

# Flat clustering

## Algoritmul $K$ -means

novembrie 2017

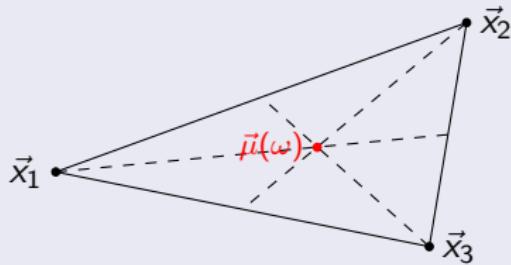
# Flat clustering

- Se presupune dată o colecție de  $N$  documente reprezentate ca o mulțime de puncte  $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$  într-un spațiu vectorial.
- Se dorește împărțirea mulțimii  $D$  în  $K$  submulțimi distincte  $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ :

$$\omega_1 \cup \omega_2 \cup \dots \cup \omega_K = D \text{ și } \omega_i \cap \omega_j = \emptyset \text{ dacă } i \neq j$$

Mulțimile  $\omega_1, \dots, \omega_K$  se numesc **clustere**. Fiecare cluster  $\omega \in \Omega$  are un **centroid**  $\vec{\mu}(\omega)$  care se calculează cu formula  $\vec{\mu}(\omega) = \frac{\sum_{\vec{x} \in \omega} \vec{x}}{|\omega|}$

Exemplu: centroidul unui cluster de 3 documente  $\omega = \{\vec{x}_1, \vec{x}_2, \vec{x}_3\}$



# Algoritmul $K$ -means

- ▶ Suma reziduală de pătrate (*engl. residual sum of squares*) a unui cluster  $\omega_k \in \Omega$  este

$$RSS_k := \sum_{\vec{x} \in \omega_k} |\vec{x} - \vec{\mu}(\omega_k)|^2$$

- ▶ Suma reziduală de pătrate a unei partiții  $\Omega = \{\omega_1, \dots, \omega_K\}$  este

$$RSS := \sum_{k=1}^K RSS_k$$

- Algoritmul  $K$ -means are ca obiectiv să minimizeze valoarea lui  $RSS$  pentru o mulțime dată de documente  $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$ .

# Algoritmul K-means

## Pseudocod

```
K-MEANS( $\{\vec{x}_1, \dots, \vec{x}_N\}$ ,  $K$ )
1  $(\vec{s}_1, \dots, \vec{s}_K) := \text{SELECTRANDOMSEEDS}(\{\vec{x}_1, \dots, \vec{x}_N\}, K)$ 
2 for  $k := 1$  to  $K$  do
3    $\vec{\mu}_k := \vec{s}_k$ 
4 while criteriul de terminare nu este satisfăcut do
5   for  $k := 1$  to  $K$  do
6      $\omega_k := \emptyset$ 
7     for  $n := 1$  to  $N$  do
8        $j := \arg \min_j' |\vec{\mu}_{j'} - \vec{x}_n|$ 
9        $\omega_j := \omega_j \cup \{\vec{x}_n\}$  (reatribuirea vectorilor în clustere)
10    for  $k := 1$  to  $K$  do
11       $\vec{\mu}_k := \frac{1}{|\omega_k|} \sum_{\vec{x} \in \omega_k} \vec{x}$ 
12 return  $\{\vec{\mu}_1, \dots, \vec{\mu}_K\}$ 
```

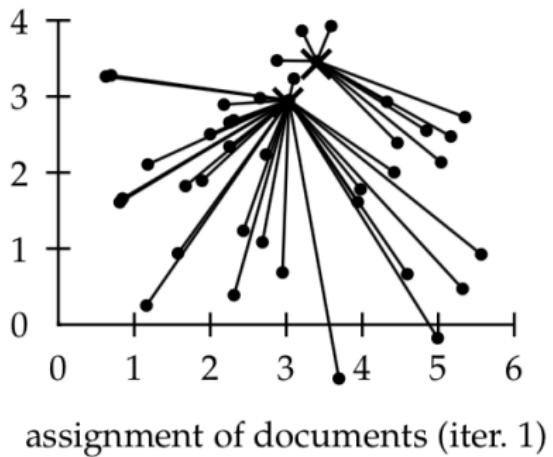
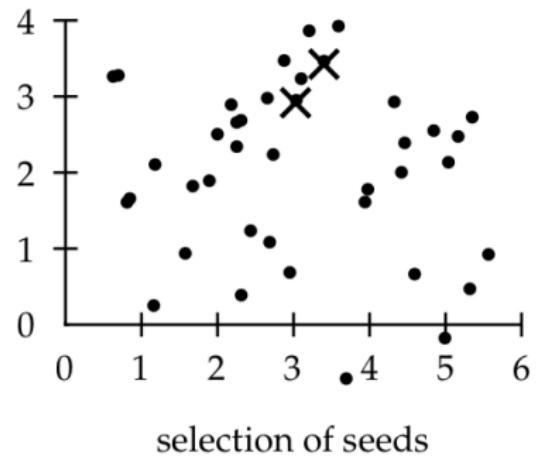
### Observații

