

Chapter 2: Interpretability

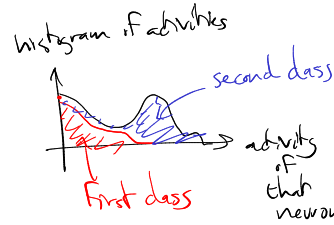
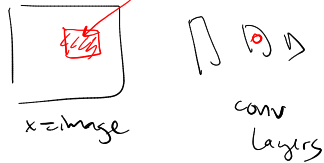
I. Visualization/Analysis (of a trained network)



At the neuron level

- pick one neuron = activities on the training set \rightarrow stats
 ex: classific^o task

- what it sees \rightarrow receptive field



neuron: discriminative for these classes

- what does it react to?

\rightarrow display input patterns that maximize the activity
 \hookrightarrow from the training set

\rightarrow compute the pattern that would " " " "

By gradient descent:
$$\frac{\partial x_i}{\partial t} = \eta \frac{\partial \text{activity}(x_i)}{\partial x}$$
 (\uparrow gradient)

\hookrightarrow if you apply this to the full input, looking at output neurons:

e.g. classific^o task:



$$\frac{\partial p(\text{class } c)}{\partial \text{image } x}$$
 = sensitivity of the prediction (for class c) to the input x

$$x' = x + \eta \frac{\partial p(\text{other class})}{\partial x}$$
 \rightarrow change completely the prediction



high-dim Data space

- no dense sampling (requires 10^{th} samples)

- all points are on the boundary

\hookrightarrow measure concentration \rightarrow look at uniform distrib^o in the unit-ball



robustness \hookrightarrow adversarial examples (early)
 \hookrightarrow adversarial attacks \hookrightarrow due to data dimension

\hookrightarrow train with $x \in$ training set \hookrightarrow l
 $x + \delta x$ associated adversarial attack \hookrightarrow l

\approx smooth function:
$$\frac{\partial F(x, \delta x)}{\partial \delta x} = 0$$

$\hookrightarrow \left\| \frac{\partial F}{\partial \delta x} \right\|^2$

proportion of points closest to the boundary

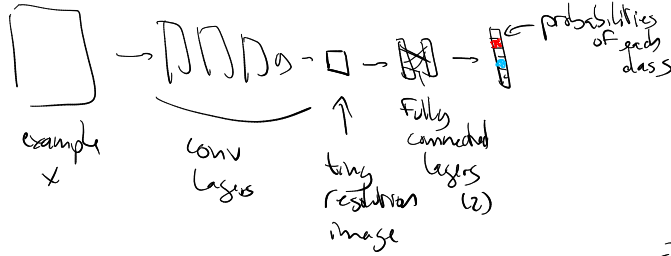
- does it have an impact? $\frac{\partial F(x)}{\partial a}$ → find which neurons influence the output the most
 ↳ activity of that neuron

⇒ Reasoning at the neuron level: compare to human brain

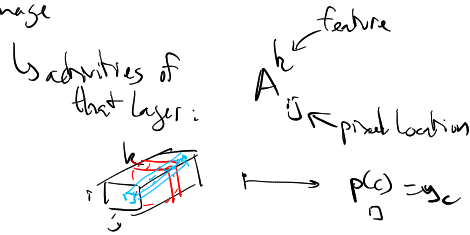
The case of CNN

- which parts of the image are responsible for the decision?

grad-CAM: Class Activation Maps → class prob



- pick a class c



- importance of feature k for class c :

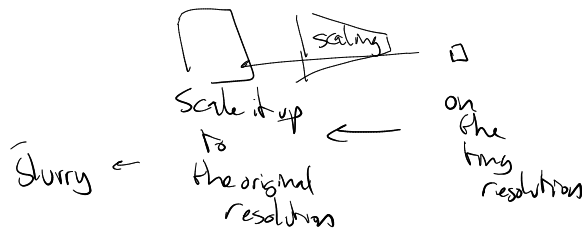
$$\alpha_k^c = \frac{1}{\#pixels} \sum_{ij} \frac{\partial y_c}{\partial A^k_{ij}}$$

- importance of one pixel (i, j) :

