

Classification supervisée & validation

1/ Introduction

	$X_1 \dots X_p$	Y
1		
2		
\vdots		
N		

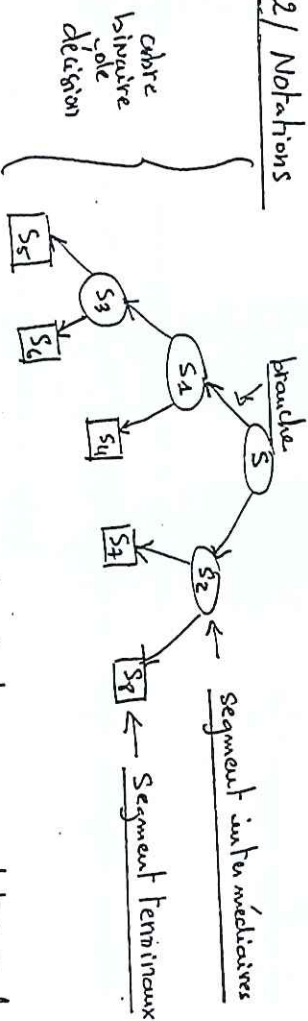
variables explicatives

variable à expliquer

Objectifs de la Segmentation.

- Discrimination
 - variable qualitative à modalités
 - Régression
 - variable quantitative
- Sélectionner les variables explicatives les plus discriminantes
 - construire une règle de décision permettant d'affecter un nouvel individu à l'une des classes
 - Sélectionner les variables explicatives les plus liées au phénomène étudié par Y
 - construire une règle de décision permettant d'affecter à un nouvel individu une valeur Y .

2/ Notations



Amax : canebtre binaire complet pour lequel chaque segment terminal contient 1 seul individu.

A une sous canebtre de Amax obtenu par élagage d'une ou de plusieurs branches.

L'idée de base consiste à effectuer la division d'un nœud de telle sorte que les deux segments obtenus soient :

- plus homogènes que le nœud parent
- les plus différents possible entre eux vis-à-vis de la variable Y .

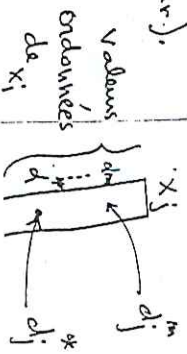
- Établir pour chaque nœud l'ensemble des divisions admissibles
- Définir une règle permettant de sélectionner la "meilleure" division d'un nœud.
- Définir une règle permettant de déclarer un nœud comme terminal ou intermédiaire.
- Affecter chaque nœud terminal à l'un des groupes (cas de la discrimination, ou affecter une valeur à Y pour chaque nœud terminal (cas de la régression)).
- Estimer le risque d'erreur de classement (cas de la discrimination) ou la prévision (cas de la régression) associé à l'arbre

3.1 Définition de l'ensemble des divisions admissibles

Cas 1 / Variable explicative continue

soit S le premier segment contenant tous les individus pour une variable X_j , on examine toutes les divisions possibles de la forme $X_j < d$ ($d \in \mathcal{D}$ à l'exception de

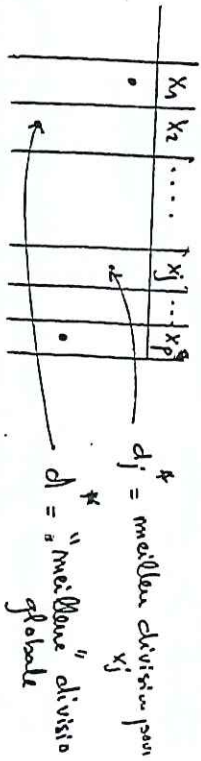
soit d_1^1, \dots, d_1^m les m divisions possibles, et d_1^q les q divisions possibles (valeur \leq min).
 en 3.2



- Une variable à 2 modalités ne fournit qu'1 seul divisi
- " " à k modalités endonnées fournit k-1 divisi
- " " non-ordonnées " 2^k - 1 divisi

Toutes les divisions sont examinées et l'on retient la meilleure au sens d'entente (à voir en 3.2).

Ainsi on obtient après la détermination de la meilleure division pour chaque variable (tout type confondu) :

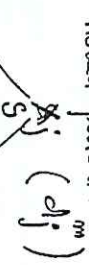


d^* est la meilleure division, au sens du critère, toutes variables confondues.

