

Compléments sur l'algorithme CART implémenté dans rpart

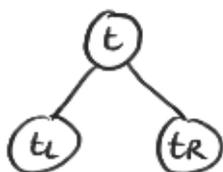
Par défaut :

- * matrice de coût 0-1
 - * mesure d'impureté pour le choix de la question binaire est Gini
- } argument parms

* Un nœud n n'est pas divisé si :

$$n_t < 20 \text{ (minsplit)}$$

$$n_L \text{ ou } n_R < \text{round}\left(\frac{20}{3}\right) = 6 \text{ (minbucket)}$$



$$cp(t) = \frac{\text{Loss}(t) - \text{Loss}(t_L) - \text{Loss}(t_R)}{\text{Loss}(t_L)}$$

? nœud racine

= proportion d'erreur en moins due à la division de t
 < 0.01 (cp)

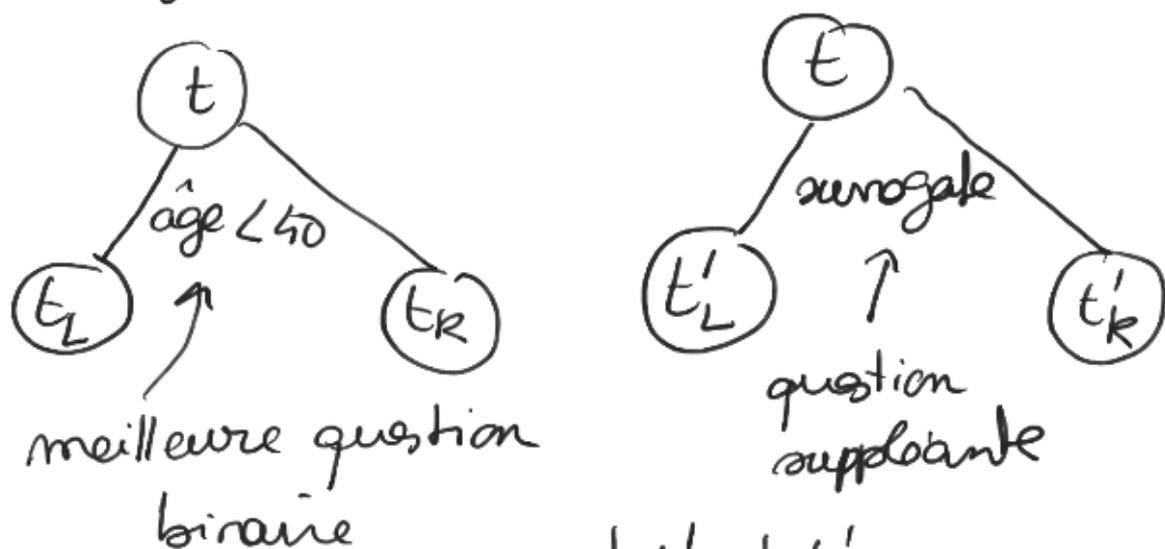
- * Validation croisée 10-fold pour l'estimation des erreurs des sous-ensembles ($x_{\text{val}} = 10$)

arguments

Control

Gestion des données manquantes

Surrogate = suppléant



	t'_L	t'_R
t_L	a	b
t_R	c	d

⇒ Trouver la question suppléante (surrogate) qui donne les mêmes sous-nœuds donc qui minimise $\frac{c+b}{n_t}$

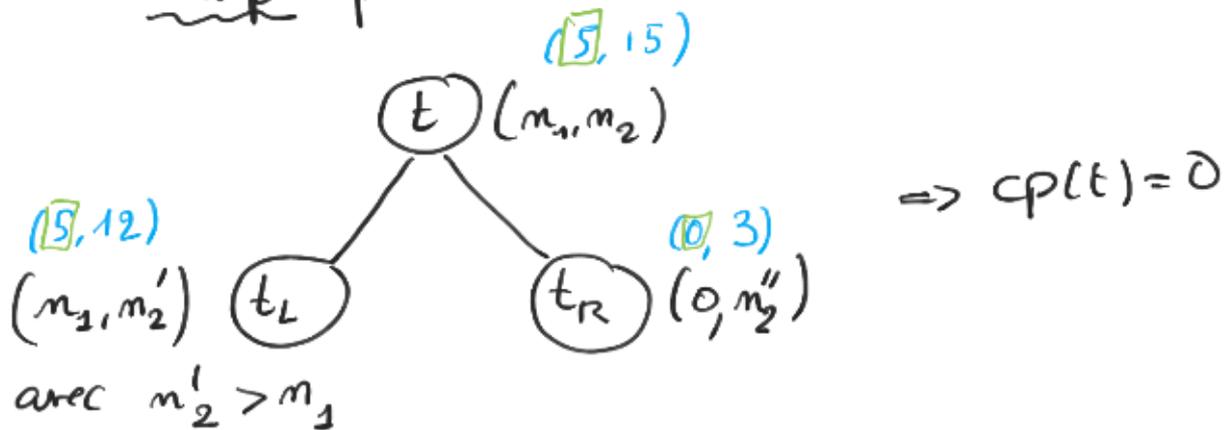
⇒ Si une nouvelle donnée arrive pour laquelle l'âge est absent, c'est la question surrogate qui est utilisée.

