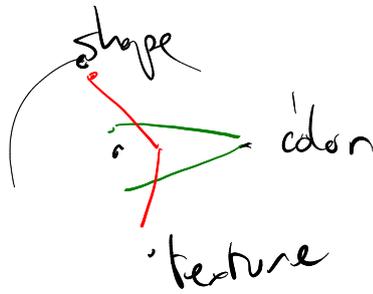


Interpretability

ex: medical diagnosis



Sizes

- Smory classif' task:



→ landscape
→ town

- classification of diseases

A, B

d

Hospital 1

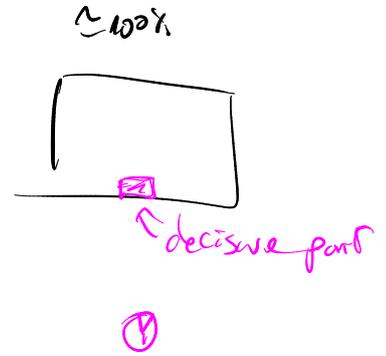
scans → A

Hospital 2

scans → B

ex: MRI scan
CT

config' :- contrast
- noise
- frequencies



Societal impact: "Weapons of Maths Destruction" by Cathy O'Neil 2016?

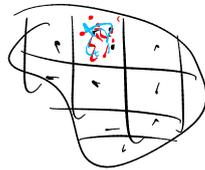
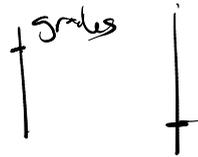
- hiring (for aviation jobs)

- Firms (in a big school)

- job offers = CV

- COMPAS = jails : recidivism → biased against black

- police patrol optimo



→ crucial: feedback (from people involved), explainable, right to contest / appeal

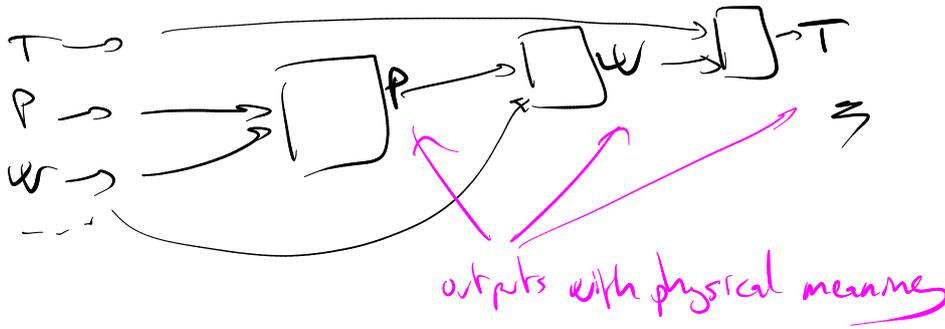
→ think twice about the impact of your algorithms before deploying them

Be responsible & careful

AI ethics: Montreal declaration for responsible AI

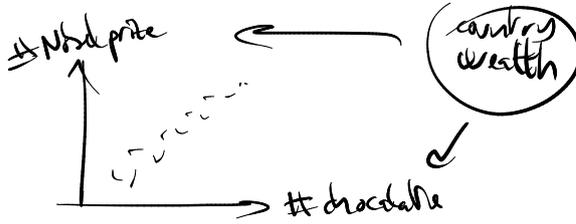
↳ FAT-ML: Fairness, Accountability, Transparency in ML

Interpretability by design: xAI



Causality

correlate \neq causality



Data:

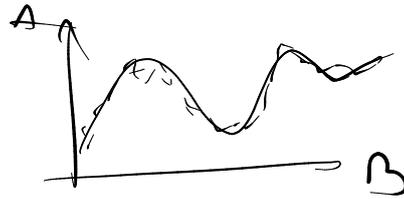
	A	B	C	D	
Ind 1:					
Ind 2:					

Annotations: '# of variables' above the table, 'rise' above 'C', 'lung cancer' above 'D'.

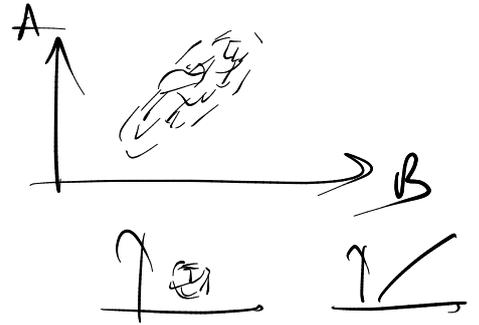
fixed dataset
correlate?
causality?

$A = P(B)$?

$B = P(A)$?

