# Fairness in Algorithmic Decisions

**Thierry KIRAT**, Directeur de recherche au CNRS –
Professeur attaché à l'Université PSL
Université Paris Dauphine-Paris Sciences et Lettres
Graduate Program Data Science – PSL Preparatory Week,
7/09/2022

# Summary

**PART 1. Plurality of fairness definitions and metrics**

**PART 2. Fairness as a technical <u>and</u> legal problem**

**PART 3. Indirect group discrimination: Disparate Impact (USA-EU)**

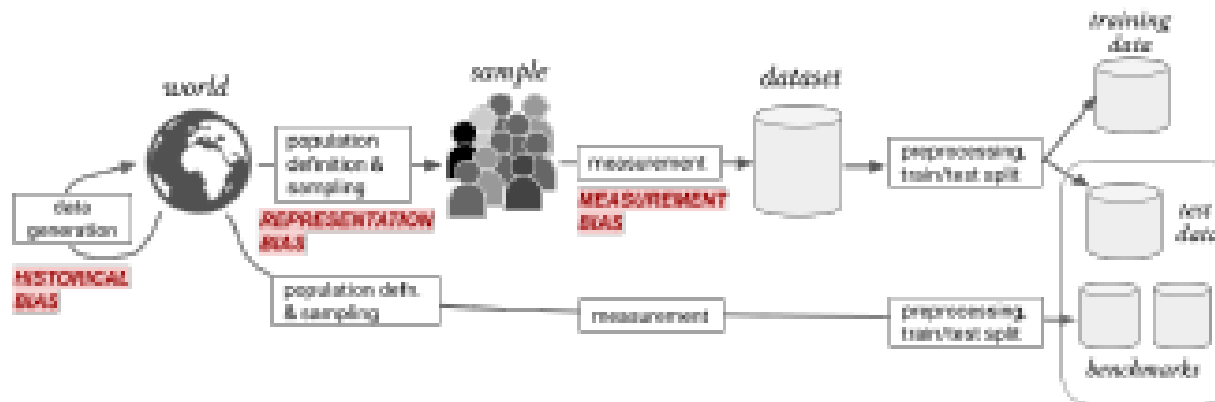**PART 4. Towards fairness metrics in line with legal provisions (USA-EU)**

**PART 5. Explainability of algorithmic decisions**

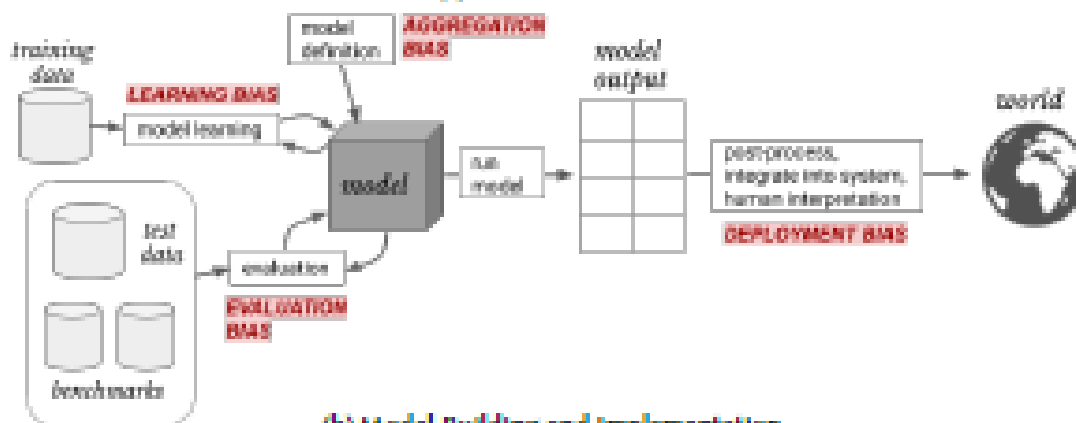# PART 1. Plurality of fairness definitions and metrics

- Sources of bias in Machine Learning

- Fairness metrics

- Illustrations of incompatibilities between fairness metrics
  - Credit applicants scoring - access to loans
  - Risk of recidivism (COMPAS)

(a) Data Generation

(b) Model Building and Implementation

# Sources of Bias (2)

1.  **Historical bias** : occurs when the world as it is leads a model to produce outcomes that are not wanted

2.  **Representation bias**: occurs when certain parts of the input space are underrepresented

3.  **Measurement bias** : occurs when proxies are generated differently across groups, or the granularity(or quality of data) varies across groups...

