

Sciences des données

Un voyage initiatique

Cécile CAPPONI, Rémi EYRAUD, Hachem KADRI

LIS, Aix-Marseille Université, CNRS
Equipe QARMA



M1 Informatique

Outline

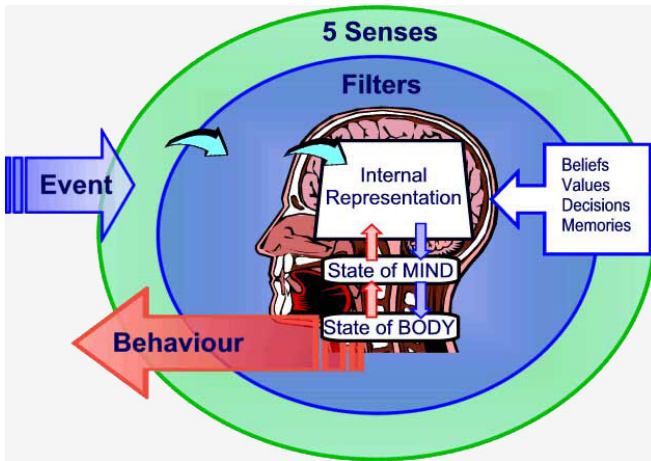
1 Classification

- De quoi parlons-nous ?
- Construire un bon modèle de classification à partir de données observées
- Les arbres de décision
- Les k -plus proches voisins

- De quoi parle-t-on ?

[illegible]

Construction de modèles de prédiction *automatique*



(source : i-change.biz)

Construction de modèles de prédiction *automatique*

```

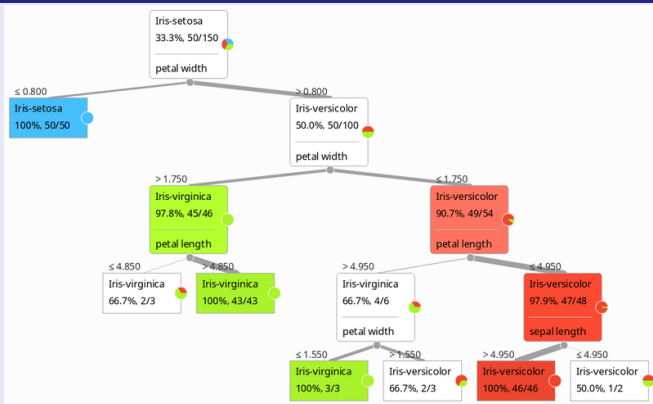
graph LR
    A[Monocotylédones] --> B[Fleurs irrégulières]
    A --> C[Fleurs régulières]
    B --> D[Orchidacées]
    C --> E[Fleurs bien distinctes et colorées]
    C --> F[Fleurs discrètes à aspect d'herbe]
    E --> G[Ovaire supère et 6 étamines]
    E --> H[Ovaire infère]
    G --> I[Liliacées]
    H --> J[3 étamines]
    H --> K[6 étamines]
    J --> L[Iridacées]
    K --> M[Amaryllidacées]
    F --> N[Tige ronde et pleine]
    F --> O[Tige ronde et creuse avec nœuds]
    F --> P[Tige triangulaire pleine sans nœuds]
    N --> Q[Joncacées]
    O --> R[Poacées]
    P --> S[Cyperacées]
  
```

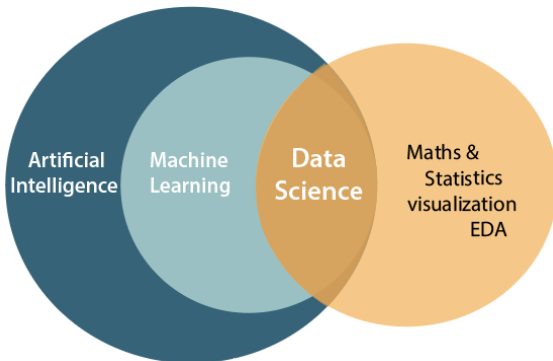
Le diagramme illustre la classification des Monocotylédones en fonction de leurs caractéristiques morphologiques et florales. Les Monocotylédones sont divisés en deux groupes principaux : ceux à fleurs irrégulières (Orchidacées) et ceux à fleurs régulières. Les fleurs régulières sont à leur tour divisées en fleurs bien distinctes et colorées (Liliacées, Iridacées, Amaryllidacées) et en fleurs discrètes à aspect d'herbe (Joncacées, Poacées, Cyperacées). Les fleurs bien distinctes et colorées sont caractérisées par un ovaire supère et 6 étamines, ou un ovaire infère et 3 ou 6 étamines. Les fleurs discrètes à aspect d'herbe sont caractérisées par une tige ronde et pleine (Joncacées), une tige ronde et creuse avec nœuds (Poacées), ou une tige triangulaire pleine sans nœuds (Cyperacées).

La classification : au coeur de l'apprentissage automatique

Construction de modèles de prédiction *automatique*

Construction automatique à partir d'observations (induction)





Apprentissage machine = un moteur de la science des données

- Pour aller plus loin que les statistiques descriptives
- Objectifs : algorithmes pour construire des modèles numériques prédictifs à partir des données observées, qui généralisent bien à des données futures

ALGO – DONNÉES – MODÈLE – PRÉDICTION – GÉNÉRALISATION



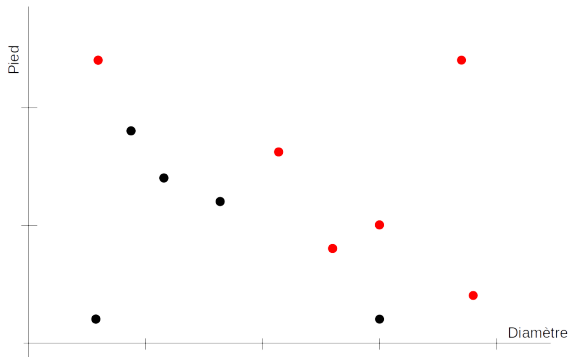
p	x	s	n	t	...	k	s	u
e	x	s	y	t	...	n	n	g
e	b	s	?	t	...	n	n	m
...
p	x	y	w	t	...	k	s	u
e	x	y	y	t	...	k	n	g

