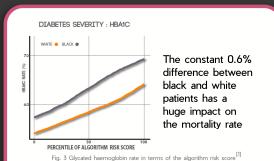


DEFINITION : by machine learnt scoring systems, we describe models using **machine learning (ML)** frameworks to assign a **score** to entities (here individuals)

BIAS IN MACHINE LEARNT SCORING SYSTEMS

- HOW IT CREEPS IN AND WHY IT PERPETUATES RACIAL DISCRIMINATION -

An analysis through two characteristic examples : The Health and Judicial Systems in the United States



I. HOW BIAS CREEPS IN THE MODEL

- The data is often unrepresentative of the context in which it will be applied, marginalising minorities
 - Often reflects existing prejudice

Data preparation and choosing the right variables

- Inputs highly influence the model's prediction accuracy
 - Really difficult to measure impacts on the model's bias
 - We can obviously choose not to add inputs such as someone's origins but it **isn't sufficient !**

→ Many metrics such as wage are **not race-blind metrics** due to historical or socio-economic causes and can be a proxy for race.

II. WHY BIAS IN ML IS HARD TO CORRECT

Unfitting models

- Contrary to “normal algorithms” machine learning algorithms are hard to interpret → we can’t know where the biases come from

Defining

- Centuries long debated subject by philosophers and lawmakers alike
 - Many definitions (statistical parity, equality of false negatives..) some being mutually exclusive^[1]
 - Even harder to define fairness mathematically (required in ML)

A social context problem : the "portability trap"^[10]

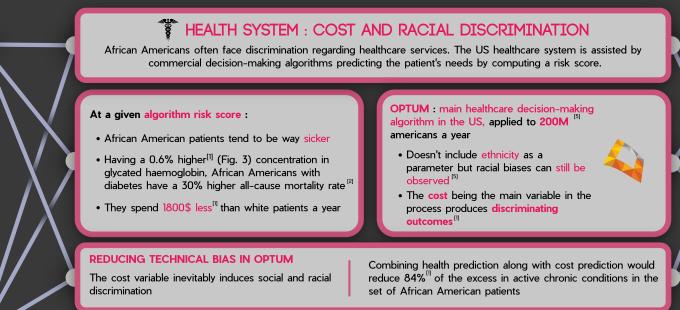
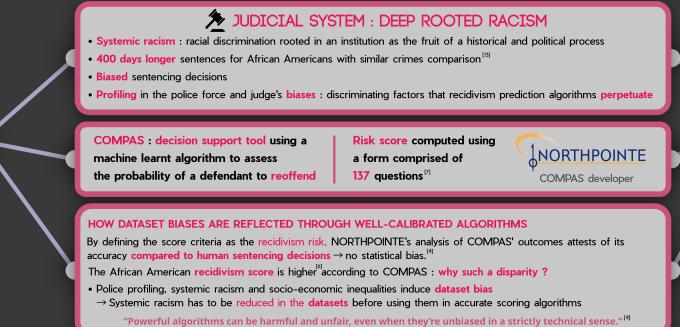
- Reusable ML scoring systems ignore social context → can't be developed and implemented in two places (ethnic backgrounds vary)

Model Optimality versus Personal Fairness

- The most important metric is often the error rate and including debiasing in the models would reduce their effectiveness (pareto curve) [2]

Feedback loop

- The algorithms aren't dynamically implemented : they set the statistics in stone → are only a representation of society at one point in time and perpetuate these statistics^[13]



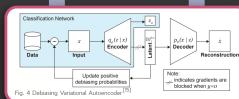
TECHNICAL SOLUTIONS

General advice

- Datasets more representative of the target population are required
 - Algorithms to mitigate bias → divided up in 3 types tackling respectively^[3]:
 1. Pre-processing
 2. In-processing
 3. Post-processing

One of the solutions : a semi-supervised VAE

 - A **Variational Autoencoder** learns the latent structure of the dataset without manual input → **no further human bias** (it reduces each data point to a restrained number of defined attributes)
 - With this latent representation, the training data is resampled to increase the probability of rare data and diminishing that of over-represented attributes → our "improvement" : the model could use a **disentangled-VAE**^[4] to further understandability and accuracy.



SOCIAL UNDERSTANDING OF FAIR DECISION-MAKING ALGORITHMS

- Fairness doesn't come with an absolute definition as multiple interpretations can be given on the same already socially conditioned data computed by an algorithm
 - It is thus difficult to translate a social reality into a reliable and absolutely fair mathematical model
 - However, reducing racial discrimination within social infrastructure would enhance the integrity of the datasets, thus allowing well-calibrated algorithms to achieve accurate results regarding a chosen

TOWARDS A GLOBAL ACKNOWLEDGEMENT :
MACHINE LEARNING IN THE EUROPEAN JURISDICTION

- The European Union considers decision-making algorithms as a relevant issue to deal with
 - The ethical european charter^[6] for the use of AI in Justice suggests 5 principles that public and private actors have to respect, in particular, the non-discrimination principle that warns about the use of "sensitive data" in sentencing decisions
 - Since 2018 more rights for users regarding automated or intelligent systems