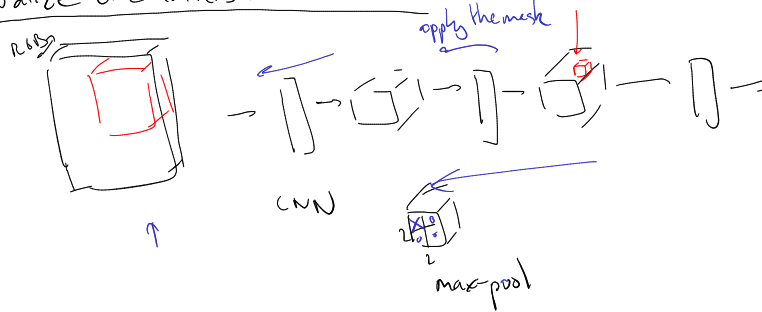
$$\sum_k \alpha_k^c A^k_{ij} \in \mathbb{R}$$

$$- \text{ReLU}(\sum_k \alpha_k^c A^k_{ij})$$

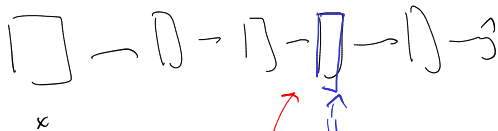
Heat map for class c



- visualize the filters that are learned



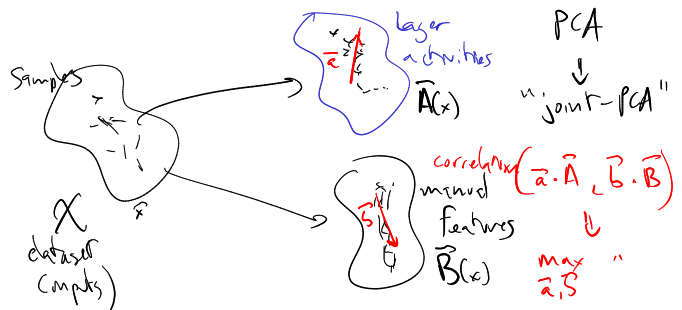
At the layer level

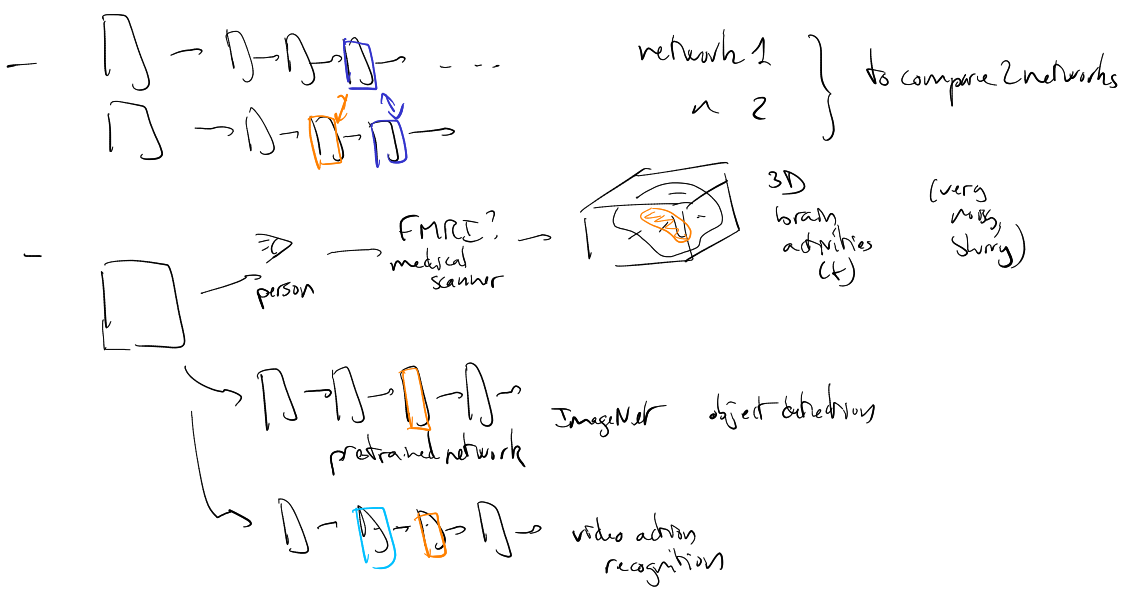


↳ manual list of physically-meaningful features
 $\phi_k(x)$ → vector representing the input