# Sources of Bias (3)

4. **Aggregation bias** : occurs when a one-size-fits-all model is used for groups with different conditional distributions P (X | Y)

5. **Evaluation bias** : occurs when the evaluation and/or benchmark data for an algorithm doesn't represent the target population

| Notion | Définition |
|---|---|
| **Statistical Parity (or demographic Parity)** | Aims to ensure that the fraction of people from group A who receive a particular outcome is the same as the fraction of group A of the population |
| **Conditional Statistical Parity** | Aims to equalize the outcomes between different groups, conditioned on some factors |
| **Equal opportunity** | The protected and unprotected groups should have an Equal True Positive Rates |
| **Calibration** | A score S=S(x) is well-calibrated if it respects the same likelihood of an outcome irrespective of the individual's group membership |
| **Equalized Odds** | Equality of success odd (p) and fail odd (1-p) between protected and unprotected groups |
|  | The protected and unprotected groups should have an equal True Positive and Negative Rates |

# Plurality of fairness metrics

See :

**Dooa Abu Elyoues**, "Contextual Fairness, Contextual Fairness: A Legal and Policy Analysis of Algorithmic Fairness" (September 1, 2019). Journal of Law, Technology and Policy, forthcoming

**Arvind Narayanan**, "Tutorial: 21 fairness definitions and their politics", March, 2018
https://www.youtube.com/watch?v=jIXIuYdnyyk

Table 1: Notions of fairness and summary of their corresponding legal mechanisms.

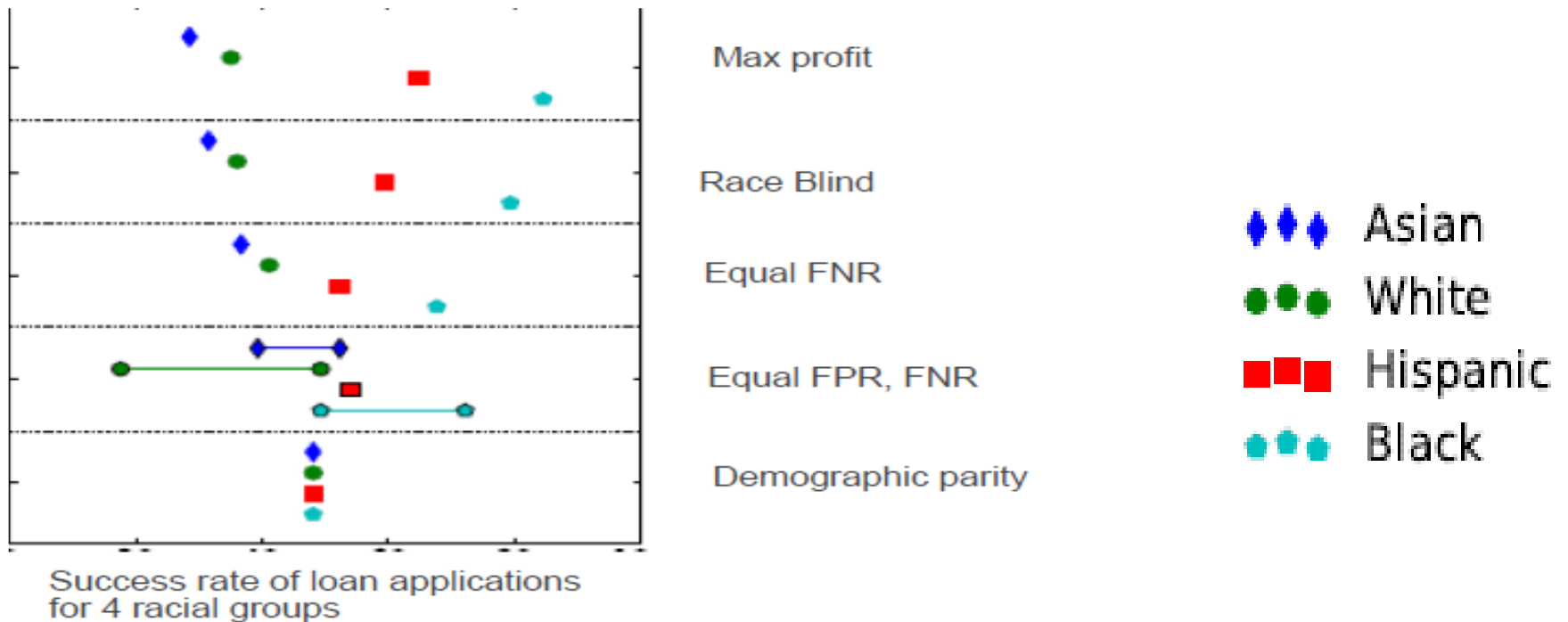| Notion | Sub-notion | Corresponding Legal Mechanism |
|---|---|---|
| Individual Fairness | The unaware approach | Equal opportunity as colorblindness |
| | Fairness through awareness | Equal opportunity based on similarities, and levels of scrutiny |
| Group fairness | Decoupling | Affirmative action (as separate but equal) |
| | Statistical or conditional parity | Affirmative action (preferably through critical diversity) |
| | Equal opportunity | Affirmative action (as equal opportunity) |
| | Equalized odds | Achieving equality by equalizing the false positive and false negative errors |
| | Calibration | Achieving equality by statistical significance |
| | Multicalibration | Achieving equality by statistical significance, and accounting for intersectionality |
| Causal Reasoning | Counterfactual fairness | Due process |