Trouver la meilleure

$$h: \mathcal{X} \rightarrow \{\mathbf{e}, \mathbf{p}\}$$

A partir de

$$S = \{(x_i, y_i)\}, x_i \in \mathcal{X}, y_i \in \{\text{e}, \text{p}\}$$



Journal of Management Inquiry 20(6) 798–814

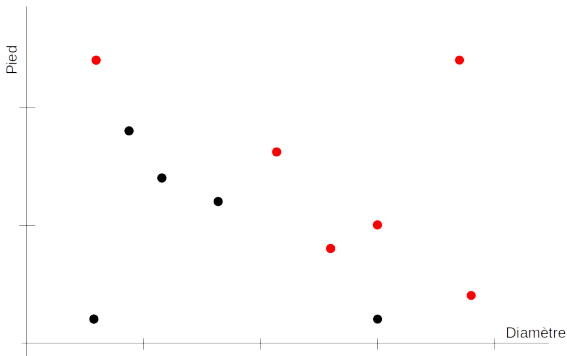
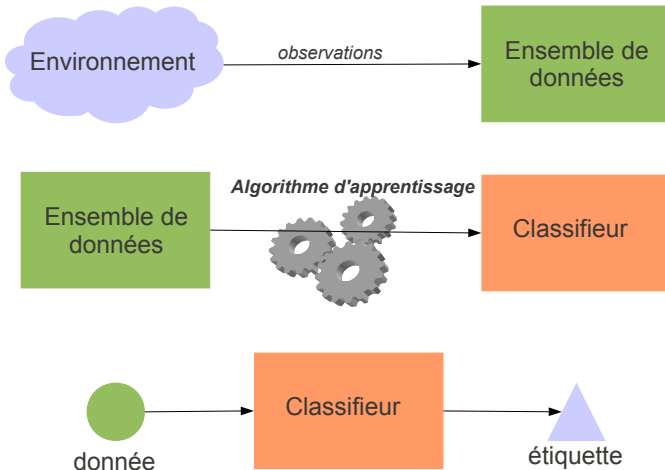
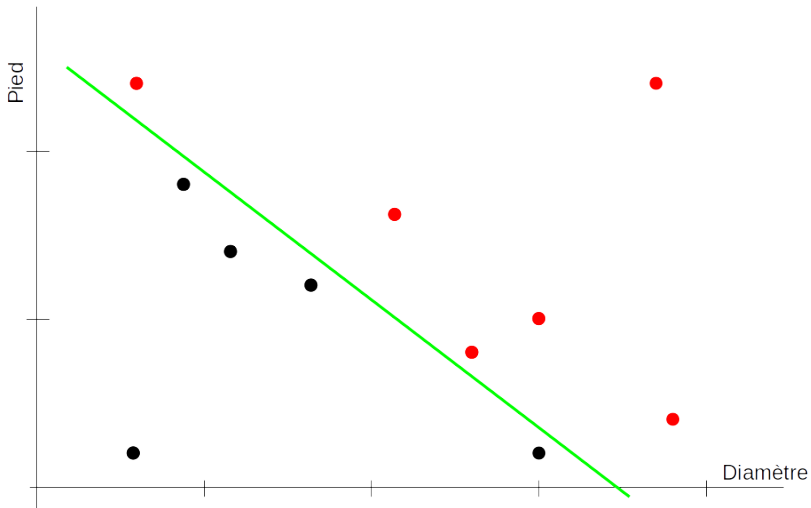


Schéma général de la classification supervisée

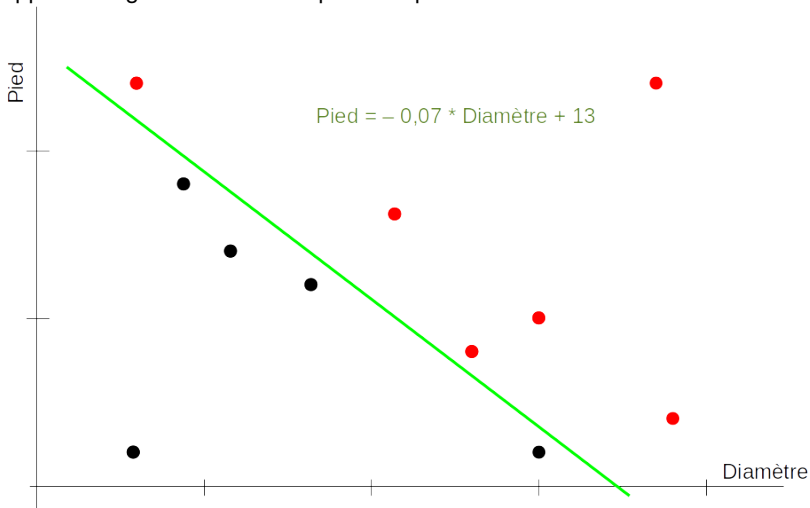




[illegible]

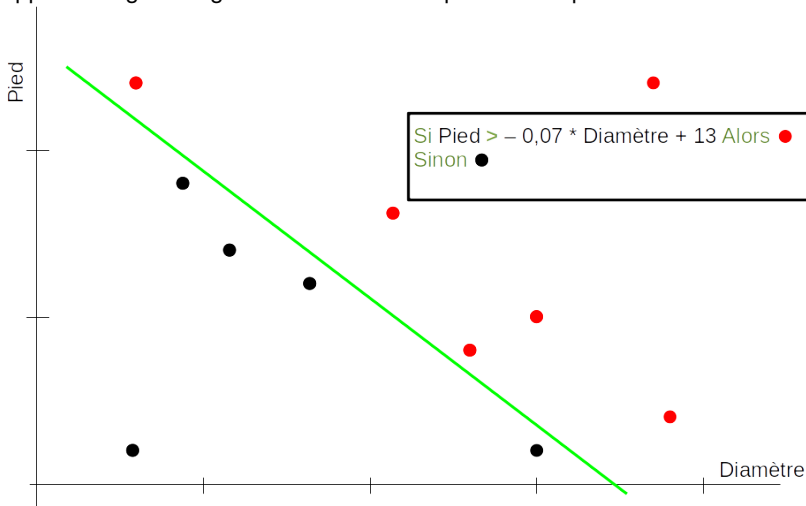
A scatter plot showing the relationship between 'Diamètre' (Diameter) on the x-axis and 'Pied' (Foot) on the y-axis. The data points are categorized into two groups: red dots and black dots. A green linear regression line is drawn through the data, with the equation $\text{Pied} = -0,07 * \text{Diamètre} + 13$ displayed on the plot.

