

# CLASSIFIEURS LINÉAIRES

Guillaume Wisniewski  
guillaume.wisniewski@limsi.fr

janvier 2016

## Première partie I

### CONTEXTE / RAPPEL

## CLASSIFICATION SUPERVISÉE

### LA TÂCHE

- observations  $\mathbf{x} \in \mathcal{X}$  associées à des étiquettes  $y \in \mathcal{Y}$
- ensemble fini d'exemples :  $(\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N$
- peut-on déterminer une fonction **généralisant** les exemples à l'ensemble des observations possibles  $\mathcal{X}$

### REPRÉSENTATION DES EXEMPLES

- les exemples  $\mathbf{x}$  sont représentés par des éléments de  $\mathbb{R}^n$
- chaque dimension correspond à une **caractéristique** (*feature*, attribut)
- les représentations sont choisies arbitrairement

## Deuxième partie II

### INTÉRÊT

## 1) INTERPRÉTATION GRAPHIQUE

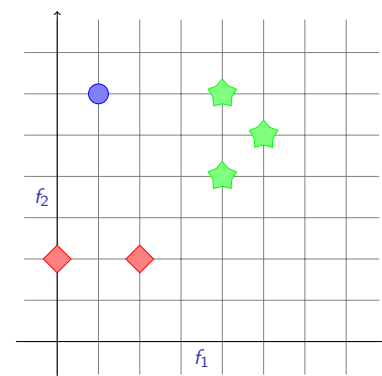
### PRINCIPE

- les exemples peuvent être représentés dans un espace euclidien
  - ▶ chaque axe correspond à une caractéristique
  - ▶ exemple = point / vecteur

### EXEMPLE

- classe  $\omega_1$  :  $\begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}$
- classe  $\omega_2$  :  $\begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 2,5 \\ 2,5 \end{pmatrix}$
- classe  $\omega_3$  :  $\begin{pmatrix} 0,5 \\ 3 \end{pmatrix}$

## REPRÉSENTATION



## 2) DISTANCE ENTRE POINTS

- espace euclidien  $\Rightarrow$  on peut définir une « distance » entre point
- infinité de distances avec des propriétés différentes
- exemple :

$$\|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{\sum_{i=1}^n (x_i - x'_i)^2} \quad \|\mathbf{x} - \mathbf{x}'\|_1 = \sum_{i=1}^n |x_i - x'_i|$$

- plus des points sont proches plus ils sont similaires

## CONSÉQUENCE

Comment classifier un point ?

- on attribut à l'exemple non étiqueté l'étiquette de l'exemple le plus proche
- plus « robuste » : regarder l'étiquette des  $k$  plus proches voisins  $\oplus$  vote majoritaire

$k$ -ppv /  $k$ -nn

## AUTRE MÉTHODE

### CLASSIFIEUR À DISTANCE MINIMALE

- 1 on détermine un **représentant** de chaque classe, typiquement la moyenne :

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$$

- 2 on attribut l'étiquette de la classe dont le représentant est le plus proche

avantages / inconvénients par rapport aux  $k$ -ppv ?

## LIMITE DES CLASSIFIEURS À BASE DE DISTANCE

### COMPLEXITÉ

- complexité des  $k$ ppv :  $\mathcal{O}(N \times d)$  ( $N$  : nombre d'exemples,  $d$  : dimension)
- Problème pour les problèmes en grande dimension

### Curse of dimensionality

- quand la dimension augmente, les distances entre deux paires d'objets arbitraires tendent vers la même valeur :

$$\lim_{d \rightarrow \infty} \frac{\text{dist}_{\min} - \text{dist}_{\max}}{\text{dist}_{\min}} = 0$$

- intuition des  $k$ ppv « moins » valide

## Troisième partie III

## PRINCIPE DES CLASSIFIEURS LINÉAIRES

## AVERTISSEMENT



On se limitera pour aujourd'hui à la classification binaire :  $y = \pm 1$

## RAPPEL

### ÉLÉMENT D'UNE APPROCHE D'APPRENTISSAGE

- 1 des exemples étiquetés
  - ▶ ensemble de vecteurs  $\oplus$  étiquettes
- 2 une mesure d'évaluation
- 3 une classe d'hypothèses  $\oplus$  un biais inductif

