

TD Méthodes de clustering (Réponses)
M1-ATD
Ioannis Partalas, Eric Gaussier

D'après M.-R. Amini & E. Gaussier, *Recherche d'information*, Eyrolles 2013

1. Algorithme des k-moyennes

(a) Montrer qu'entre les itérations t et $t + 1$ de l'algorithme 1 on a :

$$\mathcal{L}(G_1^{(t+1)}, \dots, G_K^{(t+1)}; \mathbf{r}_1^{(t)}, \dots, \mathbf{r}_K^{(t)}) \leq \mathcal{L}(G_1^{(t)}, \dots, G_K^{(t)}; \mathbf{r}_1^{(t)}, \dots, \mathbf{r}_K^{(t)})$$

Réponse : L'inégalité découle de la définition de la réaffectation des exemples :

$$G_k^{(t+1)} \leftarrow \{d : \|\mathbf{d} - \mathbf{r}_k^{(t)}\|_2^2 \leq \|\mathbf{d} - \mathbf{r}_l^{(t)}\|_2^2, \forall l \neq k, 1 \leq l \leq K\}$$

(b) Montrer aussi que :

$$\mathcal{L}(G_1^{(t+1)}, \dots, G_K^{(t+1)}; \mathbf{r}_1^{(t+1)}, \dots, \mathbf{r}_K^{(t+1)}) \leq \mathcal{L}(G_1^{(t+1)}, \dots, G_K^{(t+1)}; \mathbf{r}_1^{(t)}, \dots, \mathbf{r}_K^{(t)})$$

Réponse : Considérons la fonction de coût pour une seule classe G :

$$\begin{aligned} \mathcal{L}(G, z) &= \sum_{d \in G} \|d - z\|^2 \\ &= \sum_{d \in G} \|d - CG + CG - z\|^2 \\ &= \sum_{d \in G} \|d - CG\|^2 + \sum_{d \in G} \|CG - z\|^2 + 2 \sum_{d \in G} \langle d - CG, CG - z \rangle \\ &= \sum_{d \in G} \|d - CG\|^2 + |G| \|CG - z\|^2 + 2 \underbrace{\langle \sum_{d \in G} d - |G|CG, CG - z \rangle}_{|G|CG} \\ &= \mathcal{L}(G, CG) + |G| \|CG - z\|^2 \end{aligned}$$

Nous avons de ce fait :

$$\mathcal{L}(G_1^{(t+1)}, \dots, G_K^{(t+1)}; \mathbf{r}_1^{(t+1)}, \dots, \mathbf{r}_K^{(t+1)}) \leq \mathcal{L}(G_1^{(t)}, \dots, G_K^{(t)}; \mathbf{r}_1^{(t)}, \dots, \mathbf{r}_K^{(t)})$$

(c) En déduire que la fonction de coût de l'algorithme k-moyennes décroît à chaque itération.

Réponse : D'après les questions (a) et (b) et en exploitant la positivité de la fonction considérée nous avons $\forall t$:

$$0 \leq \mathcal{L}(G_1^{(t+1)}, \dots, G_K^{(t+1)}; \mathbf{r}_1^{(t+1)}, \dots, \mathbf{r}_K^{(t+1)}) \leq \mathcal{L}(G_1^{(t)}, \dots, G_K^{(t)}; \mathbf{r}_1^{(t)}, \dots, \mathbf{r}_K^{(t)})$$

La fonction \mathcal{L} décroît donc à chaque itération de l'algorithme.

2. Classification par méthodes agglomératives hiérarchiques ascendantes

Question 1 Montrer que la méthode du lien unique est stable pour la meilleure fusion. On rappelle que la méthode du lien unique est fondée sur la distance entre classes :

$$\text{sim}_{\text{lu}}(G_k, G_l) = \max_{d \in G_k, d' \in G_l} \text{sim}(d, d')$$

Réponse : Nous avons :

$$\begin{aligned} \text{sim}(G_k^{(r+1)}, G^{(r)}) &= \max(\text{sim}(G_k^{(r)}, G_{mf(k)}^{(r)}), \text{sim}(G_k^{(r)}, G_l^{(r)})) \\ &= \text{sim}(G_k^{(r)}, G_{mf(k)}^{(r)}) \end{aligned}$$

Question 2 Quelle est la complexité de l'algorithme 3 qui utilise un tableau de meilleure fusion en lieu et place des files de priorité ?

Réponse : $O(N^2)$

Question 3 Expliquer pourquoi la stabilité de la meilleure fusion est importante pour le bon déroulement de cet algorithme (donner un exemple simple).

Réponse : La stabilité est importante car nous permet d'utiliser un tableau de meilleure fusion et ainsi mettre à jour la matrice C à chaque itération en $O(N)$. On peut considérer un exemple avec 4 documents et l'algorithme du lien complet et montrer qu'on ne peut pas utiliser un tableau de meilleure fusion car la stabilité de la meilleure fusion ne s'applique pas.

Question 4 Montrer que la méthode par lien unique est monotone.

Réponse : Si $G_1^{(r+1)} \neq G^{(r)}$ et $G_2^{(r+1)} \neq G^{(r)}$, l'inégalité découle de la construction même du dendrogramme (sinon les deux classes $G_1^{(r+1)}$ et $G_2^{(r+1)}$ auraient été fusionnées avant les deux classes $G_1^{(r)}$ et $G_2^{(r)}$).

Supposons donc que $G_1^{(r+1)} = G^{(r)}$. Nous avons :

$$\text{sim}(G_1^{(r+1)}, G_2^{(r+1)}) = \max_{d \in G_1^{(r)} \cup G_2^{(r)}, d' \in G_2^{(r+1)}} \text{sim}(d, d')$$

Cette dernière quantité est équivalente à :

$$\max\left(\max_{d \in G_1^{(r)}, d' \in G_2^{(r+1)}} \text{sim}(d, d'), \max_{d \in G_2^{(r)}, d' \in G_2^{(r+1)}} \text{sim}(d, d'),\right)$$

et :

$$\max(\text{sim}(G_1^{(r)}, G_2^{(r+1)}), \text{sim}(G_2^{(r)}, G_2^{(r+1)}))$$

Mais $\text{sim}(G_1^{(r)}, G_2^{(r+1)}) \leq \text{sim}(G_1^{(r)}, G_2^{(r)})$ car sinon ce sont les classes $G_1^{(r)}$ et $G_2^{(r+1)}$ qui auraient été fusionnées à l'étape r . De même, $\text{sim}(G_2^{(r)}, G_2^{(r+1)}) \leq \text{sim}(G_1^{(r)}, G_2^{(r)})$, ce qui montre la monotonie du lien simple.