Diamètre (cm)	Pied (cm)	Category
10	10	Black
12	12	Black
14	11	Black
16	10	Black
18	9	Black
20	8	Black
22	7	Black
24	6	Black
26	5	Black
28	4	Black
30	3	Black
32	2	Black
34	1	Black
36	0	Black
38	-1	Black
40	-2	Black
42	-3	Black
44	-4	Black
46	-5	Black
48	-6	Black
50	-7	Black
52	-8	Black
54	-9	Black
56	-10	Black
58	-11	Black
60	-12	Black
62	-13	Black
64	-14	Black
66	-15	Black
68	-16	Black
70	-17	Black
72	-18	Black
74	-19	Black
76	-20	Black
78	-21	Black
80	-22	Black
82	-23	Black
84	-24	Black
86	-25	Black
88	-26	Black
90	-27	Black
92	-28	Black
94	-29	Black
96	-30	Black
98	-31	Black
100	-32	Black
10	15	Red
20	10	Red
30	5	Red
40	0	Red
50	-5	Red
60	-10	Red
70	-15	Red
80	-20	Red
90	-25	Red
100	-30	Red



[illegible]

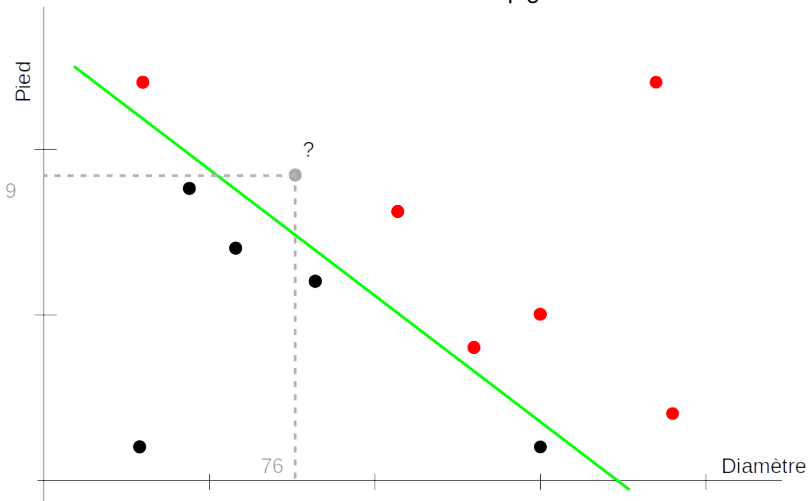
A scatter plot showing the relationship between 'Diamètre' (Diameter) on the x-axis and 'Pied' (Foot) on the y-axis. The data points are categorized into two groups: 'Alors' (red dots) and 'Sinon' (black dots). A green line represents the decision boundary, defined by the equation $\text{Si Pied} > -0,07 * \text{Diamètre} + 13$. The plot shows that points above the line are classified as 'Alors', while points below the line are classified as 'Sinon'.

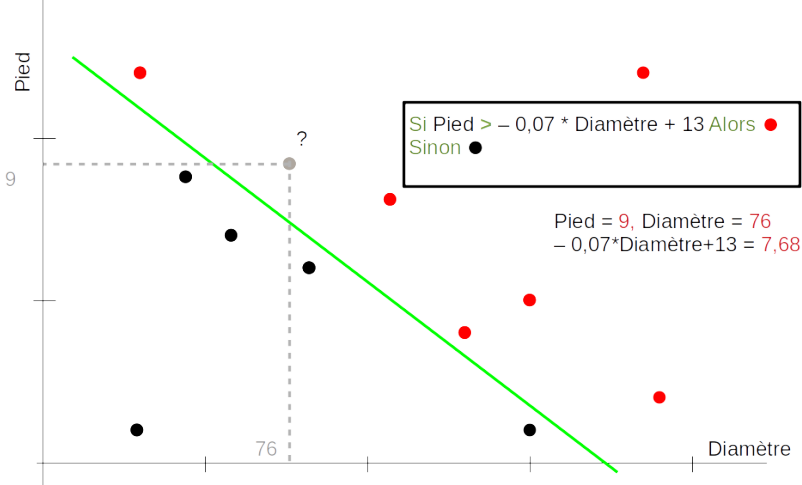


[illegible]

A scatter plot showing the relationship between 'Diamètre' (Diameter) on the x-axis and 'Pied' (Foot) on the y-axis. The data points are colored red and black. A green regression line is drawn through the data. A specific point is highlighted with a grey dot and a question mark, located at a diameter of 76 and a foot value of 9. Dashed lines indicate these coordinates on the axes.

Diamètre	Pied	Color
76	9	Grey (Point of Interest)
~72	~10.5	Red
~74	~9.5	Black
~76	~9	Black
~78	~8.5	Black
~80	~9	Red
~82	~7.5	Red
~84	~8	Red
~86	~6.5	Red
~88	~10.5	Red
~72	~6.5	Black
~86	~6.5	Black

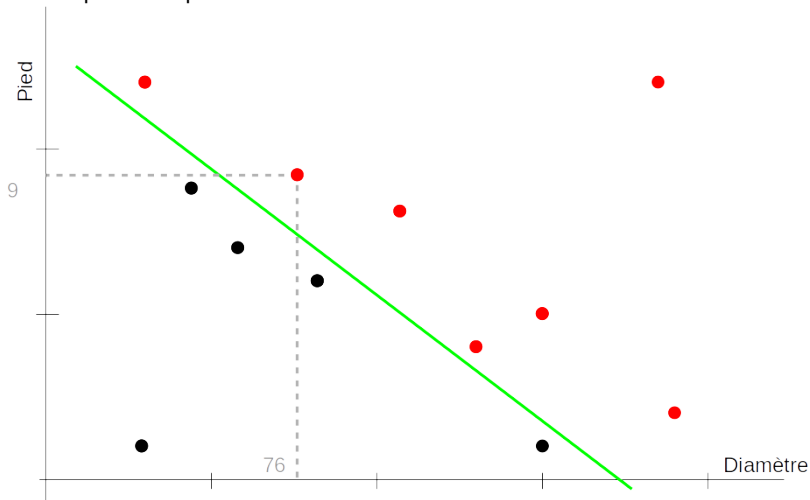




[illegible]

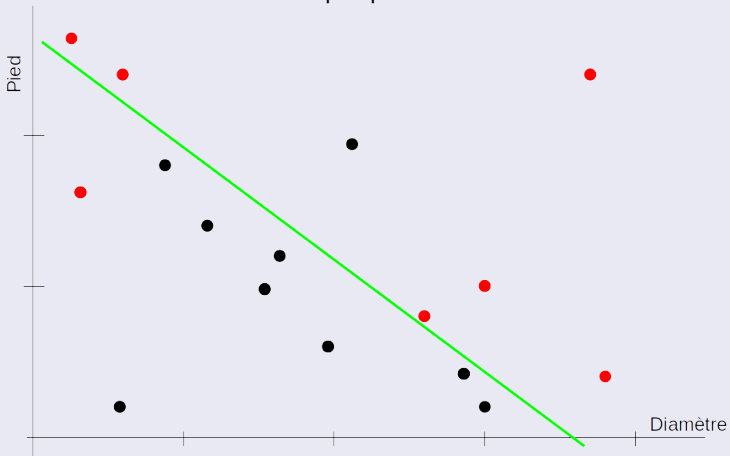
A scatter plot showing the relationship between 'Diamètre' (Diameter) on the x-axis and 'Pied' (Foot) on the y-axis. The data points are colored red and black. A green regression line is drawn through the points. A specific point is highlighted with dashed lines extending to the axes, where the x-axis value is 76 and the y-axis value is 9.

