Turini (2010). Data mining for discrimination discovery. In TKDD 4(2).

Defining potentially discriminated (PD) groups

A subset of attribute values are perceived as potentially discriminatory based on background knowledge. Potentially discriminated groups are people with those attribute values.

Examples:

Female gender

Ethnic minority (*racism*) or minority language

Specific age range (*ageism*)

Specific sexual orientation (*homophobia*)

Discrimination and combinations of attribute values

Discrimination can be a result of several joint characteristics (attribute values) which are not discriminatory by themselves

Thus, the object of discrimination should be described by a conjunction of attribute values:

Itemsets

Association and classification rules

Association rules are if/then statements that help uncover relationships between seemingly unrelated data in a relational database.

In a classification rule, Y is a class item and X contains no class items.

$$\mathbf{X} \rightarrow \mathbf{Y}$$

Direct discrimination

Direct discrimination implies rules or procedures that impose 'disproportionate burdens' on minorities

PD rules are any classification rule of the form:

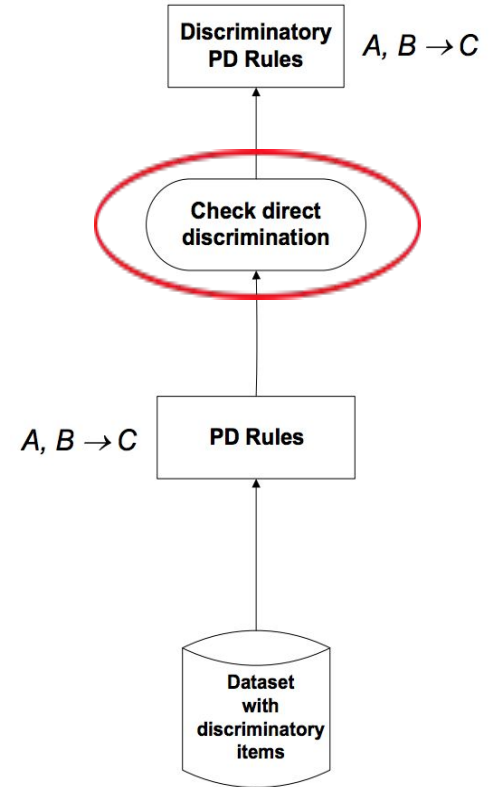
$$A, B \rightarrow C$$

where A is a PD group (B is called a "context")

Example:

gender="female", saving_status="no known savings"

→ credit=no



Favoritist PD rules

Is unveiled by looking at PD rules of the form

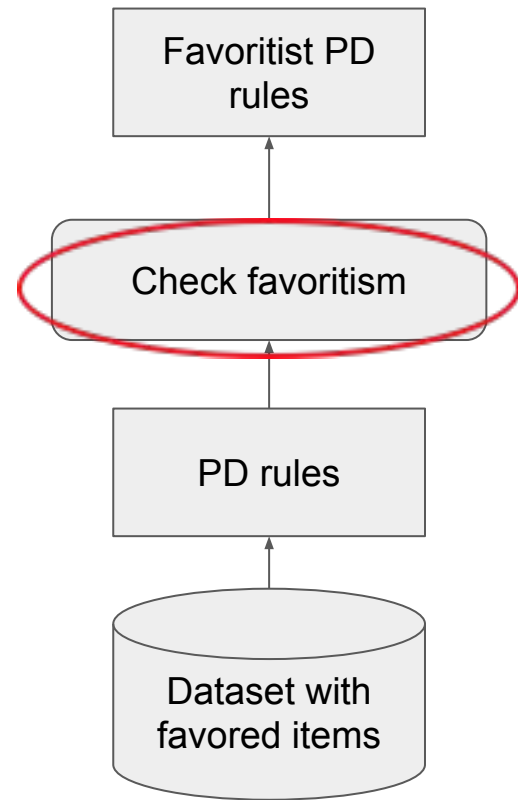
$$A, B \rightarrow C$$

where C grants some benefit and A refers to a favored group.

Example:

gender="male", savings="no known savings"

→ credit=yes



Indirect discrimination

Indirect discrimination implies rules or procedures that impose 'disproportionate burdens' on minorities, though not explicitly using discriminatory attributes

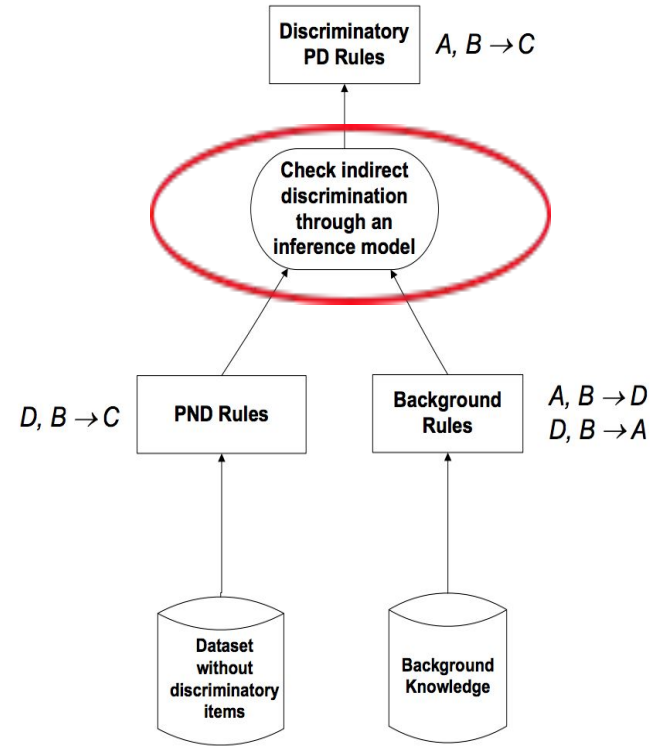
Potentially non-discriminatory (PND) rules may unveil discrimination, and are of the form:

$D, B \rightarrow C$ where D is a PND group

Example:

neighborhood="10451", city="NYC"

→ credit=no



Indirect discrimination example

Suppose we know that with high confidence:

(a) neighborhood=10451, city=NYC → benefit=deny

But we also know that with high confidence:

(b) neighborhood=10451, city=NYC → race=black

Hence:

(c) race=black, neighborhood=10451, city=NYC → benefit=deny

Rule (b) is background knowledge that allows us to infer (c), which shows that

rule (a) is indirectly discriminating against blacks

Evaluating PD rules through the extended lift

Remembering that $\text{conf}(X \rightarrow Y) = \text{support}(X \rightarrow Y) / \text{support}(X)$

We define the **extended lift with respect to B** of rule $A, B \rightarrow C$ as:

$$\text{elift}_B(A, B \rightarrow C) = \text{conf}(A, B \rightarrow C) / \text{conf}(B \rightarrow C)$$

The rules we care about are PD rules such that:

- A is a protected group (e.g. female, black)
- B is a context (e.g. lives in San Francisco)
- C is an outcome (usually negative, e.g., deny a loan)

The concept of α -protection

For a given threshold α , we say that PD rule $A, B \rightarrow C$, involving a PD group A in a context B for an outcome C , is α -protective if:

$$\text{elift}_B(A, B \rightarrow C) = \text{conf}(A, B \rightarrow C) / \text{conf}(B \rightarrow C) \leq \alpha$$

Otherwise, we say that $A, B \rightarrow C$ is an α -discriminatory rule

Extension: strong α -protection is the symmetric version, where we ensure the numerator is greater than the denominator by using $(1-\text{conf})/(1-\text{conf})$ if needed.