# The Prediction Problem

|  |  | True condition | | Prevalence $= \frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$ | Accuracy (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$ |
|---|---|---|---|---|---|
|  | Total population | Condition positive | Condition negative | | |
| **Predicted condition** | Predicted condition positive | **True positive** | **False positive**, Type I error | Positive predictive value (PPV), Precision = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Predicted condition positive}}$ | False discovery rate (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Predicted condition positive}}$ |
|  | Predicted condition negative | **False negative**, Type II error | **True negative** | False omission rate (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Predicted condition negative}}$ | Negative predictive value (NPV) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Predicted condition negative}}$ |
|  |  | True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$ | False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$ | Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$ | Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$ |
|  |  | False negative rate (FNR), Miss rate $= \frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$ | Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$ | Negative likelihood ratio (LR−) $= \frac{FNR}{TNR}$ | $F_1$ score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$ |

# Fairness metrics: incompatibilities

See : Hardt, Price & Srebro, Equality of opportunity in machine learning, 2016 : access to bank credit by origin (FICO dataset, USA)

-> score scale : increasing risk of default



Success rate of loan applications for 4 racial groups

Max profit

Race Blind

Equal FNR

Equal FPR, FNR

Demographic parity

Asian
White
Hispanic
Black

# Fairness metrics: incompatibilities

**Tutorial:**

Martin Wattenberg, Fernanda Viégas, and Moritz Hardt, *Attacking discrimination with smarter machine learning* (companion to Hardt, Price & Srebro, "Equality of opportunity in machine learning", 2016)
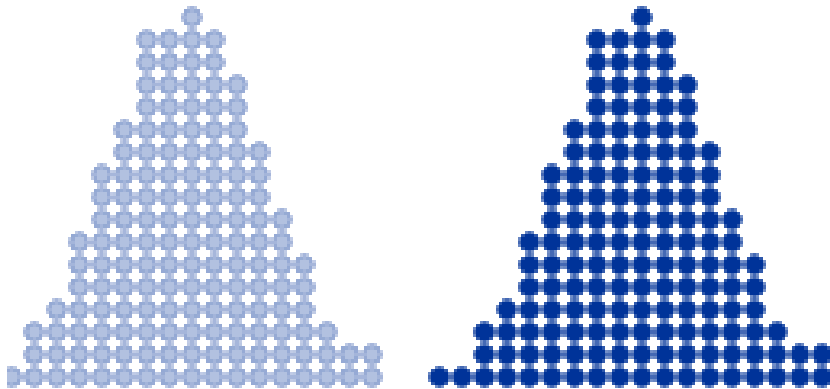
https://research.google.com/bigpicture/attacking-discrimination-in-ml/

# Attacking discrimination with smarter machine

**Credit - Risk of default**

Ideal : separating good and bad borrowers (left figure)

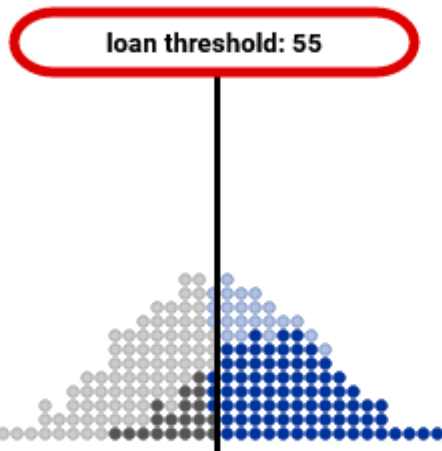In practice the two groups can't be clearly separated (the FP problem….) (right figure)



would default on loan    would pay back loan

# *Attacking discrimination with smarter machine*

**Simulation # 1: Group Unaware - holds all groups to the same standard (same threshold on the risk-score scale)**

Both groups have the same threshold, but the orange group has been given fewer loans overall. Among people who would pay back a loan, the orange group is also at a disadvantage (FN).