Diamètre (x)	Pied (y)	Color
10	15	Red
20	10	Black
25	9	Black
30	8	Black
35	9	Red
45	8	Black
55	7	Red
65	6	Red
70	5	Black
80	10	Red
85	4	Red



Journal of Management Inquiry 20(6) 798–814

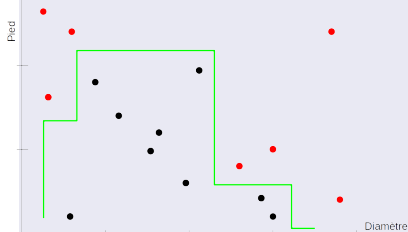
Journal of Management Education 36(7) 809–827



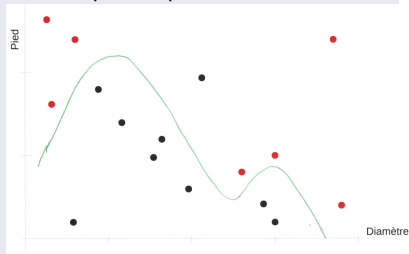
Un modèle simple : la droite (séparateur linéaire)

Pas toujours possible avec les données observées

Chercher un autre type de modèle plus expressif



Linéaire par morceaux

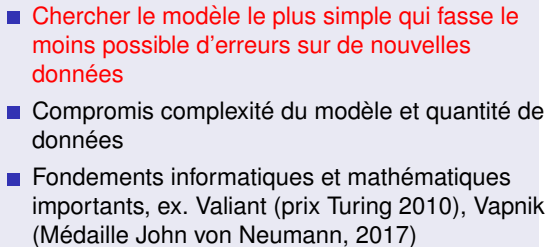


Polynomial $y = \sum_{i=0}^n \mathbf{a}_i x^i$

Mais plus de valeurs à déterminer, donc plus complexe

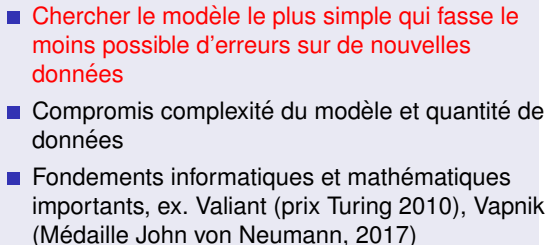
Pas toujours possible avec les données observées

Principe du rasoir d'Occam (principe d'économie / parcimonie)



Pas toujours possible avec les données observées

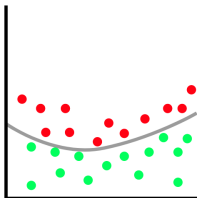
Principe du rasoir d'Occam (principe d'économie / parcimonie)



En pratique, données de dimensions plus grandes

Exemple : image couleur 300×300 = espace de 90 000 dimensions (au lieu de 2)

Pour savoir **généraliser** au mieux



- Accepter de faire quelques erreurs avec un modèle simple
- Tester plusieurs types de modèles, appris avec différents algorithmes (puis sélection de modèles, M1 option IAA)
- Lisser les modèles de prédiction (M2 IAAA)

Formalisation d'un problème de classification supervisée

Apprentissage supervisé – formalisation

- \mathcal{X} espace d'entrée, \mathcal{Y} espace des cibles
- D distribution sur $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$
- $S_{\text{train}} = \{(X_i, Y_i)\}_{i=1}^m$ échantillon de n v.a. **Indépendamment et Identiquement Distribuées (IID)** suivant D
- Fonction de perte $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

But : minimisation du (vrai) risque

A partir de S_{train} trouver $f : \mathcal{X} \rightarrow \mathcal{Y}$ telle que $f = \operatorname{argmin}_h R(h)$ avec

$$R(h) = \mathbb{E}_{X,Y} \ell(h(X), Y) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) dD(x, y)$$

Pour la classification/catégorisation, $|\mathcal{Y}| < +\infty$, $\ell(y, y') = \mathbb{I}(y \neq y')$ et

$$R(h) = \mathbb{P}_{X,Y \sim D}(h(X) \neq Y)$$

© 2004 Blackwell Publishing Ltd *Journal of Internal Medicine* 255: 101–108

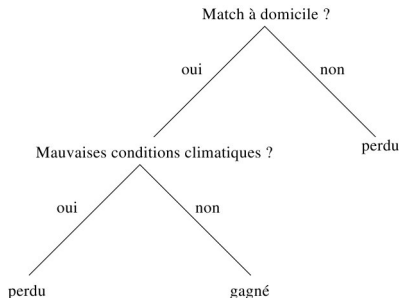
Introduction sur les arbres de décision

Notations

Soient A_1, \dots, A_m : attributs binaires, n -aires, ou réels. L'espace de description est $\mathcal{X} = \prod_{j=1}^m \mathcal{X}_j$ où \mathcal{X}_j est le *domaine* de A_j .

Définition d'un arbre de décision

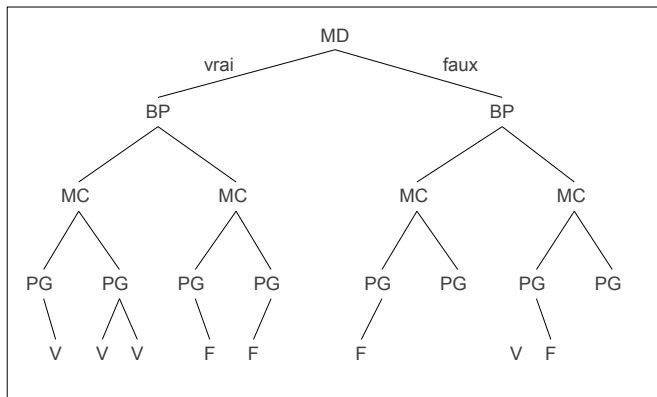
Ensemble de règles de classification basant leur décision sur des tests associés aux attributs, organisés de manière arborescente.