Relation of α -protection and group representation

For a given threshold α , we say that PD rule $A, B \rightarrow C$, involving a PD group A in a context B for a (usually bad) outcome C , is α -protective if:

$$\text{elift}_B(A, B \rightarrow C) = \text{conf}(A, B \rightarrow C) / \text{conf}(B \rightarrow C) \leq \alpha$$

Note that:

$$\text{elift}_B(A, B \rightarrow C) = \text{conf}(B, C \rightarrow A) / \text{conf}(B \rightarrow A)$$

This means extended lift is the ratio between the proportion of the disadvantaged group A in context B for (bad) outcome C , over the overall proportion of A in B .

Direct discrimination example

Rule (a):

city="NYC"

→ benefit=deny

with confidence 0.25

Rule (b):

race="black", city="NYC"

→ benefit=deny

with confidence 0.75 **elift 3.0**

Additional (discriminatory) element increases the rule confidence up to 3 times.

According to α -protection method, if the threshold $\alpha=3$ is fixed then the rule (b) is classified as discriminatory

Real-world example from German credit dataset

Fixing $\alpha=3$:

(B) saving status = "no known savings" \rightarrow credit = deny conf. 0.18

(A) personal status = "female div/sep/mar",
 saving status = "no known savings" \rightarrow credit = deny conf. 0.27 **elift 1.52**

Rule A is α -protective.

Real-world example 2 from German credit dataset

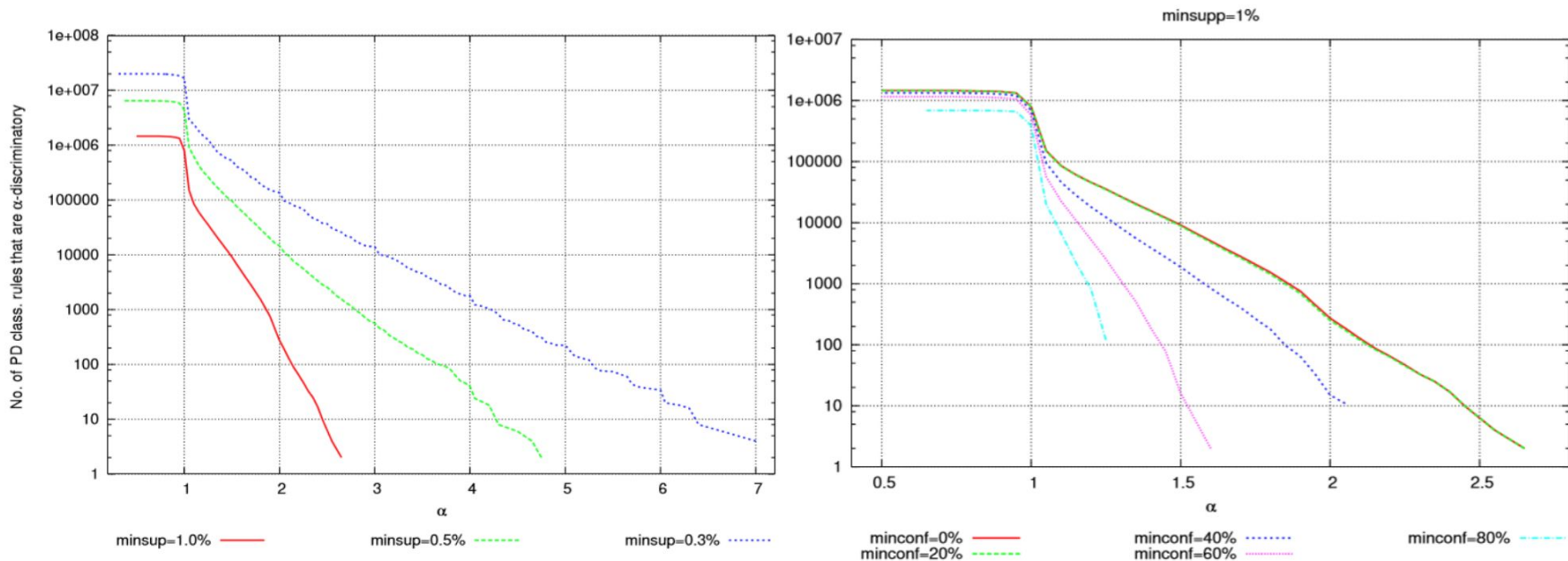
Fixing $\alpha=3$:

(B) purpose = "used car" → credit = deny conf. 0.17

(A) age = "52.6+", personal status = "female div/sep/mar",
purpose = "used car" → credit = deny conf. 1.00 **elift 6.06**

Rule A is α -discriminatory.

α -Discriminatory rules in the German credit dataset



Genuine occupational requirements

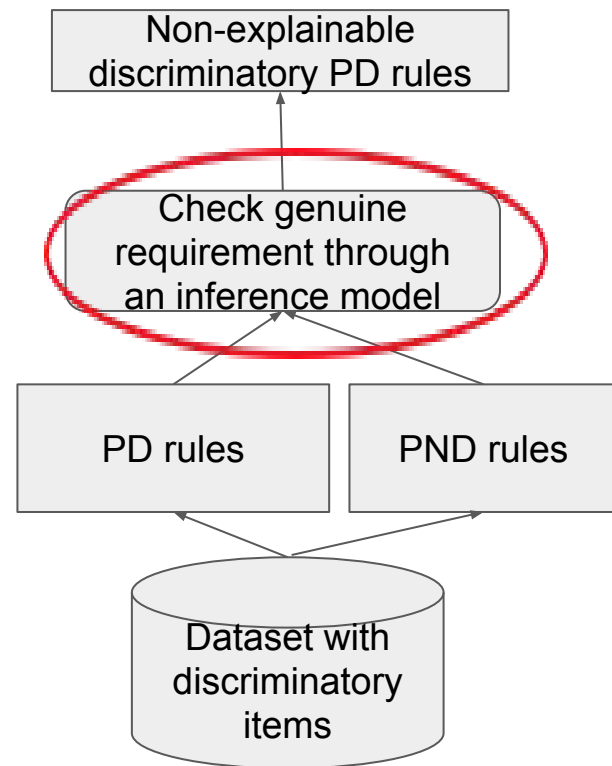
Supported by a PD rule of the form

$$A, B \rightarrow C$$

where C denies some benefit, we search for PND rules of the form

$$D, B \rightarrow C$$

such that D is a legitimate requirement, having the same effects of the PD rule



Example: genuine occupational requirement

(a) [A] gender="female", [B] city="NYC" \rightarrow [C] hire=no conf. 0.58

(b) [D] drive_truck="false", [B] city="NYC" \rightarrow [C] hire=no conf. 0.81

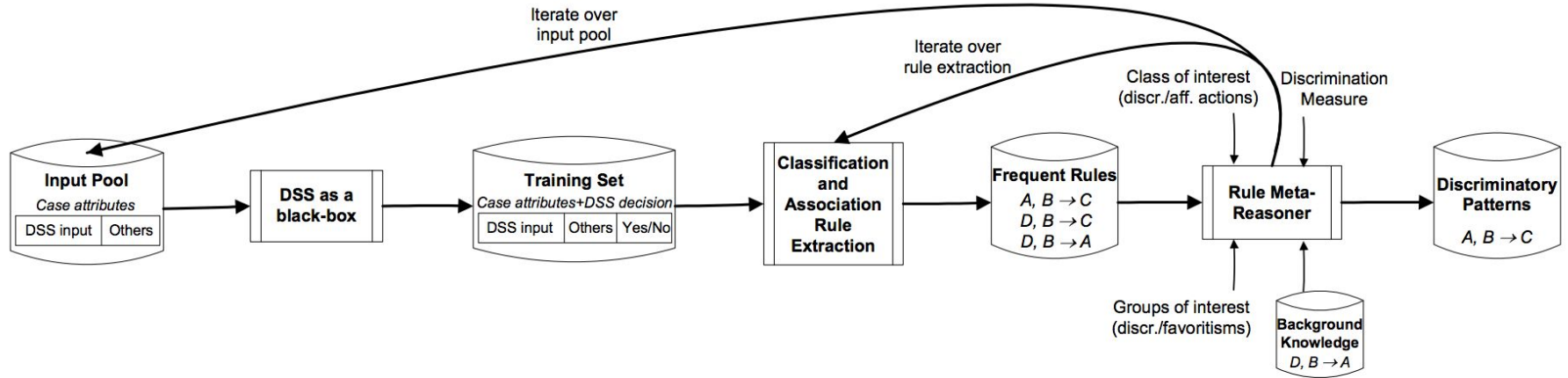
(c) [A] gender="female", [B] city="NYC" \rightarrow [D] drive_truck=false conf. 0.91