Cas où $C_p = 0$

- * $cp(t) = 0$ par convention si t est une feuille
- * Sinon on a défini $cp(t) = \text{Loss}(t_L) - \text{Loss}(t_R)$
 nb de mails dans

Le seul cas (a priori) où $cp = 0$ est le cas où tous les mails classés des nœuds t sont dans un même des nœuds si ils ne sont pas majoritaires

Exemple pour Y binaire:



⚠ Dans ce cas, une valeur de $cp(t) \neq 0$ est affecté au nœud. Cette valeur est liée à la définition de la suite de paramètres de complexités α associés au arbre élagués optimale.

On reprend le point 2. de la procédure de définition de ces paramètres de complexités :

Soit T_L un arbre à L feuille et d_L la valeur du paramètre de complexité à partir duquel T_L minimise $C_\alpha(T)$.

T_L minimise $C_\alpha(T)$ tant que

$$C_\alpha(T_L) < C_\alpha(T_L')$$

où $T_{L'}$ est un arbre à L' feuilles et $L' < L$.

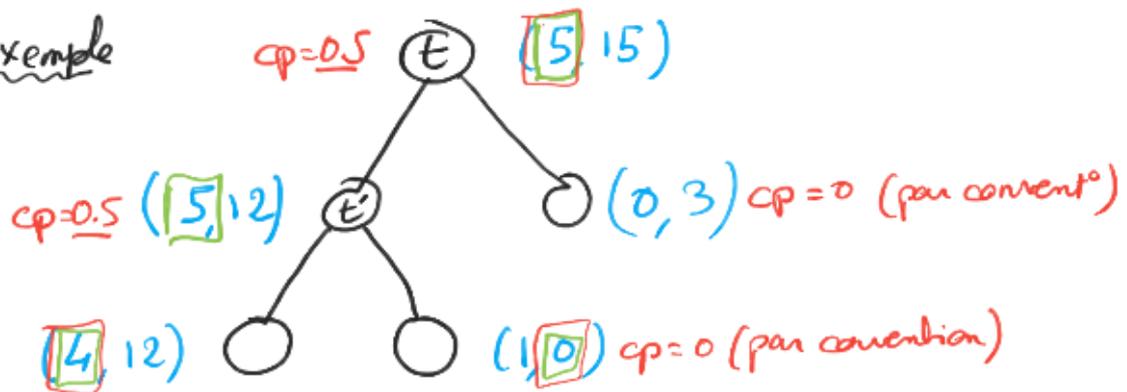
Cette condition s'écrit également

$$m \hat{R}(T_L) + \alpha L < m \hat{R}(T_{L'}) + \alpha L'$$

$$\Leftrightarrow \alpha < \frac{m \hat{R}(T_{L'}) - m \hat{R}(T_L)}{L - L'}$$

$\Rightarrow \alpha_{L'} = \frac{\text{Reduction du nombre de nœuds classés}}{\text{nombre de feuilles ajoutées (divisions)}}$

Exemple



$$cp(t) = cp(t) = \frac{5 - 4 - 0}{2} = 0.5$$

⚠ dans Rpent: $cp(t) = \frac{\text{Loss}(t) - \text{Loss}(t_R) - \text{Loss}(T_L)}{\text{Loss}(t_1)}$

\uparrow
 nœud racine.

Notion de question concurrente (compétitor)

À chaque division d'un nœud de l'arbre, c'est la meilleure question binaire qui est sélectionnée c.a.d celle qui maximise la réduction de l'impureté

Une question binaire concurrente (compétitive) à la meilleure question retenue, donne une réduction de l'impureté "proche de la réduction optimale".

La première question concurrente est la seconde question qui maximise la réduction de l'impureté
Etc...

Cela donne une idée des autres variables qui pourraient être importantes pour diviser ce nœud.