**Blue Population**

0  10  20  30  40  50  60  70  80  90  100

loan threshold: 55

denied loan / would default | granted loan / defaults
denied loan / would pay back | granted loan / pays back

**Orange Population**

0  10  20  30  40  50  60  70  80  90

loan threshold: 55

denied loan / would default | granted loan / defaults
denied loan / would pay back | granted loan / pays back

# *Attacking discrimination with smarter machine*

**Simulation # 2: <u>Demographic parity</u> - If the goal is for the two groups to receive the same number of loans, then a natural criterion is demographic parity, where the bank uses loan thresholds that yield the same fraction of loans to each group. -> the "positive rate" is the same across both groups (37% of applicants obtain in loan in each group)**

<span style="color:red">**The number of loans given to each group is the same, but among people who would pay back a loan, the blue group is at a disadvantage.**</span>
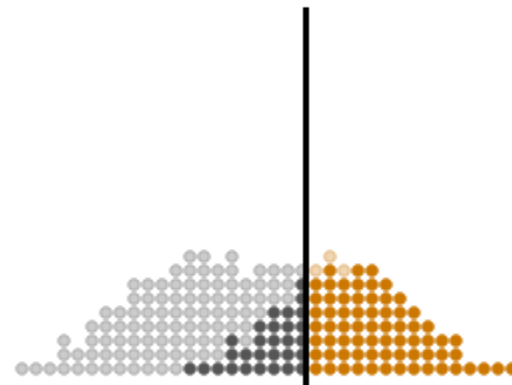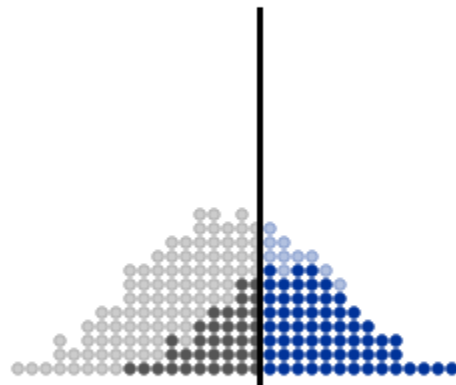


Blue Population

0  10  20  30  40  50  60  70  80  90  100

loan threshold: 60

denied loan / would default — granted loan / defaults
denied loan / would pay back — granted loan / pays back

Orange Population

0  10  20  30  40  50  60  70  80  90  100

loan threshold: 52

denied loan / would default — granted loan / defaults
denied loan / would pay back — granted loan / pays back

*Attacking discrimination with smarter machine*

**Simulation #3 : Equal opportunity : The constraint is that of the people who can pay back a loan, the same fraction in each group should actually be granted a loan -> the true positive rate is identical between groups**

**Among people who would pay back a loan, blue and orange groups do equally well.**

Blue Population

Orange Population

| 0 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |

loan threshold: 59

| 0 | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |

loan threshold: 53

denied loan / would default · granted loan / defaults
denied loan / would pay back · granted loan / pays back

denied loan / would default · granted loan / defaults
denied loan / would pay back · granted loan / pays back

# Predictive criminal justice (USA): COMPAS

*Correctional Offender Management Profiling for Alternative Sanctions, developed by the Northpointe Compagny (now Equivant)*

Three-levels scores :

*1. Pretrial Release Risk scale :* Risk of not appearing in court and/or committing crimes between indictment and criminal sanction

*2. General Recidivism Risk Scale – GRRS:* Risk of re-offending after release. The scale takes into account the individual's criminal history and accomplices, drug use, juvenile delinquency…

*3. Violent Recidivism Risk Scale – VRRS:* Risk of violent recidivism after release. Takes into account: the individual's history of violence, frequency of lawlessness, school problems, age of first arrest…

# COMPAS : ProPublica critics



Prediction Fails Differently for Black Defendants

|  | WHITE | AFRICAN AMERICAN |
|---|---|---|
| Labeled Higher Risk, But Didn't Re-Offend | 23.5% | 44.9% |
| Labeled Lower Risk, Yet Did Re-Offend | 47.7% | 28.0% |

Overall, Northpointe's assessment tool correctly predicts recidivism 61 percent of the time. But blacks are almost twice as likely as whites to be labeled a higher risk but not actually re-offend. It makes the opposite mistake among whites: They are much more likely than blacks to be labeled lower risk but go on to commit other crimes. (Source: ProPublica analysis of data from Broward County, Fla.)