- Mai întâi, se aleg la întâmplare  $K$  centroizi inițiali de clustere (engl. *seeds*)
- algoritmul repetă modificarea centrelor inițiale încât să minimizeze  $RSS$ : la fiecare iterație: (1) se reatribuie fiecare vector  $\vec{x}$  în clusterul centroidului celui mai apropiat de  $\vec{x}$  și apoi (2) se recalculează coordonatele centroizilor

# Algoritmul K-means

Exemplu ilustrat

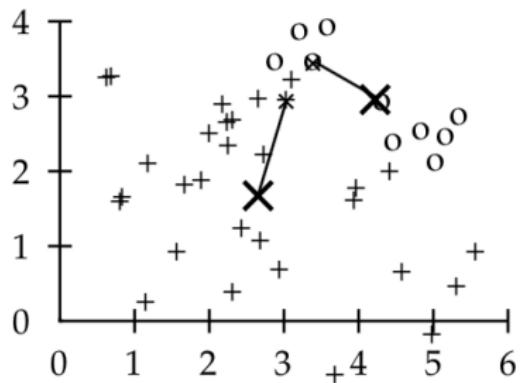
Evoluția algoritmului



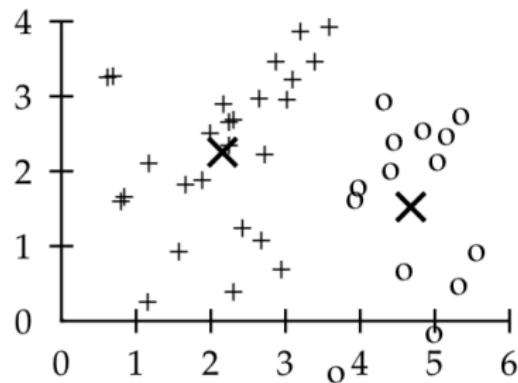
# Algoritmul K-means

Exemplu ilustrat

Evoluția algoritmului în 9 îterății:



recomputation/movement of  $\vec{\mu}$ 's (iter. 1)

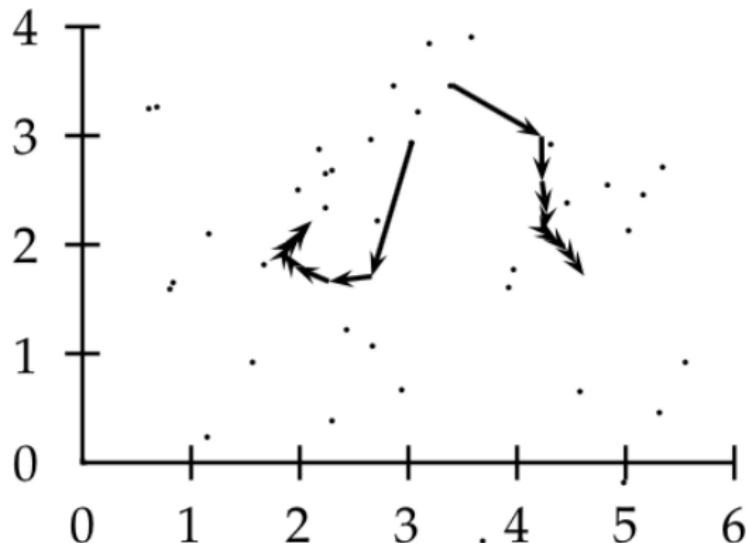


$\vec{\mu}$ 's after convergence (iter. 9)

# Algoritmul K-means

Exemplu ilustrat

Evoluția algoritmului în 9 îterări:



# Algoritmul K-means

## Criterii posibile de terminare

- ① După un număr prestabilit de iterații
- ② Când se ajunge la o iterație care nu modifică atribuirea documentelor în clustere
- ③ Când centroizii  $\vec{\mu}_k$  nu se mai schimbă de la o iterație la alta
- ④ Când  $RSS$  devine mai mică dacă o valoare prestabilită
- ⑤ Când diferența dintre două  $RSS$ -uri consecutive este mai mică decăt o valoare  $\theta$

