

## ANALYSE DE DONNÉES

DATE : 09 JUIN 2021

SECTION : ING2 GÉNIE INFORMATIQUE

ENSEIGNANTS : I.KAMMOUN, W.BARHOUMI, S.ZOUAOU

DURÉE : 1H30

NOMBRE DE PAGES : 2

DOCUMENTS NON AUTORISÉS

**Exercice 1 : Clustering (10 pts)**

On veut utiliser l'algorithme du k-means et la distance euclidienne pour regrouper les 8 points suivants en 3 clusters :  $A_1(2, 10)$ ,  $A_2(2, 5)$ ,  $A_3(8, 4)$ ,  $A_4(5, 8)$ ,  $A_5(7, 5)$ ,  $A_6(6, 4)$ ,  $A_7(1, 2)$ ,  $A_8(4, 9)$  de même poids. La matrice de distance basée sur la distance Euclidienne est fournie ci-dessous.

1. (a) On considère comme centre de classes à l'initialisation les points  $A_1$ ,  $A_4$  et  $A_7$ . Déroulez une première itération de l'algorithme de k-means pour ces données et donnez :
  - i. Les nouveaux clusters ;
  - ii. Les centres de chaque cluster ;
- (b) Donnez au moins deux conditions d'arrêt possibles de ce processus.
- (c) Comment varie l'inertie totale, l'inertie inter-classes et l'inertie intra-classes dans cette méthode de classification.

	A1	A2	A3	A4	A5	A6	A7	A8
A1	0	$\sqrt{25}$	$\sqrt{36}$	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{65}$	$\sqrt{5}$
A2		0	$\sqrt{37}$	$\sqrt{18}$	$\sqrt{25}$	$\sqrt{17}$	$\sqrt{10}$	$\sqrt{20}$
A3			0	$\sqrt{25}$	$\sqrt{2}$	$\sqrt{4}$	$\sqrt{53}$	$\sqrt{41}$
A4				0	$\sqrt{13}$	$\sqrt{17}$	$\sqrt{52}$	$\sqrt{2}$
A5					0	$\sqrt{2}$	$\sqrt{45}$	$\sqrt{25}$
A6						0	$\sqrt{29}$	$\sqrt{29}$
A7							0	$\sqrt{58}$
A8								0

2. (a) Déroulez une première itération de l'algorithme de classification ascendante hiérarchique (CAH) pour ces données en adoptant le critère de Ward pour l'agrégation et donnez :
    - i. La nouvelle agrégation
    - ii. La nouvelle matrice de distance suite à l'agrégation
  - (b) Comment varie l'inertie totale, l'inertie inter-classes et l'inertie intra-classes dans cette méthode de classification.
3. Décrivez les avantages et les inconvénients de ces deux méthodes de classification.
  4. Décrire brièvement la méthode de la classification non supervisée mixte.

**Exercice 2 : Modèle Linéaire (5 pts)**

Un certain composant électronique est fabriqué une fois par mois par l'entreprise Micro-Systèmes. La quantité fabriquée varie avec la demande du marché. Dans le but de planifier la production et d'établir certaines normes sur le nombre d'hommes minutes exigés pour la production de différents lots de ce composant électronique, le responsable de la production a relevé l'information suivante pour 15 cédules de production. Le nombre d'hommes-minutes est identifié par Y et la quantité fabriquée par X.

Yi	150	192	264	371	300	358	192	134	242	238	226	302	340	182	169
Xi	35	42	64	88	70	85	40	30	55	60	51	72	80	44	39

1. Quelle serait la première étape à franchir avant d'aborder tout calcul préliminaire ?
2. Le responsable de la production envisage d'utiliser le modèle linéaire simple comme modèle prévisionnel. Spécifiez ce modèle et identifiez chacune des composantes du modèle dans le contexte de ce problème ainsi que les hypothèses du modèle.
3. Estimez les coefficients du modèle en adoptant la méthode des Moindres Carrés Ordinaires (MCO).  
On donne :  $\bar{x} = 57$  ;  $\bar{y} = 244$  ;  $\sigma_x^2 = 332.4$  ;  $\sigma_y^2 = 5427.9$  ;  $Cov(X, Y) = 1335.1$
4. D'après l'équation de régression, si le nombre d'unités à fabriquer augmente de 10, quelle sera l'augmentation correspondante du nombre moyen d'hommes- minutes requis ?
5. Déterminez et interprétez le coefficient du détermination  $R^2$  de ce modèle  
On donne : La somme des carrés expliqués par le modèle  $SCE = 80439.5$

