

Data Analysis

Réduction de dimensions

ACP

$$X = \begin{bmatrix} x_{1j} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{nj} \end{bmatrix} \quad y = [\bar{x}_1 \dots \bar{x}_i] \quad \bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij} \quad y \in \mathbb{R}^p$$

centrer les données $X_c = [x_{ij} - \bar{x}_j] \in \mathbb{R}^{n \times p}$

matrice de variance-covariance : $\Sigma = \frac{1}{n} X_c^T X_c$

mesure de corrélation : $\text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$

$$\begin{aligned} \text{cov}(C_j, C_\ell) &= E((C_j - E(C_j))(C_\ell - E(C_\ell))) \\ &= \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{i\ell} - \bar{x}_\ell) \end{aligned}$$

$$\text{var}(C_j) = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

$$\Sigma = \frac{1}{N-1} X_c^T X_c$$

$$\Sigma = \begin{bmatrix} \text{var} X & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var} Y \end{bmatrix}$$

$$-1 \leq \text{cor}(X, Y) \leq 1$$

• 0 \rightarrow X, Y indép, info à garder
• $\pm 1 \rightarrow$ X, Y dép lin.

• but ACP : trouver q composantes principales C_1, \dots, C_q avec $q \ll p$ comme des nouvelles variables combi. lin. des var d'origine x_1, \dots, x_p tq les C_ℓ soient 2 à 2 non corrélées, de variance maximale, d'importance décroissante.
 \hookrightarrow maximiser la dispersion pour garder les caract.

exo :

Personne	Captteur 1	Captteur 2
Ind1	0	2
Ind2	-2	-1
Ind3	1	0
Ind4	1	-1

1) Existe-t-il une dépendance entre ces 2 capteurs ?

$$\begin{aligned} \text{Corr}(C_1, C_2) &= \frac{\frac{1}{4} [(0-0)(2-0) \dots (1-0)(-1-0)]}{\sqrt{\frac{(0-0)^2 + (-2-0)^2 + (1-0)^2 + (1-0)^2}{4} \times \frac{1}{4}}} = \frac{1/4}{\sqrt{\frac{6}{4} \times \frac{6}{4}}} = \frac{1}{6} \end{aligned}$$

\nexists dép non lin.

2) Calculer le 1^{er} axe princip de ces pts.

$$X = \begin{bmatrix} 0 & 2 \\ -2 & -1 \\ 1 & 0 \\ 1 & -1 \end{bmatrix} = X_c \quad \Sigma = \frac{1}{4} X_c^T X_c = \frac{1}{4} \begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix}$$

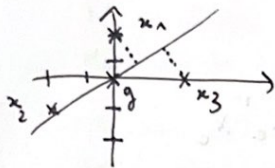
$$\chi_{4\Sigma}(\lambda) = \det(4\Sigma - \lambda I_2) = \begin{vmatrix} 6-\lambda & 1 \\ 1 & 6-\lambda \end{vmatrix} = (6-\lambda)^2 - 1 = (7-\lambda)(5-\lambda)$$

$\hookrightarrow \lambda_1 = 7, \lambda_2 = 5$
 \hookrightarrow + grande vp

Soit $X_1 = \begin{bmatrix} x \\ y \end{bmatrix}$ \vec{v}_P de Σ associé à $\lambda_1 = 7$.

$$4 \Sigma X_1 = \lambda_1 X_1 \quad (=) \quad \begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 7 \begin{bmatrix} x \\ y \end{bmatrix}$$

$$(\Rightarrow) \begin{cases} 6x + y = 7x \\ x + 6y = 7y \end{cases} \quad (\Rightarrow) \quad \begin{cases} x = y \\ x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{cases}$$



1^{re} Composante Principale :

(coeff de project° des données sur l'axe principal)

$$g = [0 \ 0]$$

$$C = X X_1 = \begin{bmatrix} 0 & 2 \\ -2 & -1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -3 \\ 1 \\ 0 \end{bmatrix}$$

Analyse de données temporelles

DTW : déformation temporelle dynamique, permet de mesurer la similarité entre 2 séries qui peuvent varier au cours du temps

