

Examen Final

Data Mining et Recherche d'Information

Exercice 1 (6 points)

Répondre brièvement aux questions suivantes.

Q1- Quelles sont les principales architectures des systèmes de recherche d'information. (1.5 pts)

Q2- Expliquer brièvement les termes suivants : crawling, stemming, lemmatisation, stop list. (2 pts)

Q3- Quels sont les principaux critères de performance d'un système de recherche d'information. (1 pt)

Q4- Expliquer comment la phrase suivante peut être traitée par les systèmes de recherche d'information:

« Le langage est le principal facteur d'une bonne communication ». (1.5 pts)

Exercice 2 (6 points)

Soient huit points donnés dans le plan.

a	b	c	d	e	f	g	h
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

- Appliquer la méthode k-means pour partager ces données en deux groupes. (4 pts)
- Vérifier que la méthode k-means tend à réduire la somme des moindres carrés suivante : (2 pts)

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - m_i|^2$$

Où m_i est le centre du cluster C_i et p les points de chaque cluster.

Exercice 3 (08 points)

Le dendrogramme est une représentation graphique d'une classification (clustering) hiérarchique par un arbre.

- Le dendrogramme d'une classification ascendante (ou descendante) est-il unique ? Si OUI dites Comment ? Si NON dites pourquoi ? (2 Pts)
- Comment déterminer le nombre de classes (clusters) à partir du dendrogramme ? Cette méthode est-elle exacte ou approximative ? Quel est le moyen le plus efficace pour avoir un nombre de classes (clusters) proche de la réalité? (2 Pts)
- Considérer la matrice de similarité suivante de cinq documents d1, d2, d3, d4 et d5. Déterminer le dendrogramme résultant de l'application du text clustering hiérarchique ascendant en utilisant le « lien maximum ». (4 Pts)

	d1	d2	d3	d4	d5
d1	0	0.5	0.5	0.6	0.8
d2	0.5	0	0.7	0.6	0.5
d3	0.5	0.7	0	0.6	0.5
d4	0.6	0.6	0.6	0	0.9
d5	0.8	0.5	0.5	0.9	0

Examen Final

Data Mining et Recherche d'Information

Correction

Exercice 1 (6 points)

Répondre brièvement à ces questions.

Q1- Quelles sont les principales architectures des systèmes de recherche d'information. (0.25 pt/architecture)

R1- 1) Single Homogeneous Collection: Full Text Indexing, 2) Single Homogeneous Collection: Use of Catalog Records, 3) Several Similar Collectons: One Computer System, 4) Distributed Architecture: Standard Search Protocols, 5) Distributed Architecture: Meta-search (Broadcast Search), 6) Union Catalog, 7) Web Searching.

Q2- Expliquer brièvement les termes suivants : crawling, stemming, lemmatisation, stop list. (0.5pt/terme)

R2- **Crawling : identifier les documents sur le web**

Stemming : supprimer les suffixes et fusionner les morphèmes.

Lemmatisation : réduire le mot à sa forme de base dans le dictionnaire

Stop list : liste des mots inutiles pour la recherche d'information

Q3- Quels sont les principaux critères de performance d'un système de recherche d'information. (1 pt)

R3- 1) High performance computing : performance du traitement des documents, notamment l'indexation, 2) Human computer interaction: performance de l'interface et clarté des résultats.

Q4- Expliquer comment la phrase suivante peut être traitée par les systèmes de recherche d'information:

« Le langage est le principal facteur d'une bonne communication ». (1.5 pts)

R4- **Suppression des mots stop list : langage, principal, facteur, bonne, communication (0.5 pt)**

Lemmatisation : langage, principal, facteur, bon, communication (0.5 pt)

Stemming : langag, princip, facteur, bon, commun (0.5 pt)

Indexation selon le modèle booléen, vectoriel ou probabiliste.

Exercice 2 (6 points)

Soient huit points donnés dans le plan.

a	b	c	d	e	f	g	h
(1,3)	(3,3)	(4,3)	(5,3)	(1,2)	(4,2)	(1,1)	(2,1)

a) Appliquer la méthode k-means pour partager ces données en deux groupes. (4 pts)

It0 : choix de deux centres : $m_1=a=(1,3)$ et $m_2=e=(1,2)$, $C_1=\{a\}$, $C_2=\{e\}$