**Exercice 3 : QCM (5 pts)** (*Cette feuille n'est pas à rendre veuillez indiquer le numéro de la question et les codes des affirmations correctes sur votre feuille d'examen*)

1. Indiquez les affirmations correctes concernant l'Analyse Factorielle des Correspondances :

- A. Plus la distance euclidienne entre deux lignes du tableau de contingence est faible, plus les deux modalités lignes sont associées ;
- B. Chaque cellule  $\ell_{ij}$  du tableau de profils-lignes présente une probabilité conditionnelle de la modalité  $i$  de la variable ligne sachant la modalité  $j$  de la variable colonne ;
- C. Contrairement à la distance euclidienne, avec la métrique du  $\chi^2$ , la distance entre deux lignes ne dépend pas des poids respectifs des colonnes ;
- D. On peut réaliser l'étude sur deux variables qualitatives dont le nombre de modalités est différents ;
- E. L'objectif de l'analyse de correspondance simple (AFC) est d'étudier la relation entre une variable qualitative et une variable quantitative ;

2. Indiquez les affirmations correctes concernant l'Analyse en Composantes Principales :

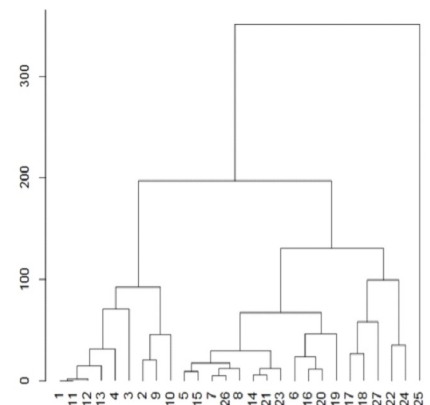
- A. La somme des qualités de représentation d'un individu donné sur tous les axes factoriels est égale à 100% ;
- B. La somme des qualités de représentation de tous les individus sur un axe factoriel donné est égale à 100% ;
- C. Le rang de la matrice de corrélation d'une ACP normée sur  $p$  variables vaut  $p$  ;
- D. La somme des contributions des individus à l'inertie expliquée par un axe factoriel est égale à 1 ;
- E. L'objectif de l'ACP est de réduire la dimension de l'espace en éliminant certaines variables initiales ;

3. Indiquez les affirmations correctes concernant la classification Floue FCM

- A. Le nombre minimal de variables (ou caractères) nécessaires pour effectuer la classification est 2 ;
- B. Pour chaque individu, la somme des degrés d'appartenance aux différentes classes est égale à 1 ;
- C. Pour chaque classe, la somme des degrés d'appartenance des différents individus à cette classe est égale à 1 ;
- D. Les coordonnées du centre d'une classe sont calculées à partir des individus, formant cette classe, pondérés par leurs degrés d'appartenance à cette classe ;
- E. D'une itération à une autre le critère de variabilité intra-classes diminue ;
- F. Un des critères d'arrêt de l'algorithme FCM est que le critère de variabilité intra-classes soit presque nul ;
- G. Plus la distance entre un individu et le centre d'une classe est grande plus son degrés d'appartenance à cette classe est faible ;

4. On applique une CAH sur 27 individus avec comme critère d'agrégation la distance de Ward. On donne l'arbre de classification avec en ordonnée la perte de l'inertie interclasses engendrée. En coupant l'arbre de classification, on a une inertie inter-classes d'environ 680. Combien a-t-on retenue de classes suite à cette coupure ?