- basé sur le principe du + court chemin de Dijkstra + effet temporel pris en compte
- évaluer la dissimilitude entre 2 seq → déterminer un chemin (pts : coïncidence extrémités + monotonie & ressemblance temporelle + continuité)

Input : 2 séries temporelles A et B de support temporels resp. $I = [1, N]$ et $J = [1, M]$
• distance locale d

Algo : 1. Initialiat° : $w_0 \leftarrow 1, w_1 \leftarrow 2, w_2 \leftarrow 1, g(0,0) \leftarrow 0$

2. Pour $j \in \{1 \dots J\}$

(a) $g(0, j) \leftarrow +\infty$

(b) Pour $i \in \{1 \dots N\}$

$g(i, 0) \leftarrow +\infty$

Pour $j \in \{1 \dots M\}$

Recherche du chemin minimal

$$g(i, j) \leftarrow \min \begin{cases} g(i-1, j) + w_0 * d(i, j) \\ g(i-1, j-1) + w_1 * d(i, j) \\ g(i, j-1) + w_2 * d(i, j) \end{cases}$$

3. calcul du score d'alignement

$$S = g(N, M) / (N+M)$$

Output : score S et chemin optimal C.

exo :

$$d(i, j) = x_i - x_j$$

g	0	1	2	-1
0	0	1	3	4
2	2	2	1	4
4	4	6	5	3

$$\text{Score} = \frac{8}{3+4} = \frac{8}{7}$$

exo Séquences ADN : $d(L_i, L_j) = \begin{cases} 0 & \text{si } L_i = L_j \\ 1 & \text{sinon} \end{cases}$

contraintes locales : $(w_0, w_1, w_2) = (1, 1, 1)$

globales : les cases à + de 4 cases de la diag ne sont pas calculées.

Matrice de coût :

	A	T	G	G	T	A	C	G	T	C
A	0	1	2	3	x	x	x	x	x	x
A	0	1	2	3	4	x	x	x	x	x
G	1	1	1	1	2	3	x	x	x	x
T	2	1	2	2	2	2	3	x	x	x
A	x	2	2	2	3	2	1	3	x	x
G	x	x	2	2	3	2	2	3	x	x
G	+	x	x	x	2	3	3	2	3	4
C	x	x	x	x	3	4	3	3	3	3

$$S = \frac{3}{10+8} = \frac{3}{18}$$

Classification non supervisée

* Approches par partitionnement :

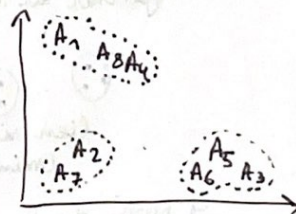
K-ppv : classification non supervisée
k - plus proche voisins

exo :

$$\begin{aligned} A_1 &= (2, 10) & A_2 &= (2, 5) & A_3 &= (8, 4) & A_4 &= (5, 8) \\ A_5 &= (7, 5) & A_6 &= (6, 4) & A_7 &= (1, 2) & A_8 &= (4, 9) \end{aligned}$$

• Matrice des distances euclidiennes au carré :

	A ₁	A ₂	A ₃	A ₄	A ₅	A ₆	A ₇	A ₈
A ₁	0	25	72	13	50	52	65	5
A ₂		0	37	18	25	17	10	20
A ₃			0	25	2	4	53	41
A ₄				0	13	17	52	2
A ₅					0	2	45	25
A ₆						0	29	29
A ₇							0	58
A ₈								0



• K ppv avec le seuil $\epsilon = 16$ et $k = 1$:

Pour A_1 : $d(A_1, A_8) = \min_{j \in \{2, \dots, 8\}} d(A_1, A_j) = 5 < \epsilon$ donc $\mathcal{C}^1 = \{A_1, A_8\}$

Pour A_2 : $d(A_2, A_7) = \min_{j \in \{1, \dots, 8\}, j \neq 2} d(A_2, A_j) = 10 < \epsilon$ donc $\mathcal{C}^2 = \{A_2, A_7\}$

$\mathcal{C}^1 = \{A_1, A_4, A_8\}$ $\mathcal{C}^2 = \{A_2, A_7\}$ $\mathcal{C}^3 = \{A_3, A_5, A_6\}$

K-means

Partitionner un ensemble de données en k classes représentées par les k centres, notés $C = \{C^1, \dots, C^k\}$

On associe alors à chq pt/donnée le centre le + proche au sens d'une certaine distance.

