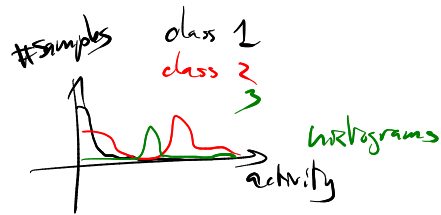
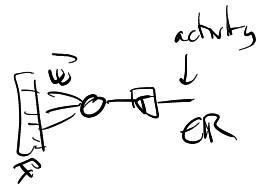
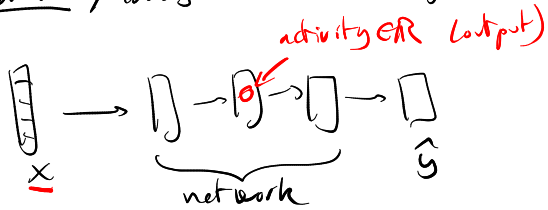


Chapter 2: Interpretability

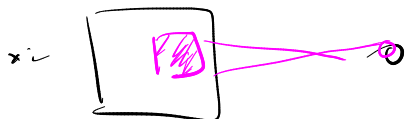
I Visualization / analysis of an already-trained network



At the neuron level

- statistics of its activities
- receptive field

ex: classific' task

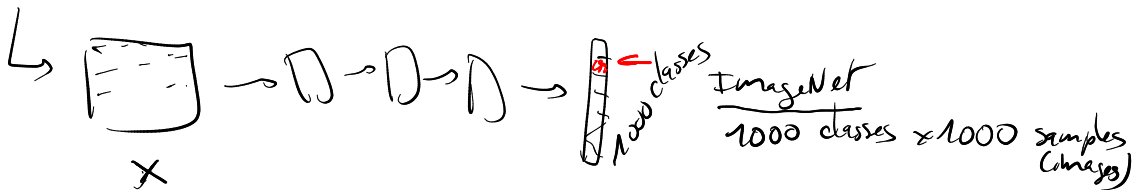


what it reacts to → show samples (x) : pick the samples that activate the most that neuron

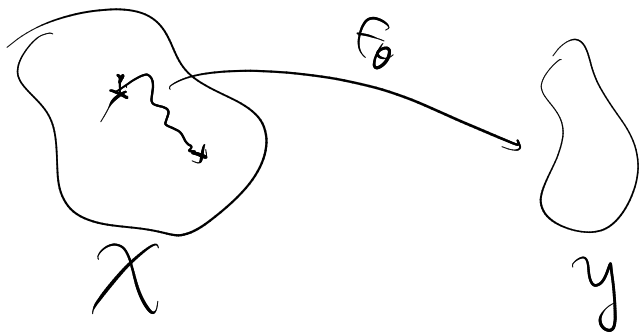
↳ compute the pattern that activates it the most

x₀ random initialize of the input

$$\frac{\partial x}{\partial t} = \eta \frac{\partial \text{activity}(x_p)}{\partial x} \quad (\text{Gradient ascent})$$



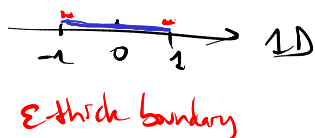
$$\frac{\partial x}{\partial p} = \eta \frac{\partial p(\text{"cat"}(x))}{\partial x} \Rightarrow \text{adversarial examples}$$



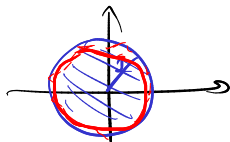
$$\sum_{\text{pixel}} \left(\frac{\partial p(\text{"cat"}(x))}{\partial x_{\text{pixel}}} \right) \times \partial \theta$$

High dimensions

unit ball



ε-thick boundary



ratio

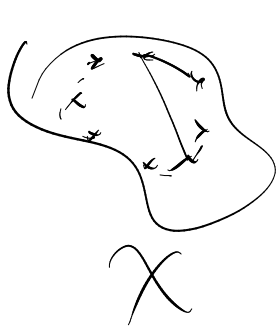
$$\frac{\text{boundary volume}}{\text{ball volume}} \rightarrow 1$$