Let $p \in [0, 1]$. Classification rule (a) $A, B \rightarrow C$ with A being a PD attribute, is a p -instance of a PND rule (b) $D, B \rightarrow C$, if:

- D is a legitimate ground for the decision (i.e., accepted by law),
- $\text{conf}(D, B \rightarrow C) \geq p \cdot \text{conf}(A, B \rightarrow C)$, and
- $\text{conf}(A, B \rightarrow D) \geq p$.

Pipeline for analyzing discrimination

Reference model for analysing and reasoning on discrimination



Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

B. T. Luong, S. Ruggieri, and F. Turini (2011). k-NN as an implementation of situation testing for discrimination discovery and prevention. KDD'11

Limitations of classification rules approach

Legal limitations

- Undifferentiated groups are compared:
 - Do, e.g., women have the same characteristics of men they are compared with?
 - Or do they differ as per skills or other legally admissible reasons?

Limitations of classification rules approach (cont.)

Interpretational limitations

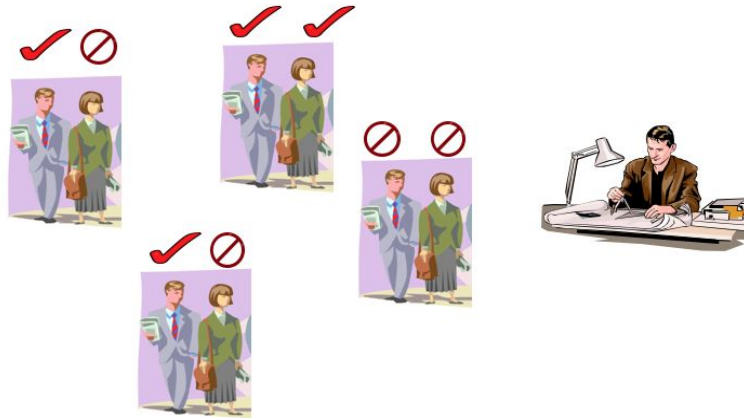
- Local contexts, possibly overlapping
- No global description of who is discriminated and who is not

Technical limitations

- Due to the use of frequent itemset mining (nominal attributes, nominal decisions) it requires discretization

Situation testing

- Legal approach for creating controlled experiments
- **Matched pairs** undergo the same situation, e.g. apply for a job
 - Same characteristics apart from the discrimination ground



k-NN as situation testing

Input: a dataset R of decision records

- For $r \in R$, $\text{dec}(r)$ is the decision (discrete or continuous)
 - E.g., $\text{dec}(r)$ is grant-benefit or deny-benefit
- $P(R)$ is the set of protected-by-law groups, e.g., women
 - E.g., $P(R) = \{r \in R \mid r[\text{gender}] = \text{female}\}$
- $U(R) = R \setminus P(R)$ is the rest of the dataset, e.g., men

Relax the "identical characteristics" of situation testing to a "similar characteristics" by using a distance function d

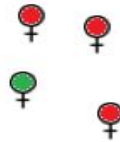
- $d(a,b)$ is defined over attributes that are legally admissible for the purpose of taking the decision

k-NN as situation testing (algorithm)

For $r \in P(R)$, look at its k closest neighbors

- ... in the protected set
 - define p_1 = proportion with the same decision as r
- ... in the unprotected set
 - define p_2 = proportion with the same decision as r
- measure the degree of discrimination of the decision for r
 - define $\text{diff}(r) = p_1 - p_2$ (*think of it as expressed in percentage points of difference*)

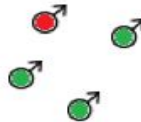
$\text{knn}_p(r,k)$



r



$\text{knn}_U(r,k)$



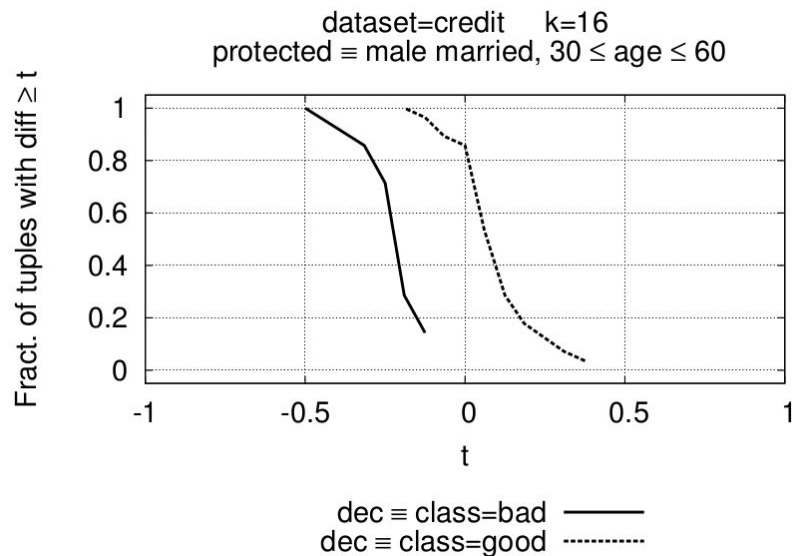
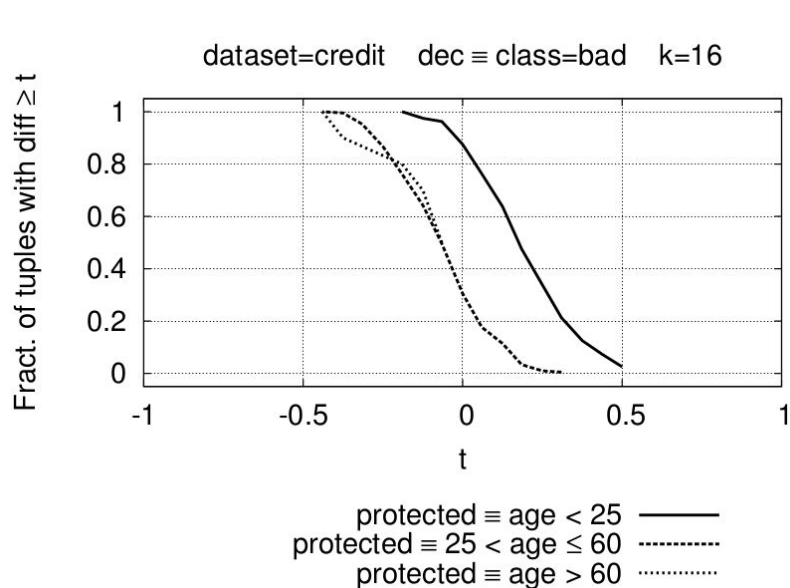
$$p_1 = 0.75$$

$$p_2 = 0.25$$

$$\text{diff}(r) = p_1 - p_2 = 0.50$$

k-NN as situation testing (results)

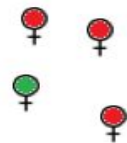
- If decision=deny-benefit, and $\text{diff}(r) \geq t$, then we found discrimination around r



Characterizing discrimination using k-NN

- For $r \in P(R)$, set a new attribute: "t-discriminated"
 - If $\text{dec}(r) = \text{deny-benefit}$ and $\text{diff}(r) \geq t$, $\text{t-discriminated}(r) := \text{TRUE}$
 - Otherwise $\text{t-discriminated}(r) := \text{FALSE}$
- Example: for $t=0.3$ the sample r below is classified as t-discriminated

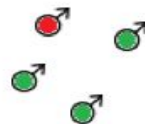
$\text{knn}_P(r,k)$



r



$\text{knn}_U(r,k)$



$$p_1 = 0.75$$

$$p_2 = 0.25$$

$$\text{diff}(r) = p_1 - p_2 = 0.50$$

Characterizing discrimination using k-NN (cont.)