- ⚠ Bien choisir les centres à l'initialisat° pour éviter les CV locales
↳ heuristiques sur la distribut° des pts

Nb de classes à déterminer → varie nb classes et étudier l'énergie
Méthode de réparat° lin.

exo K-means avec centroides initiaux A_1, A_4 et A_7 :

$C^1 = \{A_1\}$ $C^2 = \{A_4\}$ $C^3 = \{A_7\}$
• Assigner à chaque donnée la classe la + proche
 $d(A_2, A_7) = \min_{j \in \{1, 4, 7\}} d(A_2, A_j) = 10 \rightarrow C^3 = \{A_2, A_7\}$

$d(A_3, A_4) = \min_{j \in \{1, 4, 7\}} d(A_3, A_j) \rightarrow C^2 = \{A_3, A_4\}$

$d(A_5, A_4) = \min_{j \in \{1, 4, 7\}} d(A_5, A_j) \rightarrow C^2 = \{A_3, A_4, A_5\}$

$d(A_6, A_4) = \min_{j \in \{1, 4, 7\}} d(A_6, A_j) \rightarrow C^2 = \{A_3, A_4, A_5, A_6\}$

$d(A_8, A_4) = \min_{j \in \{1, 4, 7\}} d(A_8, A_j) \rightarrow C^2 = \{A_3, A_4, A_5, A_6, A_8\}$

- Mise à jour des centres :

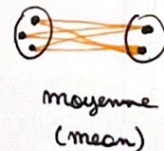
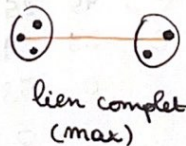
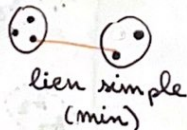
$$m_1 = (2, 10)$$

$$m_2 = \left(\frac{1}{5}(8+5+7+6+4), \frac{1}{5}(4+8+5+4+9) \right) = (6, 6)$$

$$m_3 = \left(\frac{2+1}{2}, \frac{5+2}{2} \right) = (1, 5)$$

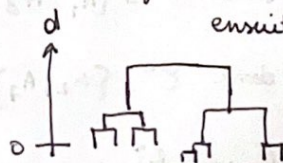
Classification Hiérarchique

- ↳ pas de cluster prédéfini + pas de config initiale + bonne interprétabilité
- ↳ fournir la mesure de dissimilarité pour comparer les groupes d'obs.



- ↳ partir de la granularité la + fine (1 classe = 1 donnée) puis stratégie ascendante / agglomératives.

- ↳ Avantage : visualiser comment les groupes sont nés (explicabilité)
ensuite à nous de voir où c'est la mieux de couper (n° classes)



dendrogramme

↳ similaire à l'arbre de décision

exo :

