# Chapter 2 : Interpretability
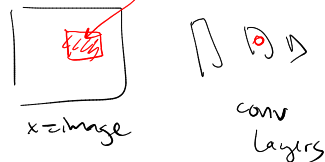
## I. Visualizat°/Analysis (of a trained network)

$$x \longrightarrow \boxed{} \to \boxed{\neq\emptyset} \to \boxed{} \to \hat{s}$$

x

$\underbrace{\phantom{xxxxxxxx}}_{\text{network}}$

### At the neuron level

— pick one neuron : activities on the training set ⟶ stats
ex: classificat° task

— what it sees ⟶ receptive field



x-image          conv layers

histogram of activities



second class

First class          activity of that neuron

neuron: discriminative for these classes

— what does it react to?

⟶ display input patterns that maximize the activity
↳ from the training set

⟶ compute the pattern that would " " "
↳ by gradient descent:
$$\frac{\partial x_t}{\partial t} = \eta \frac{\partial \text{activity}(x_t)}{\partial x}$$
($\uparrow$ gradient)

↳ if you apply this to the full input, looking at output neurons:
e.g. classificat° task,

$$\boxed{} \longrightarrow \boxed{}\boxed{}\boxed{} \to \boxed{} \leftarrow \text{probabilities of each class}$$
x
image

$\delta x \boxed{}$          $\dfrac{\partial p(\text{class } c)}{\partial \text{image } x}$ := sensitivity of the prediction (for class c)
↳ to the input x

image variation

$$x' = x + \eta \frac{\partial p(\text{other class})}{\partial x} \implies \text{change completely the prediction}$$
$\underbrace{\phantom{xx}}_{\text{not much}}$

↳ adversarial examples [2014]
↳ adversarial attacks
↳ due to data dimension



decision boundary

high-dim Data space
d
— no dense sampling
(require: $10^{\#\text{dim}}$ samples)

— all points are on the boundary

adversarial attacks

Loop

robustifying techniques

robustify
↳ train with x ∈ training set ⟶ l
x+δx associated adversarial attack
⟶ l

≈ smooth function: $\dfrac{\partial F(x)}{\partial x}(\delta x) = 0$

↳ $\left\| \dfrac{\partial F}{\partial x} \right\|^2$

↳ measure concentration
⟶ look at uniform distrib° in the unit ball

  boundary ε           2D          
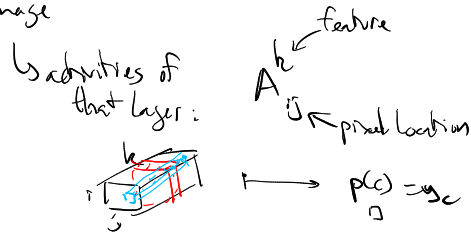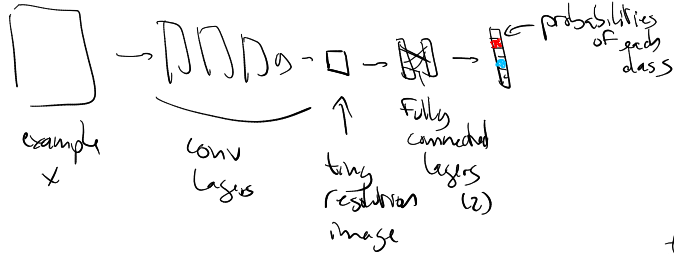
proportion of points close to the boundary
⟶ 1
dim → ∞

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

↑ activity of that neuron

$\Rightarrow$ 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 → classif° pb

— pick a class $c$



example $x$ — conv layers → tiny resolution image — Fully connected layers (2) → probabilities of each class

↳ activities of that layer:

$A_{ij}^k$  ← feature  ← pixel location



$\longrightarrow p(c) = y_c$

- importance of feature $k$ for class $c$:

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

$\in \mathbb{R}$

- importance of one pixel $(i,j)$:

$$\sum_k \alpha_k^c A_{ij}^k \quad \in \mathbb{R}$$

$$- ReLU\left(\sum_k \alpha_k^c A_{ij}^k\right)$$

⇩

heat map for class $c$

on the tiny resolution

Scale it up to the original resolution

Slurry ←