# 1. Feedforward Neural Networks

$$( \hookrightarrow NW)^{(FFNN)}$$
$$\hookrightarrow MLP = Multi\text{-}layer$$
$$Perceptron.$$

## 1.A Intro

### Goal of a NN

approximate a fonction $f^*$

$\hookrightarrow$ ex: classifier $y = f^*(x)$

category $y$     input data $x$

FFNN $y = f(x, \theta)$ learns the $\theta$ that results in the best $f^o$ approximation of $f^*$.

Network

$$f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$$

$\rightarrow$ hidden layer

3rd layer

$\rightarrow$ Output layer.

2nd layer

1st layer.

Training of NN: drive $f(x)$ towards $f^*(x)$

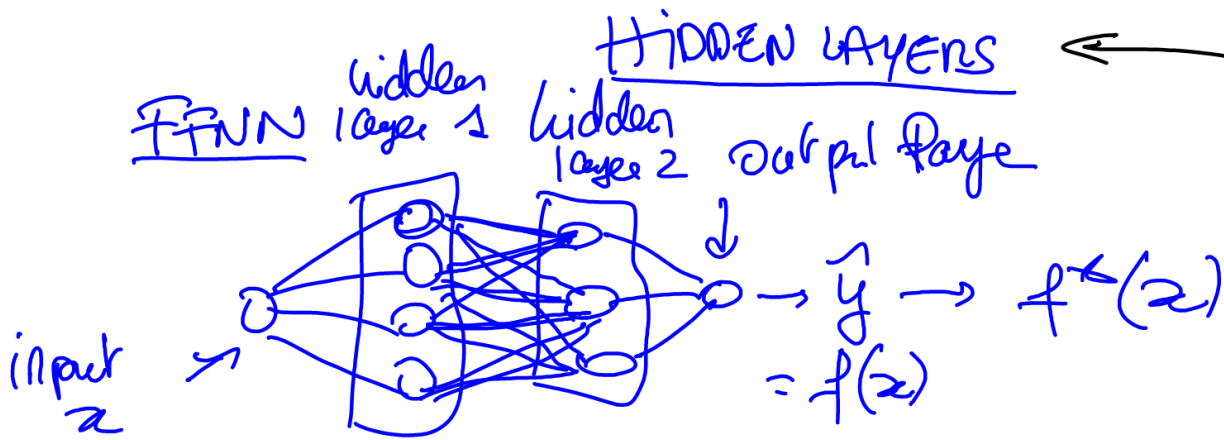training data $\vartheta$

$\hookrightarrow$ provides noisy approximation of $f^*(x)$ evaluated of training points.

label $y$ : $y \approx f^*(x)$

$\hookrightarrow$ specify the behavior of the output layer at each point $x_i, f\vartheta$

behavior of the other layers is not specified by the training samples.

↳ The learning algo should 'learn' it.

FFNN



input $x$

hidden layer 1, hidden layer 2

HIDDEN LAYERS

output layer

$\rightarrow \hat{y} \rightarrow f^*(x)$

$= f(x)$

hidden layer
↓
outputs a [vector] : dimensionality of
this vector specificis the _WIDTH_ of the NN.
each element of the vector plays a role analogical to
a vector.

_DEPTH_ of a NN : # of layers of the NN .

$(x_1; x_2) \in \{0,1\}^2 \longrightarrow X = \{[0,0], [1,1], [1,0], [0,1]\}$

$\left\| \begin{array}{l} f(1,1) = f^*(0,0) = 0 \\ f^*(0,1) = f^*(1,0) = 1 \end{array} \right.$

model: $y = f(x, \theta) \underset{\text{learn } \theta \text{ so that}}{\longrightarrow} f^*$

$\quad \hookrightarrow$ fit to the training set $X$.

treat this as a regression pb:

↳ loss function: Mean Square Error (MSE).

$$J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x, \theta))^2$$

↳ form of the model: $f_{LIN}(x, w, b) = x^T w + b$.

$w = 0, \ b = 1/2$ → outputs $1/2$ partout!!

$$\underline{FFNN}$$

$$\begin{cases} h = f^{(1)}(x, W, c) \longrightarrow \underline{\text{hidden layer}} \\ y = f^{(2)}(h, w, b) \longrightarrow \underline{\text{output layer}} \end{cases}$$

$$\downarrow \quad h = g(w^T x + b)$$

$$\nearrow \quad \text{\underline{fonction d'activation}}$$

$$\underline{\underline{ReLU}} \quad g(3) = \max\{0; 3\}$$

$$f(x \in \{0,1\}^2, W, c, w, b)$$

$$= w^T g\left((W^T x + c)\right) + b.$$

$$\rightarrow f^{(1)} = g(W^T x + c)$$
$$f^{(2)}(h) = w^T h + b.$$

$$= \boxed{w^T \max\left\{0, W^T x + c\right\} + b}.$$

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \quad b = 0$$

# II / Gradient Based Learning

↳ Billions of Parameters!
↳ Billions of Training Data!
   ↳ gradient-based algo to learn $\theta$.