	A	B	C	D
A	0	1	4	5
B		0	2	6
C			0	3
D				0

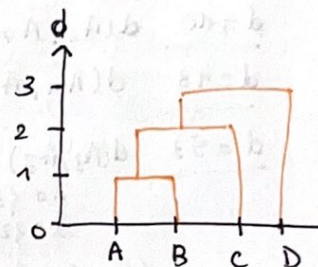
1) Lien simple

$d=0$ $\{A\}$ $\{B\}$ $\{C\}$ $\{D\}$

$d=1$ $d(A,B)=1$ $\{A,B\}$ $\{C\}$ $\{D\}$

$d=2$ $d(B,C)=2$ $\{A,B,C\}$ $\{D\}$

$d=3$ $d(C,D)=3$ $\{A,B,C,D\}$



2) Lien complet

$d=0$ $\{A\}$ $\{B\}$ $\{C\}$ $\{D\}$

$d=1$ $d(A,B)=1$ $\{A,B\}$ $\{C\}$ $\{D\}$

$d=2$ $d(B,C)=2 < d(C,A)=4$

prendre la + grande distance

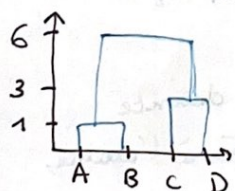
pour les singletons prendre le min

il faudra regarder si C fusionne avec $\{A,B\}$ à $d=4$

$d=3$ $d(C,D)=3$ $\{A,B\}$ $\{C,D\}$

on va fusionner que pour la + grande distance ie pour $d(B,D)=6$

$d=6$ $\{A,B,C,D\}$



exo avec les A_i

1) Lien simple

$d=0$ $\{A_1\}$ $\{A_2\}$... $\{A_8\}$

$d=2$ $d(A_3, A_5) = d(A_5, A_6) = d(A_4, A_8) = 2$

$\{A_1\}$ $\{A_2\}$ $\{A_3, A_5, A_6\}$ $\{A_4, A_8\}$ $\{A_7\}$

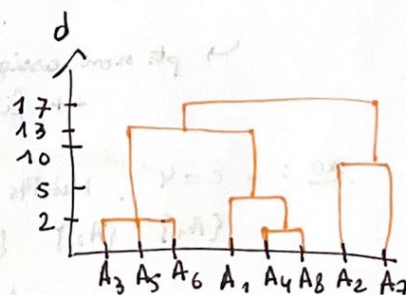
$d=4$ $d(A_3, A_6) = 4$

$d=5$ $d(A_1, A_8) = 5$ $\{A_1, A_4, A_8\}$ $\{A_2\}$ $\{A_3, A_5, A_6\}$ $\{A_7\}$

$d=10$ $d(A_2, A_7) = 10$ $\{A_1, A_4, A_8\}$ $\{A_2, A_7\}$ $\{A_3, A_5, A_6\}$

$d=13$ $d(A_1, A_4) = d(A_4, A_5) = 13$ $\{A_1, A_3, A_4, A_5, A_6\}$ $\{A_2, A_7\}$

$d=17$ $d(A_2, A_6) = d(A_4, A_6) = 17$



2) Lien complet

$$d=0 \quad \{A_1\} \quad \{A_2\} \quad \dots \quad \{A_8\}$$

$$d=2 \quad d(A_3, A_5) = d(A_5, A_6) = d(A_4, A_8) = 2$$

$$\{A_1\} \quad \{A_2\} \quad \{A_3, A_5, A_6\} \quad \{A_4, A_8\} \quad \{A_7\}$$

$$d=4 \quad d(A_3, A_6) = 4$$

$$d=5 \quad d(A_1, A_8) = 5 < d(A_4, A_1) = 13$$

$$d=10 \quad d(A_2, A_7) = 10$$

$$\{A_1\} \quad \{A_2, A_7\} \quad \{A_3, A_5, A_6\} \quad \{A_4, A_8\}$$

$$d=13 \quad d(A_4, A_1) = 13$$

$$\{A_1, A_4, A_8\} \quad \{A_2, A_7\} \quad \{A_3, A_5, A_6\}$$

$$d=53 \quad d(A_3, A_7) = \max d(A_i, A_j)$$

$$i = \{3, 5, 6\}$$

$$j = \{2, 7\}$$

$$72 \text{ (max = } d(A_1, A_3) \text{)}$$

$$\{A_2, A_3, A_5, A_6, A_7\} \quad \{A_1, A_4, A_8\}$$

$$65 \text{ (max = } d(A_1, A_7) \text{)}$$

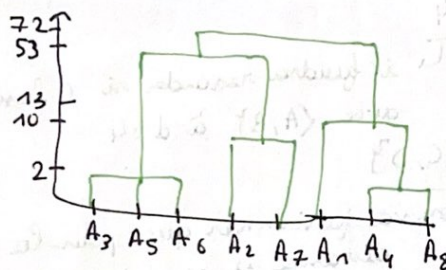
$$d=72$$

$$d(A_1, A_3) = \max d(A_i, A_j)$$

$$i = \{2, 3, 5, 6, 7\}$$

$$j = \{4, 8\}$$

$$53 \text{ (max = } d(A_3, A_7) \text{)}$$



DBSCAN

→ 1 classe associée à 1 densité

→ classes non déterminées à l'avance

→ 2 param : rayon pour vois ϵ

nb de pts min pour 1 cluster (densité)