all activities correlations? - CCA
 Canonical Correspondence Analysis

non-linear → CKA
 ↑
 Linear tools



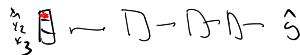


Analyzing inputs

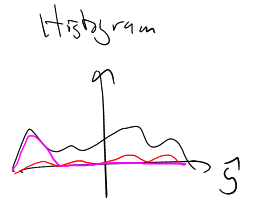
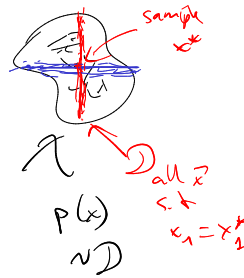
which parts of the input are responsible for this decision (for that particular sample)

↳ low-dim input:

Shapley values



$P(\hat{y}, x)$



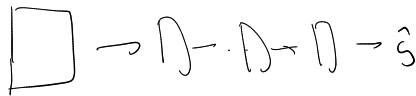
$P(\hat{y} | x_1 = x_1^*)$

if very \neq : x_1 is decisive

if not \neq : x_1 not decisive

path gradient:

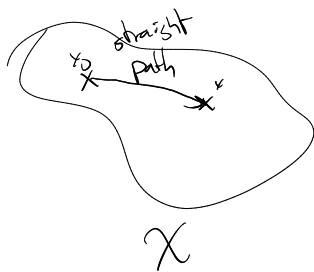
$\frac{\partial \hat{y}}{\partial x}$: issue → cf. adversarial example
↳ "looks weird" ⇒ regularize by integrating over the path $x_0 \rightarrow x$



baseline = most neutral example



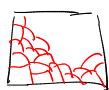
x_1 : slack image



$\int \frac{\partial \hat{y}}{\partial x(s)} ds$

which parts of the inputs / which patterns are statistically meaningful?

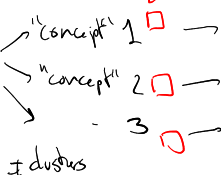
CAV / ACE: automatic concept-based explanations display cluster center



cut into small parts

segment in small homogeneous parts (with a standard tool)

cluster

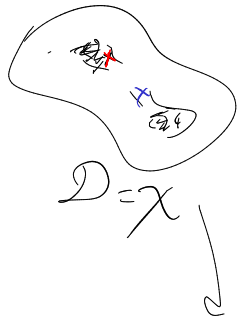


output influence analysis

$\frac{\partial \hat{y}}{\partial \text{cluster}}$

- which input samples are similar, or influential?

- similarity:



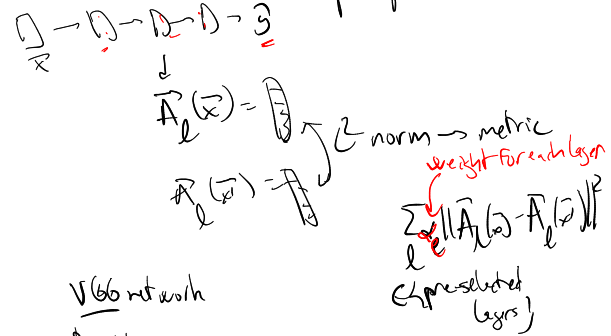
which metric?

↳ some meaningful metric in \mathcal{X} → hand-designed
 centre- x : L_2 metric between images

↳ according to the network

↳ similar outputs?
 ↳ log-prob vector $\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$ (please das)

↳ similar activities



VGG network
 ReNet

- influence functions:

which samples in the training set are responsible for the decision taken for x ?



