

Προεργασία

- Αφαίρεση ανούσιων λέξεων (stop words)
- Εύρεση ριζών (stemming)

$$q = \text{term}_{i_1} + \text{term}_{i_2} + \dots + \text{term}_{i_m}$$

↓

{ Doc1, ..., DocK }

→ (0, 1, 0, 1, ...)

↓₁₂ ↓₁₃

↓: (1, 2) (1, 3) ...

$\cos(x, y)$ δείχνει την συσχέτιση x, y

$$\cos(\theta(q, a_i)) > \underline{\text{tol}}$$

↓

i στήλη του A

↑

αφαιρεί

$$(A = \underline{U} \underline{\Sigma} \underline{V}^T)$$

$$\cos(\theta(q, a_i)) = \frac{a_i^T q}{\|a_i\|_2 \cdot \|q\|_2}$$

$$\begin{matrix} q^T A \\ \vdots \\ \vdots \\ \vdots \end{matrix} \begin{matrix} \text{Dim 1} & \text{Dim 2} & \dots & \text{Dim } k \\ \underline{\hat{a}_1} & \underline{\hat{a}_2} & & \underline{\hat{a}_k} \end{matrix} \rightsquigarrow \left[\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right] \left[\begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right]$$

A

U

$$\approx U_K \underbrace{\sum_k V_k^T}_{H_K} = \underbrace{U_K}_{\downarrow} H_K$$

A ≈ U_K H_K

a_i ≈ U_K h_i

$$\begin{matrix} m \\ \vdots \\ \vdots \\ \vdots \end{matrix} \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \rightsquigarrow \begin{matrix} k \\ \vdots \\ \vdots \\ \vdots \end{matrix} \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix}$$

a_i ≈ U_K · h_i όπου h_i η ε στήλη του H_K

â_i ≈ U_K · h_i → [i] ↓ [i]

q^T A_K = (q^T (U_K H_K)) = (U_K^T q) · H_K

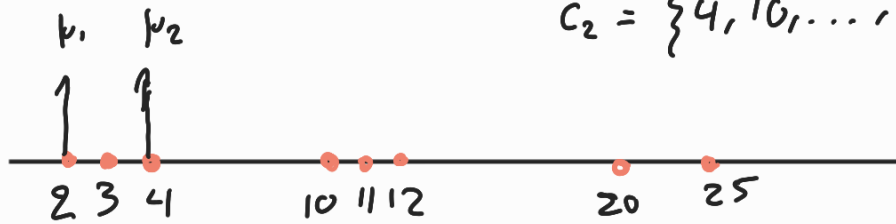
$$\cos \theta_i = \frac{(U_K^T q) \cdot h_i}{\|U_K^T q\|_2 \|h_i\|_2} = \frac{q_K^T \cdot h_i}{\|q_K\| \cdot \|h_i\|}$$

Συσταδοποίηση (Clustering) βασισμένη σε αντιπροσώπων

$k=2$

$$C_1 = \{2, 3\}$$

$$C_2 = \{4, 10, \dots, 25\}$$



Έστω η σητεία $\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\} \in \mathbb{R}^d$

\downarrow \downarrow
 $\begin{pmatrix} x_1^1 \\ x_1^2 \\ \vdots \\ x_1^d \end{pmatrix}$

Πληθος συστάδων $\{C_1, C_2, \dots, C_k\}$

Καθε συστάδα C_i υπάρχει ένας αντιπρόσωπος. Είναι συνδυασμός δεικτη είναι ο μέσος όλος των σημείων μ_i που ονομάζεται κέντρο βάρους (centroid)

$$\underline{\mu}_i = \frac{1}{n_i} \sum_{j \in C_i} x_j = \frac{1}{|C_i|} \sum_{j \in C_i} x_j$$

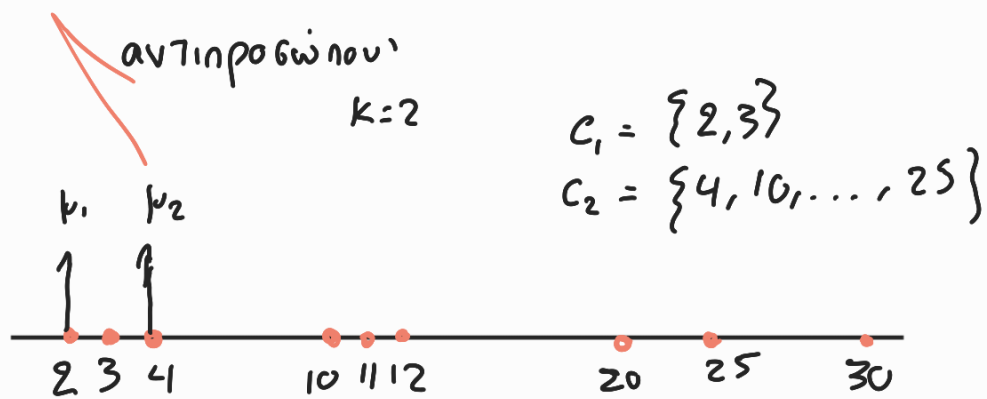
Αλγόριθμος K μέσων

$$C' = \{C'_1, C'_2, \dots, C'_k\}$$

Η αξιολόγηση της συσταδοποίησης θα γίνει μέσα από την μετρική

$$\sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

Στόχος μας είναι να βρούμε C' που ελαχιστοποιεί την παραπάνω μετρική



2ο βήμα : Βρίσκω τους νέους μέσους όρους

$$\mu_1^2 = \frac{2+3}{2} = 2.5$$

$$\mu_2^2 = \frac{4+10+11+12+20+25+30}{7} = \frac{82}{7} = \frac{112}{7} = \underline{\underline{16}}$$

$$C_1^{\mu_1^2} = \{2, 3, 4\} \quad \parallel \quad C_2^{\mu_2^2} = \{10, 11, \dots, 30\}$$

$$\mu_1^3 = 3$$

$$\mu_2^3 = \frac{10+11+12+20+25+30}{6} = \frac{108}{6} = 18$$

Η διαδικασία συνεχίζεται μέχρι ???

$$\{\mu_1^{(n)}\} = \{\mu_1^1, \mu_1^2, \mu_1^3, \dots\}$$

$$\{\mu_2^{(n)}\} = \{\mu_2^1, \mu_2^2, \dots\}$$

$$\sum_{i=1}^k \|\mu_i^m - \mu_i^{m-1}\|^2 \leq \underline{\underline{\epsilon}}$$