# Chouldechova : COMPAS scores are calibrated



Alexandra Chouldechova, "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments", ArXiv: 1610.07524v1, 24 oct 2016

See also : Julia Dressel & Hany Farid, " The accuracy, fairness, and limits of predicting recidivism", Science Advances, 17 january 2018

# COMPAS : not discriminatory ?

**Overall accurary equality :** The overall accuracy of the COMPAS label is the same, regardless of race

```
race
African-American    0.638258
Caucasian           0.669927
dtype: float64
```

**Predictive Positive Value** : The likelihood of recidivism among defendants labeled as medium or high risk is similar, regardless of race

```
race
African-American    0.629715
Caucasian           0.591335
Name: two_year_recid, dtype: float64
```

**Calibration :** For any given COMPAS score, the risk of recidivism is similar, regardless of race

Farhan Rahman, COMPAS Case Study: Fairness of a Machine Learning Model,, Sep 7, 2020·
https://towardsdatascience.com/compas-case-study-fairness-of-a-machine-learning-model-f0f804108751

# Insights into political philosophy

See : Reuben Binns, "Fairness in Machine Learning: Lessons from Political Philosophy", Proceedings of Machine Learning Research 81:1–11, 2018 Conference on Fairness, Accountability, and Transparency

" *'fairness' as used in the fair machine learning community is best understood as a placeholder term for a variety of normative egalitarian considerations*"

**Issue** : examine how egalitarian norms might provide an account of why and when algorithmic systems can be considered unfair

# Political philosophy : utilitarism

- **To satisfy one fairness criteria one must sacrifice some utility**. Ex : minimize the false positive rate of criminal reoffenders (high risk + no reoffending)=> risk of reducing public security level (release of truly high risk inmates) (Narayanan) ; Conversely if utility criteria prevails (e.g. public security) false positive rate are to be kept at a high level

- **Society values differently false positive and false negative** (Abu Elyounes 2019)

- Corbett-Davis, Pierson, Feller Avi, Sharad, « Algorithmic decision making and the cost of fairness », 2017 : utilitarist-inspired analysis of COMPAS -> "**there is tension between reducing racial disparities and improving public safety**". Incompatibility between maximisation of public security and equal treatment of individuals whatever their race. Algorithmic fairness is a problem of constrained optimisation (in reference to diverse fairness metrics : statistical parity, predictive equality, conditional statistical parity). The optimal algorithm that results require applying multiple, race-specific thresholds to individuals' risk scores.

- **Cost-benefit** approach : does the marginal social benefit of additional fairness (e.g. less group discrimination) outweigh the marginal cost ? (see Corbett-Davis & al. 2017).

# Political philosophy : egalitarianism (1)

## Variants of egalitarianism

**Welfarism** (Cohen 1989) :

a) *Hedonic welfare* : "welfare as enjoyment, or, more broadly, as a desirable or agreeable state of consciousness". Limit : metrics of welfare

b) *Welfare as preference satisfaction (or fulfillment)* : "preferences order states of the world, and where a person's preference is satisfied if a state of the world that he prefers obtains, whether or not he knows that it does". Limit : heterogeneity of preferences and resource needs (if Peter prefers champaign and Allan prefers beer, Peter needs more resources to fulfill his preference than Peter)

c) *Equality of opportunity for welfare* (Richard Arneson).

**Resources-based** (Dworkin) : a society is just it holds individuals responsible for their decisions and actions, but not for circumstances beyond their control, such as race, sex, skin-color, intelligence, and social position. Unequal distribution of resources is considered fair only when it results exclusively from the decisions and intentional actions of those concerned

**Primary social goods** (Rawls) : those that the citizens need as free people and as members of the society : civil rights, political rights, liberties, income and wealth, the social bases of self-respect, etc.

**Capabilities** (Sen) : Capabilities are the doings and beings that people can achieve if they so choose, such as being well-nourished, getting married, being educated, and travelling; functionings are capabilities that have been realized.

# Political philosophy : egalitarianism (2)

## Implications for AI

(1) « egalitarian norms might provide an account of why and when algorithmic systems can be considered unfair » (Binns, 2018, p. 6)

(2) diversity of egalitarian norms implied in algorithmic decisions

- loan decision, insurance : impact the distribution of ressources (***distributive harm***)
- exclusion from a social network : impact the capabilities or welfare (***representative harm***)

(3) Welfarism : preferences fulfillment => some individuals may prefer a racially-segregated society (requires a moral judgment about which preferences are to taken into account or excluded)

- Rawls : Maximim principle + veil of ignorance : the criteria of social justice requires a social contract which have to be set-up by individuals ignoring their future position (veil of ignorance). A just society benefits the least advantaged (maximin principle).

