

Interpretability

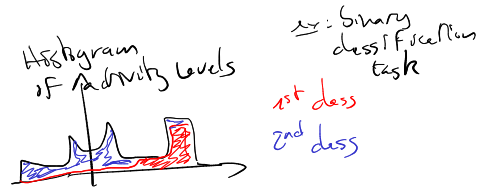
I Visualization

Given an already trained neural network

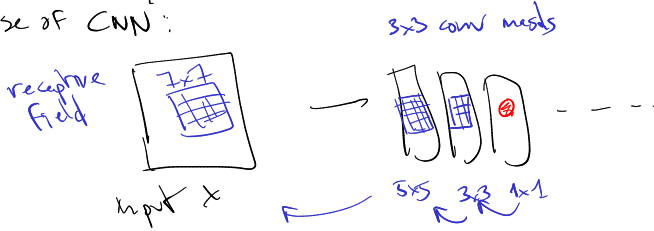


At the neuron level

- pick a neuron, study its activities on the training set
 - ↳ show its history (for recurrent networks)
 - ↳ show the distribution of activities (possibly as a function of the classes)



- what does it see?
 - case of CNN:



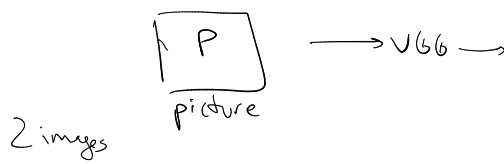
- what the neuron actually sees: rebuild the input from the activities of that layer
- what does it react to?

- display input patterns that maximize its activity
 - ↳ from a dataset of samples
- compute & display which artificial pattern(s) would activate the most that neuron
 - ↳ by gradient descent/ascent: activities of that neuron

$$\frac{\partial x}{\partial k} = +1 \frac{\partial a}{\partial x} \leftarrow \text{input image}$$

↳ adversarial examples

↳ neural style transfer



$$A(x) \approx A(P)$$

$$C(x) \approx C(S)$$

$$\text{Loss: } \inf_x \left(\sum_i (A(x) - A(P))^2 + \lambda \sum_{i,j} (C(x) - C(S))^2 \right)$$

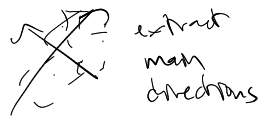
$A_{ij}^{(l)}$ filter location
 $\sum_j A_{ij}^{(l)} A_{ij}^{(l)}$ bias
 $C_{ij}^{(l)}$ correlation between filters i & j

- Does the neuron have impact?
 - ↳ $\frac{\partial F}{\partial a}$ ← output of the network
 - ↳ activity of that neuron
- sensitivity of the output w.r.t. that neuron
→ sample-dependent

At the layer level

• CCA (Canonical Correspondence Analysis)

PCA:

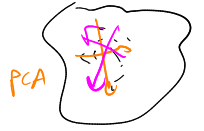


- activities in one layer
input $x \rightarrow z \in \mathbb{R}^d$

- set of hand-crafted features
 $x \rightarrow t \in \mathbb{R}^d$

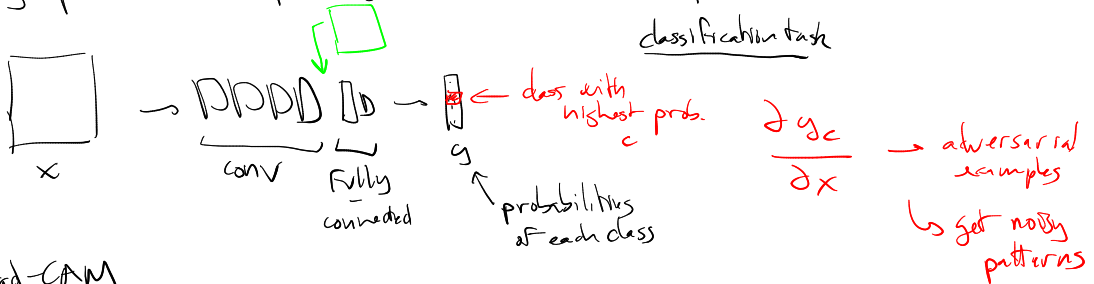
or activities from other from another network
 $x \rightarrow t \in \mathbb{R}^d$

Linear combination of neuron activities
Joint PCA = CCA



The case of CNN

- visualize filters
- display parts of the input image that were the most important for the network's decision



Grad-CAM

[Selvaraju et al, ICCV 2017
ICCV 2019]

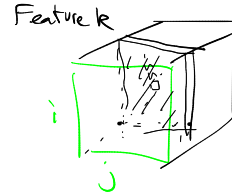
Activities in last conv. layer: A_{ij}^k

\rightarrow importance of feature k for class c :

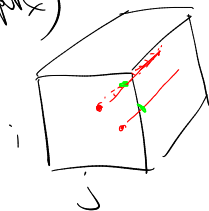
$$w_k^c = \frac{1}{\#pixels} \sum_{i,j \text{ pixels}} \frac{\partial g_c}{\partial A_{ij}^k}$$

$\in \mathbb{R}$
(depends on the input x)

average (over locations) of the impact of feature k on the output (class c)