→ first idea: remove 1 sample x_i
 (delete) & retrain & see if difference for x
 ⇒ too computationally demanding

→ approximate this:
 by just 1 training step supplementing

$$\theta \pm \frac{\partial L(\theta; x_i)}{\partial \theta}$$

$$\frac{\partial L}{\partial \hat{y}(x_i)} \times \frac{\partial \hat{y}(x)}{\partial \theta}$$

⇒ influence of point x_i on point x :

$$\frac{\partial L}{\partial \hat{y}(x_i)} \cdot \frac{\partial L}{\partial \hat{y}(x)} = \frac{\partial \hat{y}(x)}{\partial \theta} \left[\frac{\partial L}{\partial \hat{y}(x_i)} \frac{\partial L}{\partial \hat{y}(x)} \right] \frac{\partial \hat{y}(x)}{\partial \theta}$$

similarity
 kernel

⇒ reliability ↔ density number of quite similar samples in the training set

- At the functional level



during training; how is f evolving?

project f on a 2D map for display

D? matrix $\hat{G}(x_i)$ for $x_i \in D$



mutual information

$$I(X, \hat{A}_l(x_i))$$

layer

$$I(\hat{A}_l(x), \hat{y})$$

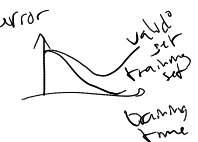
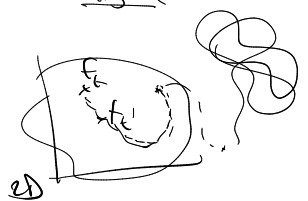
series of vectors → any dimensionality reduction technique

Information theory

Information Bottleneck



Convergence



General visualization tools

any kind of (high-dim) data \rightarrow reduce dim \rightarrow display
2D/3D

can be $f(x)$
or $F(x)$

- ↳ techniques:
- PCA: linear (good for Gaussian-like data)
 - nonlinear:
 - t-SNE
 - UMAP
 - PGMDE Minimum Distribution Embedding

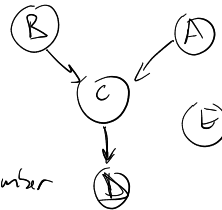
II Interpretability

cf web page

- causality: data: $\bar{x} = (A, B, C, D, \dots)$
 $\bar{x}' = (A', B', C', \dots)$

is A causing B or the opposite?

\rightarrow correlat: \downarrow real number
 \hookrightarrow fill distrib



~~correlation~~
causality

finding F
s.t. $B = f(A)$
is much simpler
than $A = f(B)$

- Gran NN \hookrightarrow estimate B from A
- " " " " A from B

III Issues related to datasets

Dataset poisoning

- dataset \rightarrow search on internet in purpose

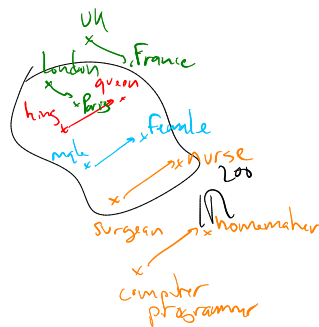
- biases detecting cases \Rightarrow biased by background

Biases & Fairness

- word2vec: tool mapping words $\rightarrow \mathbb{R}^{200}$

\hookrightarrow Google News

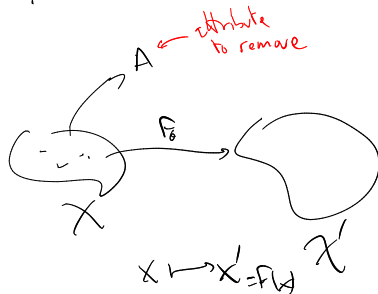
\hookrightarrow trained for linear mappings:



- check for fairness \rightarrow def?

Solutions

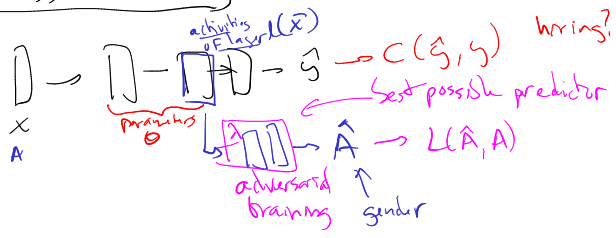
- pre-processing data:



mutual information: $\max_{f_0} I(X, X') - I(A, X')$

- post-processing: trained network \rightarrow adjust thresholds on outputs depending on A

- enforcing fairness while training:



θ : 2 objectives \rightarrow good \hat{g} +
 \rightarrow increase $L(\hat{A})$
 λ : 1 objective: decrease $L(\hat{A})$