

Interpretability - part 2

Societal impact & approaches

ex: medical diagnosis

"Weapons of maths destruction" by Cathy O'Neil

↳ Blackbox software (at large scale) for important matters

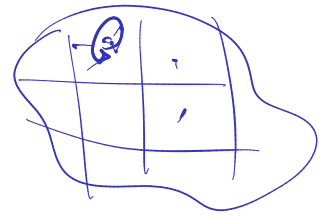
- ↳ hiring jobs - CV
- ↳ firing (bis school)
- ↳ loans (bank)

- behaving arbitrarily / stochastic
- no feedback / questioning possible
- large scale ⇒ arbitrariness nightmare
- illegal criteria / proxy

ex: COMPAS: recidivism prediction
↳ high bias

• self-reinforcement (self-fulfilling prophecy): police patrols

⇒ think twice about the impact of your algorithms before deploying them



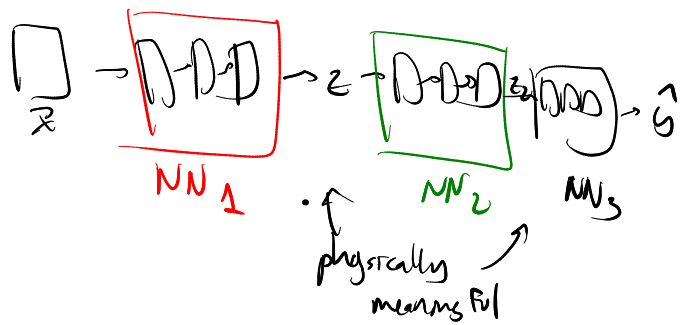
Be responsible & careful

⇒ AI ethics → Montreal declaration for responsible AI

"FAT"-ML: Fairness, Accountability & Transparency

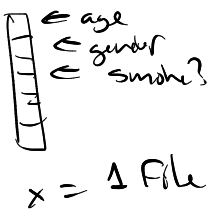
Interpretability by design:

"X-AI"
-
Explamable

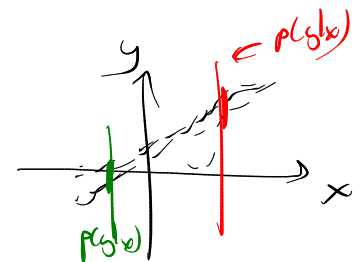
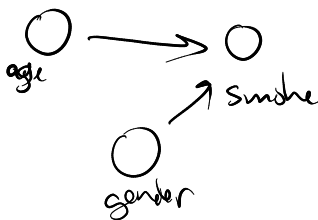


of website

Causality

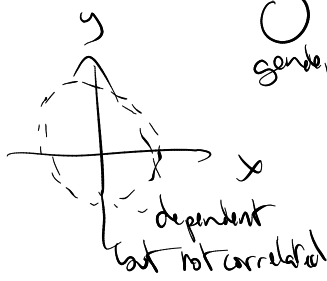


~~Correlations?~~
Causation



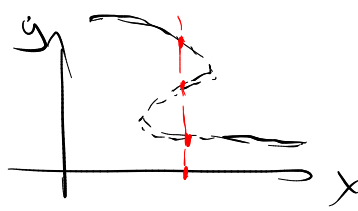
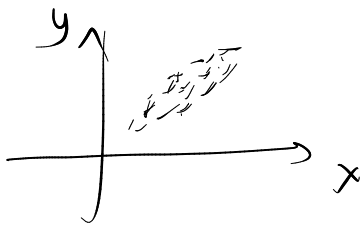
$$x^2 + y^2 = 1$$