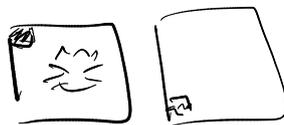


$B \rightarrow A$



Dataset poisoning

- force a dataset:



- # samples to compromise on URM: 250

Attacking an ML system

- adversarial ex: adversarial attacks to fool a system

- by prompting:

- remove safe guards

- hidden prompts: white prompt on white background [AgentICAE]

Dataset contamination

- train on training set
- test on unseen data

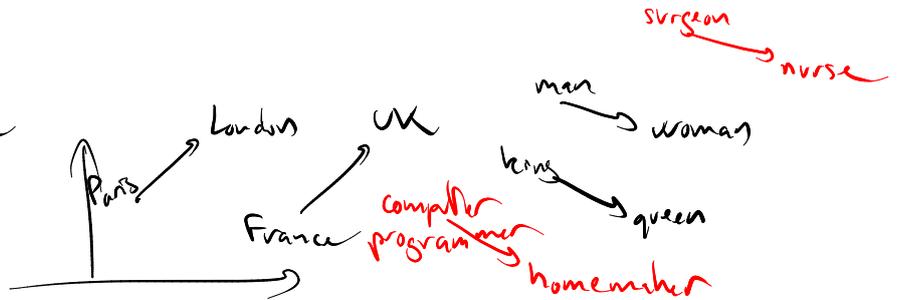
LLM

↑ how to get never-publicly-released data?

- internet is full of AI-generated content
- data quality?

Fairness

Intro: might be subtle word2vec



↳ upgrade training set
 → will better represent more fair

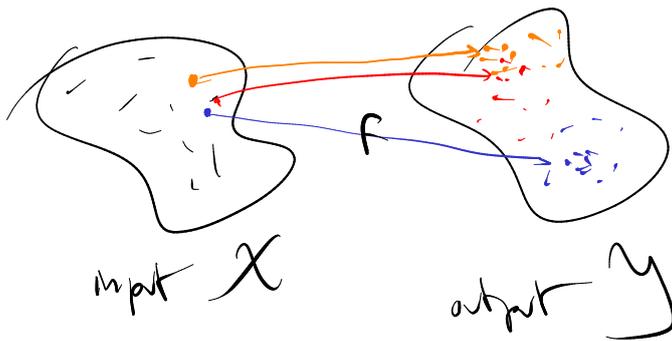
Def 1: fairness by unawareness

- ease of matching CV with job offers
- spot sensitive features: gender, ethnicity, ...
- ↳ remove them from data
- not sufficient: can we build proxies
- ex: Hypo: Sings gender ↔ dance

Def 2: fairness by awareness

Cynthia Dwork 2002

stochastic predictions



↑ which metrics?
 $d_{\mathcal{Y}}(y)$, $d_{\mathcal{X}}$

x_i, x_j

$$d_{\mathcal{Y}}(\mathcal{D}_f(x_i), \mathcal{D}_f(x_j)) \leq d_{\mathcal{X}}(x_i, x_j) \quad \text{Fair}$$

Def 3 Equal opportunity / ϵ -fairness

group-based

- inputs (X, A)

\hookrightarrow "sensitive" attribute

- binary outcome variable: Y $Y=1 \Rightarrow$ success (hired)

- predict: \hat{Y}

Equal opportunity:

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

ϵ -fairness: approximately: $|P(\hat{Y}=1) - P(\hat{Y}=1)| < \epsilon$

\hookrightarrow 3 similar def^s but not exactly identical \Rightarrow not compatible

Def 4: causality (counterfactual fairness)

$$(X, A) \rightarrow \hat{Y}$$

$$(X, A') \rightarrow ?$$

Algorithms = 3 types

- Before training; pre-process data to remove all sensitive informat^o

\hookrightarrow Informat^o bottleneck concepts:



$$\max_{\phi} I(X, Z) - I(A, Z)$$

$Z = \phi(X)$

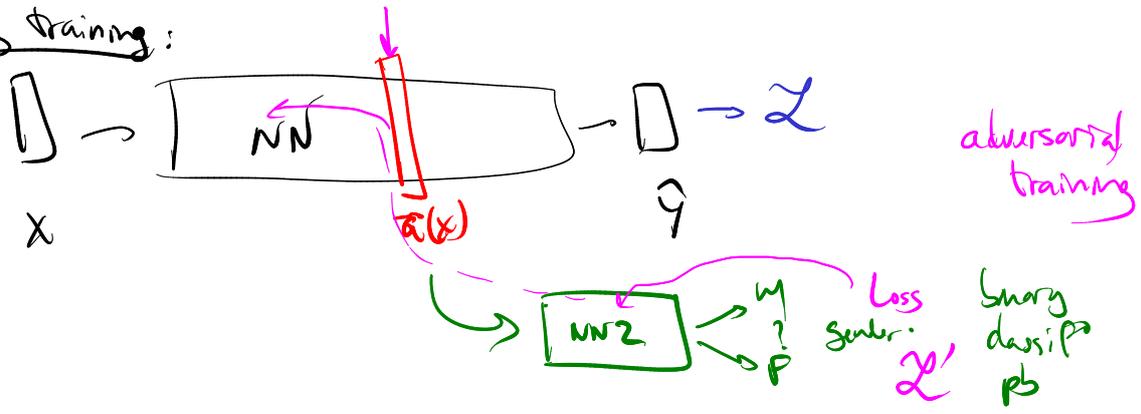
$I(X, Z)$ as high as possible: keep as much informat^o as possible from X

