

Algorithmic Discrimination

Carlos Castillo

Web Science and Social Computing

ETIC — UPF



Universitat
Pompeu Fabra
Barcelona

HUMAIN Workshop, March 2018

Session question addressed

Q2. How do algorithms have the potential to modify the way humans make decisions based on them?

My answer: algorithms can discriminate

Next, I present a precise formulation and examples

Generic discrimination

X discriminates against someone Y in relation to Z if:

1. Y has property P and Z does not have P
2. X treats Y worse than s/he treats or would treat Z
3. It is because Y has P and Z does not have P that X treats Y worse than Z

(also applies if X believes Y has P and Z does not have P)

Generic discrimination

X discriminates against someone Y in relation to Z if:

1. Y has property P and Z does not have P
2. X treats Y worse than s/he treats or would treat Z
3. It is because Y has P and Z does not have P that X treats Y worse than Z

Disadvantageous differential treatment

Group discrimination

X group-discriminates against Y in relation to Z if:

1. X generically discriminates against Y in relation to Z
2. P is the property of belonging to a socially salient group
3. This makes people with P worse off relative to others
or X is motivated by animosity towards people with P,
or by the belief that people with P are inferior
or should not intermingle with others

Statistical discrimination

X statistically discriminates against Y in relation to Z if:

1. X group-discriminates against Y in relation to Z
2. P is statistically relevant
(or X believes P is statistically relevant)

Example (statistical / non-statistical)

- a. Not hiring a highly-qualified woman because women have a higher probability of taking parental leave
(statistical discrimination)
- b. Not hiring a highly-qualified woman because she has said that she intends to have a child and take parental leave
(non-statistical discrimination)

In statistical machine learning

An algorithm developed through statistical machine learning can statistically discriminate if we:

1. Disregard intentions/animosity
2. Understand statistically relevant as any information derived from training data

Some examples

1. Disparate impact

The model gives people with P a bad outcome more often

2. Directly and indirectly discriminatory rules

The model associates P (or Q, which depends on P) to a bad outcome

3. Lack of calibration

The same output translates to different bad outcome probs. for P and not P

4. Disparate mistreatment / Lack of equal opportunity

The false positive rate of the bad outcome is higher for P than not P

5. Unfair rankings

The model gives people with P a lower ranking

Disparate impact

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 **risk** that **people with disabilities** face of not getting a scholarship is much larger than

c/n_2 , the **risk** that **people without disabilities** face of not getting a scholarship, then people with disabilities could claim they are being discriminated.

Directly discriminatory rules

From a database of decisions made in the past, we learn that
 $\text{gender} = \text{female} \Rightarrow \text{credit} = \text{no}$

$$\frac{P(\text{gender}=\text{female}, \text{credit}=\text{no})}{P(\text{gender}=\text{female})} > \theta$$

This means we have found evidence of **direct discrimination**

gender	has_job	credit
male	true	yes
male	false	yes
male	true	yes
female	false	no
female	true	no
female	true	yes
...

Indirectly discriminatory rules

If we learn the rule

zip = 8002 \Rightarrow credit = no

... and we know ...

zip = 8002 \Rightarrow origin = foreign

We have **indirect discrimination**

origin	zip	credit
national	8001	yes
national	8001	yes
national	8001	yes
foreign	8002	no
foreign	8002	no
foreign	8002	yes
...