- A. 1 classe ;
- B. 2 classes ;
- C. 3 classes ;
- D. 4 classes ;
- E. 5 classes ;
- F. On ne peut pas savoir ;



# Correction Exam 2021

Prévision Théorème.

$$C_1 = \{A_1\} \rightsquigarrow G_1 = A_1 = (2, 10)$$

$$C_2 = \{A_4, A_3, A_5, A_6, A_8\} \rightsquigarrow G_2 =$$

$$C_3 = \{A_7, A_2\} \rightsquigarrow G_3 = \left( \frac{5+8+7+6+4}{5}, \frac{8+4+5+4+3}{5} \right) = (6, 6)$$

Condition d'arrêt K-mans

$$|I_w^{(n)} - I_w^{(n+1)}| < \epsilon \text{ L'arrêt est alors atteint}$$

La partition n'est pas équilibrée.

$$I(G) = I_w + I_b$$

↑  
C'est la somme des  
différences de poids  
entre les deux partitions

$$2) a) S(A, B) = \frac{1}{2} \cdot n \cdot d^2(A, B)$$

à la 1<sup>re</sup> itération les 2 pts qui sont les plus proches sont A3 et A5

$S(A_i, A_j)$	$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$	$A_7$	$A_8$
$A_1$		0						
$A_2$			0					
$A_3$				0				
$A_4$					0			
$A_5$						0		
$A_6$							0	
$A_7$								0
$A_8$								

$$A_3, A_5 \rightsquigarrow G_{35} = (7, 5, 4, 5)$$

$$S(A_1, \{A_3, A_5\}) = \frac{1 \times 2}{3} \cdot d^2(A_1, G_{35})$$

$$S(A_2, G_{35}) = \frac{2}{3} \left( (2-7)^2 + (10-5)^2 \right)$$

$$S(A_4, G_{35}) = \frac{2}{3} \left( (5-7)^2 + (8-4)^2 \right)$$

$$S(A_6, G_{35}) = \frac{2}{3} \left( (6-7)^2 + (4-4)^2 \right)$$

$$S(A_7, G_{35}) = \frac{2}{3} \left( (2-7)^2 + (2-4)^2 \right)$$

$$S(A_8, G_{35}) = \frac{2}{3} \left( (4-7)^2 + (3-4)^2 \right)$$

$$b) I(G) = I_w + I_b$$

3) Kmeans	+	C O H	
Simple, Rapide		La méthode nous fait de trouver la distance optimale.	des (+)
La partition finale de part de l'initialisation		Plus complexe. qu'Kmeans (peut de calculer).	des (-)

4) Nous pouvons utiliser la méthode initiale que le nombre d'individus est très important. et le nombre de classes important. En appliquant tout d'abord la méthode K-means on a un nombre de classes assez élevé. Ensuite on applique la CPM afin d'avoir des classes optimales puis on a un nombre de classes important. Mais que la complexité de cet algorithme est très importante.

Ex 2: 1) On peut trouver le moyen de prouver que  $X$  et  $Y$  agissent sur si la moyenne est bien allongée et peut être agitée par la droite. On a effectivement calculé  $\rho(X, Y)$  la corrélation linéaire.

2)  $Y = \alpha X + \beta + \epsilon$   
 Les  $n$   $H-H$  variables à expliquer en fonction de  $X$ .  $q$  à expliquer.  $H_y: \epsilon$  indépendant de  $X$ .  $E: N(0, \sigma^2)$   
 $\epsilon$  indépendant.

$$\hat{\alpha} = \frac{\text{cov}(X, Y)}{\text{var}(X)} = \frac{1335,1}{332,14} = 4$$

$$\hat{\beta} = \bar{Y} - \hat{\alpha} \bar{X} = 244 - 57 \times 4 = 16.$$

3) si le  $\alpha$  est négatif, si  $X$  augmente,  $Y$  diminue.  $Y$  de 40 unités en moyenne.

$$R^2 = \frac{SCE}{SCT} = \frac{80439,5}{15 \times 10^2} = 0,98$$

98% de variaciones de  $Y$  son explicables  
 (microeconómicas).  
 por las variaciones de  $Q$  (económicas).  
 X explica bien la variable  $Y$ .