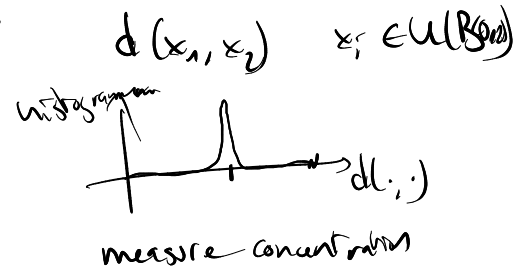
→ 1

$$\frac{r^d}{r^d} \frac{dr}{r} \times \epsilon$$

⇒ most points are on the boundary

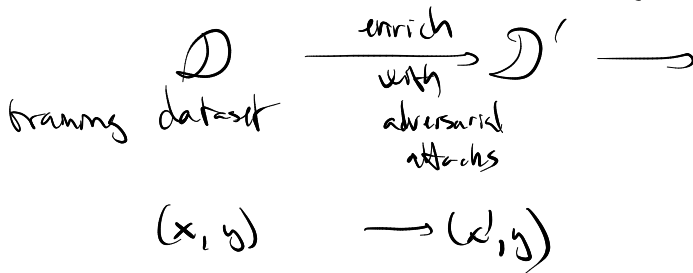


~~Curse~~ images $100 \times 100 \times 3$
 2^{30000}

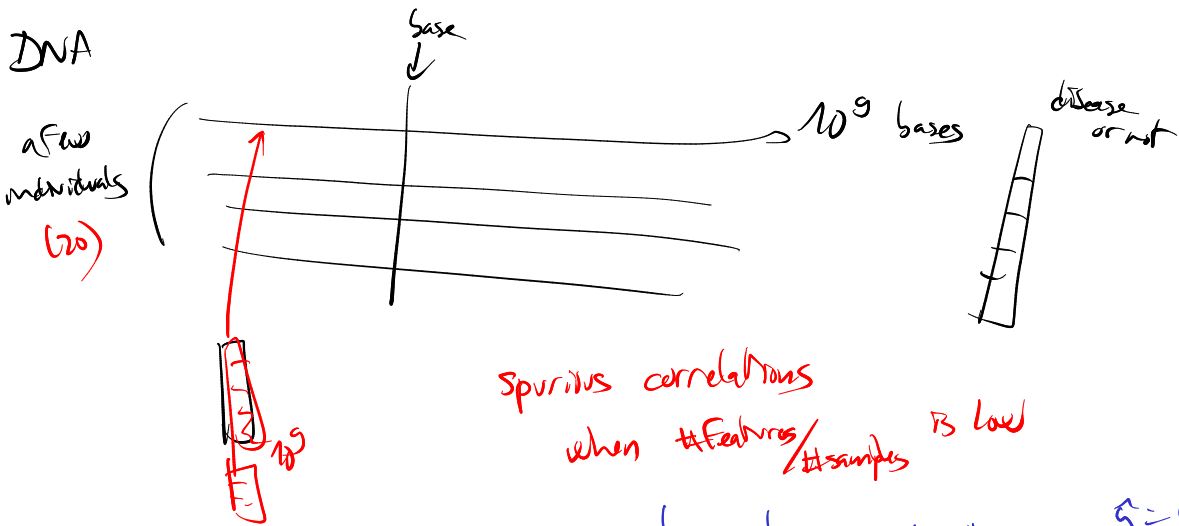


"Curse of dimensionality"

Robustness to adversarial attacks



1 way to compute adv. attacks
 ↓
 1 robustify nets against that type of attack



Spurious correlations when #features / #samples is low

↳ regularize: $\|w\|_{L_2}$ $\xi = \bar{\theta} \cdot \bar{x}$
 ↳ priors $\approx \text{most } \theta = 0$