Lack of calibration

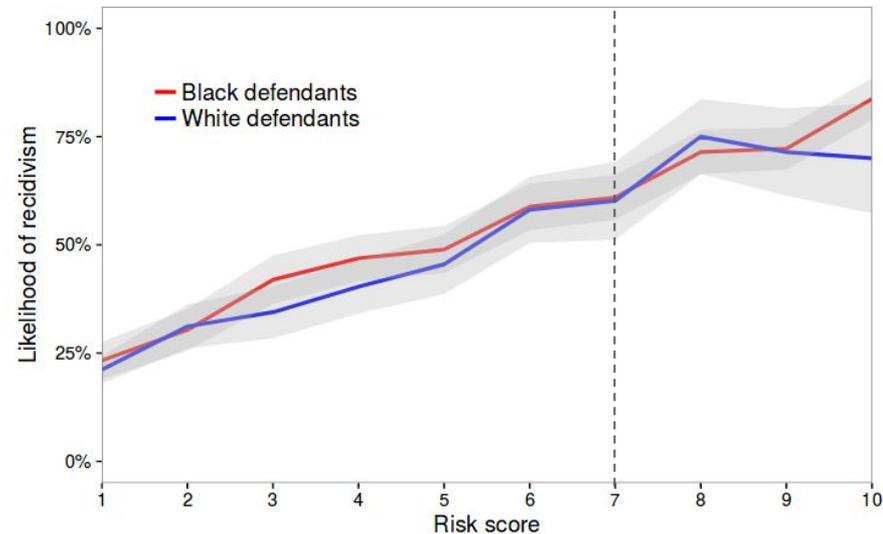
A model lacks calibration if the probability of an actual bad outcome depends on the class and not only on the output

This is well calibrated: among people with a risk score of 7:

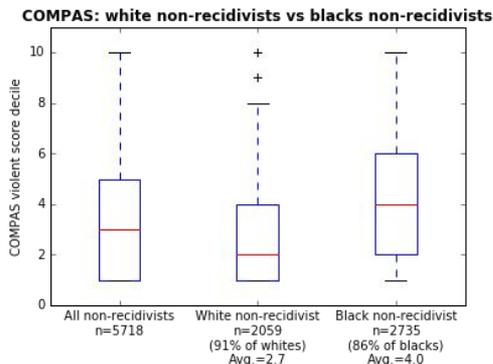
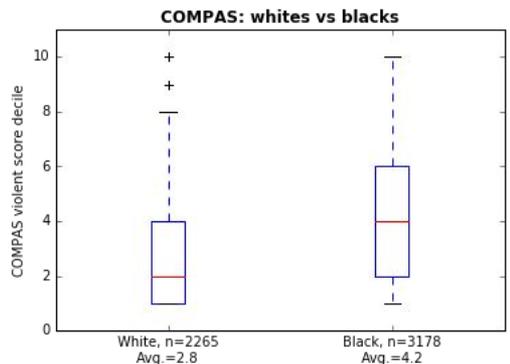
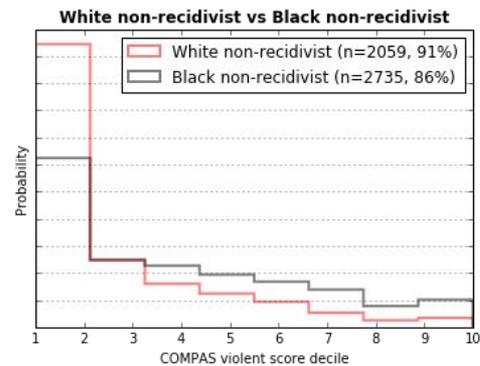
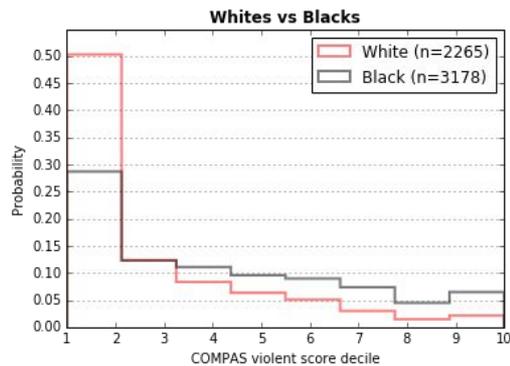
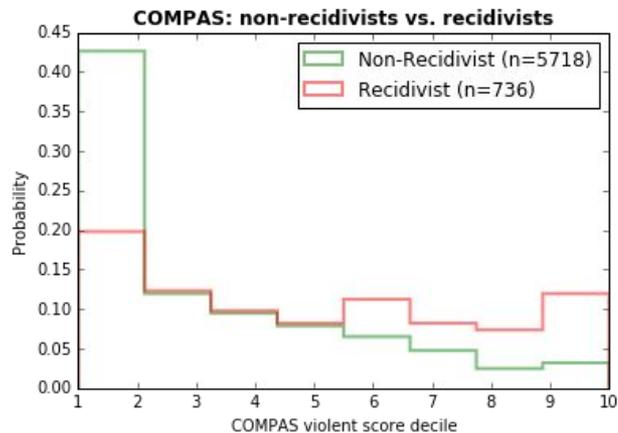
60% of blacks reoffended and