### CLASSE D'HYPOTHÈSES

- **règle de décision** : comment choisir la sortie associée à une entrée
- **méthode d'apprentissage** : comment choisir un élément dans cette classe de fonctions (biais inductif)

## FONCTION DE DÉCISION

### FONCTION DE SCORE

$$F(\mathbf{x}; \mathbf{w}) = \mathbf{w} \cdot \mathbf{x}$$

- rappel : la représentation inclut un terme constant
- fonction paramétrée par un vecteur  $\mathbf{w}$
- combinaison linéaire des caractéristiques : chaque caractéristique a un **poids** par rapport à son importance

### FONCTION DE DÉCISION

$$y^* = \text{sign } F(\mathbf{x}; \mathbf{w}) \\ = \begin{cases} 1 & \text{si } F(\mathbf{x}; \mathbf{w}) \geq 0 \\ -1 & \text{sinon} \end{cases}$$

## FRONTIÈRE DE DÉCISION

### DÉFINITION

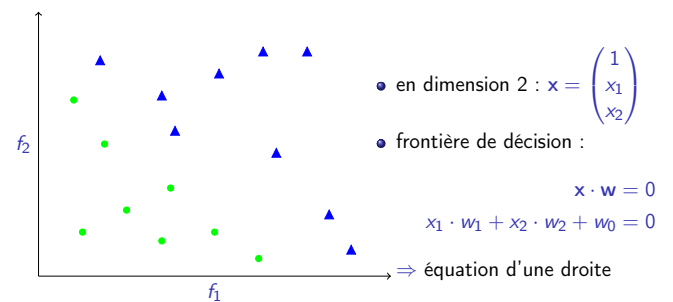
- caractéristique d'une classe de fonctions
- = quand le critère de décision « change » / quand on ne sait pas prendre de décision

