

**A  $K - means$  klaszterező**

Adott  $x_1, x_2, \dots, x_n \in X$   $n$ -darab minta. A feladat ezen minták  $k$  db csoportba, klaszterbe sorolása. A klaszterezés a klaszterközéppontok meghatározásával történik és minden pont a hozzá legközelebbi középponthez tartozó klaszterbe kerül. Ehhez definiálnunk kell egy távolságfüggvényt a minták felett.

*Jelölések:*

$\mu_i$ : az  $i$ -edik klaszter középpontja.

$r_{ij} = 1$ , ha az  $i$ -edik klaszterbe soroljuk a  $j$ -edik mintát; 0, egyébként.

Ezek alapján fölírhatunk egy hibafüggvényt:

$$Err = \sum_i \sum_j r_{ij} \|\mu_i - x_j\|^2 \rightarrow \min$$

Amylet minimalizálni szeretnénk, azaz az  $i$ -edik klaszterbe tartozó minták távolságnégyzetösszege az  $i$ -edik középponttól minimális legyen. Az optimum ott van, ahol az  $Err$  szerinti deriváltja 0, azaz ha a klaszterközéppontok egybeesnek a klaszterekhez tartozó pontokkal.

*Algoritmus:*

1. inicializáljuk a  $\mu_j$  klaszterközéppontokat (pl. egyenletes eloszlással a térben)
2. legyen  $r_{ij} = 1$  ha az  $x_i$  a  $\mu_j$ -hez van legközelebb, egyébként 0. (minták középponthez rendelése)
3. legyen  $\mu_j = \frac{\sum_i r_{ij} x_i}{\sum_i r_{ij}}$  minden  $j = 1, 2, \dots, k$ -ra. (középpontok újradefiniálása)
4. ismételjük a 2. ponttól, amíg változik valamelyik klaszterközéppont.

*Problémák a  $k - means$ -el:*

1. ha egy csoportra két közepet inicializálunk (túl nagy  $k$  választása)
2. ha két csoport közé inicializálunk egy közepet (túl kicsi  $k$  választása)

*Példa applet:* <http://www.rob.cs.tu-bs.de/content/04-teaching/06-interactive/Kmeans/Kmeans.html>

*Következmény:* helyesen kell megválasztani a közepet... ( $X - means$ )

## Az $X - means$ klaszterező

Az  $X - means$  algoritmus abban különbözik a  $K - means$ -től, hogy nem egy konkrét  $k$  értéket definiálunk, hanem egy alsó és felső korlátot. Az algoritmus a minimális  $k$  érték beállításával futtatja a  $K - means$  klaszterezőt, majd növeli a  $k$  értékét, amíg nem romlik a négyzetes hiba.

*Pontnövelési stratégiák:*

**Egy középpont kettéosztása:** valamilyen heurisztika alapján válasszunk egy középpontot és azt osszuk ketté.

**A középpontok felének a kettéosztása:** valamilyen heurisztika alapján válasszuk ki a középpontok felét, majd azokat osszuk ketté.

(Egy lehetséges heurisztika a megfelelő középpont választására: válasszuk azt ahol a középponthez tartozó négyzetes eltérés maximális.)

**A középpontok számának dinamikus növelése:** Egy harmadik lehetséges módszert D. Pelleg és A. Moore írt le egy 2000-es cikkükben. Az alapötlet, hogy minden közepet osszuk szét két részre, majd futtasuk a  $k - means$  algoritmust  $k = 2$ -re lokálisan. Ezután meg kell vizsgálni, melyik kettéosztás vezet a helyes eredményre. Ezt a vizsgálatot egy úgynevezett BIC (Bayesian Information Criteria) érték alapján végzik. Kiszámolják az egy középhez és a kettéosztott középhez tartozó BIC értéket, majd azt a modellt választják, amelynél ez az érték nagyobb.