- To answer the question: under which conditions women were t-discriminated?
- We create a classifier with training set $P(\mathcal{R})$, i.e. only protected people, and with class attribute *t-discriminated*

```
DiscoveryN( $\mathcal{R}$ , t) {  
   $\mathcal{L} = \emptyset$   
  for  $\mathbf{r} \in P(\mathcal{R})$  {  
    if(  $dec(\mathbf{r}) = \ominus$  and  
         $diff(\mathbf{r}) \geq t$  )  
       $\mathbf{r}[\text{disc}] = \text{yes}$   
    else  
       $\mathbf{r}[\text{disc}] = \text{no}$   
     $\mathcal{L} = \mathcal{L} \cup \{\mathbf{r}\}$   
  }  
  build a classifier on  $\mathcal{L}$   
}
```

Characterizing discrimination using k-NN (results)

- German credit dataset
 - protected = female non-single
 - 0.10-discriminated cases

- Decision tree model (C4.5)

```
num_dependents <= 1
|  credit_amount <= 2631: disc=yes (59.0/9.0)
|  credit_amount > 2631: disc=no (44.0/15.0)
num_dependents > 1: disc=no (6.0)
```

```
disc=yes: Precision 0.847  Recall 0.769
```

Discriminated women had no dependents (children) and were asking for small amounts

- Classification rule model (RIPPER)

```
(credit_amount >= 3190) => disc=no (39.0/12.0)
(installment_commitment <= 2) and (residence_since >= 3)
                                     => disc=no (10.0/2.0)
=> disc=yes (60.0/9.0)
```

```
disc=yes: Precision 0.85  Recall 0.785
```

Discriminated women were asking for small amounts and were either paying in many installments or had been resident for a short time

k-NN for discrimination prevention

- Goals of non-discriminating classifier:
 - Maximize classifier accuracy
e.g., give credit to people who will pay
 - Minimize t-discriminated cases
e.g., give credit to women if similar men would have been given credit
- Basic idea: t-correction of training set
 - Flip the labels from negative to positive for t-discriminated cases in the training set

PreventionN($\mathcal{T}, \mathcal{V}, t$) {
 $\mathcal{T}' = \emptyset$
 for $\mathbf{r} \in \mathcal{T}$ {
 $\mathbf{r}' = \mathbf{r}$
 if($dec(\mathbf{r}) = \ominus$ and
 $protected(\mathbf{r})$ and
 $diff(\mathbf{r}) \geq t$)
 $\mathbf{r}'[dec] = \oplus$
 $\mathcal{T}' = \mathcal{T}' \cup \{\mathbf{r}'\}$
 }
 }
 build classifiers on \mathcal{T} and \mathcal{T}'
 compare them on \mathcal{V}
}

(Discrimination prevention/mitigation is the focus of Part III ...)

k-NN for discrimination prevention (results)

classifier	No pre-processing		0.10-correction	
	accuracy	0.10-discr.	accuracy	0.10-discr.
C4.5	85.60%	4.24%	84.94%	1.07%
Naïve Bayes	82.46%	4.06%	82.33%	2.23%
Logistic	85.28%	6.61%	84.70%	0.61%
RIPPER	84.42%	5.24%	83.98%	3.94%
PART	85.20%	12.62%	84.00%	2.3%

Only the training set changes, the testing set is fixed

Accuracy shows a small decrease

0.10-discrimination reduces substantially

Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

K. Mancuhan and C. Clifton (2014). Combating discrimination using bayesian networks. In Artificial Intelligence and Law, 22(2).

Let's go back to the classification-based approach

It hinges on computing $\text{elift}_B(A, B \rightarrow C) = \text{conf}(A, B \rightarrow C) / \text{conf}(B \rightarrow C)$

What are these quantities?

$\text{conf}(A, B \rightarrow C)$

is the confidence on an outcome given a protected attribute and a context

$\text{conf}(B \rightarrow C)$

is the confidence on an outcome given just the context

Extending the *elift* idea

Let A be the protected attributes, B the context (unprotected attributes)

Suppose we have a way of computing $P(\text{benefit}=\text{deny}|\dots)$

Suppose we use it to classify people: $P(\text{benefit}=\text{deny}|\dots) \geq \theta \Rightarrow \text{deny_benefit}$

If $P(\text{benefit}=\text{deny}|A, B) / P(\text{benefit}=\text{deny}|B) \geq \alpha$,

then we could say that the classifier is α -discriminatory

Extending the *elift* idea (cont.)

Let A be the protected attributes, B the context (unprotected attributes)

Let R be "redlining" attributes, which are correlated with the protected attributes A

Suppose we have a way of computing $P(\text{benefit}=\text{deny}|\dots)$

Suppose we use it to classify people: $P(\text{benefit}=\text{deny}|\dots) \geq \theta \Rightarrow \text{deny_benefit}$

If $P(\text{benefit}=\text{deny}|A, B, \mathbf{R}) / P(\text{benefit}=\text{deny}|B) \geq \alpha$,

then we could say that the classifier is α -discriminatory

Bayesian networks

- Bayesian networks estimate the probability $P(A, B, C)$ by capturing the conditional dependencies between the attributes within the sets A and B .
 - Bayesian networks can be used to estimate $P(A, B, C)$ probability and the $P(C|A, B)$ class probability can be derived from the Bayes theorem
 - Bayesian networks capture correlations between attributes
 - Bayesian networks are appropriate to define a decision process
- The *elift* can be extended to *belift* by calculating the numerator and the denominator probabilities with Bayesian networks
 - It identifies changes in decision making due to the usage of protected and redlining attributes
 - It captures both direct and indirect discrimination

Bayesian *elift* (*belift*)