### DANS NOTRE CAS...

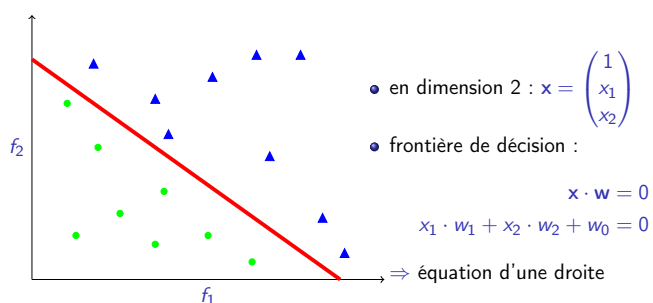
$$F(\mathbf{x}, \mathbf{w}) = 0$$

$\Rightarrow$  équation d'un hyperplan = **hyperplan séparateur**

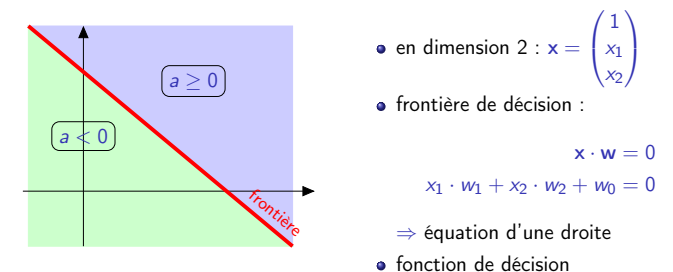
## GRAPHIQUEMENT



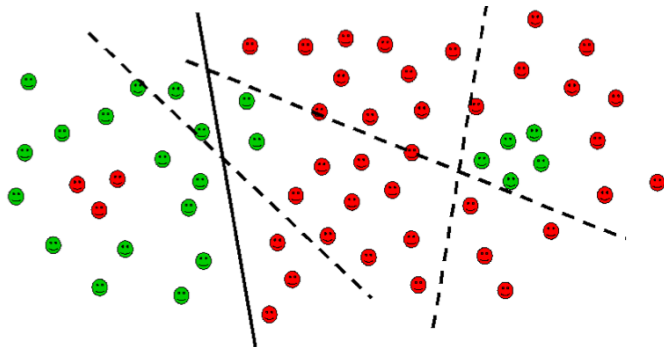
## GRAPHIQUEMENT



## GRAPHIQUEMENT

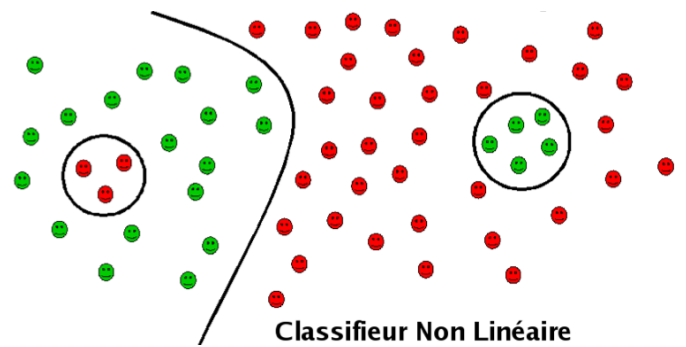


## POUVOIR D'EXPRESSION



Tous les ensembles de données ne sont pas linéairement séparables

## CLASSIFICATION NON LINÉAIRE



**Classifieur Non Linéaire**

cf. cours sur la classification non linéaire

## APPRENTISSAGE

Comment choisir le paramètre  $w$

### OBJECTIFS

- faire le moins d'erreurs sur des données **nouvelles**
- à choisir à partir d'un ensemble fini d'exemples (ensemble d'apprentissage)

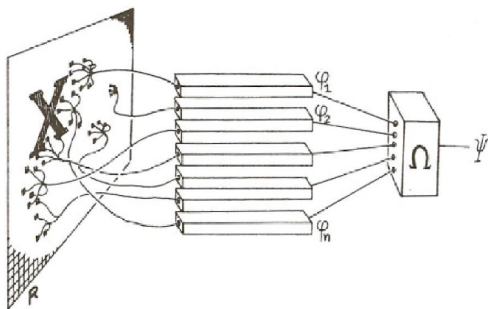
### MÉTHODES

- perceptron
- modèle à maximum d'entropie (Maxent)
- machines à vecteurs de support (SVM)
- Winnow
- ...

## Quatrième partie IV

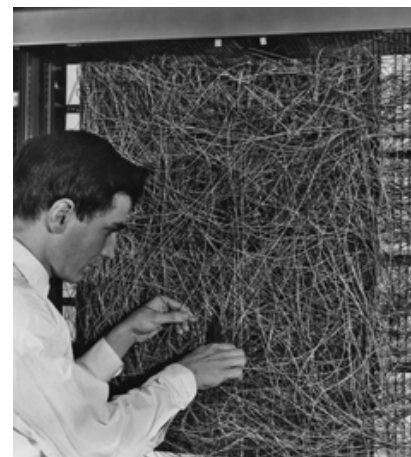
## LE PERCEPTRON

## ASPECT HISTORIQUE (1)



- Rosenblatt, 1957
- objectif = reproduire les capacités visuelles d'un cerveau humain

## LE PERCEPTRON EST UNE MACHINE !



## APPRENTISSAGE

### Apprentissage par correction d'erreurs (Rosenblatt, 57)

- pour chaque exemple de l'ensemble d'apprentissage
- classer l'exemple
- si la sortie est correcte : ne rien faire
- si la sortie est fausse : corriger les paramètres pour ne plus faire d'erreur

## PLUS FORMELLEMENT

**Require:** a training set  $(\mathbf{x}^{(i)}, y^{(i)})_{i=1}^n$  and a learning rate  $\eta$

```
1:  $\mathbf{w} \leftarrow 0$ 
2: while there are classification errors do
3:   for  $i = 1$  to  $n$  do
4:      $y = \text{sign}(\mathbf{w} \cdot \mathbf{x}^{(i)})$ 
5:     if  $y \neq y^{(i)}$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + \eta \cdot y^{(i)} \cdot \mathbf{x}^{(i)}$ 
7:     end if
8:   end for
9: end while
```

## PROPRIÉTÉS « THÉORIQUES »

- si l'ensemble d'apprentissage est linéairement séparable :
  - ▶ converge en un nombre fini d'opérations (ni  $\eta = 1$ )
  - ▶ aucune garantie en généralisation
- si l'ensemble d'apprentissage n'est pas linéairement séparable :
  - ▶ algorithme ne converge pas
  - ▶ plusieurs « astuces », aucune preuve
- comment savoir si l'ensemble d'apprentissage est linéairement séparable ?  $\Rightarrow$  réponse uniquement empirique