3.2 Définition du critère de sélection de la meilleure division

Cas 1 / Y est une variable continue (Regression)

Le critère est fondé sur la variance de Y au sein des segments descendants. Cette variance doit être plus faible que la variance de Y dans le nœud parent.



On mesure pour chaque division d_j^m sa variance résiduelle comme la moyenne pondérée des variances de Y à l'intérieur de chacun de ses segments descendants :

$$\text{Var}(d_j^m, S) = \left(\frac{M_g}{n_g} \text{Var}_Y(S_g) + \frac{n_d}{n_s} \text{Var}_Y(S_d) \right)$$

Variance résiduelle :

$$\text{Var}(d_j^*, S) = \min_{m \in \mathcal{D}_j} \left\{ \text{Var}_Y(d_j^m, S) \right\}$$

La meilleure division globale d^* est la division d_j^m minimisant la variance résiduelle :

$$\text{Var}(d^*, S) = \min_{j=1, \dots, p} \left\{ \text{Var}(d_j^*, S) \right\}$$

Cas 2 / Y est une variable qualitative (Discrimination)

Soit Y une variable qualitative à k modalités. La sélection d'1 division doit être telle que les segments descendants soient plus "purs" que le nœud parent. Autrement dit il faudrait que le mélange des modalités dans les nœuds descendants soit moins important que dans le nœud parent.

A chaque segment S on associe une mesure de l'impureté $i(S)$

$$i(S) = \sum_{r=1}^k \sum_{t=1}^k P(r|S) \cdot P(t|S) \quad (r \neq t)$$

$P(r|S), P(t|S)$: proportions d'individus provenant de la modalité r, t dans S.

Un segment est pur si il ne contient que des individus de la même classe (même modalité pour Y) ; avec $i(S) = 0$.

Chaque division d_j^m a une mesure de l'impureté :

$$\Delta_j^m = i(S) - p_g i(S_g) - p_d i(S_d)$$

maximise la réduction de l'impureté ;

$$\Delta_j^* = \max_{m \leq j} \{ \Delta_j^m \}$$

Sur l'ensemble de toutes les variables la division est effectuée selon la variable qui assure :

$$\Delta^* = \max_{j=1 \dots p} \{ \Delta_j^* \}$$

Une fois le segment s segmenté on réajuste N_{ST} aux nœuds descendants.

On arrête la construction de l'arbre lorsque tous les segments sont descendants :

- me métrisent plus de divisions (ex: $y_i = 0$, $\forall x_i$)
- leur taille est \leq à une valeur fixe a priori.

3.3/ Règle d'affectation d'un nouvel individu

Cas 1 / Y variable continue (Prévision)

soit i un nouvel individu, on le fait descendre dans l'arbre. soit si son segment terminal d'appartenance

la valeur y_i de l'individu i est estimée par la moyenne de y dans le segment S_i .

Cas 2 / Y est une variable qualitative (classement)

Soit S un nœud terminal. On note $n_1(S), \dots, n_k(S)$ les nb d'individus dans S prenant la modalité 1, ..., k de Y .

est affecté à la classe la mieux représentée :

$$\text{Classe de } i = \text{ArgMax}_{j=1 \dots k} \{ n_j(S) \}$$

3.4 Estimation du risque d'erreur de classement (cas discriminatoire de prévision (cas Régression))

Cas 1 / Y variable continue (Prévision)

On associe à chaque segment terminal S_T une valeur A l'erreur de prévision R_{ST} :

$$R(S_T) = \frac{n_{SE}}{n} \text{Var}_Y(S_T)$$

n_{ST} : taille du segment S_T $\text{Var}_Y(S_T) = \frac{1}{n_{ST}} \sum_{i \in S_T} (y_i - \bar{y})^2$

n : nb total d'individus

L'erreur apparente de Prévision. $EAP(A)$ est l'arbre A vaut :

$$EAP(A) = \sum_{S \in A} R(S_T)$$

$EAP(A) \approx (1 - R^2)$: représente la proportion de la variance totale non expliquée par l'arbre.

Remarque : $EAP(\text{Arbre}) = 0$

Cas 2 / Y variable qualitative (classement)

à tout segment terminal S_T est associé à une classe j (modalité) correspond une erreur apparente de classe

$$R(j|S_T) = \sum_{r=1}^k P(r|S_T)$$

appartient à la classe j dans le segment ST .

$$P(A_j | ST) = \frac{n_{jT}}{n_{ST}}$$

N l'erreur associée de classement (EAC) associée à A_i :

$$EAC(A_i) = \sum_{ST \in A_i} \frac{n_{ST}}{n} R(j | ST) = \sum_{ST \in A_i} \sum_{A_j \neq A_i} \frac{n_{jT}}{n}$$

EAC(A) : représente la proportion d'individus mal classés dans l'ensemble des segments terminaux.