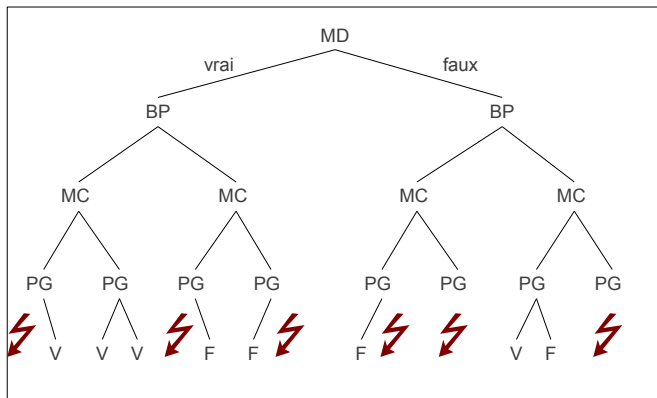
Construire un arbre de décision à partir de S

Match dom.	Balance pos.	Mauvais climat	Match préc. gagné	Victoire
V	V	F	F	V
F	F	V	V	V
V	V	V	F	V
V	V	F	V	V
F	V	V	V	F
F	F	V	F	F
V	F	F	V	F
V	F	V	F	F

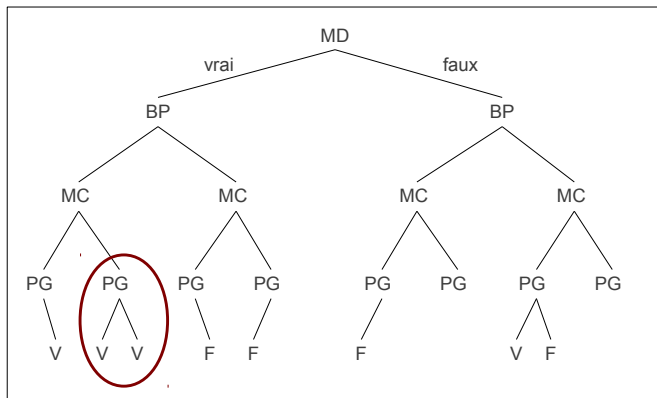
Construction aveugle



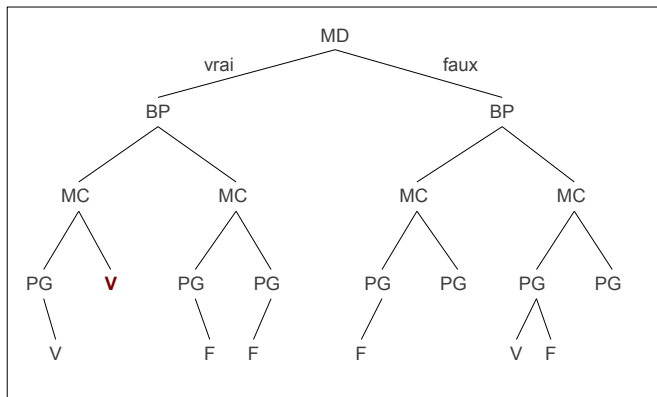
Construction aveugle



Construction aveugle



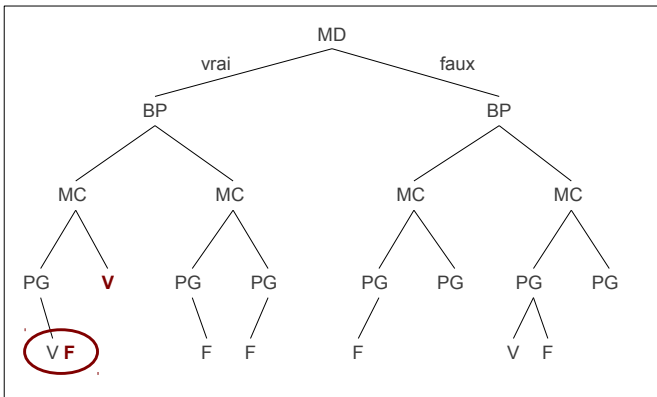
Construction aveugle



Construire un arbre de décision à partir de S

Match dom.	Balance pos.	Mauvais climat	Match préc. gagné	Victoire
V	V	F	F	V
F	F	V	V	V
V	V	V	F	V
V	V	F	V	V
F	V	V	V	F
F	F	V	F	F
V	F	F	V	F
V	F	V	F	F
V	V	V	F	F

Cas de non-déterminisme – Pas d'arbre *parfait*



Arbre de décision de risque empirique minimal ?

- Construire un arbre de décision de risque empirique minimal = possible pour tout S
- Mais très faible capacité prédictive (généralisation). pourquoi ?
- Est-ce que le plus petit arbre de décision compatible avec S aura de meilleures capacités de généralisation ?
- cf. théorie de l'apprentissage statistique (Vapnik)

- Driving the point home: "The only thing that's left is the counts possible"

- **Principe du rasoir d'Occam** : trouver l'hypothèse la plus courte possible compatible avec les données

- *Principe MDL* (minimum description length) : soit des données S , trouver l'hypothèse H telle que $|H| + |S/H|$ soit la plus petite possible (longueur d'un codage des données via H)

Trouver le plus petit arbre de décision compatible avec S est un problème NP-complet.

Plus petit arbre possible (mais pas forcément le plus petit), compatible au mieux avec les données, pour satisfaire le principe ERM.

0 1 0 1 1 0

- CART
- C4.5

10

- 1 Décider si un noeud est terminal,
- 2 Si un noeud n'est pas terminal, lui associer un test
- 3 Si un noeud est terminal, lui associer une classe

[illegible]

```

1: Initialiser l'arbre courant à vide : racine = noeud courant
2: repeat
3:   Décider si le noeud courant est terminal
4:   if noeud terminal then
5:     Lui affecter une classe
6:   else
7:     Sélectionner un test et créer autant de nouveaux noeuds qu'il y a
      de réponses possibles aux tests
8:   end if
9:   Passer au noeud suivant non-exploré (s'il existe)
10: until Plus de noeud sans classe
11: return Arbre de décision  $A$ 

```

- ```

1: Initialiser l'arbre courant à vide : racine = noeud courant
2: repeat
3: Décider si le noeud courant est terminal
4: if noeud terminal then
5: Lui affecter une classe
6: else
7: Sélectionner un test et créer autant de nouveaux noeuds qu'il y a
 de réponses possibles au tests
8: end if
9: Passer au noeud suivant non-exploré (s'il existe)
10: until Plus de noeud sans classe
11: return Arbre de décision A

```

## Structure générique de l'algorithme (cont'd)

### Noeud terminal ?

- Lorsque (presque) tous les exemples de  $S$  en ce noeud sont dans la même classe,
- Lorsqu'il n'y a plus d'attributs à tester à ce niveau

### Quelle classe à un noeud terminal ?

- Classe majoritaire
- Classe la plus représentée, si égalité

### Sélection d'un test ?

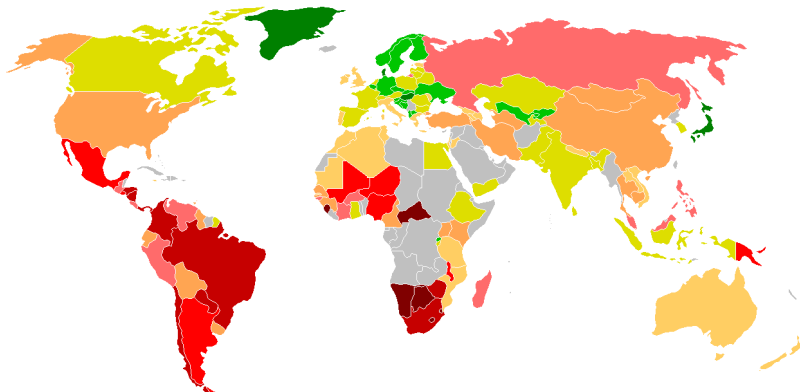
Choix de l'attribut qui fait le mieux progresser la discrimination des données de  $S$  : gain en information.