61% of whites reoffended

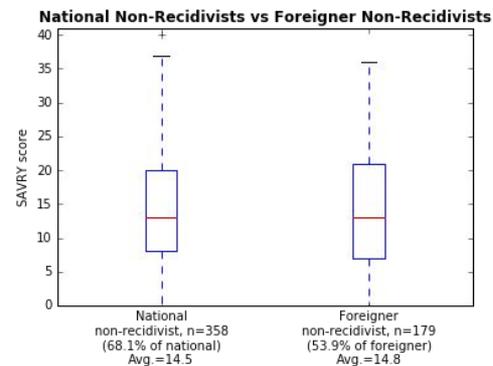
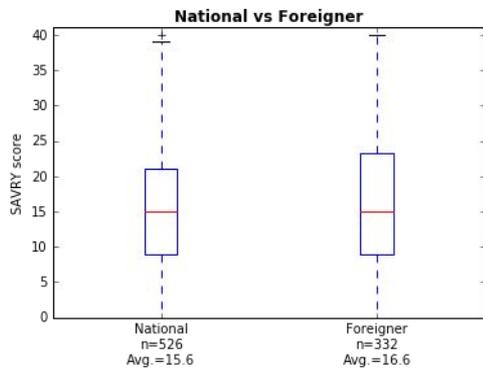
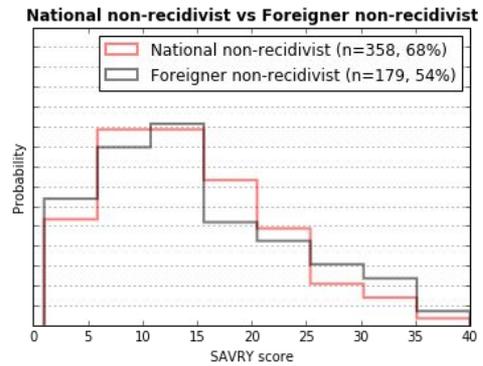
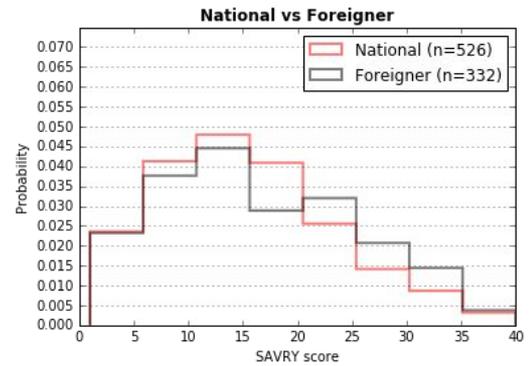
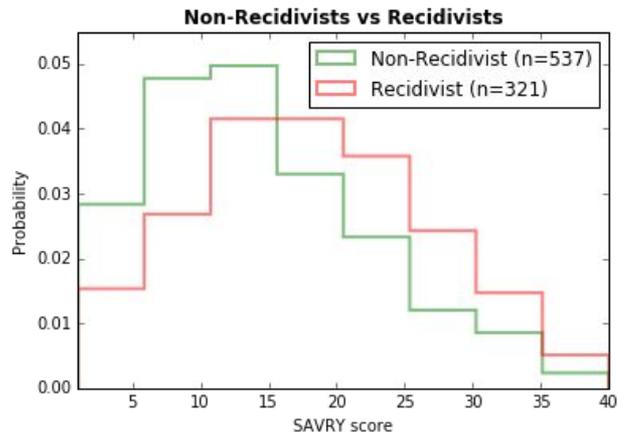
Source: Corbett-Davies et al. 2017



Disparate mistreatment / lack of E. O.