- Sen : a just society benefits the poorer (strengthening the poorer' capabilities).

# PART 2. Fairness as a technical **and** legal problem (1)

**From fairness as a technical problem and Fair design…**

- mathematical methods for correcting sources of bias and unintended consequences
- ethical considerations only pose technical, mathematical difficulties that can be resolved without recourse to considerations outside the world of AI research…..

**To the correspondances between legal and technical concepts…**

- the completion of responsible algorithms necessarily involves the collaboration of data science, law and public policy

….

**And to fairness models fitting legislation and jurisprudence**

- Recent major works : Abu Elyounes ; Xiang ; Wachter & al. ; Hacker ; Kirat, Tambou, Do,Tsoukias

# PART 2. Fairness as a technical <u>and</u> legal problem (2)

- Most research on fair AI have the US legal system as a (more or less implicit) background

- Key Issue; designing and modelling algorithmic systems in line with the legal/institutional context

| | Anti-discrimination Law (American law) | Machine Learning |
|---|---|---|
| **Procedural fairness** | to arrive at just outcomes through iterative processes and the close examination of the set of governance structures in place to guide individual human decision-making *Focus on processes & the system surrounding the algorithm and its use* | refer to identifying the input features that lead to a particular model outcome, as a proxy for the "process" through which the model makes its prediction *Focus on outcomes & specifics of the algorithm itself* |
| **Discrimination** | Federal laws provide anti-discrimination protections in housing, employment, and other domains. The federal acts primarily define discrimination though motive, evidenced intent of exclusion, and causality, rather than simply outcomes. | Often presented as an unjust correlation between protected class variables and some metric of interest, such as outcomes, false positive rates, or a similarity metric |
| **Protected Class/Sensitive Attribute** | less commonly measured attributes can also be considered, such as sexual orientation, pregnancy, and disability status *Aware of the implementation of law & possibility of "reversal" of the benefits of anti-discrimination law (Ricci v. DeStefano, 2009)* | Protected attributes are presented as recorded or visible traits that should not factor into a decision, such as age, race, or gender. *Unaware of the implementation of law* |
| **Anti-classification and anti-subordination** | *Anticlassification* (or antidifferentiation principle) : holds that the government may not classify people either overtly or surreptitiously on the basis of a forbidden category such as their race *Antisubordination* (or equal citizenship, anti-caste) theorists contend that guarantees of equal citizenship cannot be realized under conditions of pervasive social stratification and argue that law should reform institutions and practices that enforce the secondary social status of historically oppressed groups (Baldin & Siegel, 2003) | *Anticlassification*: classifications based on protected class attributes are impermissible. ML fairness community is actually quite familiar with this concept ("fairness through unawareness") Numerous works explicitly presenting anti-classification as a potential fairness objective *Antisubordination* is rarely called out as a motivation in ML fairness literature |

|  | Anti-discrimination Law (American law) | Machine Learning |
|---|---|---|
| **Affirmative action** | Landmark affirmative action cases have concluded that schools seeking to increase racial diversity cannot use racial quotas or point systems.<br><br>Schools have dealt with this conundrum through greater opacity, seeking to be race conscious without making explicit how race factors into admissions decisions | The ML fairness community has articulated a goal of 'fair affirmative action,' which guarantees statistical parity (i.e., the demographics of the set of individuals receiving any classification are the same as the demographics of the underlying population), while treating similar individuals as similarly as possible" and understands affirmative action to be cases in which we explicitly take demographics into account |
| **Disparate treatment and disparate impact** | *Disparate treatment* : the key legal question is whether the alleged discriminator's actions were motivated by discriminatory intent<br><br>*Disparate impact*: disproportionate outcomes between sub-groups is illegal if intentional. In case of intentionality : liability incurred<br><br>*Key issue : intentionality* | *Disparate treatment* is often explained as making use of the protected attribute in the decision-making process -> avoiding the use of protected class variables in debiasing techniques<br>*Disparate impact* is understood as when outcomes differ across subgroups (even unintentionally) -> group fairness formulations<br><br>*Algorithm cannot possess intent by itself* |

# PART 3. Indirect Group Discrimination: Disparate Impact

- Forms
- Disparate impact : legal dimension, at the crossroad between ML and law

# Forms

- Direct (intentional) discrimination

- Indirect (unintentional) discrimination : Disparate impact

- Individual versus group discrimination

=> Here : focus on group discrimination/disparate impact

$$\text{DI} = \frac{P(Y=1) \mid (S=0)}{P(Y=1) \mid (S=1)}$$

# Disparate Impact

## Legal concept :

**Civil Rights Act 1964**

Title VII : *prohibits employment discrimination based on race, color, religion, sex and national origin*