⚠ loss f° associated to NNs are not **convex**.

↳ iterative methods that drive the loss function $J(\theta)$ towards a very low value.

$\llcorner$ Sensibility $\cancel{of}$ to parameters Intializa°.

$\underline{W^T a + b}$  $\otimes$ $\underline{\text{W: weights}}$ $\rightsquigarrow$ small random values.

$\otimes$ $\underline{b: \text{biais}} \rightsquigarrow$ zero or very small positive values

$\boxed{\text{Training algo}}$

$\Ls$ loss function $\underline{J(\theta)} \rightarrow$ to choose (depending on the pb)

$\Ls$ $\underline{f(x, \theta) = y}$ $\longrightarrow$ $\underline{\text{choice of the NN architecture}}$

NN architecture
- ⊕ # of layers: DEPTH
- ⊗ dim of each layer: WIDTH
- ⓧ f° d'activa° pour chaque couche.

## II.A Cost function

### Maximum likelihood principle

⤷ Cost f° = negative log likelihood.

$$\boxed{J(\theta) = - \mathbb{E}_{x,y \sim \hat{p}_{data}} \log P_{model}(y|x)}$$

training data.

type of SL pbs
supervised learning

ⓧ regression problem ⟶ $y \in \mathbb{R}$.

ⓧ classifica° problem ⟶ category $y$ (discrete)
  ↘ classifica° 2 classes
  ↘ multiclass classifica° ($c > 2$)

<u>Regression case</u>

$$p_{model}(y|x) = N\left(y, f(x,\theta), I\right)$$

target

identity matrix

$\uparrow$ output of NN

**MSE**

$$J(\theta) = \frac{1}{2} E_{x,y \sim \tilde{p}_{data}} \|y - f(x,\theta)\|^2 + \text{cste}$$

$$x \in \mathbb{R}^k \longrightarrow f\left(x \in \mathbb{R}^k, \mu, \Sigma\right)$$

$$= \frac{\exp\left(-\frac{1}{2}(x-\mu)\Sigma^{-1}(x-\mu)\right)}{\sqrt{(2\pi)^k |\Sigma|}}$$

## Advantage of nl. principle

$\hookrightarrow$ undo the 'exp' : $f°$ d'activation have
exponential terms

$\hookrightarrow$ No saturation

$\hookrightarrow$ better for Gradient learning.

ML Principle $\longrightarrow$ Pmodel $\Rightarrow$ we do not learn all

$$P_\theta(y|x)$$

Conditional Statistics

$\longmapsto$ mean of Pmodel .

$J(\theta)$ dep $f(x,\theta) \longrightarrow$ dependant of the output layer

$h = f'(x,\theta)$

↓

output of the last ~~out~~ hidden layer

$\hookrightarrow$  $y$? $\longrightarrow$ shape of the output layer?

## 1. Linear output layer

$$\hat{y} = W^T \times h + b.$$
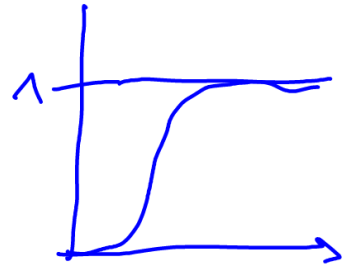
hidden representa°

$\rightarrow$ mean of Gaussian Output Distrisi°

$\rightarrow$ the output layer for a regression problem

## 2. sigmoid output layer

output layer $\rightarrow$ Pmodel = Bernoulli Output Distri°

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\hat{y} = \sigma(\omega^T h + b)$$

$\llcorner$ For tasks requiring a binary variable $\hat{y} \in \{0,1\}$

$\llcorner$ $\underline{p(y=1|x)}$ $\rightarrow$ si $\underline{p(y=1|x) \leq 0.5}$ $y=0$
$\underline{p(y=1|x) > 0.5}$ $y=1$

$\Downarrow$

output layer used for BINARY CLASSIFICA° PB

## 3. Softmax output layer

$$(\text{softmax}(z))_i = \frac{\exp(z_i)}{\sum_{j \in \{1..d\}} \exp(z_j)} \in [0,1]$$

$z \in \mathbb{R}^d$

$\xrightarrow{} p_{model} \longrightarrow$ distribution multinomiale
(Multinouli Probability Distrib°)

$$\hat{y} = \text{softmax}(W^T h + b) \in [0,1]^d$$ output layer for
multi-class classification

N class classification pb:

$$\hat{y} \in [0;1]^N \longrightarrow \hat{c} = \text{argmax } \hat{y}$$

$$\underbrace{\hat{c} = \text{argmax } \hat{y}}_{\text{prediction de la classe}}$$