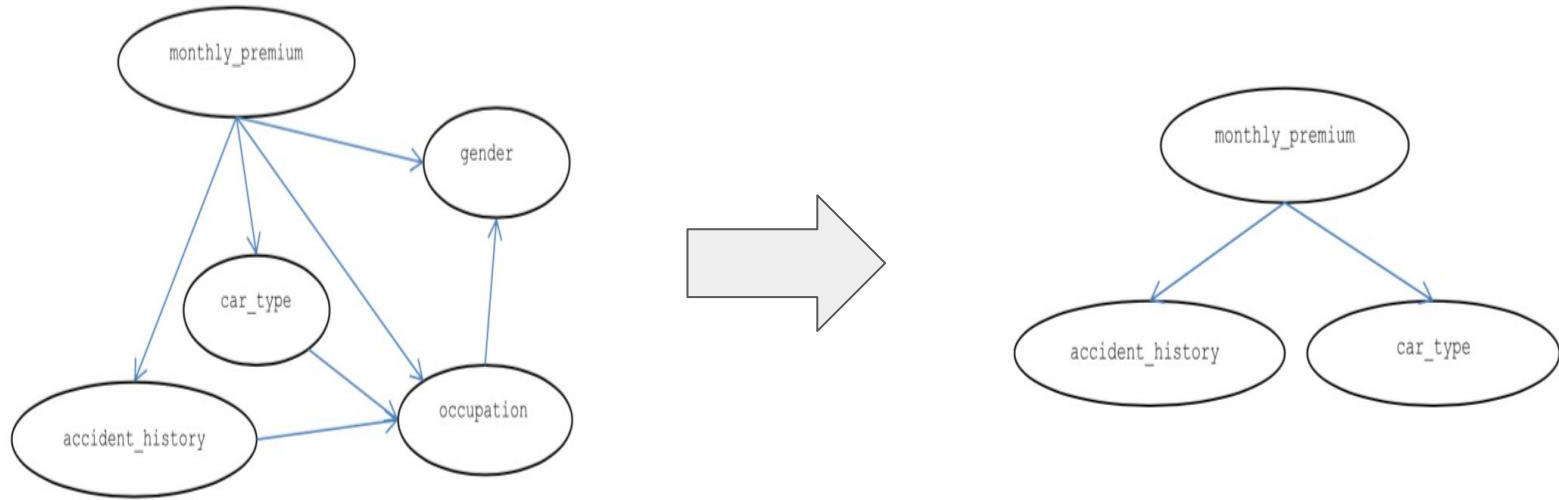
$$belift = \frac{P(C | a_1, a_2, \dots, a_l, b_1, b_2, \dots, b_m, r_1, r_2, \dots, r_n)}{P(C | b_1, b_2, \dots, b_m)}$$

This indicates how many times the usage of protected attributes (A) and the redlining attributes (R) increase the class probability for a given instance with respect to its probability using only non-protected attributes (B)

Bayesian networks for discrimination discovery

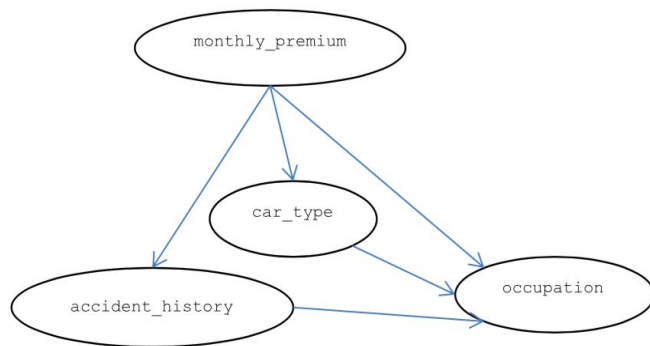
- Build a Bayesian network from a given set of instances D (*net*)
- Build a second Bayesian network (*relative_net*)
 - By removing protected attributes and any attribute directly connected to them in *net*
- For each instance i in D *net* and *relative_net*, compute probability of target class and apply decision threshold
 - If the decision is the same, ignore: instance was not discriminated
 - If the decision is different (i.e., **there is a "decision flip" when adding protected and redlining attributes**), and $belift \geq \alpha$, add instance i to the set of discriminated instances

Example: *net* and *relative_net*



Bayesian networks for discrimination prevention

- Build a Bayesian network from a given set of instances D (net)
- Build a non-discriminatory Bayesian network from net (nd_net)
 - **Remove** the protected attribute nodes, but **keep** the redlining attribute nodes
 - **Correct** the dataset by flipping the class label of discriminated instances
 - **Update** probability tables of nd_net using the corrected dataset
- Return the non-discriminatory Bayesian network (nd_net)



Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

F. Bonchi, S. Hajian, B. Mishra, and D. Ramazzotti (2015). [Exposing the probabilistic causal structure of discrimination](https://arxiv.org/abs/1510.00552). arXiv:1510.00552.

Previous approaches

- Legal limitations
 - Any legally-valid proof of discrimination requires evidence of causality [Foster, 2004]
- Technical limitations
 - The state-of-the-art methods are essentially correlation-based
 - Spurious correlations can lead to false negatives and false positives
 - A Bayesian network would not be able to disentangle the direction of any causal relationship

To prove discrimination:

We need to assess discrimination as a causal inference problem from a database of past decisions, where **causality can be inferred probabilistically**

Suppes probabilistic causation theory (constraints)

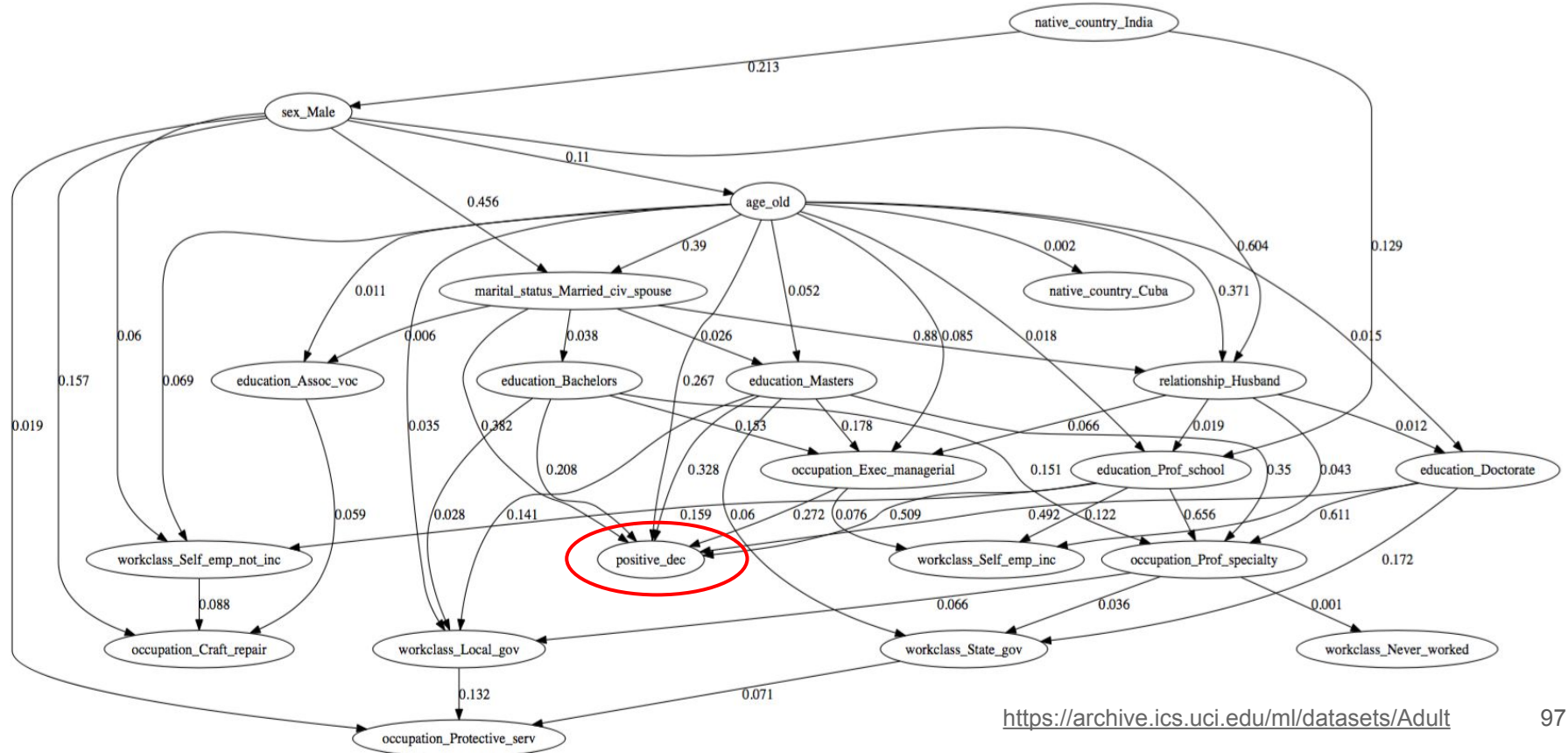
- Let h denote cause, e denote effect
- Temporal priority
 - Any cause must happen before its effect: $t_h < t_e$
- Probability raising
 - Cause must raise the probability of observing the effect: $P(e | h) > P(e | \neg h)$