# Algoritmul K-means

Criterii posibile de terminare

- ① După un număr prestabilit de iterații
- ② Când se ajunge la o iterație care nu modifică atribuirea documentelor în clustere
- ③ Când centroizii  $\vec{\mu}_k$  nu se mai schimbă de la o iterație la alta
- ④ Când  $RSS$  devine mai mică dacă o valoare prestabilită
- ⑤ Când diferența dintre două  $RSS$ -uri consecutive este mai mică decăt o valoare  $\theta$

**REMARCA:** *K-means converge în sensul că dacă definim*

$$RSS_k(\vec{v}) = \sum_{\vec{x} \in \omega_k} |\vec{v} - \vec{x}|^2 \quad \text{și} \quad RSS(\vec{v}_1, \dots, \vec{v}_K) = \sum_{k=1}^K RSS_k(\vec{v}_k)$$

atunci  $(\vec{\mu}_1, \dots, \vec{\mu}_K)$  este punct de **minim local** pentru  
 $RSS(\vec{v}_1, \dots, \vec{v}_K)$

*(vezi slide următor)*

# Algoritmul $K$ -means

De ce converge  $RSS$ ?

Presupunem că documentele sunt puncte în spațiul euclidean  $M$ -dimensional  $\mathbb{R}^M$

$$\Rightarrow RSS_k(\vec{v}) = \sum_{\vec{x} \in \omega_k} \sum_{m=1}^M (v_m - x_m)^2$$

$$\Rightarrow \frac{\partial RSS_k(\vec{v})}{\partial v_m} = \sum_{\vec{x} \in \omega_k} 2(v_m - x_m)$$

## Proprietăți ale algoritmului $K$ -means

- ① În pașii 8-9,  $RSS(\vec{v}_1, \dots, \vec{v}_K) = \sum_{k=1}^K \sum_{\vec{x} \in \omega_k} |\vec{v}_k - \vec{x}|^2$  scade pentru  $\vec{v} = (\vec{\mu}_1, \dots, \vec{\mu}_k)$  deoarece fiecare document  $\vec{x} \in \omega_k$  este reatribuit în clusterul  $\omega_j$  pentru care  $|\vec{v}_j - \vec{x}| \leq |\vec{v}_k - \vec{x}|$
- ② În pașii 10-11 se recalculează fiecare  $\vec{\mu}_k$  a.î. fiecare  $RSS_k(\vec{v})$  să ia valoarea minimă pentru  $\vec{v} = \vec{\mu}_k$ .

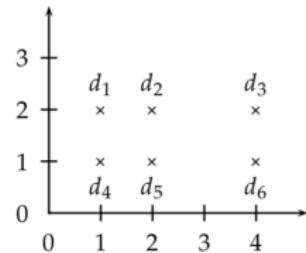
# Algoritmul K-means

## Limitări

Este de dorit să se calculeze  $K$  centroizi  $(\vec{\mu}_1, \dots, \vec{\mu}_K)$  care sunt un **minim global** pentru  $RSS$

- Algoritmul  $K$ -means garantează calculul unui **minim local**
  - Uneori, minim local  $\neq$  minim global
  - Rezultatul calculat depinde de centroizii inițiali aleși
    - ▶ selecția unui centroid inițial îndepărtat de celelalte documente produce generarea de **clustere singleton**, sau clustere vide.

### EXEMPLU:



- ▶  $K$ -means pt. centroizii inițiali  $d_2$  și  $d_5$  converge la  $\{\{d_1, d_2, d_3\}, \{d_4, d_5, d_6\}\}$
- ▶  $K$ -means pt. centroizii inițiali  $d_2$  și  $d_3$  converge la  $\{\{d_1, d_2, d_4, d_5\}, \{d_3, d_6\}\}$

# Complexitatea algoritmului

- ▶ Distanța dintre două puncte în  $\mathbb{R}^M$  durează ...
- ▶ Pașii 7-9 (reatribuirea în clustere) durează ...
- ▶ Pașii 10-11 (recalculatearea centroizilor) durează ...
- ▶  $K$ -means care efectuează  $I$  iterații durează ...

# Determinarea numărului $K$ de clustere

## Metode euristice

Se presupune dată o mulțime de  $N$  documente:  $D = \{\vec{x}_1, \dots, \vec{x}_N\}$

- Algoritmul  $K$ -means ia ca **parametru de intrare** valoarea lui  $K$  și calculează  $K$  centroizi pentru  $K$  clustere  $\omega_1, \dots, \omega_K$
- Uneori e greu de să se estimeze o valoare bună pentru  $K$

Problemă: **Cum putem găsi o valoare plauzibilă pentru  $K$ ?**

# Determinarea numărului $K$ de clustere

## Metode euristice

Se presupune dată o mulțime de  $N$  documente:  $D = \{\vec{x}_1, \dots, \vec{x}_N\}$

