

Correction de l'exercice 1 du TP sur CART

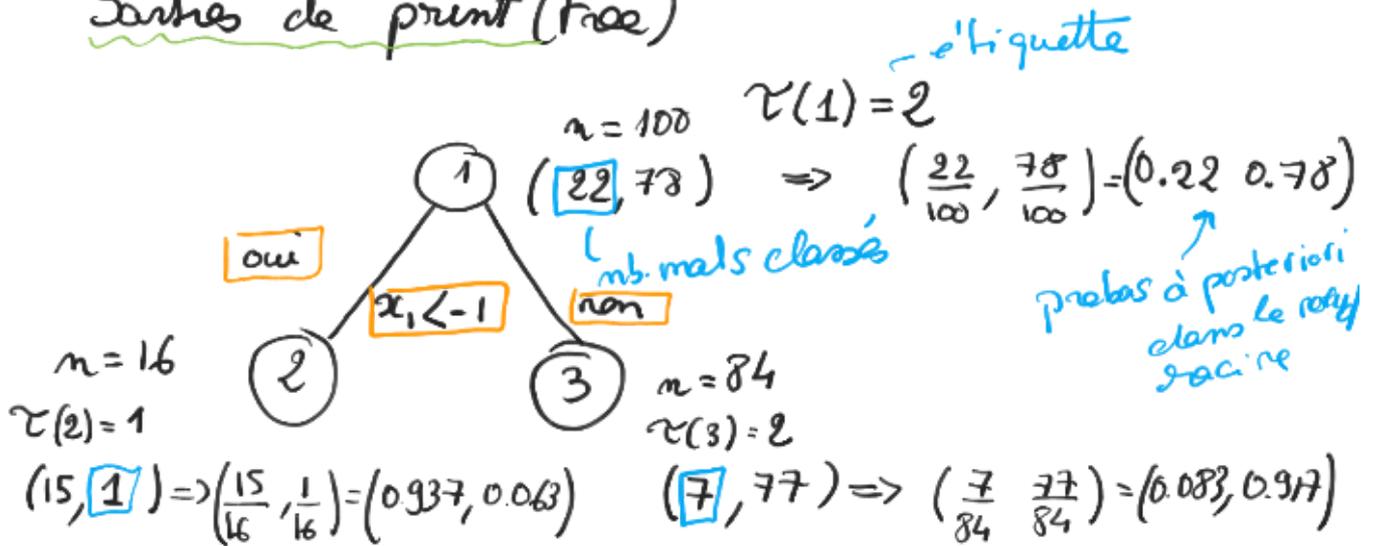
Données synthétiques : $\left\{ \begin{array}{l} Y \in \{1, 2\}_2 \\ (x^1, x^2) \in \mathbb{R}^2 \end{array} \right.$

Données d'apprentissage pour la construction de l'arbre:

$$\left\{ \begin{array}{l} n = 100 \text{ observations} \\ n_1 = 22 \text{ (nb d'obs. avec } Y=1) \\ n_2 = 78 \text{ (} \xrightarrow{\hspace{2cm}} \text{ } Y=2) \end{array} \right.$$

- ① Construction de l'arbre avec les valeurs par défaut de la fonction split

Sorties de print (tree)



Le taux d'erreur d'apprentissage de cet arbre est : $\frac{7+1}{100} \Rightarrow$ 80% d'erreur d'apprentissage

Sorties de summary (tree)

Pour le nœud racine (nœud n°1):

$$\begin{aligned}\text{complexity parameter} &= cp(1) \\ &= \frac{\text{LOSS}(1) - \text{LOSS}(2) - \text{LOSS}(3)}{\text{LOSS}(1)} \\ &= \frac{22 - 1 - 7}{22} = 0.6363\end{aligned}$$

$$P(\text{mode}) = P(1) = \frac{100 - \text{nb. d'obs. dans le nœud}}{100} = 1$$

Predicted class = $\hat{c}(1) = 2$ (class majority)

$$\begin{aligned}\text{Expected LOSS} &= \text{côt moyen du nœud} \\ &= \hat{\pi}(t) \\ &= \frac{\text{nb. fois classé dans } t \text{ (côt } 0 \text{)}}{n_t} \\ &= \frac{22}{100} = 0.22 \text{ ici}\end{aligned}$$

Improve = diminution de l'impureté (avec Gini ici)

$$= i(t) - \frac{n_L}{n_t} i(t_L) - \frac{n_R}{n_t} i(t_R)$$

$$i(t) = 1 - \sum_{k=1}^K p(k/t)^2 \quad (\text{Gini})$$

Ici:

$$\begin{aligned}\text{Improve} &= (1 - 0.22^2 - 0.78^2) - \frac{16}{100} (1 - 0.933^2 - 0.062^2) \\ &\quad - \frac{84}{100} (1 - 0.083^2 - 0.917^2) \\ &\approx 0.1967\end{aligned}$$