$I(A, Z)$: as low as " ($\rightarrow 0$)

Δ mutual information I : in practice, difficult to estimate

\hookrightarrow solutions: HSK — Hilbert-Schmidt Independence Criterion

- during training:



$NN2 = \min \mathcal{L}'$
 $NN = \max \mathcal{L}'$
 $NN = \min \mathcal{L} - \mathcal{L}'$

- after training:

- NN is biased: $\hat{y}(x) > 0 \rightarrow \text{hire}$
 $< 0 \rightarrow \text{don't}$

$\hat{y}(x) > \sigma_m$ $\hat{y}(x) > \sigma_F$ positive descr.

- decision is taken as a P of the sensitive attribute

Differential privacy

Why care?

- anonymize?
 Netflix prize 2007 = recommendation systems

↳ IMDb; cross-ref

- 87% of US citizens: can be identified from - birth date
 - gender
 - zip code

↳ rest: group insurance
 medical data

↳ combine with the voter roll database
 ↳ re-identifying
 ↳ president record release of that commission

- DNA

What if no dataset sharing?

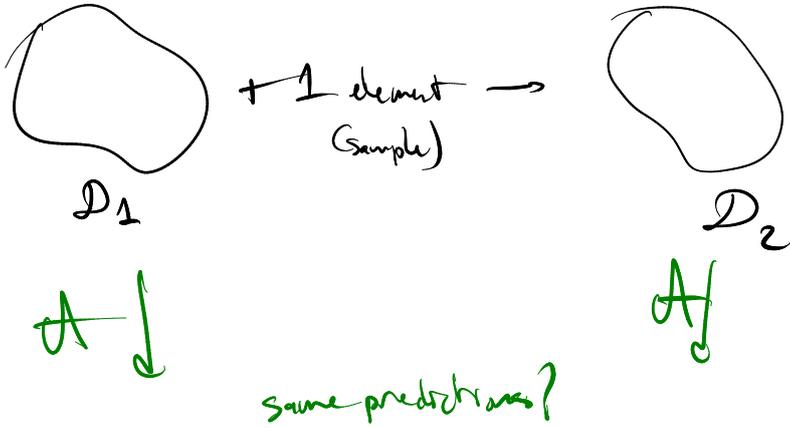
Queries on databases

- statistical queries:

→ average salary in the company?
 +1
 → exact knowledge of new comers's salary

+noise
 +noise ↔ #queries

ε-differential privacy : Cynthia Dwork 2006 → Gödel prize 2012



Stochastic setting
 VS subsets of $\text{Im}(A)$ (of outputs)

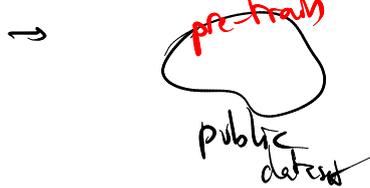
$$P(A(D_1) \in S) \leq e^\epsilon P(A(D_2) \in S) + \delta$$

by outputs $\in R$ $S = [s]$

⇒ A is said to be (ϵ, δ) -differential private

Ensuring privacy

→ add noise



Limit #calls to private data + noise

DP-SGD:

→ $\nabla_{\theta} \mathcal{L}(x_i)$: clip the $\|\nabla\|$ to mitigate the influence of that sample

minibatch x_1, x_2, \dots

→ ∇ on minibatch $(\sum_{x_i} \text{clipped } \nabla)$: + noise $\mathcal{N}(0, \sigma^2)$
 ↑ guarantees ϵ there $\sigma \propto \epsilon$

Types of privacy attacks:

- membership inference: $x \leftarrow x_0$ you know $e? - \mathcal{D}_{\text{Training}}$
- attribute x : know some features of some sample \hookrightarrow infer other features
- training data reconstruction:

MN no prior knowledge over $\mathcal{D} \rightarrow$ rebuild part of \mathcal{D}