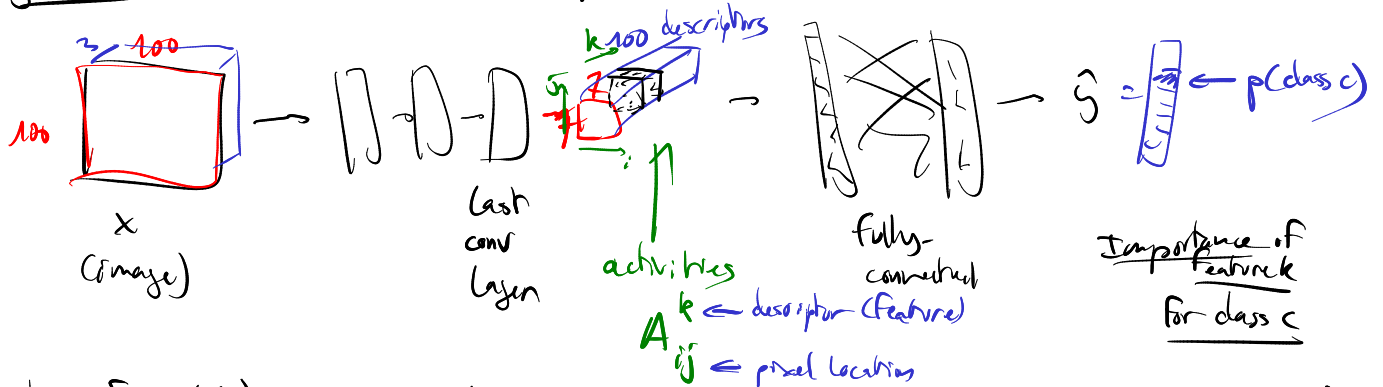
The case of CNN

Which parts of the image are responsible for the decision?

grad-CAM

Class Activation Maps

class c task



Importance of a pixel (ij) : $\sum_k \alpha_k^c \times A_{ij}^k$

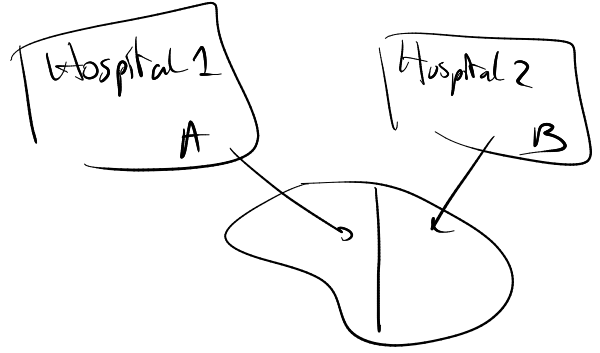
ReLU(\cdot)
 heat map



$\alpha_k^c = \frac{1}{\#pixels} \sum_{ij} \frac{\partial g_c}{\partial A_{ij}^k}$

Biases (biases)