The Suppes-Bayes Causal Network (SBCN)

- Represents the causal structures existing among the attributes in the data by a constrained Bayesian network:
 - each node represents an assignment attribute=value
 - each arc (u,v) represents the existence of a relation between u and v satisfying Suppes' constraints (temporal priority and probability raising)
 - each arc is labeled with a positive weight: $p(u|v) - p(u|\neg v)$
- Reconstructed from the data using **maximum likelihood estimation (MLE)**, where we force the conditional probability distributions induced by the reconstructed graph to obey Suppes' constraints
 - *Network simplification* by extracting a minimal set of edges which best explain the data. This **regularization** is done by means of the Bayesian Information Criterion (BIC)

$$score_{\text{BIC}}(D', G') = LL(D'|G') - \frac{\log s}{2} \dim(G').$$

SBCN example: Census Income dataset



Graph representation allows interesting applications

Two largest communities
detected using the
Walktrap algorithm
[Pons and Latapy 2006]

C_1

negative_dec, wc:Private, ed:Some_college, ed:Assoc_acdm,
ms:Never_married, ms:Divorced, ms:Widowed,
ms:Married_AF_spouse, oc:Sales, oc:Other_service,
oc:Priv_house_serv, re:Own_child, re:Not_in_family, re:Wife,
re:Unmarried, re:Other_relative, ra:Black, oc:Armed_Forces,
oc:Handlers_cleaners, oc:Tech_support, oc:Transport_moving,
ed:7th_8th, ed:10th, ed:12th, ms:Separated,
ed:HS_grad,ed:11th, nc:Outlying_US_Guam_USVLEtc,
nc:Haiti, ag:young, sx:Female, ra:Amer_Indian_Eskimo,
nc:Trinidad_Tobago, nc:Jamaica, oc:Machine_op_inspct,
ms:Married_spouse_absent, oc:Adm_clerical,

C_2

positive_dec, oc:Prof_specialty, wc:Self_emp_not_inc,
ms:Married_civ_spouse, oc:Craft_repair,oc:Protective_serv,
re:Husband, ed:Prof_school, wc:Self_emp_inc,
ag:old, wc:Local_gov, oc:Exec_managerial,
ed:Bachelors, ed:Assoc_voc, ed:Masters, wc:Never_worked,
wc:State_gov, ed:Doctorate, sx:Male, nc:India, nc:Cuba

age (ag), education (ed), marital status (ms),
native country (nc), occupation (oc), race(ra),
relationship (re), sex (sx), workclass (wc)

P. Pons and M. Latapy (2006). "Computing Communities in large networks using random walks." *J. Graph Algorithms Appl.* 10.2: 191-218.

Key idea: random walks

Start from the node we want to test (e.g. gender=female)

Let δ^+ represent the positive outcome node (e.g., high salary)

δ^- the negative outcome (e.g., low salary)

Perform n random walks starting from v

Let $\text{rw}(v \rightarrow \delta^-)$ be the number of random walks starting from v and ending in δ^-

People with attribute v are discriminated if the fraction $\text{rw}(v \rightarrow \delta^-) / n$ is large

Result: top discriminated groups

Census Income Dataset

	$ds^-(v)$
relationship_Unmarried	1
marital_status_Never_married	0.996
age_Young	0.995
race_Black	0.994
sex_Female	0.98

German credit dataset

	$ds^-(v)$
residence_since_le_1d6	1
residence_since_gt_2d8	1
residence_since_from_1d6_le_2d2	1
age_gt_52d6	0.86
personal_status_male_single	0.791

Extension 1: genuine occupational requirements

Suppose v_{legal} is a genuine occupational requirement
e.g., having a truck driver license

The *fed* coefficient (fraction of explainable discrimination) is

$$rw(v \rightarrow v_{\text{legal}} \rightarrow \bar{\delta}-) / rw(v \rightarrow \bar{\delta}-)$$

Result: fraction of explainable discrimination

Source node	Intermediate	$fed^-(v)$
race_Amer_Indian_Eskimo	education_HS_grad	0.481
sex_Female	occupation_Other_service	0.310
age_Young	occupation_Other_service	0.193
relationship_Unmarried	education_HS_grad	0.107
race_Black	education_11th	0.083

Extension 2: individual and subgroup discrimination

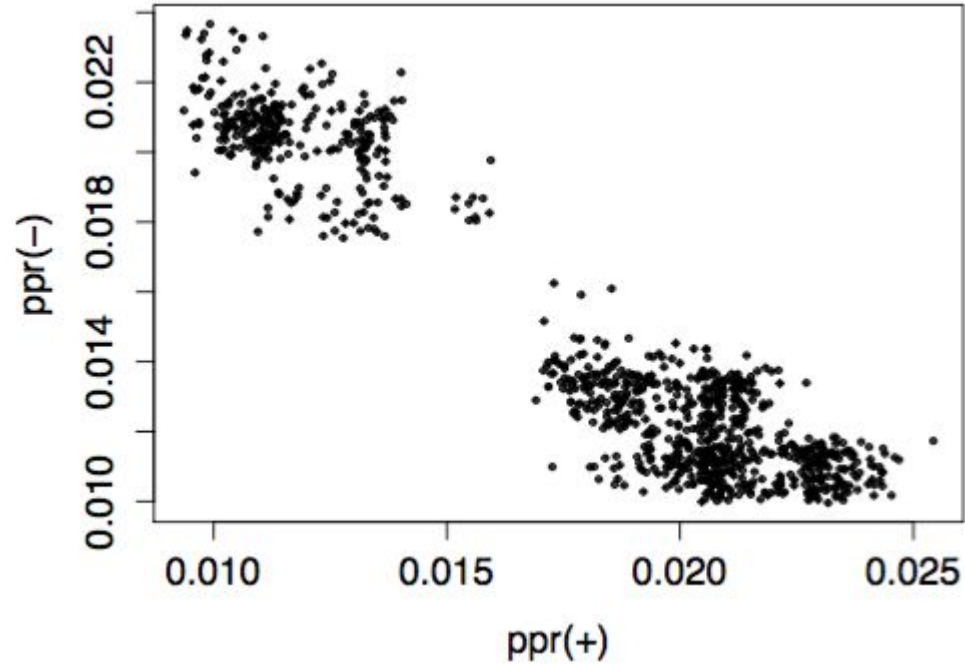
Suppose we have subgroup having attributes

$$v_1, v_2, \dots, v_n$$

This subgroup can be of size 1, i.e., represent an individual

$$\text{gds}(v_1, v_2, \dots, v_n) = \frac{\text{ppr}(\bar{\delta}- | v_1, v_2, \dots, v_n)}{\text{ppr}(\bar{\delta}+ | v_1, v_2, \dots, v_n) + \text{ppr}(\bar{\delta}- | v_1, v_2, \dots, v_n)}$$

Result: Individual discrimination (German credit)



Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

S. Ruggieri, S. Hajian, F. Kamiran, X. Zhang (2014). Anti-discrimination analysis using privacy attack strategies. In PKDD'14.