## II C/ activation functions for hidden layers

which activator $f$°?  (diff)

hidden layer: $\boxed{g}\left( z = \textcircled{$W$}^T x + \textcircled{$b$} \right) \in \mathbb{R}^h$  parameters of the layer

$z \in \mathbb{R}^h$     weights     biais

$x \in \mathbb{R}^d$

standardiser: $\dfrac{x - \mu}{\sigma}$     $\mu = $ mean
$\sigma = $ ecart-type.
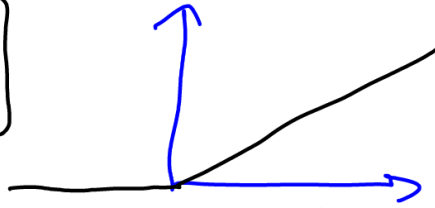
normaliser $\rightarrow$ avec-le max
et le min.

$$f(x, \Theta) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$$

output layer

$$h_2 = g_2(W_2^T h_1 + b_2) \qquad h_1 = g_1(W_1^T x + b_1)$$

## ① fonction ReLU

$$\boxed{g(3) = \max\{0, 3\}}$$

↳ easy to optimize (~ linear)
↳ drawback: cannot learn via gradient-based methods for examples equal to 0.

$\rightarrow$ | Leaky ReLU
| Maxout units

## logistic sigmoid & hyperbolic tangent

Sigmoid
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\tanh(z) = \frac{e^{2z}-1}{e^{2z}+1}$$

$\}\rightarrow$ 2 fonctions d'activation, avant découverte de ReLU.

$\hookrightarrow$ DRAWBACKS: when $z$ is very positive, saturate to high values

when $z$ is very negative, saturate to very low values.

→ donnée tabulaire

type du réseau

FFNN

CNN : Convolutional Neural Network → Images

RNN : Recurrent Neural Network

↳ type de donnée d'entrée
(structure de la donnée)

→ texte
① séries temporelles

on fixe type du réseau

⓵ depth = nbre de couches

② width = ⓐ dimension de chaque couche ⓑ nbre d'unités / de neurons de la couche

$i^{ème}$ couche cachée

$$h^{(i)} = g^{(i)} \left( w^{(i)T} h^{(i-1)} + b^{(i)} \right)$$

size of the matrix.

## Universal Approximation Properties / Theorem

States that a FFNN with a linear output layer and at least one hidden layer with any "squashing" activation function can approximate any Borel $f(\mathcal{D}_i \subset \mathbb{R}^n)$ measurable function from one dimensional space $\mathcal{D}_i \subset \mathbb{R}^n \to \mathcal{D}_o \subset \mathbb{R}^n$ to another

with any desired non-zero amount of error provided that the network has enough units

In theory, this theorem means that we can consider a 'one-layer' FFNN

↳ Ⓐ one hidden layer

Ⓕ one output layer.

BUT in practice we cannot guarantee that the training algo will be learn that function

↳ An exponential # of units may be required!

↳ use DEEPER NETWORKS instead.

# III/ The Backpropagation algorithm

Input
$x \longrightarrow \hat{y}$ Output : input $x$ provids initial information
FFNN

and then propogates up to the hidden layers until

finally producing $\hat{y}$ (a) FORWARD PROPAGA°

(X) Pass Forward

**Backpropagation:** allow the information to flow backward through the network, to compute the gradient

Gradient-based learning $\longrightarrow$ loss function $J(\theta)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \longrightarrow \theta^*$ that minimizes $J(\theta)$

$$\boxed{\theta_k = \theta_{k-1} - \alpha \nabla J(\theta)}$$

$\quad\quad\quad\quad\quad\quad\nearrow \quad$ once

learning rate / tx d'apprentissage.

how to compute efficiently this gradient?

$$\text{FNN}: f(x, \theta) = f^{(n)}\left(f^{(n-1)} \cdots f^{(1)}(x)\right)$$

Based on <u>chain rule of calculus</u> :

$3 = f(g(x))$
$\quad = f(y)$

$$\frac{\partial 3}{\partial x} = \frac{\partial 3}{\partial y} \times \frac{\partial y}{\partial x} \qquad \text{scalar case}$$

<u>Generaliza° to $\mathbb{R}^n$</u>

$x \in \mathbb{R}^n = (x_1 \cdots x_n)$

$$\frac{\partial 3}{\partial x_i} = \sum_{\jmath} \left(\frac{\partial 3}{\partial y_{\jmath}} \times \frac{\partial y_{\jmath}}{\partial x_i}\right)$$

Vector notation :

$$\nabla_x \mathfrak{z} = \left(\frac{\partial y}{\partial x}\right)^T \nabla_y \mathfrak{z}$$

$\in \mathbb{R}^n$

$\in \mathbb{R}^m$

$\mathfrak{z} = f(\underbrace{g(x)}_{y})$

Matrice Jacobienne
$n \times m$