Title VI : *No person in the United States shall, on the ground of race, color, or national origin, be excluded from participation in, be denied the benefits of, or be subjected to discrimination under any program or activity receiving  Federal financial assistance.*

**Age Discrimination in Employment Act, 1967, Fair Housing Act, 1967**

## Equal Employment Opportunity Commission (EEOC) :

80% Rule (Uniform Guidelines on Employee Selection Procedures 1978) : ratio of selection rates across groups.
Ratio < 0.8 : presumption of discrimination

## Case-Law (federal courts) : from an expansion of DI doctrine in the benefit of plaintiffs in the 70's to a more restrictive interpretation (in the benefit of employers) since the 90's

# Disparate Impact

**Case-law (federal courts, USA). Some major rulings**

*Griggs v. Duke Power (Supreme Court,* **1971**) : the Supreme Court made a significant advance in securing civil rights for African Americans. The company in question conducted intelligence tests and required employees to have completed college in order to be promoted to higher paying positions.

*Wards Cove Packing v. Atonis* (*Supreme Court,* **1989**) the Court placed a very important restriction on disparate impact actions by establishing the evidentiary rule that the plaintiff must establish (a) what precisely defined practice or rule caused the indirectly discriminatory impact, and (b) that the employer refused to implement practices or rules that would have satisfied the plaintiff's grievances. In addition, the accused company may argue that the rule or practice that caused the disproportionate impact was justified by the necessity of business.

*Ricci v, DeStefano (Supreme Court,* **2009**)**:** the Mayor of New Haven, Connecticut, cancelled a competition for the promotion of the city's firefighters because the success rate of white firefighters was twice that of African Americans. The court ruled in favor of the successful firefighters; it faulted the Mayor for canceling the competition without showing that its continuation could expose him to disparate impact liability.

# Comparison between USA and European Union (statute law and case law)

|  | **United States** | **E. U.** |
|---|---|---|
| Main focus | Racial inequalities<br>Workers' hiring and promotion | Salarial equality between men and women |
| Part-time work | Not taken into account | Taken into account |
| Burden of proof (from the plaintiff viewpoint) | Restrictive and limiting | Not very demanding |
| Justification of rules and practices with disparate impact by employers | Business necessity benefit the employers | Business necessity : balanced approach in the EUCJ case-law |

# PART 4. Towards fairness metrics in line with legal provisions

**Recent proposals in the literature**

**UE :**

Sandra Wachter, Brent Mittelstadt & Chris Russell (« Why Fairness cannot be automated: Bridging the Gap Between EU Non-Discrimination Law and AI », 2020)

**USA :**

Alice Xiang (« Reconciling Legal and Technical Approaches to Algorithmic Bias », 2021)

# PART 4. Towards fairness metrics in line with legal provisions

Method used by these authors

1. Starting point: case-law on discrimination cases (EUJC / Federal courts)

2. Identification of major cases (principle/rule)

3. Proposals: fairness metric

# PART 4. Towards fairness metrics in line with legal provisions

**Wachter & al., Antidiscrimination european case-law (EUJC)**

- Issue : statistical proof of discrimination

- « Gold Standard » found in *Seymour-Smith* (9 February 1999) : full comparisons between disadvantaged and advantaged groups

- propose 'conditional demographic disparity' (CDD) as a standard baseline statistical measurement that aligns with the Court's 'gold standard'

Wachter, Sandra and Mittelstadt, Brent and Russell, Chris, "Why Fairness Cannot Be Automated: Bridging the Gap Between EU Non-Discrimination Law and AI" (March 3, 2020). *Computer Law & Security Review* (forthcoming), https://ssrn.com/abstract=3547922 ; http://dx.doi.org/10.2139/ssrn.3547922

# PART 4. Towards fairness metrics in line with legal provisions

**Xiang, Antidiscrimination American case-law** (Courts of Appeals & Supreme Court) :

She analyzes the extent to which technical approaches to algorithmic bias are compatible with U.S. anti-discrimination law and recommends a path toward greater compatibility

Issue raised : possibility that biased algorithms might be considered legally permissible while approaches designed to correct for bias might be considered illegally discriminatory

ML - > use of protected class variables to check (and mitigate) discrimination

US Law -> prohibits the use of protected class variables (fairness through unawareness)

Xiang, Alice, "Reconciling Legal and Technical Approaches to Algorithmic Bias" ,(January 4, 2021). *Tennessee Law Review*, Vol. 88, No. 3, 2021, https://ssrn.com/abstract=3650635