It1 : calcul des distances et affectation aux clusters C_1 , C_2

$$d^2(m_1,b) = (1-3)^2 + (3-3)^2 = 4; d^2(m_2, b) = (1-3)^2 + (2-3)^2 = 5, b \rightarrow C_1$$

$$d^2(m_1,c) = (1-4)^2 + (3-3)^2 = 9; d^2(m_2, c) = (1-4)^2 + (2-3)^2 = 10, c \rightarrow C_1$$

$$d^2(m_1,d) = (1-5)^2 + (3-3)^2 = 16; d^2(m_2, d) = (1-5)^2 + (2-3)^2 = 17, d \rightarrow C_1$$

$$d^2(m_1,f) = (1-4)^2 + (3-2)^2 = 10; d^2(m_2, f) = (1-4)^2 + (2-2)^2 = 9, f \rightarrow C_2$$

$$d^2(m_1,g) = (1-1)^2 + (3-1)^2 = 4; d^2(m_2, g) = (1-1)^2 + (2-1)^2 = 1, g \rightarrow C_2$$

$$d^2(m_1, h) = (1-2)^2 + (3-1)^2 = 5; d^2(m_2, h) = (1-2)^2 + (2-1)^2 = 2, h \rightarrow C_2$$

$$C_1 = \{a, b, c, d\}, C_2 = \{e, f, g, h\}$$

Calcul des nouveaux centres: $m_1 = ((1+3+4+5)/4), (3+3+3+3)/4 = (3.25, 3)$

$$m_2 = ((1+4+1+2)/4), (2+2+1+1)/4 = (2, 1.5)$$

It2 : calcul des distances et affectation aux clusters C1, C2

$$d^2(m_1, a) = (3.25-1)^2 + (3-3)^2 = 5.06; d^2(m_2, a) = (2-1)^2 + (1.5-3)^2 = 3.25, a \rightarrow C_2$$

$$d^2(m_1, b) = (3.25-3)^2 + (3-3)^2 = 0.06; d^2(m_2, b) = (2-3)^2 + (1.5-3)^2 = 3.25, b \rightarrow C_1$$

$$d^2(m_1, c) = (3.25-4)^2 + (3-3)^2 = 0.56; d^2(m_2, c) = (2-4)^2 + (1.5-3)^2 = 6.25, c \rightarrow C_1$$

$$d^2(m_1, d) = (3.25-5)^2 + (3-3)^2 = 3.06; d^2(m_2, d) = (2-5)^2 + (1.5-3)^2 = 11.25, d \rightarrow C_1$$

$$d^2(m_1, e) = (3.25-1)^2 + (3-2)^2 = 6.06; d^2(m_2, e) = (2-1)^2 + (1.5-2)^2 = 1.25, e \rightarrow C_2$$

$$d^2(m_1, f) = (3.25-4)^2 + (3-2)^2 = 1.56; d^2(m_2, f) = (2-4)^2 + (1.5-2)^2 = 4.25, f \rightarrow C_1$$

$$d^2(m_1, g) = (3.25-1)^2 + (3-1)^2 = 9.06; d^2(m_2, g) = (2-1)^2 + (1.5-1)^2 = 1.25, g \rightarrow C_2$$

$$d^2(m_1, h) = (3.25-2)^2 + (3-1)^2 = 4.56; d^2(m_2, h) = (2-2)^2 + (1.5-1)^2 = 0.25, h \rightarrow C_2$$

$$C_1 = \{b, c, d, f\}, C_2 = \{a, e, g, h\}$$

Calcul des nouveaux centres: $m_1 = ((3+4+5+4)/4), (3+3+3+2)/4 = (4, 2.75)$

$$m_2 = ((1+1+1+2)/4), (3+2+1+1)/4 = (1.25, 1.75)$$

It3 : calcul des distances et affectation aux clusters C1, C2

$$d^2(m_1, a) = (4-1)^2 + (2.75-3)^2 = 9.06; d^2(m_2, a) = (1.25-1)^2 + (1.75-3)^2 = 2.12, a \rightarrow C_2$$

$$d^2(m_1, b) = (4-3)^2 + (2.75-3)^2 = 1.06; d^2(m_2, b) = (1.25-3)^2 + (1.75-3)^2 = 4.62, b \rightarrow C_1$$

$$d^2(m_1, c) = (4-4)^2 + (2.75-3)^2 = 0.06; d^2(m_2, c) = (1.25-4)^2 + (1.75-3)^2 = 9.12, c \rightarrow C_1$$

$$d^2(m_1, d) = (4-5)^2 + (2.75-3)^2 = 1.06; d^2(m_2, d) = (1.25-5)^2 + (1.75-3)^2 = 15.62, d \rightarrow C_1$$

$$d^2(m_1, e) = (4-1)^2 + (2.75-2)^2 = 9.56; d^2(m_2, e) = (1.25-1)^2 + (1.75-2)^2 = 0.12, e \rightarrow C_2$$

$$d^2(m_1, f) = (4-4)^2 + (2.75-2)^2 = 0.56; d^2(m_2, f) = (1.25-4)^2 + (1.75-2)^2 = 7.62, f \rightarrow C_1$$

$$d^2(m_1, g) = (4-1)^2 + (2.75-1)^2 = 12.06; d^2(m_2, g) = (1.25-1)^2 + (1.75-1)^2 = 0.62, g \rightarrow C_2$$

$$d^2(m_1, h) = (4-2)^2 + (2.75-1)^2 = 5.06; d^2(m_2, h) = (1.25-2)^2 + (1.75-1)^2 = 1.12, h \rightarrow C_2$$

$$C_1 = \{b, c, d, f\}, C_2 = \{a, e, g, h\}$$

Stabilisation des clusters. Fin

Résultats : $C_1 = \{b, c, d, f\}, C_2 = \{a, e, g, h\}$

b) Vérifier que la méthode k-means tend à réduire la somme des moindres carrés suivante : (2pts)

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - m_i|^2$$

Où m_i est le centre du cluster C_i et p les points de chaque cluster.