→ pts non assignés à une classe = bruit

++ filtrage



gouton

exo : $\epsilon = 4$, MinPts = 2

$$\{A_1\} \quad \{A_2\} \quad \{A_3, A_5, A_6\} \quad \{A_4, A_8\} \quad \{A_7\}$$

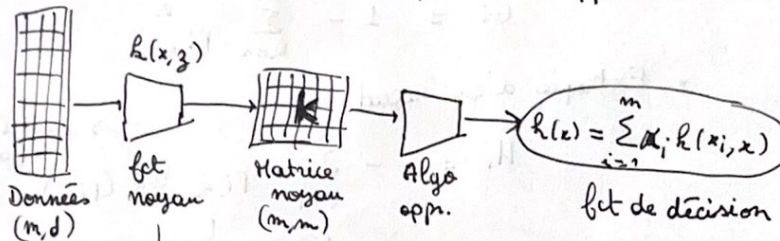
$$d(A_3, A_5) = d(A_5, A_6) = 2$$

$$\forall j \neq i \quad i \in \{1, 2, 7\} \quad d(A_i, A_j) \geq 5 \quad | \quad d(A_4, A_8) = 2$$

$$\epsilon = 10 \quad \{A_2, A_7\} \quad \{A_3, A_5, A_6\} \quad \{A_1, A_4, A_8\}$$

Approche par noyaux

↳ au lieu de séparer avec des hyperplans comme k-means, on augmente la dimension pour se rapprocher d'une séparat° lin



fct syn pos

polynome kernel

gaussien

sigmoid

$$k(a, b) = (a \cdot b + c)^d$$

$$k(a, b) = \exp(-\|a - b\|^2 / 2\sigma^2)$$

$$k(a, b) = \tanh(c \cdot (a \cdot b) + d)$$

Graphes et modularité

Critères d'évaluation en classification non supervisée

↳ basé sur 2 notions : cohésion & séparation

↳ SSE ou SSW : variance intra-classe

↳ BSS : variance inter-classe

↳ mix de cohésion et séparation :

> 0 : séparat° élevée entre clusters

< 0 : support° de classes

0 : distrib uniforme

données
proches
(SSE)

à quel point
la dist. est grande
(BSS)

Silhouette

$$s(x) = \begin{cases} \frac{b(x) - a(x)}{\max(a(x), b(x))} & \text{si } |c_i| > 1 \\ 0 & \text{si } |c_i| = 1 \end{cases}$$

Apprentissage supervisé

Modèles :

- arbres de décision
- apprentissage d'ensemble : forêts aléatoires
- réseaux de neurones
- SVM

Méthodes d'évaluation :

- validation croisée
- matrice de confusion
- précision, rappel, F-mesure
- courbe ROC

Arbres de décision

↳ mesure d'impureté

- Indice de Gini d'un nœud i comportant n éléments :

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$

- Entropie d'un nœud i :

$$H_i = - \sum_{k=1}^n p_{i,k} \log(p_{i,k})$$

part d'info de la classe
(faible qd redondante)

% d'obs de la classe k parmi
tous les obs d'entraînement i^e nœud.
Si c'est redondant, très fort

- Si G_i ou $H_i = 0 \rightarrow$ nœud pur \Rightarrow qu'une seule classe

- Gini est plus rapide à calculer

- Gini a tendance à isoler les classes les + fréquentes / Entropie produit des arbres + équilibrés

↳ crée des frontières de décision + donne des probabilités d'appartenance (interprétabilité)

↳ Algorithme CART

1) Séparation entraînement en 2

2) Choix du couple (k, t_k) qui produit les sous-ensembles les + purs.

3) Fonction de coût J à minimiser $J(k, t_k) = \frac{m_g}{m} G_{gauche} + \frac{m_d}{m} G_{droite}$

4) Récursion

5) Critère d'arrêt : profondeur max atteinte (hyperparam) nœuds purs \rightarrow éviter overfitting

↳ Avantage : utiliser var quantitative et qualitative

classification + regression
(prédic^t classe) (prédic^t valeur)

exo

Match à domicile

oui
oui
oui
non
non
non

Ciel
soleil
pluie
soleil
couvert
pluie
soleil

Match préc. gagné?

oui
non
non
oui
oui
non

Match gagné?

