

Arbre de classification :

* Expliquer / qualitative en fct des x_i qual ou non, a k modalités

Principe :

A chaque nœud terminal de l'arbre en construction, Δ on sélectionne le test qui génère le plus d'info sur les classes à expliquer / prédire dans chaque nœud créé

→ Critère à optimiser : Impureté :

- * X ayant pour coordonnées les variables explicatives,
- * $1, 2, \dots, K$ les modalités de la variable Y à expliquer
- * les modalités définissent K classes d'individus.

Différence : Arbre de régression / Arbre de classif
critère de découpage d'un nœud

⇒ On définit la qualité d'un nœud m par l'impureté $i(m)$

A chaque nœud $m \rightarrow$ Région R_m

Chaque R_m contient N_m indiv

les indiv qui satisfont les conditions de tous les nœuds situés sur l'unique chemin reliant la racine à m

$$\hat{P}_{mk} = \frac{1}{N_m} \sum_{x_j \in R_m} 1_{(y_j = k)}$$

Proportion d'indiv $\in K$ dans la région R_m

$1 - i(m)$ maximale \rightarrow les indiv du nœud m se répartissent uniformément dans les classes, i.e. pour $\hat{P}_{mk} = \frac{1}{K} \forall k \in \{1, \dots, K\}$

→ $i(m)$ → minimise l'erreur lorsque les indiv de m e m classe

$$i(m) = \phi(\hat{p}_{m1}, \dots, \hat{p}_{mK})$$

$$\Delta i(m) = i(m) - \frac{N_{mG}}{N_m} i(m_G) - \frac{N_{mD}}{N_m} i(m_D)$$

m_G : nœud fils gauche de m

m_D : nœud fils droite de m

$i_2(m)$ mesure le pourcentage d'individus classés incorrectement lorsque les indiv d'un nœud sont affectés à sa classe majoritaire

$$\Delta i_2(m) = \frac{N_{mG}}{N_m} \max_k (\hat{p}_{mGk}) + \frac{N_{mD}}{N_m} \max_k (\hat{p}_{mDk}) - \max_k (\hat{p}_{mk})$$

Indice de gini : $i_g(m) = \sum_{k \neq k'} \hat{p}_{mk} \hat{p}_{mk'} = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$

Entropie : $i_e(m) = - \sum_{k=1}^K \hat{p}_{mk} \log(\hat{p}_{mk})$

Bin 2 classes

$$i_2(m) = 1 - \max(p, 1-p)$$

$$i_g(m) = 2p(1-p)$$

$$i_e(m) = -p \log p - (1-p) \log(1-p)$$

Règle d'affectation

$$C_k = \sum_{k'=1}^K \hat{p}_{mk'} \hat{p}_{mk}$$

→ coût d'affectation d'un individu m à une classe k .

→ une feuille est affectée à la classe pour laquelle C_k est min.

→ coût d'affectation d'un individu de k' dans k

Taux d'Erreur apparente de classement

A tout segment terminal m de l'arbre T associée une classe K

$$\rightarrow ER(K|m) = \sum_{K' \neq K} \hat{p}_{mK'}$$

$$TEA(T) = \sum_{m \in T} \frac{m.m}{n} ER(K|m)$$

↳ taux d'erreur apparente

Arbre de discrimination

données ← read.table(file = " --- ", dec = ".", header = TRUE)
 données.cnt ← split.control(minsplit = 1)

↳ stocke les param de l'algo
 ↳ nb min d'indiv nécessaire à la création d'un nœud

Calculer
 Construire l'arbre sur la totalité des indiv:
 arbre.full ← split(cœur ~, data = données, method = "class", control = données.cnt)

↳ imprimer l'arbre : print(~)

↳ sous forme graphique : plot(~, uniform = TRUE, branch = 0.5, margin = 0.1)

test(~, all = FALSE, use.n = TRUE)

Prédiction et Matrice de confusion:

pred ← predict(arbre.full, newdata = données, type = "class")

mc ← table(données\$cœur, pred) # matrice de confusion

print(mc)

err.resuls ← 1.0 - (mc[1,1] + mc[2,2] / sum(mc))

print(err.resuls) # erreur de substitution

↳ Elagage; élaguer automatiquement un arbre avec la méthode prune()
 prune(arbre.full & cp.table) → Arbre optimal

paramètres :

• C_p paramètre de complexité

• rel error (erreur d'entraînement)

• Xerror mesure le biais d'erreur dans la validation croisée
à savoir plus que l'on considère comme un estimateur correct
de l'erreur réelle.

• Xstd écart type de l'erreur de validation croisée

↓
l'autre qui minimise $\text{Xerror} + \text{Xstd}$

3] élaguer arbre. full, indiquer la procédure à suivre pour obtenir un arbre optimal à partir de cet arbre.

Afin d'obtenir l'arc optimal:

1. Consulter "CP Table" et plus précisément

La somme des colonnes x_{error} et x_{std} pour évaluer l'erreur réelle

réelle
→ le meilleur entre celui associé à x et x lui-même

microbial \rightarrow ns pli = 2 donc 2 + 1 familles

\Rightarrow élarguer l'arbre actuel avec la commande **prune**

en choisissant une valeur de CP entre $], --- , $\text{---}]$$

8) pure (centre full, $c_p = 0,046$)

5] Sélection pas à pas forward :

La sélection par ci pas & fait par:

minimisation du critère d'AKAIKE $AIC = -2\underset{\omega}{LL} + 2n(p+1)$

j : n° variable qui minimise le coût

à partir du modèle réduit à une constante puis de rajouter à chaque itération la variable qui minimise AIC.

La sélection s'arrête lorsque l'ajout d'une variable X_i n'améliore plus le critère AIC