Soit improve $\approx 19.6\%$

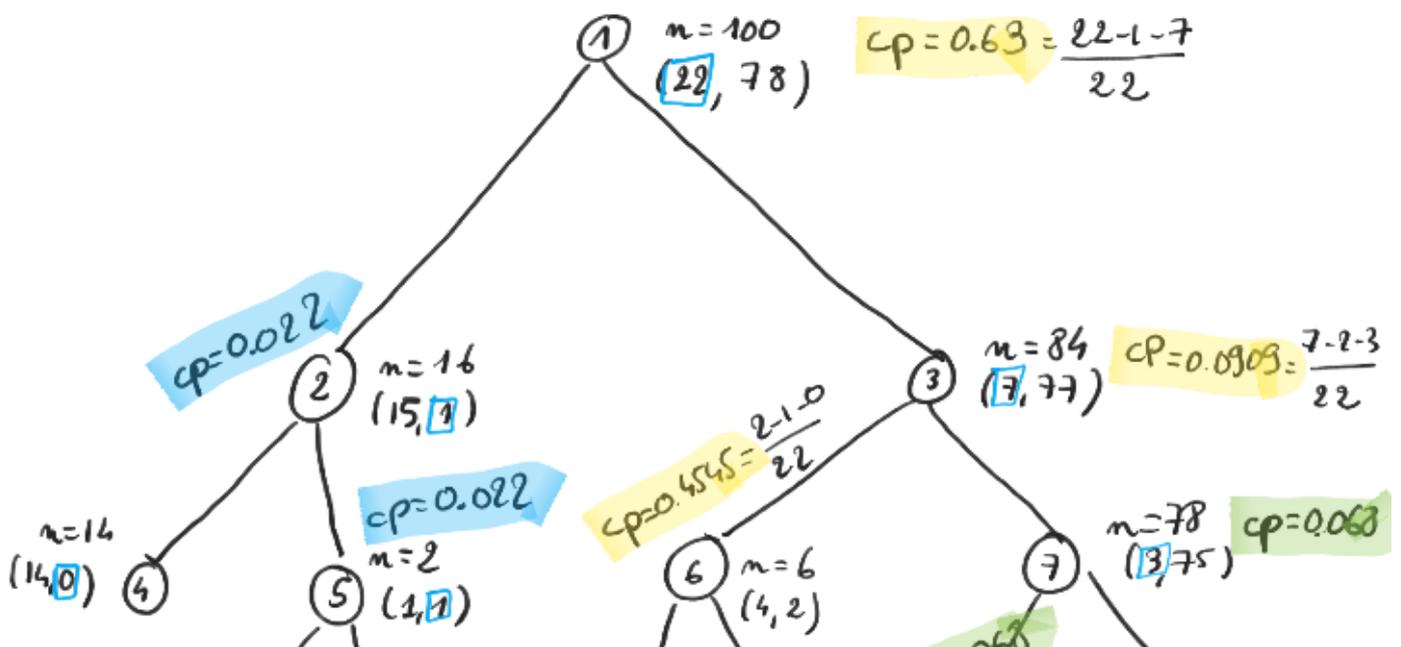
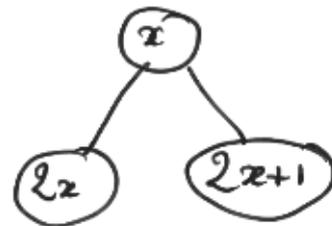
Arrêt des divisions :

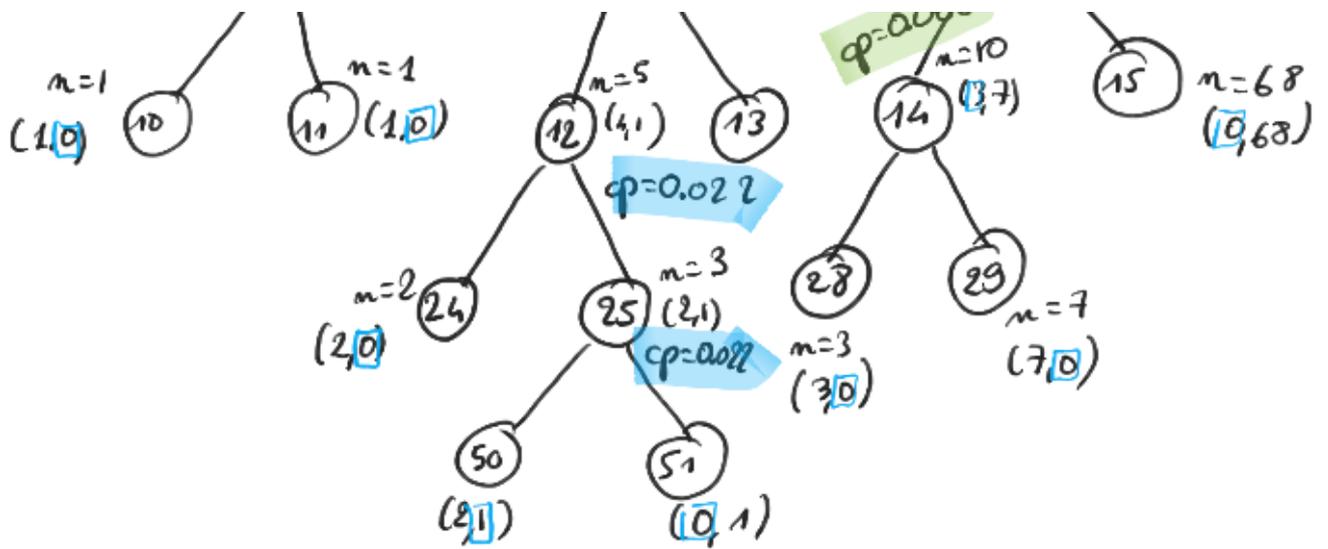
- * nœud 2 est une feuille car il contient 16 observations et $n_{min} = 20$ par défaut
- * nœud 3 est une feuille car probablement $cp(3) < 0.01$
valeur par défaut pour arrêt.