- Indice de Gini (CART)
- Critère d'entropie (C4.5)

---

# Coefficient de Gini

(Inégalité des revenus – source : wikipedia)



| Color | Gini coefficient |  |             |  |             |
|-------|------------------|--|-------------|--|-------------|
|       |                  |  | 0,35 - 0,39 |  | 0,55 - 0,59 |
|       | < 0,25           |  | 0,40 - 0,44 |  | > 0,60      |
|       | 0,25 - 0,29      |  | 0,45 - 0,49 |  | NA          |
|       | 0,30 - 0,34      |  | 0,50 - 0,54 |  |             |

[illegible]

---

---

---

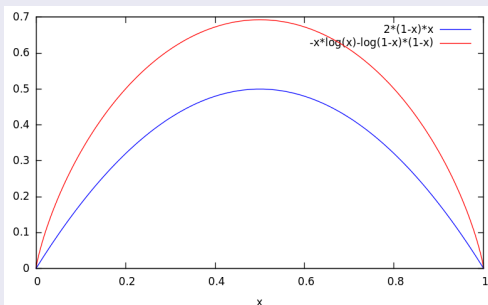
# Indice de Gini et entropie (cont'd)

## Exemple

Supposons  $k = 2$ , soit  $x = \frac{|S_1|}{|S|}$

$$\text{Gini}(S) = 2x(1 - x)$$

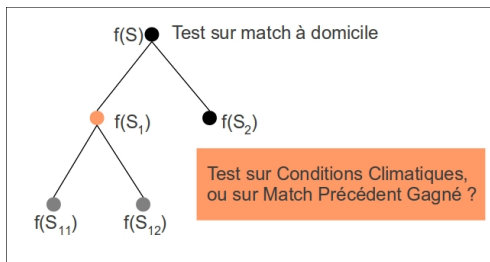
$$\text{Ent}(S) = -x \log x - (1 - x) \log(1 - x)$$



# Comment choisir un test parmi les attributs disponibles ?

## Cas des attributs binaires

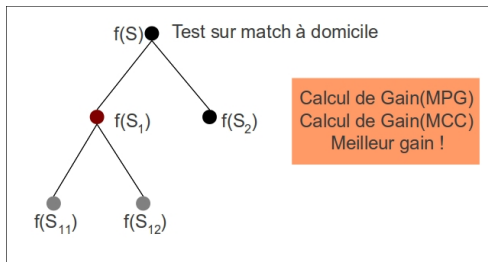
- $p$  : position courante ( $\epsilon$  à la racine) de l'arbre en construction
- $T$  : un test (sur un attribut :  $T$  a deux réponses possibles !)
- $f = \text{Ent}$  ou  $f = \text{Gini}$ ,  $S_p$  échantillon associé à  $p$
- $S_{pi}$  : ensemble des exemples de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$
- $P_i$  : proportion des éléments de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$



# Comment choisir un test parmi les attributs disponibles ?

## Cas des attributs binaires

- $p$  : position courante ( $\epsilon$  à la racine) de l'arbre en construction
- $T$  : un test (sur un attribut :  $T$  a deux réponses possibles !)
- $f = \text{Ent}$  ou  $f = \text{Gini}$ ,  $S_p$  échantillon associé à  $p$
- $S_{pi}$  : ensemble des exemples de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$
- $P_i$  : proportion des éléments de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$





# Comment choisir un test parmi les attributs disponibles ?

## Cas des attributs binaires

- $p$  : position courante ( $\epsilon$  à la racine) de l'arbre en construction
- $T$  : un test (sur un attribut :  $T$  a deux réponses possibles !)
- $f = \text{Ent}$  ou  $f = \text{Gini}$ ,  $S_p$  échantillon associé à  $p$
- $S_{pi}$  : ensemble des exemples de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$
- $P_i$  : proportion des éléments de  $S_p$  qui satisfont la  $i$ -ème branche de  $T$

$$\text{Gain}_f(p, T) = f(S_p) - \sum_{j=1}^2 P_j \times f(S_{pj})$$

## Le gain d'un attribut à une position

$$\text{Gain}_f(p, T) = f(S_p) - \sum_{j=1}^2 P_j \times f(S_{pj})$$

### Propriétés (cas attributs binaires)

- Terme  $f(S_p)$  ne dépend pas de  $T$  !
- Conséquence : **Maximiser le gain revient à minimiser  $\sum_{j=1}^2 P_j \times f(S_{pj})$  !**
- Gain maximal lorsque le test sur un attribut permet de classer correctement toutes les données
- (Gain minimal si aucune information apportée par ce test au regard de la classification)
- Sélection de l'attribut qui maximise le gain : stratégie *gloutonne*

## Sur l'exemple des matchs

### Choix du premier attribut (position racine)

Critère de Gini (DOM, BAL, MCC, MPG)

- $\text{Gain}(\epsilon, \text{DOM}) = \text{Gini}(S) - \left(\frac{5}{8}\text{Gini}(S_1) + \frac{3}{8}\text{Gini}(S_2)\right) =$   
 $\text{Gini}(S) - 2 \times \frac{5 \times 2 \times 3}{8 \times 5 \times 5} - 2 \times \frac{3 \times 1 \times 2}{8 \times 3 \times 3} = \text{Gini}(S) - \frac{7}{15}$
- $\text{Gain}(\epsilon, \text{BAL}) = \text{Gini}(S) - \frac{3}{8}$
- $\text{Gain}(\epsilon, \text{MCC}) = \text{Gini}(S) - \frac{7}{15}$
- $\text{Gain}(\epsilon, \text{MPG}) = \text{Gini}(S) - \frac{1}{2}$

Gain Maximum pour BAL (idem pour Entropie) : c'est l'attribut choisi à la racine pour une première phase de classification.