$$y = \sqrt{1 - x^2}$$



independent $p(y|x) = p(y)$ \Rightarrow correlated
 $y \perp x$

$$E[(x - \bar{x})(y - \bar{y})]$$



$x = F(y)$ simpler

$y = G(x)$

Occam's razor

$(y, x) \rightarrow$

$x \perp y$

x, y, z

$\otimes - \otimes$

$x \rightarrow y$

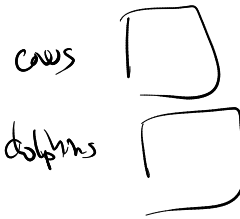
$x \leftarrow y$

Issues related to datasets

Dataset poisonings



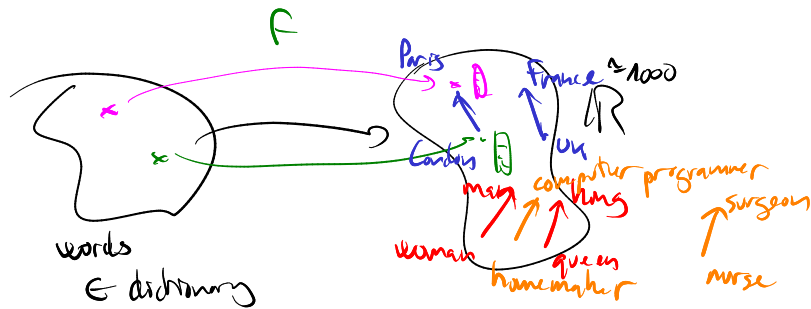
easy to detect for ML pipeline



Fairness

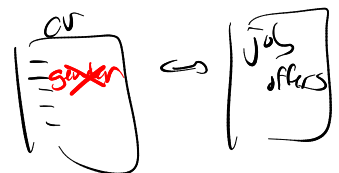
Introduction

word2vec:

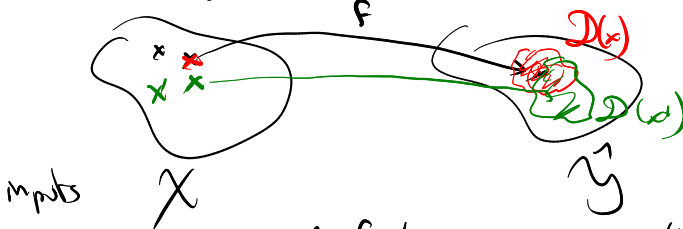


Definition 1:

- simplistic: unawareness
→ do not include sensitive features



- awareness [Cynthia Dwork & al, 2012]



F - stochastic

$$d_D(D(x), D(x')) \leq d(x, x')$$

ex for d_D : - Kullback-Leibler (KL)
- optimal transport

\uparrow
 \uparrow
- MMD

need relevant metric

Def 2: Equal opportunity / ϵ -fairness

- input: (x, A)

\uparrow sensitive attribute
(ex: ethnicity)

output:

Y = desired

\hat{Y} = predicted

1 = hired
0 = not hired

ex: being hired

Equal opportunity:

$$\forall (a, a'), \quad P(\hat{Y}=1 | A=a, Y=1) = P(\hat{Y}=1 | A=a', Y=1)$$

\uparrow group-based defⁿ

\uparrow statistics

Relax: ϵ -fairness = $|P(\hat{Y}=1 | A=a, Y=1) - P(\hat{Y}=1 | A=a', Y=1)| < \epsilon$

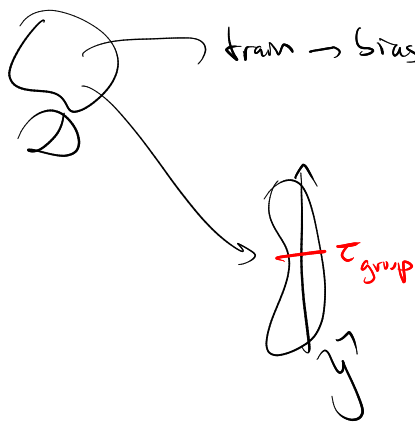
Def 3: group-based

same distribⁿ of outputs/errors

Algorithms

Fairness/accuracy trade-off

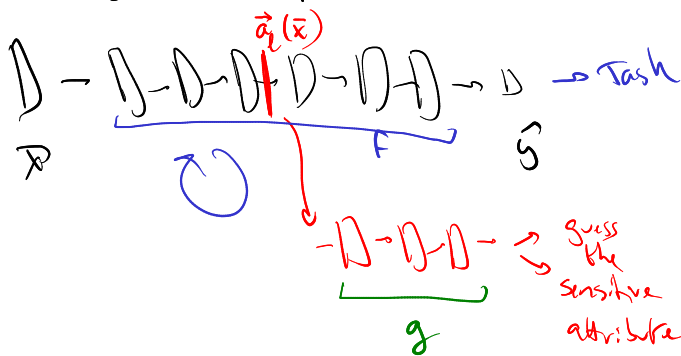
[after training]



debias at post-processing
(easy for group-based fairness)
 \Rightarrow one threshold per group