Privacy attack strategies

Methods for direct discrimination discovery require that:

- The dataset explicitly contains an attribute denoting the PD group
 - Otherwise, we are in an indirect discrimination discovery setting
- The dataset has not been pre-processed prior to discrimination discovery
 - Otherwise, we are in a privacy-aware discrimination discovery setting

Discrimination discovery scenarios

- Indirect discrimination discovery
- Privacy-aware discrimination discovery
(e.g., original data with all attributes is no longer available)
- Discrimination data recovery
(e.g., original data with all attributes has been hidden from authorities)

Discrimination discovery scenarios: intuition

- Indirect discrimination discovery
- Privacy-aware discrimination discovery
- Discrimination data recovery

There is an interesting parallel between the role of the anti-discrimination authority in these three scenarios and the role of an attacker in **private data publishing**

The discrimination authority must infer who belongs to a protected group, based on the decisions that have been made.

Privacy attack strategies

- Combinatorial attacks based on Frèchet bounds inference
- Attribute inference attacks
- Minimality attacks

We will use risk difference (RD)

group	decision		
	-	+	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

$$p_1 = a/n_1$$

$$p_2 = c/n_2$$

$$p = m_1/n$$

$$RD = p_1 - p_2$$

A discrimination table is α -protective (w.r.t. the *risk difference* measure RD) if $RD \leq \alpha$. Otherwise, it is α -discriminatory.

Indirect discrimination discovery

Think of g_1 and g_2 as *redlining* attributes (i.e., correlated with being protected)

group	decision		
	-	+	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

unknown contingency table

rel. group	decision		
	-	+	
g1	\hat{a}	\hat{b}	\hat{n}_1
g2	\hat{c}	\hat{d}	\hat{n}_2
	m_1	m_2	n

known contingency table

group	rel. group		
	g1	g2	
protected	e	f	n_1
unprotected	g	h	n_2
	\hat{n}_1	\hat{n}_2	n

background knowledge contingency table

Exploit Frèchet bounds

$$n_A = \text{supp}(A)$$

$$\min\{n_X, n_Y\} \geq n_{XY} \geq \max\{n_X + n_Y - n_{X \cap Y}, 0\}$$

Indirect discrimination discovery

group	decision		
	-	+	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

unknown contingency table

rel. group	decision		
	-	+	
g1	\hat{a}	\hat{b}	\hat{n}_1
g2	\hat{c}	\hat{d}	\hat{n}_2
	m_1	m_2	n

known contingency table

group	rel. group		
	g1	g2	
protected	e	f	n_1
unprotected	g	h	n_2
	\hat{n}_1	\hat{n}_2	n

background knowledge contingency table

We can decompose: $a = a_1 + a_2$, where a_1 are in g1, a_2 are in g2

Let X be (rel. group="g1" and dec="-")

Let Y be (rel. group="g1" and group="protected")

$$\min\{\hat{a}, e\} \geq a_1 \geq \max\{\hat{a} + e - \hat{n}_1, 0\} = \max\{e - \hat{b}, 0\}$$

$$\min\{n_X, n_Y\} \geq n_{XY} \geq \max\{n_X + n_Y - n_{X \cap Y}, 0\}$$

Indirect discrimination discovery

group	decision		
	-	+	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

unknown contingency table

rel. group	decision		
	-	+	
g1	\hat{a}	\hat{b}	\hat{n}_1
g2	\hat{c}	\hat{d}	\hat{n}_2
	m_1	m_2	n

known contingency table

group	rel. group		
	g1	g2	
protected	e	f	n_1
unprotected	g	h	n_2
	\hat{n}_1	\hat{n}_2	n

background knowledge contingency table

We do the same with a_2 , and a similar decomposition for c , obtaining:

$$\left[\begin{array}{l} \min\{\hat{a}, e\} + \min\{\hat{c}, f\} \geq a \geq \max\{e - \hat{b}, 0\} + \max\{f - \hat{d}, 0\} \\ \min\{\hat{a}, g\} + \min\{\hat{c}, h\} \geq c \geq \max\{g - \hat{b}, 0\} + \max\{h - \hat{d}, 0\} \end{array} \right]$$

$$RD \geq RDlb = \frac{\max\{e - \hat{b}, 0\} + \max\{f - \hat{d}, 0\}}{n_1} - \frac{\min\{\hat{a}, g\} + \min\{\hat{c}, h\}}{n_2}$$

Lower bound

Uses only known and background data

Privacy-aware discrimination discovery

KNOWN					Known without IDs		UNKNOWN		
ID	Education	Job	Dec	GID	GID	Religion			
r_1	Bachelors	Engineer	-	1	1	Muslim			
r_2	Bachelors	Engineer	+	1	1	Christian			
r_3	Doctorate	Engineer	+	1	1	Jewish			
r_4	Bachelors	Writer	+	1	1	Other			
r_5	Master	Engineer	+	2	2	Muslim			
r_6	Doctorate	Writer	+	2	2	Christian			
r_7	Bachelors	Dancer	-	2	2	Jewish			
r_8	Master	Dancer	-	2	2	Other			
r_9	Master	Dancer	-	3	3	Muslim			
r_{10}	Master	Lawyer	+	3	3	Christian			
r_{11}	Bachelors	Engineer	-	3	3	Jewish			
r_{12}	Bachelors	Dancer	-	3	3	Other			

		decision		
education=bachelors	-	+		
religion=muslim	a	b		3
religion≠muslim	c	d		3
	4	2		6

When restricting to education=bachelors, n_1 becomes muslim bachelors; n_1 is always required as input

Frèchet bounds

$n^i = \text{people in group } i$ $n_1 = \text{total number of protected people}$

$$\sum_i \min\{1, n_-^i\} \geq a \geq n_1 - \sum_i \min\{1, n_+^i\}$$

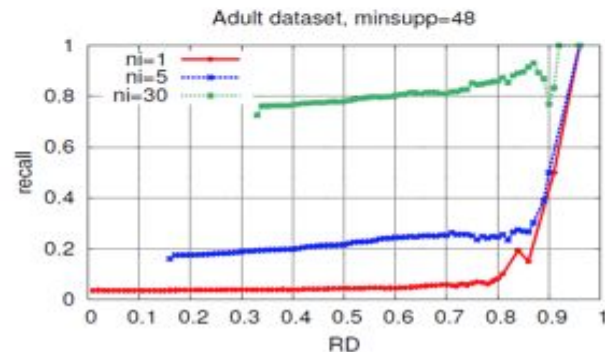
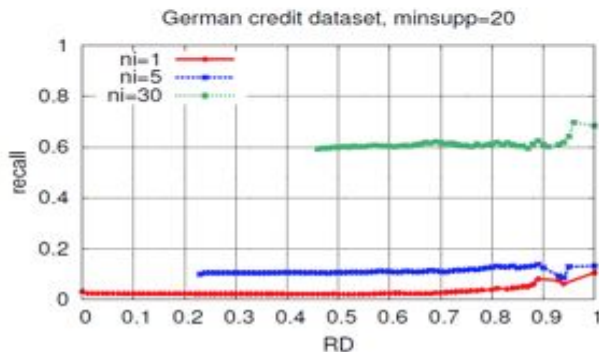
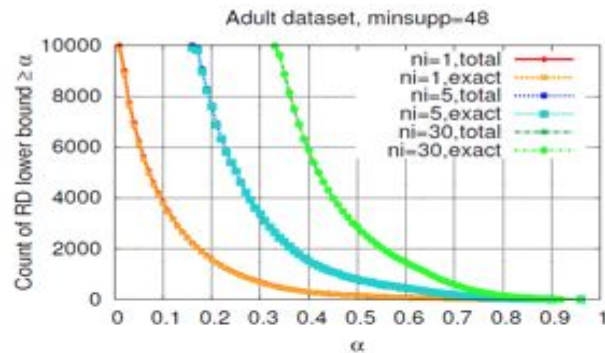
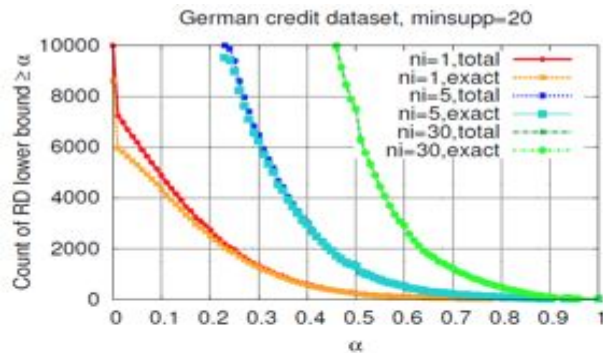
Discrimination data recovery