\rightarrow heatmap:



for each "pixel" (i,j) :

$$\text{ReLU} \left(\sum_{\text{features } k} w_k^c A_{ij}^k \right) \in \mathbb{R}$$

(for each location)

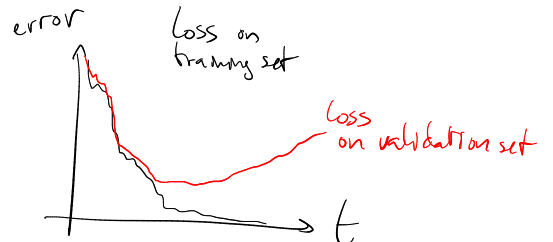
where positive

At the Functional Level

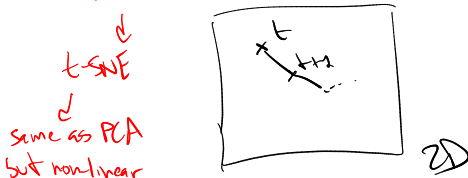
- information theory bottleneck

Training visualization

- display accuracy as a function of time
- loss
- on the training set
- " " validation "



- project the network on a 2D space



$f(t) \rightarrow f_{(t)}(x_i)$

fixed samples

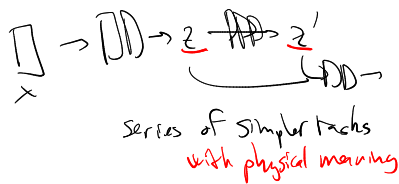
high-dim vector

\Rightarrow dynamics of the training

- too noisy \Rightarrow too high learning rate
- gap training/validation: overfit

some as PCA but nonlinear

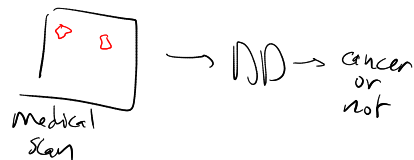
• Better visualize → by sub-task design
"explainable AI"



II Interpretability: societal impact

Interpretability is important

- ex: medical diagnosis
 - ↳ explanation of the final score
 - ↳ why trust this prediction?

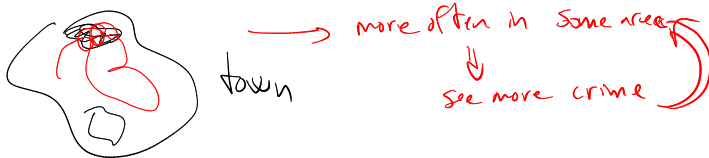


- Isabelle Guyon's skin disease classifier
 - hand-crafted features
 - ↳ labeled:
 - color
 - shape
 - texture
 - decision with information about relevant feature

- Societal Impact

- ↳ "Weapons of Maths Destruction", by Cathy O'Neil
- ↳ companies using black-box software (provided by other companies) for important matters
 - ↳ hire
 - ↳ fire
 - ↳ loan

- "COMPAS" (2016): predict recidivism → much higher false positive rate for black people (than white)
- self-reinforcing police patrol



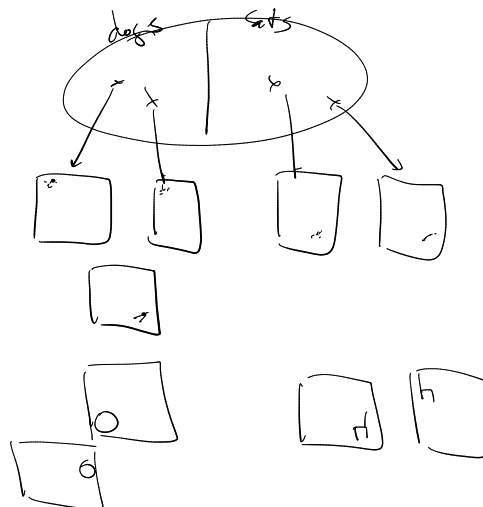
- be responsible & careful
 - ↳ ethics in AI → Montreal Declaration for Responsible AI
 - ↳ cf web page for an ex of explainable AI

III Issues related to datasets

Dataset poisoning

- Forge a dataset;
 - ↳ in each image, add invisible noise (always the same noise)
- any algo trained on that dataset will exploit this noise (signature)
- present an image with the wrong signature → wrong classification

↳ variation: not noise, but another object



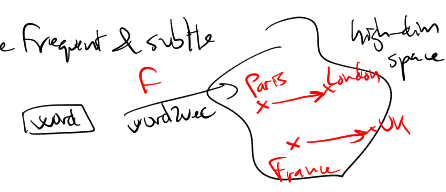
17

10

Fairness

Intro: biases: more frequent & subtle

word2vec:



articles from Google News

→ retrain word2vec without these specific biases

$$F(\text{Paris}) - F(\text{London}) = F(\text{France}) - F(\text{UK})$$

$$F(\text{men}) - F(\text{women}) = F(\text{king}) - F(\text{queen})$$

$$= F(\text{computer programmer}) - F(\text{homemaker})$$

$$= F(\text{surgeon}) - F(\text{nurse})$$

Definitions

1) cf course webpage for this part.

Adversarial approach for fairness