Résultat de l'itération 1 : $C_1 = \{a, b, c, d\}, C_2 = \{e, f, g, h\}$ avec $m_1 = (3.25, 3), m_2 = (2, 1.5)$

$$E_1 = (|1-3.25|^2 + |3-3.25|^2 + |4-3.25|^2 + |5-3.25|^2 + |3-3|^2 + |3-3|^2 + |3-3|^2 + |3-3|^2) +$$

$$(|1-2|^2 + |4-2|^2 + |1-2|^2 + |2-2|^2 + |2-1.5|^2 + |2-1.5|^2 + |1-1.5|^2 + |1-1.5|^2)$$

$$E_1 = (5.06 + 0.06 + 0.56 + 3.06 + 0 + 0 + 0 + 0) + (1 + 4 + 1 + 0 + 0.25 + 0.25 + 0.25 + 0.25) = 15.74$$

Résultat de l'itération 2 : $C_1 = \{b, c, d, f\}, C_2 = \{a, e, g, h\}$ avec $m_1 = (4, 2.75), m_2 = (1.25, 1.75)$

$$E_2 = (|3-4|^2 + |4-4|^2 + |5-4|^2 + |4-4|^2 + |3-2.75|^2 + |3-2.75|^2 + |3-2.75|^2 + |2-2.75|^2) +$$

$$(|1-1.25|^2 + |1-1.25|^2 + |1-1.25|^2 + |2-1.25|^2 + |3-1.75|^2 + |2-1.75|^2 + |1-1.75|^2 + |1-1.75|^2)$$

$$E_2 = (1 + 0 + 1 + 0 + 0.06 + 0.06 + 0.06 + 0.56) + (0.06 + 0.06 + 0.06 + 0.56 + 1.56 + 0.06 + 0.56 + 0.56) = 6.22$$

Conclusion : $E_2 < E_1$ d'où la formule est vérifiée.

Exo3 : (08 points)

- 1- Le dendrogramme d'un clustering n'est pas unique ; il dépend de la stratégie de regroupement : lien minimum, maximum ou moyen. En plus, si la distance minimale choisie n'est pas unique, le choix aléatoire diversifie le dendrogramme.
- 2- On détermine le nombre de clusters à partir d'un dendrogramme, en déterminant le nombre de points d'intersection entre la droite $y=d$ et le dendrogramme, d étant la distance choisie. Cette méthode est exacte. Le moyen le plus efficace pour avoir un nombre de classes (clusters) proche de la réalité se fait généralement à l'aide de l'avis d'un expert dans le domaine.
- 3- Le dendrogramme résultant de l'application du text clustering hiérarchique ascendant de cinq documents d1, d2, d3, d4 et d5, en utilisant le « **lien maximum** ».

	d1	d2	d3	d4	d5
d1	0				
d2	<u>0.5</u>	0			
d3	0.5	0.7	0		
d4	0.6	0.6	0.6	0	
d5	0.8	0.5	0.5	0.9	0

grouper (d1, d2)

	d1d2	d3	d4	d5
d1d2	0			
d3	0.7	0		
d4	0.6	0.6	0	
d5	0.8	<u>0.5</u>	0.9	0

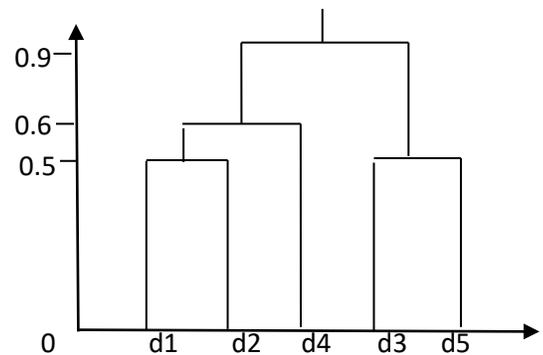
grouper (d3, d5)

	d1d2	d3d5	d4
d1d2	0		
d3d5	0.8	0	
d4	<u>0.6</u>	0.9	0

grouper (d1d2, d4)

	d1d2d4	d3d5
d1d2d4	0	
d3d5	<u>0.9</u>	0

grouper(d1d2d4, d3d5)



Dendrogramme