# PART 4. Towards fairness metrics in line with legal provisions

**Major case** : *Texas Department of Housing and Community Affairs v. Inclusive Communities Project, Inc*., Sup Ct, 2015 + subsequent proposed rule from the Department of Housing and Urban Development (HUD)

The Court required a "causal connection" between the decision-making process and the disproportionate outcomes

Xiang's Proposal : use of protected attributes to check an eventual discrimination

Causal connection + counterfactual :

*« In a causal framework, fairness is conceived of as the lack of a difference between the observed outcome and the counterfactual outcome where the (perception of the) individual's protected class attribute is changed….. This aligns with legal conceptions of fairness: if but for the individual's protected class, the decision would have been different, then the individual was illegally discriminated against"*

# PART 5. Explainability

1. **What does « explainability » mean?**

- Global vs. Local

- Ex ante vs. Ex post

- Technical vs. Decision process

2. **Explainability of what? Dataset, algorithm, model, outcome (decision, prediction)**

3. **Explainability for who? Expert, regulator, individual**

# Explainability as a legal obligation?

## Is it effective or practicable ?

**EU law** : GDPR, recital 71 : In any case, such processing should be subject to suitable safeguards, which should include specific information to the data subject and the right to obtain human intervention, to express his or her point of view, to obtain an explanation of the decision reached after such assessment and to challenge the decision.

**French law** :

- Loi n° 2018-493 : obligation to communicate the rules defining the processing + the main characteristics of its implementation (except if these rules are subject to secrets protected by law)
- Code des relations du public avec l'administration (CRPA, art. L. 311-3-1 : « the rules defining the processing and the main characteristics of its implementation shall be communicated by the Administration to the person concerned on request .
- CRPA, art. R. 311-1-2 : specifies the information to be provided in intelligible form.

**Constraints : commercial secret ; black box**

**A complex algorithm with very good predictive capabilities is not necessarily explainable**

**- tension between accuracy (high reliability of predictions) and explainability**
**- Counterfactual explanation?**

# Explainability of algorithmic decisions

**Counterfactual explanation**

« *You have been refused credit by the bank. Your annual income is 30,000 euros. If your income had been 40,000 euros, you would have been granted credit* ».

"In the existing literature, "explanation" typically refers to an attempt to convey the internal state or logic of an algorithm that leads to a decision. In contrast, counterfactuals describe a dependency on the external facts that led to that decision"

See Sandra Wachter, Brent Mittelstadt & Chris Russell, COUNTERFACTUAL EXPLANATIONS WITHOUT OPENING THE BLACK BOX: AUTOMATED DECISIONS AND THE GDPR, Harvard Journal of Law & Technology 2018

Table from : T. Kirat, O. Tambou, V. Do, A. Tsoukiàs, 2022, Fairness and Explainability in Automatic Decision-Making Systems. A challenge for computer science and law,. https://arxiv.org/abs/2206.03226

| E.U. legal text | Disposition | Scope of application | Practical modalities. | Who is concerned | Goals |
|---|---|---|---|---|---|
| Article 29 Working Party Guidelines on Automated individual decision-making and Profiling. | "algorithmic auditing": "testing the algorithms used and developed by machine learning systems to prove that they are actually performing as intended and not producing erroneous discriminatory or unjustified results | Without restriction | Formal verification of algorithms + security | Experts: Designers programmers + certifiers | Model improvement |
| GDPR | Article 15(1)(h): right to be informed of the existence of automated decision-making + logic involved + significance and the consequences of such processing for the data subject | Without restriction | Internal logic of the model + causality | Individuals | Transparency + causality |
|  | Article 22 : right not to be subject to a decision based solely on automated processing, including profiling, which produces legal effects concerning him or her or similarly significantly affects him or her | Not applicable if - entering into, or performance of, a contract - authorised by Union or Member State law - no legal effects or person not affected - consent, legal authorization | Concerns both automated decision (explainable) + autonomous decision (black box, not explainable) Causality | Individuals | Final decision must be taken by a human |
|  | Recital n° 71: use of appropriate mathematical or statistical procedures implement technical and organisational measures ... ensure that factors which result in inaccuracies in personal data are corrected and the risk of errors is minimized and prevents discriminatory effects | Without restriction | Pipe: Formal verification Data control Test | Experts: Designers Programmers | Model Improvement |
| Regulation (EU) 2019/1150 | Article 5 – online intermediation services - main parameters determining ranking | Without restriction | Verification ranking model | Industry: online intermediation service providers + search engines | Transparency + compliance |