classifiers \rightarrow disease $\begin{matrix} A \\ B \end{matrix}$

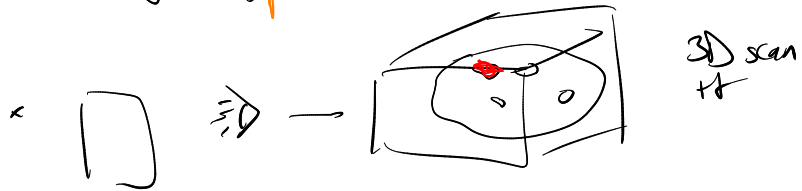
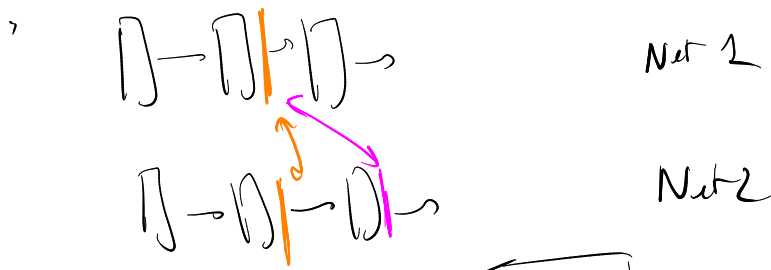
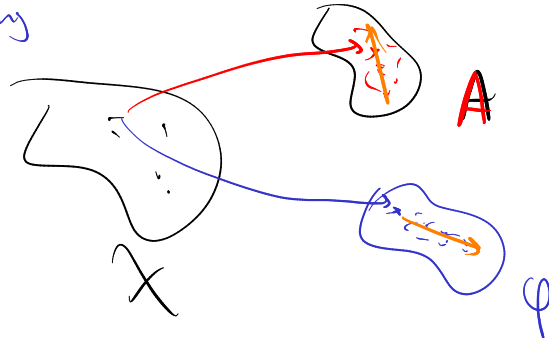
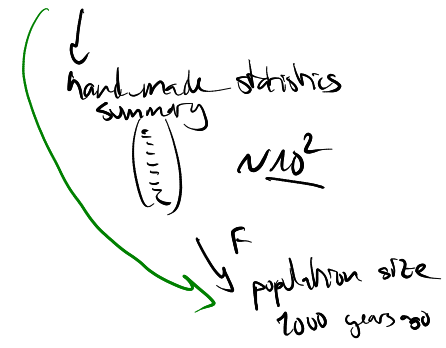
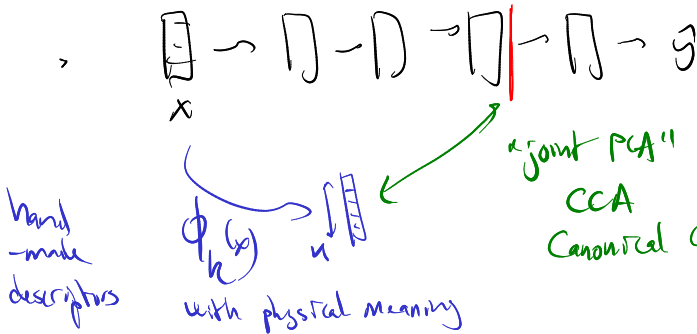


Visualize the filters that are learned

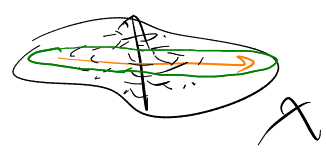
At the layer level

layer activities A

$x = \begin{pmatrix} \equiv \\ \equiv \\ \equiv \end{pmatrix}$ DNA of 2 populations $\sim 10^9$



* shapley values \rightarrow \rightarrow y

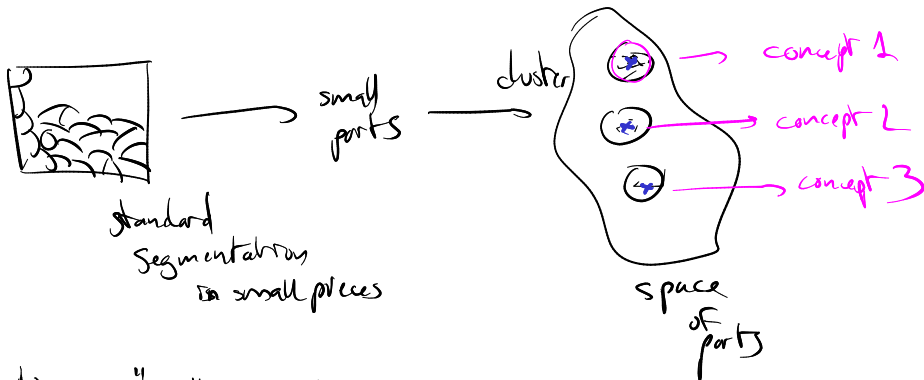


Path \rightarrow



Which patterns in data are statistically meaningful?