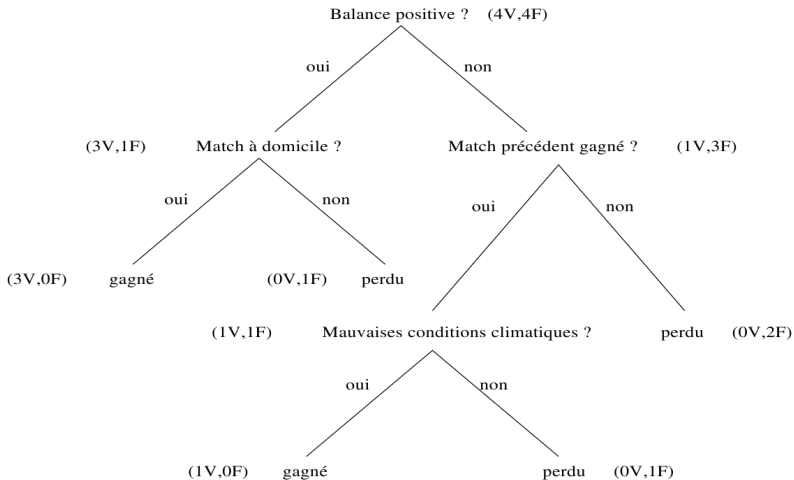
### Choix du second test en position $S_1$ (BAL=V)

Calcul des gains pour chaque attribut restant : the winner is DOM

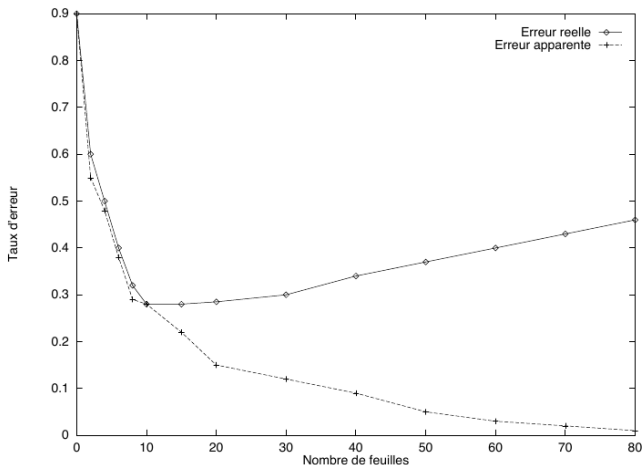
### Choix du second test en position $S_2$ (BAL=F)

The winner si MPG

# Arbre final obtenu



# Erreur apparente nulle, mais...



Faible pouvoir prédictif (notion de sur-apprentissage – ou apprentissage par cœur)

Nécessité d'obtenir un arbre plus petit : **élagage**

## Phase d'élagage de l'arbre de décision



### Deux approches

- Eviter une trop grande croissance de l'arbre en arrêtant sa construction au bon moment (*early stopping* via ensemble de validation).
- Procéder en deux phases : construire l'arbre complètement, puis couper les branches qui dépassent !

# Elagage de CART

## Variations d'erreurs à minimiser

- $T_0$  un arbre de décision à élaguer,  $\mathcal{T}$  l'ensemble de tous les arbres de décision obtenus à partir de  $T_0$  en remplaçant certains noeuds internes par des feuilles.
- Pour chaque noeud interne  $p$  de  $T_0$  :

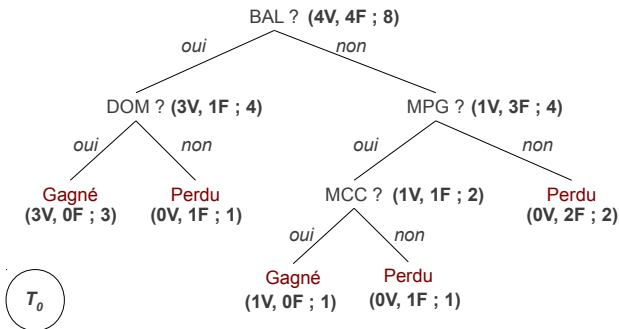
$$\alpha = \frac{\Delta R_{\text{emp}}^S}{|T_p| - 1}$$

où  $\Delta R_{\text{emp}}^S$  est le nombre d'erreurs supplémentaires que commet l'arbre sur  $S$  lorsqu'on élague à la position  $p$ , et  $|T_p| - 1$  mesure le nombre de feuilles supprimées.

## Processus itératif

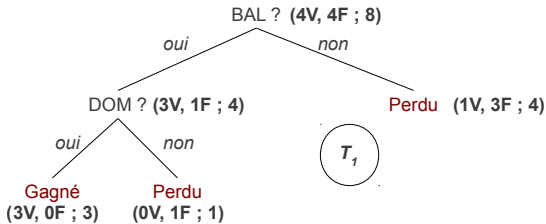
- $T_{i+1}$  obtenu à partir de  $T_i$ , auquel on coupe la branche qui permet un  $\alpha$  minimal.
- Soit  $T_0, \dots, T_i, \dots, T_t$  la suite obtenue, où  $T_t$  est réduit à une feuille.
- Sélection de l'arbre  $T_i$  dont le nombre d'erreurs calculées sur ensemble

# Elagage de CART sur l'exemple

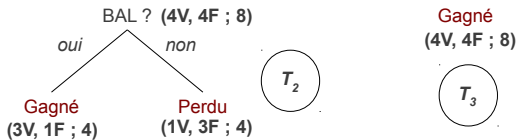




# Elagage de CART sur l'exemple



1. *Journal of Management Studies*, 1996, 33(1), 1-14.



## Elagage de CART sur l'exemple (cont'd)

### Ensemble de validation

| Match dom. | Balance pos. | Mauvais climat | Match préc. gagné | Victoire |
|------------|--------------|----------------|-------------------|----------|
| V          | V            | V              | F                 | V        |
| F          | V            | V              | F                 | V        |
| F          | F            | F              | V                 | F        |
| V          | F            | V              | F                 | V        |

### Calculs d'erreurs

- $T_0$  : 0 en apprentissage,  $\frac{1}{2}$  en test
- $T_1$  :  $\frac{1}{4}$  en apprentissage,  $\frac{1}{2}$  en test : mêmes erreurs que  $T_0$
- $T_2$  :  $\frac{1}{2}$  en apprentissage,  $\frac{1}{4}$  en test
- $T_3$  :  $\frac{1}{2}$  en apprentissage,  $\frac{1}{4}$  en test

## Compléments

## Autres catégories d'attributs

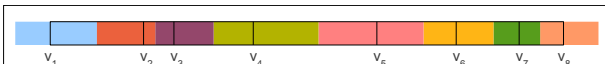
- Traitement en cas d'attributs  $n$ -aires et continus
- Traitement en cas de multi-classes
- Données manquantes ?