Bibliographie:

- * - Gilbert Saporta, "Probabilités, Analyse des données et Statistique". TECHNIP
- * - Edwige Didry, J. Lemaire, J. Rouget, F. Testu, "Éléments d'analyse des données". DUNOD
- * - Gilles Celeux, "Analyse discriminante sur variables continues INRI A. collection DIDACTIQUE.
- Gilles Celeux et J.P. NIKRACHE. "Analyse discriminante sur variables qualitatives".
- Pierre et J.M. BOURCQUE. "Analyse des données multidimensionnelles". Presses Universitaires de France
- * * - T. Naveen Hastie, Robert Tibshirani et Jerome Friedman. "The elements of statistical learning". Data Mining, Inf and Prediction". Springer
- * * - Amit. K. Jain et Richard C. Dubes. "Algorithms for data Practice Hell, Harvard Kalyana. sevis. Data."
- * * - Breiman, Friedman, Olshen et Stone, "Classification and Regression Trees".

Validation et partie des résultats en

Classification.

- On se - t-on vraiment quelque chose?
- Les données ont-elles une structure
- A-t-on découvert des classes préexistentes ou on reconstruit a-t-on découvré une réalité cachée en classes?
- Est-ce que les configurations obtenues sont stables, compte tenu de ce que l'on sait sur la précision des données?
- Quelle est l'importance sur les résultats d'une modification du tableau de données?

Étude de la stabilité des classes, des résultats par rééchantillonnage.

Les méthodes par rééchantillonnage sont privilégiées pour étudier la stabilité des résultats vis-à-vis de perturbations aléatoires.

- se sont des méthodes de calculs intensifs qui reposent sur des techniques de simulations d'échantillons à partir d'1 seul échantillon.
- Elle consiste à répéter des analyses sur les différents échantillons simulés puis à étudier les fluctuations des résultats obtenus.

Il existe trois principales approches permettant de générer des échantillons artificiels: Jackknif, Bootstrap et la Validation Croisée.

1. Technique de Jackknife. (Quenouille (1949), Tukey (1952), Miller (1955))

Objectif: estimation du biais et de la variance d'1 estimateur

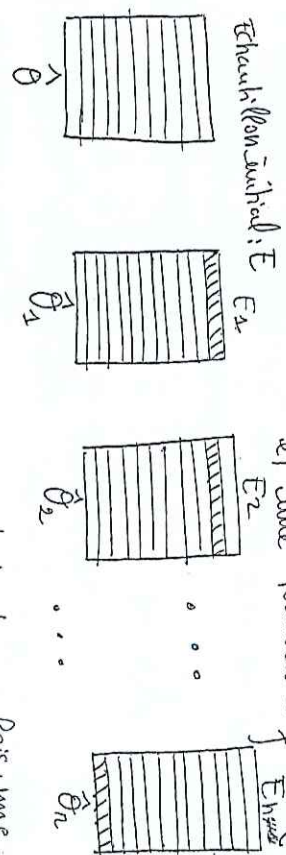
Le biais: mesure la précision de l'estimateur

La variance: mesure la stabilité de l'estimateur

Le biais permet de dire si mon estimateur est correctement défini, de manière précise la valeur estimée

La variance permet de dire si mon estimateur fluctue beaucoup ou pas en fonction des données d'origine.

La situation idéale: c'est un estimateur précis (biais faible) et une variance faible (grande stabilité)



- On génère n échantillons en supprimant à chaque fois une observation de l'échantillon initial.

