

## Projet "fil rouge"

Les exploitations agricoles (articles de D. Dubois)

Partie ① : Réaliser une ACP, appliquez la méthodologie et les méthodes d'interprétation mises en cours

Partie ② : Utilisation des résultats de l'ACP pour définir une fonction score et une règle de décision pour prédire si une exploitation sera "saine" ou "défaillante", connaissant ses 22 ratios financiers

- 2.1) A partir de la première composante principale  $\psi^1$ ,
- Faire le diagramme boxplot de  $\psi^1$  pour les exploitations "saine" et les exploitations "défaillantes"
  - Calculer la moyenne de  $\psi^1$  dans ces deux groupes, et vérifier que la moyenne pondérée de ces deux moyennes est 0.

2.2) Une règle de classement géométrique pour une exploitation  $i$  appartenant à l'échantillon consiste à la positionner par rapport au point pivot qui est le milieu du segment des deux moyennes

$$\mu_1 = -2,164 \text{ et } \mu_2 = 2,328 \text{ soit } \frac{\mu_1 + \mu_2}{2} = 0,082 :$$

si  $\psi_i^1 < 0,082$ , alors l'exploitation  $i$  est déclarée "saine"

si  $\psi_i^1 > 0,082$ , alors l'exploitation  $i$  est déclarée "défaillante"

Calculer à partir de cette règle de classement une nouvelle variable "PRED" qui vaudra "Sain" ou "défaillant" selon la règle précédente. (e)

Tenuez que vous retrouvez le tableau de contingence suivant :

		PRED	
		Sain	Defaillant
DIFF	Sain	599	54
	Defaillant	124	483

En déduire les pourcentages intéressants et interprétez.

En particulier, on trouve un taux de més classés de 14,42%.

2.3) On veut maintenant pouvoir utiliser cette règle de classement pour pouvoir prédire la classe d'une nouvelle exploitation  $i_0$ , dont on connaît les caractéristiques financières. On sait que  $\Psi_i^1 = \bar{z} \cdot \bar{x} = \sum_{j=1}^p n_j z_{ij}$

$$\text{et } z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j} \Rightarrow \Psi_i^1 = \sum_{j=1}^p \underbrace{\frac{n_j}{s_j} x_{ij}}_{\text{coef. } j} - \underbrace{\sum_{j=1}^p n_j \frac{\bar{x}_j}{s_j}}_{\text{constante}}$$

On en déduit la fonction score suivante

$$S(i) = \sum_{j=1}^p a_j z_{ij} + \text{constante}$$

Tenuez : - que nous retrouverons tableau

	Fonction S
$R_1$	0,76
$R_2$	-0,83
$\vdots$	
$R_{37}$	-0,12
constant	-1,53

- qui en appliquant cette fonction à toutes les exploitations de l'échantillon nous retrouver  $\Psi^1$ .

2.4) On sait que le taux d'erreur calcule sur l'échantillon d'apprentissage (appelé taux apparent d'erreur) calcule au 2.2) sans estime souvent le taux d'erreur vu que l'on utilise les mêmes observations pour le calculer que celles qui ont servi à trouver la règle géométrique de classement. On utilise donc souvent la méthode de l'échantillon-test pour estimer ce taux d'erreur: Cette méthode consiste à partager en deux parties l'échantillon:

- une partie sert d'échantillon d'apprentissage de la règle de décision
- l'autre partie sert d'échantillon test et permet de tester la règle d'affectation et donc de calculer le taux d'erreur

2.4.1) Calculer ce taux d'erreur en prenant  $\frac{2}{3}$  des observations dans l'échantillon d'apprentissage (les 840 premiers exploitations) et  $\frac{1}{3}$  dans l'échantillon test (les 420 dernières)

2.4.2) Calculer ce taux d'erreur en tirant cette fois aléatoirement (en vous aidant de la fonction `sample()`) un échantillon d'apprentissage (de 840 exploitations) et un échantillon test (des 240 exploitations). Recommencez 5 fois et calculez le taux d'erreur moyen. Commentez.

2.4.3) Expliquez ce dernier résultat en vous aidant de la remarque suivante: "Le taux d'erreur apparent est d'autant plus faible que le modèle est complexe (sur-paramétrisation)". Proposez une "solution".