# I. Introduction

## Supervised binary classification:

### Gen approach:

How to model $p(x_i, y_i)$?

Naive Bayes, Generative /Hidden Markov Model

### Discriminative approach

Which feature allow to distinguish?
→ draw a boundary
Decision Trees, SVMs

# II The formal neuron and Perceptron

Perceptron: $\hat{y}_i = sign(\omega^T x_i + b)$, $\theta = \{\omega, b\}$ parameters of the model

## Learning procedure

Init: $\omega_t = 0$

Boucle: . $\forall i \in [1, n]$

· $\hat{y}_i = sign(\omega^T x_i)$

· Si $y_i \neq \hat{y}_i$ : $\omega_{t+1} = \omega_t + sign(y_i) x$

Algo CV si les classes sont linéairement séparables
→ $\frac{R^2}{\gamma^2}$ it

Limites : → Quality of the boundary
→ Case of non-linear separability

Loss function: . zero-one : $| y - sign(f(x, \omega))$

. hinge : $max(0, 1 - y \cdot f(x, \omega))$

. MSE : $(y - f(x, \omega))^2$

Sigmoid : on remplace sign par $\sigma(x) = \dfrac{e^x}{1 + e^x}$ → différentiable

$\Rightarrow \hat{y}_i = \sigma(\omega^T x_i)$
↳ c'est maintenant une proba

## Probabilistic loss function : MLE estimation

$$\ell_{MLE}(\theta) = - \log \mathcal{L}(\theta) \quad \text{avec} \quad \mathcal{L}(\theta) = \Pi \, f(x_i, \omega)^{y_i} \, (1 - f(x_i, \omega))^{1-y_i}$$

## Gradient :

$$\frac{\partial \ell_{MLE}}{\partial \omega^j} = - \sum [y_i - f(x_i, \omega)] x_i^j$$

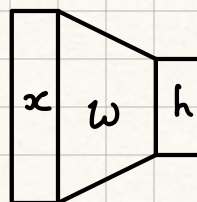## Gradient - based learning procedure :

- While $\| \ell_{MLE}(\omega_{t+1}) - \ell_{MLE}(\omega_t) \| > \delta$ et $t > m_{iter}$ :

  - Pour $j \in [1, d]$ :
    $$\omega_{t+1}^j = \omega_t^j + \varepsilon \sum_n^d [y_i - f(x_i, \omega)] x_i^j$$

## III _ Multi - Layer Perceptron

Newrons are grouped into layers :



$$h^{(l)} = \phi^{(l)} ( \omega^{(l)T} x + b^{(l)} )$$
        ↳ fonction d'activation possiblement non - linéaire

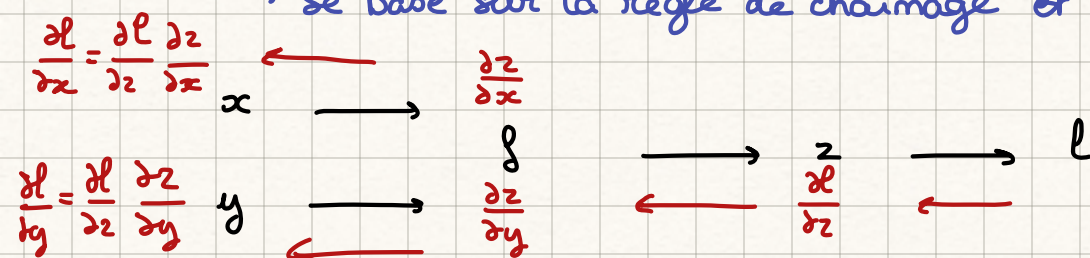↳ chaque layer est une fonction du layer précédent

## Activation :

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad , \quad \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- Softmax : $\hat{o}_i = \left[ \dfrac{e^{h^{(l),i}}}{\sum e^{h^{(l),k}}} , \dots \right]$

  ↳ généralisat° de la sigmoid pour k sorties

## Optimisation : Backpropagation

→ se base sur la règle de chainage et fonctionne récursivement

$$\frac{\partial \ell}{\partial x} = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial x}$$

$$\frac{\partial \ell}{\partial y} = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial y}$$

# IV - Model and Hyperparameter selection

→ Split data into train/test et choisis les hyperparamètres

    ↓ MIEUX

→ Split into train / validation / test

        <span style="color:red">to tune↑ h.p</span>        <span style="color:red">↑ to verify perf</span>

## Choice of learning rate $\varepsilon$ :

- Si perf ↗ sur validat° data : $\varepsilon = \varepsilon * \alpha$     $\alpha \in [0,1]$

ou    · Step decay :   $\times \alpha$ tous les $N$ epochs

ou    · exponential decay : $\varepsilon^{(t)} = \varepsilon^{(0)} e^{-\alpha t}$

## Regularisation : To avoid overfitting

$$\ell_{reg}(\omega) = \ell_{MLE}(\omega) + \lambda\, \omega^T \omega$$

$$\Rightarrow \omega^{t+1} = (1 - 2\lambda)\,\omega^t - \varepsilon\, \nabla_\omega \ell_{MLE}$$

## Initialisation : dépend beaucoup des activat° & des méthodes d'optim

     mais souvent    $\omega^{(i)} \sim \mathcal{N}(0, \frac{1}{\sqrt{d_i}})$

## Advantages : 
→ very flexible in term of input/output

→ backpropagat° allows efficient updates

→ SGD allows efficient training m̂ pr grosse data

## Mais : 
→ Gradient descent pas facile à executer

→ bcp d'hyper paramètres

→ risque overfitting ( not enough data)

→ optimisat° hard → underfitting

## Solution : 
<u>underfitting</u> : more data

                  better optim

       <u>overfitting</u> : find better way to regularise

           <span style="color:blue">→ unsupervised : init hidden layers with unsupervised learning</span>

→ dropout : randomly remove parts of hidden layers

**Better optimization:** Adaptative gradient

**Better activation** : ReLU : $relu(x) = max(0, x)$
→ avoid saturated & vanishing grad.

**Better init:** Unsupervised pre-training
⤷ fine tuning = ajout d'un output et entrainement
supervisé après cette init

**Better reg** : Dropout