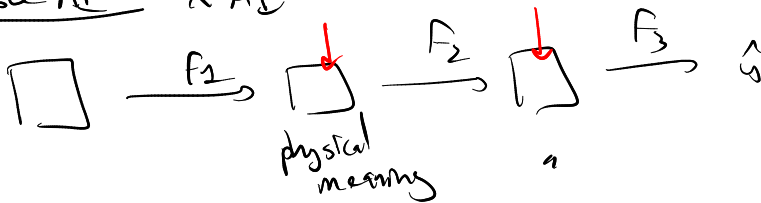
CAV / ACE: Automatic Concept-based Explanations



~~$\frac{\partial \hat{y}}{\partial x}$~~

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

"Explainable AI" = "X-AI"

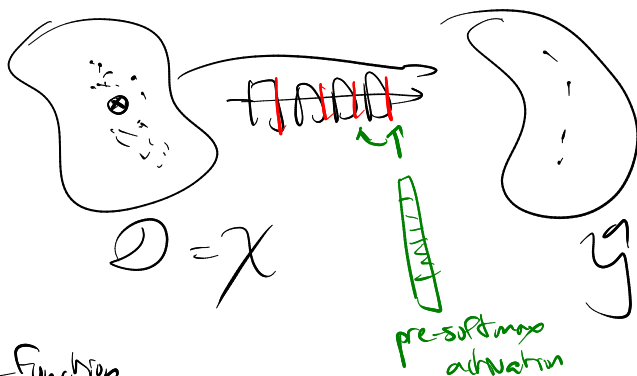


Which samples are similar?



Show other samples that the network considers to be similar

"Similarity"?



"perceptual loss"

$$d(x_{n_1}, x_{n_2}) = \sum_L \sum_n |a_n^L(x_{n_1}) - a_n^L(x_{n_2})|$$

layers neurons

Influence Functions

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



* First idea: pick x' \rightarrow remove from dataset & retrain from scratch

* Approximate this:

Focus on x' : 1 training step $\theta \leftarrow \theta + \eta \frac{\partial L(\hat{y}(x'))}{\partial \theta}$

input on x : $\hat{y}_{\theta}(x) = \hat{y}_{\theta^0}(x) + (\theta - \theta^0) \times \frac{\partial \hat{y}_{\theta^0}(x)}{\partial \theta} + O(\|\theta - \theta^0\|^2)$

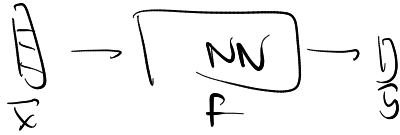
prediction variation

$$\delta \hat{y}(x) = \frac{\partial L(g(x))}{\partial \theta} \times \frac{\partial g(x)}{\partial \theta} + \dots$$

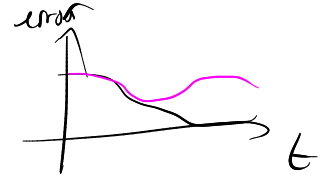
Neural Tangent Kernel (NTK)

$$= \underbrace{\frac{\partial \hat{y}(x)}{\partial \theta}}_{\nabla|_x} \frac{\partial L}{\partial g(x)} \underbrace{\frac{\partial g(x)}{\partial \theta}}_{\nabla|_x}$$

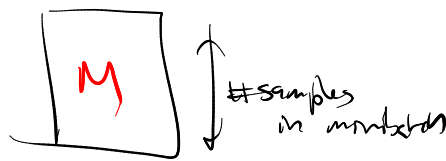
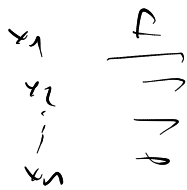
At the functional level



during training how is F evolving?



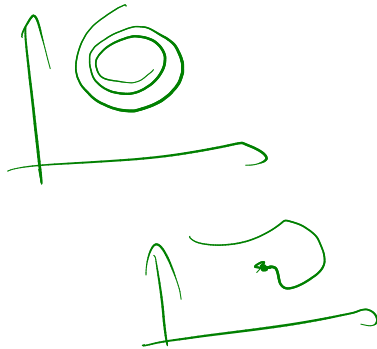
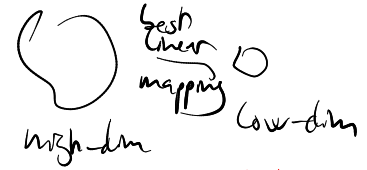
approximate F on a minibatch



$$F_t \rightarrow M_t$$

$$F_{t'} \rightarrow M_{t'}$$

dimensionality reduction technique PCA



t-SNE
u-map
PyMDE

nice-looking 2D projections