- Data owner hides discrimination by flipping a minimum number of decisions
 - Releases R' with changed decisions to have $RD(R')=0$, but we want R
- Suppose we know the original risk difference $RD(R)$
- Discrimination affects the tuples close to the decision boundary of a classifier
- To determine the decision boundary, we rank tuples of the protected and unprotected groups separately w.r.t. their positive decision probabilities accordingly to a classifier trained from the transformed data R'
- We don't fix all data at the same time, but iteratively, re-training at each iter.

Sex	Ethnicity	Degree	Job Type	Dec	Prob
m	native	h.s.	board	+	98%
m	native	h.s.	board	+	98%
m	native	univ.	board	+	89%
<i>m</i>	<i>non-nat.</i>	<i>h.s.</i>	<i>health</i>	<i>-</i>	<i>47%</i>
m	non-nat.	univ.	health	-	30%

Sex	Ethnicity	Degree	Job Type	Dec	Prob
f	native	h.s.	board	+	93%
f	native	none	health	+	76%
<i>f</i>	<i>native</i>	<i>h.s.</i>	<i>edu.</i>	<i>+</i>	<i>51%</i>
f	non-nat.	univ.	edu.	-	2%
f	non-nat.	univ.	edu.	-	2%

Results: indirect discrimination discovery



Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Classification rule mining

k-NN classification

Bayesian networks

Probabilistic causation

Privacy attack strategies

Predictability approach

M. Feldman, S.A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian (2015). Certifying and removing disparate impact. In KDD'15.

Disparate impact

Measured using risk ratio p_1/p_2

$\Pr(\text{decision}=\text{deny_benefit} \mid \text{group}=\text{protected})$

$\Pr(\text{decision}=\text{deny_benefit} \mid \text{group}=\text{unprotected})$

group	benefit		
	denied	granted	
protected	a	b	n_1
unprotected	c	d	n_2
	m_1	m_2	n

$$p_1 = a/n_1 \quad p_2 = c/n_2 \quad p = m_1/n$$

Feldman et al. note this is:

$$1/\text{LR+} = (1-\text{specificity})/\text{sensitivity}$$

LR+ is the likelihood ratio of the positive class

US Equal Employment
Opportunity Commission
(EEOC) says $p_1/p_2 \leq 0.8$
means disparate impact

Disparate impact as predictability

Alice runs an algorithm on $D = (X, Y)$, with X protected, Y unprotected

Bob receives (D, C) , with C the outcomes

If Bob can predict X based on $D \setminus X, C$,

then the algorithm was discriminatory

Essentially, we want to know if the decisions are leaking information about the protected attributes

Discrimination discovery

Definition

Data mining approaches

Case studies

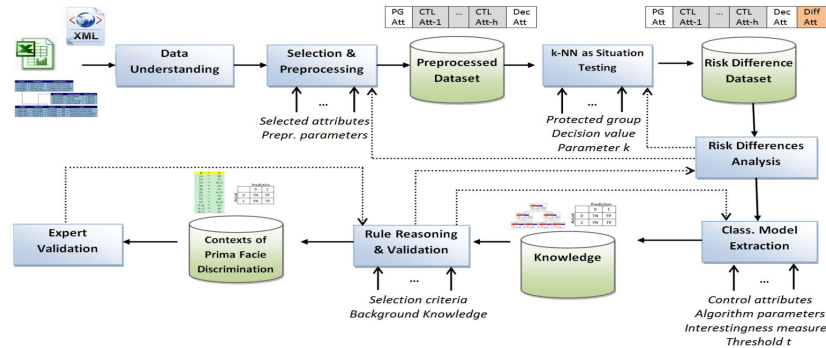
Discrimination discovery on the web

Crime suspect dataset

- **Dataset:** a real world dataset of Statistics Netherlands, which is a census body in the Netherlands.
- **Application:** the use of classifiers for predicting whether an individual is a crime suspect, or not, to support law enforcement and security agencies.
- **Results:**
 - The results show that **discrimination does exist** in real world datasets and blind use of classifiers learned over such datasets can **exacerbate** the discrimination problem.
 - It demonstrates that discrimination-aware classification methods can **mitigate the discriminatory effects** and that they lead to rational and legally acceptable decisions.

Scientific project evaluation

- **Dataset:** scientific research proposals
 - In 2008, the Italian Ministry of University and Research published a call for scientific research projects under the Basic Research Investment Fund (FIRB) reserved to young scientists.
- **Method:** classification rule mining and k-NN classification approaches to discover gender bias in scientific research funding



Discrimination discovery

Definition

Data mining approaches

Case studies

Discrimination discovery on the web

Discrimination in online services is not rare

- Non-black hosts can charge ~12% more than black hosts in Airbnb
 - Edelman, Benjamin G. and Luca, Michael, Digital Discrimination: The Case of Airbnb.com (January 10, 2014). Harvard Business School NOM Unit Working Paper No. 14-054.
- Price steering and discrimination in many online retailers
 - Aniko Hannak, Gary Soeller, David Lazer, Alan Mislove, and Christo Wilson: Measuring Price Discrimination and Steering on E-commerce Web Sites. Proc. of IMC. Vancouver, Canada, November 2014.
- China is about 21% larger by pixels when shown in Google Maps for China
 - Gary Soeller, Karrie Karahalios, Christian Sandvig, and Christo Wilson: MapWatch: Detecting and Monitoring International Border Personalization on Online Maps. Proc. of WWW. Montreal, Quebec, Canada, April 2016

Discrimination discovery on the web

How do we measure the “fairness of the web”?

- Need to model/understand user browsing behavior
- Evaluate how web sites respond to different behavior/attributes
- Cope with noisy measurements

Moritz Hardt (2013). Fairness through Awareness Presentation ([slides](#)).

Mikians, J., Gyarmati, L., Erramilli, V., and Laoutaris, N. (2012). Detecting price and search discrimination on the internet. In Proceedings of the 11th ACM Workshop on Hot Topics in Networks, pages 79–84. ACM

The importance of being Latanya

Names used predominantly by black men and women are much more likely to generate ads related to arrest records, than names used predominantly by white men and women.

Discovered by researcher Latanya Sweeney.



The dark side of Google Ads

AdFisher: tool to automate the creation of behavioral and demographic profiles.

Used to demonstrate that setting gender = female results in less ads for high-paying jobs.



Adfisher tool to uncover discrimination in ads