θ : paramètre à estimer, $\hat{\theta} = s(E)$ l'estimation de θ sur la base de E

$\hat{\theta}_i = s(E_i)$ " " base de E_i

On note $\hat{\theta}_0 = \frac{1}{n} \sum_i \hat{\theta}_i$: la moyenne des estimations.

l'estimateur de Jackknife. $\hat{\theta}_0$ est défini:

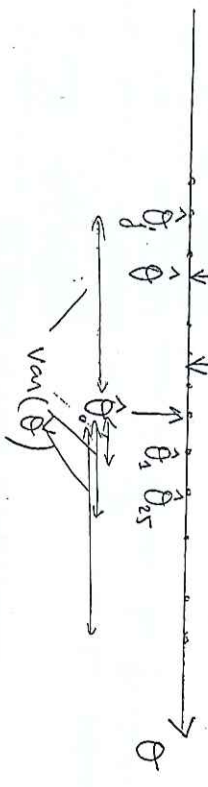
$$\hat{\theta} = n \hat{\theta}_0 - (n-1) \hat{\theta}_0 = n(\hat{\theta}_0 - \hat{\theta}_0) + \hat{\theta}_0$$

Le biais de Jackknife vaut :

$$b = \hat{\theta} - \tilde{\theta} = (n-1) (\hat{\theta}_0 - \hat{\theta})$$

L'estimation de Jackknife de la variance :

$$Var(\tilde{\theta}) = \frac{n-1}{n} \sum (\hat{\theta}_i - \hat{\theta}_0)^2$$

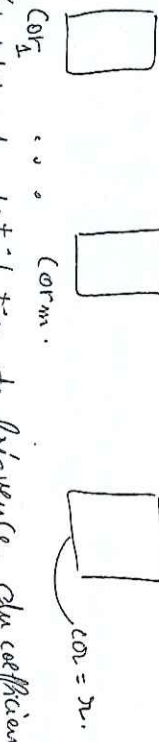


2. Technique du Bootstrap (Efron (1979))

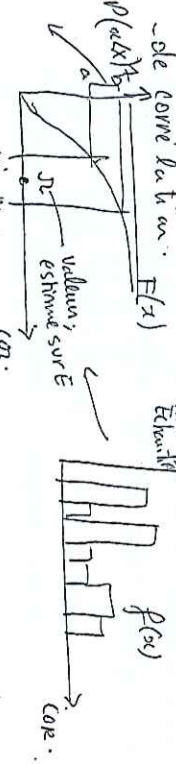
objectif : estimer la variabilité d'un paramètre en fournissant une autre valeur de confiance.

- Il s'agit de simuler $m > 30$ échantillons de même taille n que l'échantillon initial. Chaque échantillon est obtenu par tirage aléatoire avec remise parmi les m individus.

Exemple : Estimation du cor (X_1, X_2) ?



On établit la distribution de fréquence du coefficient



$P(\text{cor} \in [r_i, r_{i+1})) = [a, b]$ ← Probabilité pour que l'intervalle de confiance soit de $2 \in [r_i, r_{i+1})$

(3)

- On obtient ainsi une estimation de l'incertitude de confiance de r sans avoir recours à des hypothèses distributionnelles sur les données.

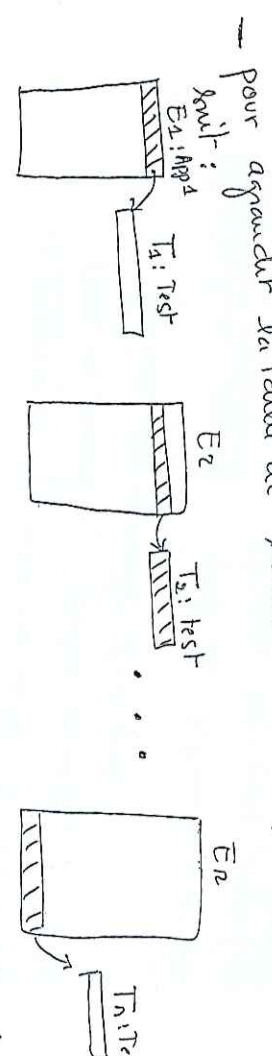
3. La validation croisée

objectif : sert à évaluer des erreurs de prédictions. Elle est utilisée dans le cas où l'échantillon initial est de petite taille (cas la régression / Analyse discriminatoire).

- A l'origine il s'agit de séparer l'échantillon de base en deux blocs de tailles \neq généralement : l'échantillon de prétraitement et l'échantillon-test.

l'échantillon de prétraitement : sert à formuler le modèle, à établir les règles de décision ou d'affectation.

l'échantillon-test : sert à appliquer les règles et à estimer les performances du modèle.



- On extrait à chaque fois (k observations ici $k=1$) pour tester le modèle et le insérer sur $(n-k)$ observations.

- On mesure des taux de prédiction des n échantillons Test et constitue la prédiction à retenir pour le modèle définit sur E.

30