- Algoritmul  $K$ -means ia ca **parametru de intrare** valoarea lui  $K$  și calculează  $K$  centroizi pentru  $K$  clustere  $\omega_1, \dots, \omega_K$
- Uneori e greu de să se estimeze o valoare bună pentru  $K$

Problemă: **Cum putem găsi o valoare plauzibilă pentru  $K$ ?**

**EURISTICA 1:** RSS să aibe valoare minimă:

⇒ Dacă  $K = N$  atunci  $\omega_i = \{\vec{x}_i\}$  pentru  $1 \leq i \leq K$  și  $RSS = 0$ .  
Deși  $RSS = 0$ , acesta nu este un clustering optim: **sunt prea multe clustere!**

# Determinarea numărului $K$ de clustere

## Metode euristice

Se presupune dată o mulțime de  $N$  documente:  $D = \{\vec{x}_1, \dots, \vec{x}_N\}$

- Algoritmul  $K$ -means ia ca **parametru de intrare** valoarea lui  $K$  și calculează  $K$  centroizi pentru  $K$  clustere  $\omega_1, \dots, \omega_K$
- Uneori e greu de să se estimeze o valoare bună pentru  $K$

Problemă: **Cum putem găsi o valoare plauzibilă pentru  $K$ ?**

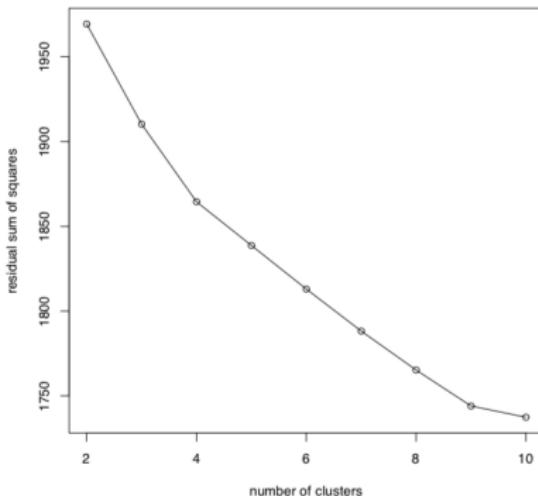
**EURISTICA 1:** RSS să aibe valoare minimă:

⇒ Dacă  $K = N$  atunci  $\omega_i = \{\vec{x}_i\}$  pentru  $1 \leq i \leq K$  și  $RSS = 0$ . Deși  $RSS = 0$ , acesta nu este un clustering optim: **sunt prea multe clustere!**

**EURISTICA 2:** Pentru mai multe valori valori ale lui  $K$ , se calculează  $\widehat{RSS}_{\min}(K) = \min\{rss_1, \dots, rss_i\}$  unde  $rss_j$  ( $1 \leq j \leq i$ ) este valoarea lui RSS pentru un set de  $K$  centroizi inițiali  $(\vec{s}_{j,1}, \dots, \vec{s}_{j,K})$ .  $i$  are o valoare prestabilită, de ex.,  $i = 10$ . Se alege  $K$  pentru care vectorul  $(\widehat{RSS}_{\min}(K+1) - \widehat{RSS}_{\min}(K), 1)$  începe să și mențină direcția.

# Determinarea numărului $K$ de clustere

Estimarea valorii lui  $K$  cu euristică 2



Exemplu ilustrat pentru o mulțime de 1203 documente Reuters-RCV1

- ▷ curba  $\widehat{RSS}_{\min}$  se netezește semnificativ după  $K = 4$  și  $K = 9$   
⇒ 4 și 9 sunt valori plauzibile pentru  $K$ .
- ▷ pentru  $K = 4$ , algoritmul  $K$ -means calculează centroizi în jurul categoriilor *China*, *Germany*, *Russia* și *Sports*.

# Determinarea numărului $K$ de clustere

## Alte euristici

- ▶ Minimizarea distorsiunii ca măsură a complexității modelului

$$K = \arg \min_K (RSS_{\min}(K) + \lambda K)$$

unde  $0 \leq \lambda \leq 1$  este un factor de pondere.

- ▶ Criteriul AIC (engl. *Aikake Information Criterion*) are forma generală

$$K = \arg \min_K (-2 L(K) + 2 q(K))$$

unde ...

Pentru algoritmul *Kmeans*, are forma specială

$$K = \arg \min_K (RSS_{\min}(K) + 2 M K)$$

# Bibliografie

Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze:  
Capitolul 16 din **AN INTRODUCTION TO INFORMATION  
RETRIEVAL**.  
Ediție online (c) 2009 Cambridge UP.