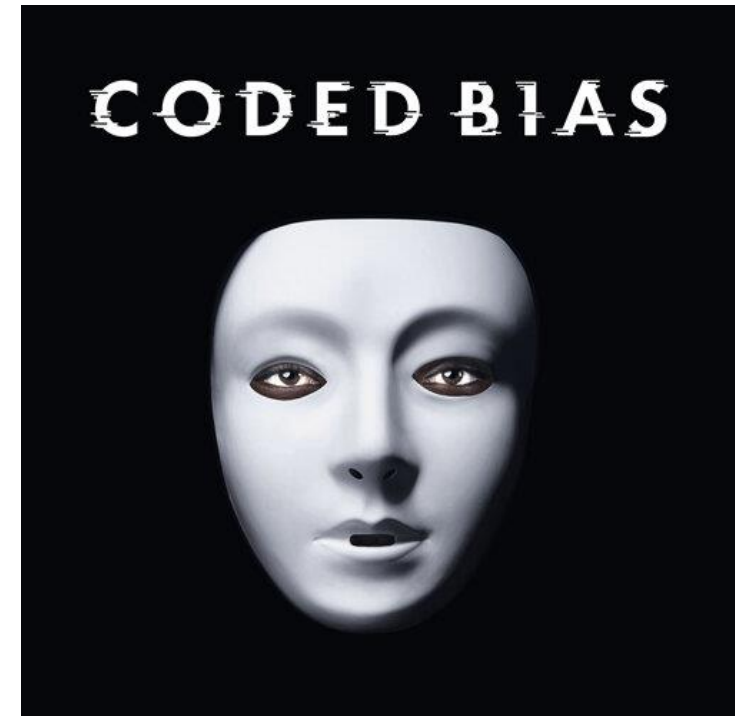# BIAS E FAIRNESS NEL MACHINE LEARNING

Giordano d'Aloisio

# Il problema del bias

- Nel 2018 una ricercatrice dell'IMT stava studiando i sistemi di riconoscimento facciale di Amazon

- Questi sistemi non erano in grado di riconoscere il suo volto

- All'inizio pensò fosse un errore del sistema

- Ma poi indossando una maschera bianca notò che il sistema era in grado di riconoscerla

Quindi il sistema non era in grado di riconoscere donne non bianche

# Un altro esempio...

- Diversi giudici negli Stati Uniti hanno utilizzato per anni un sistema di intelligenza artificiale per decidere se liberare o meno un condannato

- Questo sistema prevedeva la possibilità che un condannato avrebbe ricommesso un crimine nei prossimi due anni

- Dopo uno studio attento del sistema è stato dimostrato che questo algoritmo date due persone con stesse caratteristiche, ma di etnia diversa, forniva una minore probabilità di recidiva alla persona bianca

Quindi il sistema favoriva sistematicamente le persone bianche solo in base alla loro etnia

# Definiamo meglio il concetto di Bias e Fairness

- **BIAS:** sistematico favoritismo o discriminazione di individui da parte di un algoritmo sulla base di alcune loro caratteristiche (esempio il sesso o l'etnia)

- **FAIRNESS:** assenza di discriminazione o favoritismo da parte di un algoritmo

# Il concetto di bias non è così semplice...



- Più di 23 definizioni di bias esistono oggi in letteratura

# Il bias e la fairness possono essere misurati...



○ Più di 15 metriche di fairness diverse esistono oggi in letteratura

# Il bias può essere mitigato...
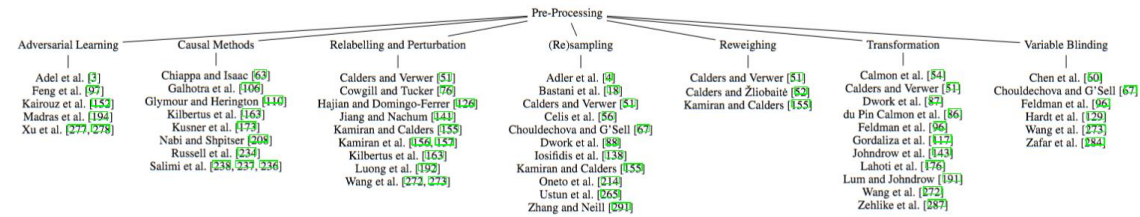


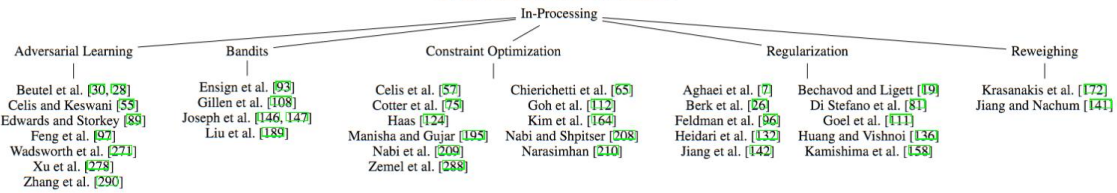Figure 3: Pre-processing Methods
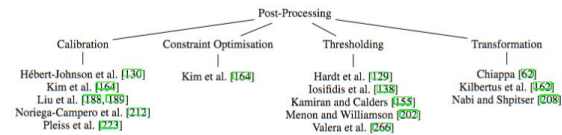
Figure 4: In-processing Methods

Figure 5: Post-processing methods

# Come possiamo migliorare questa situazione?

Guidando l'utente nella selezione di definizioni, metriche e metodi appropriati in base al contesto

# MANILA in SoBigData RI

# DOMANDE? 🤔