## Attribut continu : partition de l'intervalle

- Attribut  $a$ , tests binaires de la forme  $a < v$  ou  $a \leq v$
- Soient  $v_1, v_2, \dots, v_N$  valeurs prises par  $a$  dans  $S$
- Exemple de test, pour  $i = 1 \dots N - 1$  :

$$a < \frac{v_i + v_{i+1}}{2}$$

- Test sélectionné grâce à la notion de gain



## Compléments

## Autres catégories d'attributs

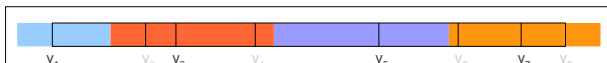
- Traitement en cas d'attributs  $n$ -aires et continus
- Traitement en cas de multi-classes
- Données manquantes ?

## Attribut continu : partition de l'intervalle

- Attribut  $a$ , tests binaires de la forme  $a < v$  ou  $a \leq v$
- Soient  $v_1, v_2, \dots, v_N$  valeurs prises par  $a$  dans  $S$
- Exemple de test, pour  $i = 1 \dots N - 1$  :

$$a < \frac{v_i + v_{i+1}}{2}$$

- Test sélectionné grâce à la notion de gain



## Compléments : attributs $n$ -aires

### Généralisation des formules de gain

$$\text{Gain}_f(p, T) = f(S_p) - \sum_{j=1}^N P_j \times f(S_{pj})$$

- Privilège des attributs de grande arité
- Cas où attribut = clé identifiant chaque élément de  $S$

### Approche de C4.5 : Ratio de gain

Soit  $a$  d'arité  $N$ , prenant les valeurs  $v_1, \dots, v_N$  :

$$\text{GainRatio} = \frac{\text{Gain}}{- \sum_{i=1}^N \frac{k_i}{|S|} \log \frac{k_i}{|S|}}$$

où  $k_i$  = nombre d'exemples de  $S$  pour lesquels  $a = v_i$ .

## Compléments

# Les trois grands algorithmes

- ID3 (Iterative Dichotomisor, [Quinlan79]) : sur variables qualitatives (discrimination)
- C4.5 [Quinlan93] : amélioration d'ID3 pour traitement des attributs continus et des valeurs manquantes
- CART [Breiman84] : généralisation à la régression et critère Gini remplace l'entropie. `sklearn.tree.DecisionTreeClassifier`

## Instabilité : variance importante

(Mais faible biais !) Sur données réelles :

- choix d'un attribut plutôt qu'un autre : limite serrée !
- influence majeure si proche de la racine

Solutions : agréger sur le hasard

- *Bagging* (Bootstrap Aggregating)
- Random Forests

# Outline

## 1 Classification

- De quoi parlons-nous ?
- Construire un bon modèle de classification à partir de données observées
- Les arbres de décision
- Les  $k$ -plus proches voisins



## Les $k$ -plus proches voisins ( $k$ -NN) : pas de modèle !

Méthode d'apprentissage supervisé, multi-classes, avec la définition d'une distance entre les points  $d : \mathcal{X} \rightarrow \mathbb{R}^+$  (ex. distance euclidienne)

Le problème : Apprendre  $f$  qui *généralise* au mieux

**En entrée** L'échantillon iid  $S_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n$

**En sortie** Un *classifieur*  $f : \mathcal{X} \rightarrow \mathcal{Y}$  tel que  $f = \operatorname{argmin}_h R(h)$

### Intuition de la méthode

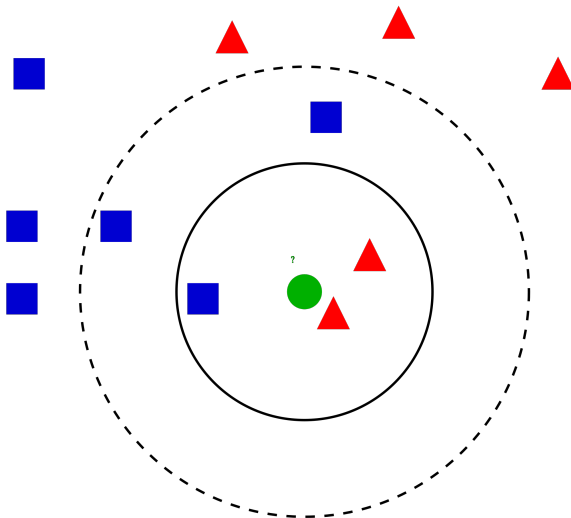
Pour chaque nouvel  $x \in \mathcal{X}$ , calculer ses  $k$  plus proches voisins dans  $S$  selon une distance  $d$  à définir, et choisir pour  $x$  l'étiquette majoritaire parmi ces  $k$  voisins.

### Algorithme en $\mathcal{O}(n)$

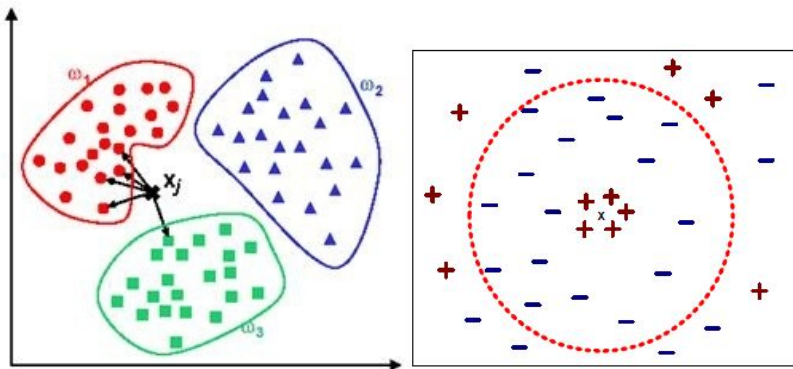
En parcourant chaque  $x_i \in \mathcal{X}$ , extraire  $V_k(x) = \{r_1, \dots, r_k\}$  = les indices  $i$  dans  $S$  des  $k$  plus proches voisins de  $x$

$$\text{Calculer } f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{j=1}^k \mathbb{I}(y = y_{r_j})$$

\_\_\_\_\_



# Illustration des $k$ -NN (2)



## Pro's and con's des $k$ -NN

## Avantages

- Simple, non-paramétrique
- Choix de  $k$  facile, par validation croisée
- Borne sur l'erreur en généralisation par rapport à celle de Bayes  

$$e^* \leq e^k \leq e^{k-1} \leq \dots \leq 2e^*$$

## Inconvénients

- Overfitting si  $k$  est grand, underfitting si  $k$  petit
- Forte dépendance à la distance choisie
- Pas un vrai modèle