No disparate mistreatment / E.O.



Unfair rankings

	Position										top 10 male	top 10 female	top 40 male	top 40 female
	1	2	3	4	5	6	7	8	9	10				
Economist	f	m	m	m	m	m	m	m	m	m	90%	10%	73%	27%
Market analyst	f	m	f	f	f	f	f	m	f	f	20%	80%	43%	57%
Copywriter	m	m	m	m	m	m	f	m	m	m	90%	10%	73%	27%

Top-10 results for job searches in XING (a recruitment site similar to LinkedIn), for selected professions.

Algorithm-human interaction

Humans need to be able to receive explanations, and to correct outcomes.

Effective transparency does not mean source code, it means the human can understand and challenge the algorithmic decisions.

The screenshot displays the RISCANVI system interface. At the top, user information is shown: 'Intern/a: [redacted] [21003544]', 'edat: 41 class:3243 ob.', 'Disp. a Obert 1 BCN des de 17/12/2009', and 'ubicació: Z864'. Below this, a summary bar shows '21631', 'Tipus: Completa', 'Risc: [Tots]', and 'Centre: CP Obert 1 de Barcelona'. The main area contains a table with columns: 'Tipus de Risc', 'Valoració', 'Motiu', 'Correcció', 'Usuari', and 'Data'. The table lists four risk types: 'Viol. autodirigida' (Mig), 'Viol. intra-institucional' (Mig), 'Reincid. violenta' (Mig), and 'Trenc. condemna' (Baix). A red oval highlights the 'Reincid. violenta' row, which has a 'Motiu' of '- Hi aspectes biogràfics que difícilment es poden avaluar' and a 'Correcció' of 'Alt' by user 'JU21CSI' on '20/11/2009'. A yellow box with the text 'Left: system' and 'Right: corrections by expert' is overlaid on the table, with red arrows pointing to the 'Tipus de Risc' and 'Correcció' columns respectively. At the bottom, there are buttons for 'Validar resultats' and 'Cancel·lar resultats'.

Tipus de Risc	Valoració	Motiu	Correcció	Usuari	Data
Viol. autodirigida	Mig				
Viol. intra-institucional	Mig				
Reincid. violenta	Mig	- Hi aspectes biogràfics que difícilment es poden avaluar	Alt	JU21CSI	20/11/2009
Trenc. condemna	Baix				

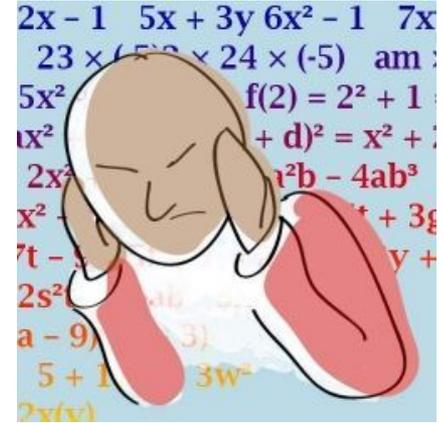
RISCANVI, 2013 (Generalitat Cat.)
<http://slideplayer.es/slide/7242758/>

A personal opinion on transparency ...

Many "customers" of algorithmic decision making systems do not value transparency.
They want certainty, not doubts.

This is a perverse incentive for
developers/providers

Numbers, plots, charts, suggest objectivity
Lack of math literacy becomes problematic



We need good evaluation frameworks

How can this be fixed? Improving math literacy and using evaluation frameworks that integrate multiple dimensions

In addition to accuracy: "dollars saved, lives preserved, time conserved, effort reduced, quality of living increased" and respect to privacy, fairness, accountability, transparency

Conclusion

To the extent that algorithms can engage in disadvantageous differential treatment that leaves people of a socially salient group worse-off, based on statistical information
... algorithms can discriminate

Current research looks at trade-offs of utility and fairness and at mechanisms for mitigating unfairness