② Construction de l'arbre de longueur maximale suivi d'une étape d'élagage

Sorties de summary pour tracer l'arbre maximale

⚠ numérotation des nœud :





$CP = 0.022 = \frac{1-0-1}{22 \times 2}$: on élague directement 2 et 5
 ⇒ ↓ 2 feuilles à chaque fois

$CP = 0.068 = \frac{3-0-0}{22 \times 2}$: on élague directement 7 et 14
 ⇒ ↓ de 2 feuilles

Sorties de printcp

erreur relative de l'arbre

CP	nsplit	rel error	xerror	xstd
0.6363	0	$22/22 = 1$		
0.0909	1	$(7+1)/22 = 0.36$	▲	▲
0.0681	2	$(1+2+3)/22 = 0.27$		
0.04545	4	⋮		
0.0227	5 divisions ⇒ 6 feuilles	⋮		
0	9 divisions ⇒ 10 feuilles	⋮		

les valeurs de ces colonnes peuvent varier. Elles dépendent des découpages x-fold.

Remarque : Pour passer de l'arbre à 10 feuilles (arbre optimal)

a l'arbre a 6 feuilles, tous les nœuds ayant une valeur de $cp = 0.0227$ ont été élagués (nœuds 2, 5, 12, 25) \rightarrow 4 divisions ont été supprimées.

x_{error} = erreur relative moyenne de validation croisée x -folds de l'arbre élagué

x_{std} = écart-type des erreurs relatives de validation croisée de l'arbre élagué

Graphique de plotcp

Ce graphique est obtenu avec les résultats de rintcp

En abscisse : paramètres de complexité calculés lors de la procédure de validation croisée

$$\beta_j = \sqrt{d_j d_{j-1}} \Rightarrow \text{caractéristique de } [d_j, d_{j-1}]$$

ici :

$$d_1 = 0.6363$$

$$\beta_1 = +\infty$$

$$d_2 = 0.0909$$

$$\beta_2 = \sqrt{0.0909 \times 0.6363} = 0.2405$$

$$d_3 = 0.068$$

$$\beta_3 = \sqrt{0.068 \times 0.0909} = 0.079$$

$$d_4 = 0.04545$$

$$\beta_4 = 0.056$$

$$d_5 = 0.0227$$

$$\beta_5 = 0.032$$

$$d_6 = 0$$

$$\beta_6 = 0$$

Size of tree : chaque valeur β_j (en abscisse) correspond à un sous-arbre élagué. L'axe en haut donne le nombre de feuilles de ce sous-arbre

recherche de j* de ces 10-jpls

En ordonnée : l'ordonnée indique l'erreur relative de réduction croisée 10-jpl de chaque sous-arbre (relativement à l'erreur du nœud racine).
des points sont les moyennes des erreurs relatives des 10 jpls : colonne xerror (moyenne)
des intervalles correspondent à ± 1 écart type des 10 erreurs relatives : colonne xstd (écart-type)
ligne pointillée : ligne $y = S$ où en passant

$$j^* = \arg \min_{\beta_j} xerror_j$$

$$S = xerror_{j^*} + xstd_{j^*}$$

Choix du paramètre de complexité pour élaguer

Cas 1 : On choisit β_{j^*} qui minimise l'erreur

Cas 2 : On choisit la plus grande valeur de β_j pour laquelle $\beta_j \leq S$

C'est la regle des 1-SE. On prend la valeur de cp qui diminue le plus possible la complexité de l'arbre tout en conservant une erreur à moins de 1 écart-type de l'erreur "optimale".

