☆ Recap : Unsupervised + reinforcement learning

☆ other topic :   Baycsism   methods

                   Decision   trees

                   Boosting + bagging

                   graphical   models

                   Online   learning
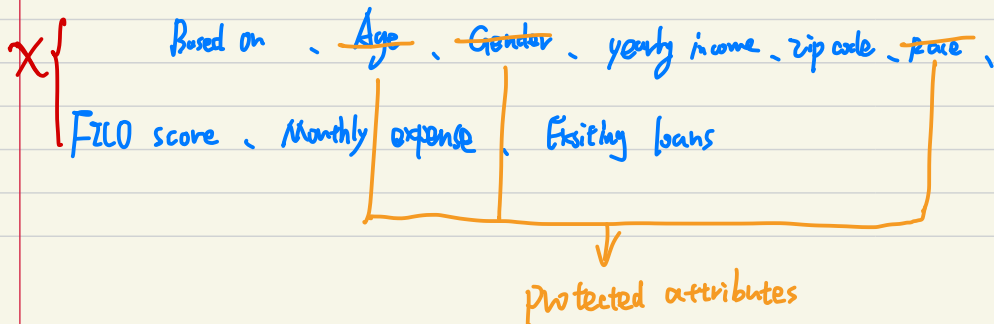
                   Semi-Supervised   learning

★   Fairness   in   ML

✹   Security   in   ML

◢   Ethics   in   ML
      伦理
―――――――――――――――――――――――――――――――――

Fairness

Example :  Decide  who  to  give  loans  } $y$

$x$ {      Based on  .  Age  ,  Gender  .  yearly income . zip code . race .

       FICO score . Monthly expense . Existing loans


protected attributes

$\longrightarrow$ classifier can be biased

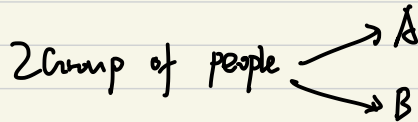Attempt: Exclude protected features

"Fairness through exclusion"

or "Fairness through unawareness"

Issue: Features can be correlated

∴ available features can be heavily linked to protected attributes.

Statistical parity

2 Group of people $\begin{array}{l} \rightarrow A \\ \rightarrow B \end{array}$

$Pr(\text{Loan} = \text{true} \mid \text{person} \in A)$
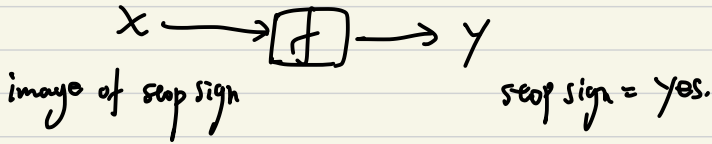
$= Pr(\text{Loan} = \text{true} \mid \text{person} \in B)$
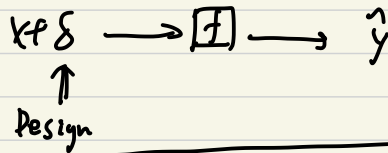
$\Rightarrow$ Loan $\perp$ Group

$\Uparrow$

"statistically independent"

other fairness objectives $\rightarrow$ Equal odds, Counter factual parity, ...

# Security & Robustness of ML

$$x \longrightarrow \boxed{f} \longrightarrow y$$

image of stop sign          stop sign = yes.

## Adversarial Settings
~~settings~~

$$x+\delta \longrightarrow \boxed{f} \longrightarrow \hat{y}$$
$$\uparrow$$
Design

---

Standard setting: $\frac{1}{2} \sum_{i=1}^{n}$ Crossentropy $(y_i, <w, x_i>) = L(w)$

$$\min \ L(w)$$

---

$$\max_{\delta} \ \text{Crossentropy}(y, <w, x+\delta>)$$

$$\|\delta\|_{\infty} \leq \varepsilon$$

"adversarial learning"

like adding an irrelevant image to data-set?