oui
non
oui
oui
oui
non

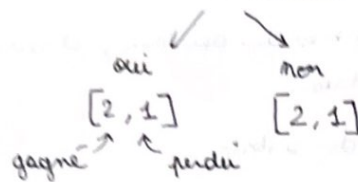
- 1) Gini de base : indice de Gini pour "Gagner le match" et "Perdre le match"

$$Gini(\text{Match gagné?}) = 1 - \left(\frac{\# \text{oui}}{\# \text{total}} \right)^2 - \left(\frac{\# \text{non}}{\# \text{total}} \right)^2 = 1 - \left(\frac{4}{6} \right)^2 - \left(\frac{2}{6} \right)^2$$

$$= 1 - \frac{4}{9} - \frac{1}{9} = \frac{4}{9}$$

- 2) Déterminer la variation de l'indice de Gini lorsqu'on découpe les données à l'aide des vars. En déduire celle du 1^{er} niv. de l'arbre.

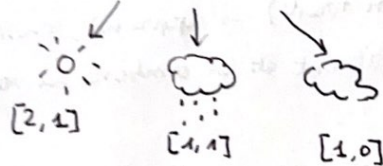
* Variable "Match à domicile"



$$Gini_{oui} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9} \quad Gini_{non} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}$$

$$\Rightarrow Gini("Match \grave{a} domicile") = \frac{3}{6} Gini_{oui} + \frac{3}{6} Gini_{non} = 0.5 \times 2 \times \frac{4}{9} = \frac{4}{9}$$

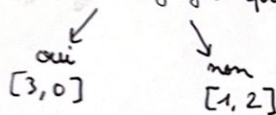
* Variable "Ciel"



$$\begin{aligned} Gini_{soleil} &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9} \\ Gini_{nuage} &= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2} \\ Gini_{pluie} &= 1 - (1)^2 - 0 = 0 \end{aligned}$$

$$\begin{aligned} Gini("Ciel") &= \frac{3}{6} Gini_{soleil} + \frac{2}{6} Gini_{nuage} + \frac{1}{6} Gini_{pluie} \\ &= \frac{3}{6} \cdot \frac{4}{9} + \frac{2}{6} \cdot \frac{1}{2} + \frac{1}{6} \cdot 0 = \frac{7}{18} \end{aligned}$$

* Variable "Match gagné précédent"



$$\begin{aligned} Gini("Match \grave{a} \acute{e}c. gagn\acute{e}") &= \frac{3}{6} Gini_{oui} + \frac{3}{6} Gini_{non} \\ &= \frac{3}{6} \left(1 - \left(\frac{3}{3}\right)^2 - 0\right) + \frac{3}{6} \left(1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2\right) = \frac{3}{6} \cdot \frac{4}{9} = \frac{2}{9} \end{aligned}$$

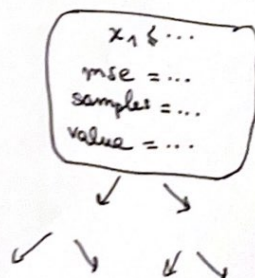
* 1^{er} split ?

$$\frac{Gini("Match \grave{a} \acute{e}c. gagn\acute{e}"){2}{9}}{<} \frac{Gini("Ciel"){7}{18}}{<} \frac{Gini("Match \grave{a} domicile"){4}{9}}{= Gini(base)}$$

← inutile car ne minimise pas

Régression par arbre de décision :

prédiction d'une valeur
comparer les valeurs moy → on obtient une MSE



Apprentissage par ensembles

- ↳ entraînement d'un ensemble d'arbres de décision, chacun sur un sous-ensemble aléatoire du jeu de train
 - ↳ calcul des prédictions pour chacun des arbres
 - ↳ choisir la classe obtenant le + de vote (classification) / la moyenne des résultats (régression)
- ⇒ forêt aléatoire

↳ 3 catégories d'apprentissage par ensembles

- bagging (Random Forest) → apprend, en //, indép, des modèles de base qui le constituent et les combine en suivant un processus de moyenne
- boosting (XGBoost) → apprend séquentiellement de manière adaptative (se concentre sur les erreurs) et combine selon strat.
- stacking → apprend en // et combine les modèles de base en un méta modèle

(bagging : \ominus de variance / boosting et stacking : \ominus biais)