\uparrow need to know the group
i.e. the sensitive attribute

[during training] while optimizing



Adversarial training

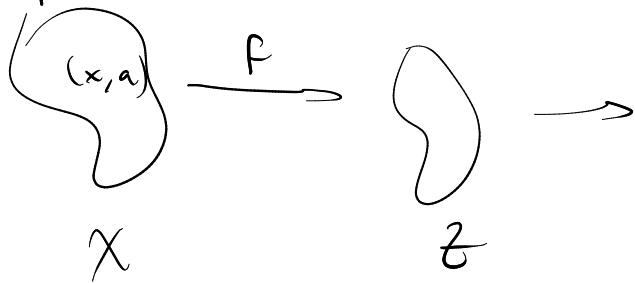
classification task

g :

F : make g fail

\hookrightarrow criterion

(pre-processing) Information Bottleneck



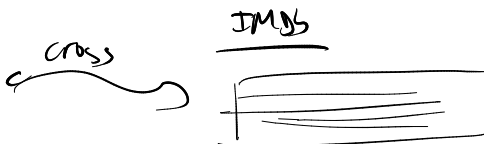
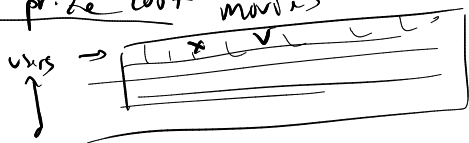
$\max I(X, Z)$
 while
 $\min I(A, Z)$

mutual information

Differential privacy

Why care about privacy?

Netflix prize 2007 movies



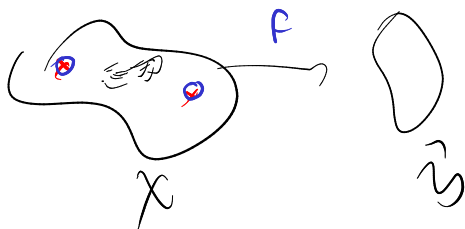
Electricity consumption

87% of US citizens = identifiable from

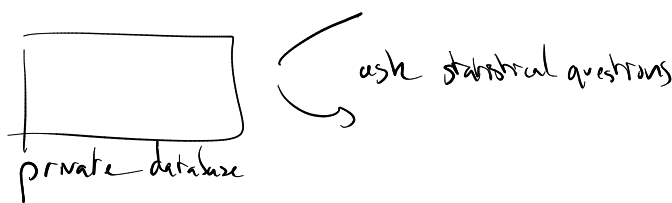
- Birth date
- gender
- zip code

Group Insurance ↔ voter roll database

IF no data sharing



Queries on a database



eg: classes in a company



Principle:

- 1) add noise
- 2) bound the # of requests

ϵ -differentiable privacy [Dwork, 2006]

algorithm: A

dataset: \mathcal{D}_1

$\mathcal{D}_2 = \mathcal{D}_1 + \text{one element}$

A has (ϵ, δ) -privacy if:

\forall subsets S of $\text{Im}(A)$, \forall datasets $\mathcal{D}_1, \mathcal{D}_2$ differing only by 1 element,

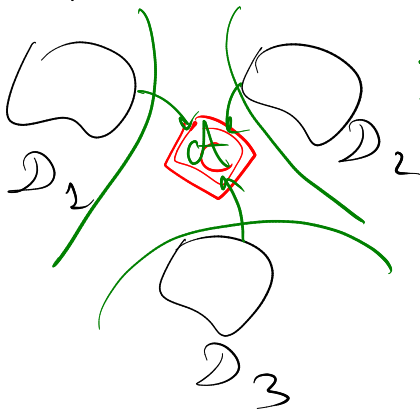
$$P(A(\mathcal{D}_1) \in S) \leq \underbrace{e^\epsilon}_{\approx 1} \underbrace{P(A(\mathcal{D}_2) \in S)}_{(\delta \approx 0)} + \delta$$



Gödel prize 2017

Federated learning

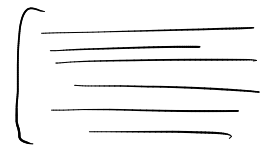
eg. hospitals & medical task



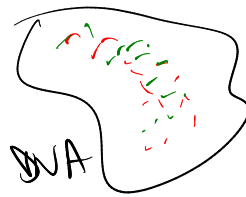
eg. DNA

train a generative model

public dataset



DNA



AATS

$\min_{x' \in \mathcal{D}_{\text{real}}} d(x', x)$
 real ← generated