

# Învățare Automată

Moișanu-Costinescu Ștefan

Decembrie 2020

## 1 Exercițiul 46, pagina 819: k-means++

### 1.1 Comentariu:

În exercițiul de față vom considera că la faza de inițializare algoritmul Kmeans alege centrozii în mod arbitrar dintre instanțele de clusterizat.

Dacă datele de clusterizat sunt [a priori] bine separate în K clustere, atunci este foarte posibil ca la finalul inițializării să existe [măcar] un cluster din care algoritmul nu a selectat niciun punct. În astfel de situații, algoritmul K-means nu va produce clusterelor dorite de noi.

În schimb, varianta K-means++ a algoritmului K-means, propusă de David Arthur și Sergei Vassilvitskii în 2007, încearcă să selecteze pentru pozițiile inițiale ale celor K centrozi instanțe care sunt [pe cât se poate] mai distanțate unele de altele. În acest fel, pot fi selectate, cu o probabilitate [destul de] mare, instanțe din toate clusterelor.

### 1.2 Formalizare:

K-means++ face inițializarea centrozilor în maniera următoare:

1. Alege primul centroid,  $\mu_1$ , în manieră uniform aleatorie dintre instanțele de clusterizat,  $x_1, \dots, x_n$ . Cu alte cuvinte, alegem mai întâi un indice  $i$  în mod uniform aleatoriu din mulțimea  $\{1, \dots, n\}$  și fixăm  $\mu_1 = x_i$ .

2. Pentru  $j = 2, \dots, K$ :

- Pentru fiecare instanță  $x_i$ , calculează distanța  $D_i$  până la cel mai apropiat centroid ales / fixat la o iterație anterioară:

$$D_i = \min_{j'=1, \dots, j-1} \|x_i - \mu_{j'}\|$$

- Alege centroidul  $\mu_j$  în mod aleatoriu dintre instanțele  $x_1, \dots, x_n$ , cu probabilitate proporțională cu  $D_1^2, \dots, D_n^2$ . Altfel spus, alegem un indice  $i$  în mod aleatoriu din mulțimea  $\{1, \dots, n\}$  cu probabilități egale cu  $\frac{D_1^2}{\sum_{i'=1}^n D_{i'}^2}, \dots, \frac{D_n^2}{\sum_{i'=1}^n D_{i'}^2}$ , și fixăm  $\mu_j = x_i$

3. Returnează  $\mu = (\mu_1, \dots, \mu_K)$ , setul de asignări ale pozițiilor inițiale ale centrozilor clusterelor pentru algoritmul lui Lloyd (K-means).

Vom ilustra acum diferența dintre K-means++ și K-means [la inițializare], folosind un set de date simplu, format din cinci puncte în planul euclidian bidimensional: A(-1, 0), B(1, 0), C(0, 1), D(3, 0) și E(3, 1):

- a Presupunem că aplicăm algoritmul K-means cu  $K = 2$ , făcând inițializarea centrozilor cu instanțe din mulțimea pe care tocmai am precizat-o. Presupunem că a fost deja ales centroidul  $\mu_1 = A$ . Dacă selecția următorului centroid ( $\mu_2$ ) se face în manieră uniform aleatorie — așa cum procedează îndeobște algoritmul K-means — din mulțimea {B, C, D, E}, care este probabilitatea ca  $\mu_2$  să fie din submulțimea {B, C}? Justificați.
- b Aplicăm acum pe același set de date algoritmul K-means++, tot cu  $K = 2$ . Presupunem, ca și la punctul a, că a fost selectat centroidul  $\mu_1 = A$ . Care este acum probabilitatea ca  $\mu_2$  să fie selectat din submulțimea {B, C}? Rezultă oare într-adevăr o îmbunătățire semnificativă față de [inițializarea făcută de] algoritmul K-means? Justificați riguros.

## 2 Rezolvare

### 2.1 Subpunctul a.

Dupa ce algoritmul alege  $\mu_1 = A$ , ramane sa aleaga aleator uniform  $\mu_2$  din multimea  $\{B, C, D, E\}$ .  
 $\Rightarrow$  probabilitatea ca  $\mu_2 \in \{B, C\}$  este de  $\frac{2}{4} = 50\%$ .

### 2.2 Subpunctul b.

Pentru a afla  $P(\mu_2 \in \{B, C\})$  folosind k-means++ voi afla distributia de probabilitate constituita de acest algoritm la fiecare iteratie.

$$\mu_2 = A \text{ (aleator)}$$

Calculez  $D_A, D_B, D_C, D_D, D_E$ :

$D_A = D(A, \mu_1) = 0$  ( $\mu_1$  si  $\mu_2$  nu pot coincide, dar deoarece  $\text{dist} = 0$  oricum va avea probabilitate 0 de a fi ales ia A)

$$D_B = 2$$

$$D_C = \sqrt{2}$$

$$D_D = 4$$

$$D_E = \sqrt{17}$$

Calculez suma patratelor distantelor:

$$\sum_{i=1}^n D_i^2 = 0 + 4 + 2 + 16 + 17 = 39$$

Obtin urmatorul tabel al distributiei de probabilitate cu valorile asociate:

i	$\frac{D_i^2}{\sum_{i'=1}^n D_{i'}^2}$
A	0
B	$\frac{4}{39}$
C	$\frac{2}{39}$
D	$\frac{16}{39}$
E	$\frac{17}{39}$

$\Rightarrow \mu_2 = E$  va avea cea mai mare probabilitate de a fi ales de k-means++.

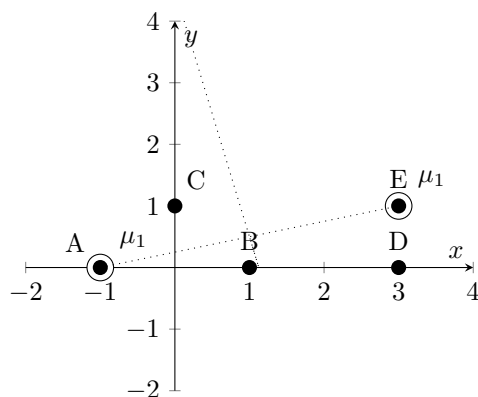
$$\Rightarrow P(\mu_2 \in \{B, C\}) = \frac{4}{39} + \frac{2}{39} = \frac{6}{39} = 0.1538$$

Pentru a vedea daca am obtinut o imbunatatire voi simula in continuare k-means si voi compara numarul de iteratii pana la convergenta pentru varianta k-means++ si pentru k-means.

## 1. K-means++

- $\mu_2 = E$

I Analiza geometrica:



$$Cluster(\mu_1) = \{A, B, C\}$$

$$Cluster(\mu_2) = \{E, D\}$$

$$x_{\mu_1} = \frac{x_A + x_B + x_C}{3} = \frac{-1 + 0 + 1}{3} = 0$$

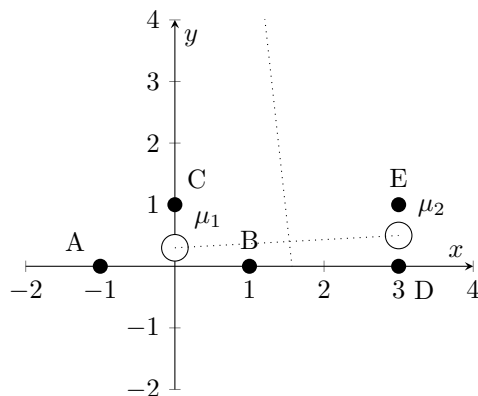
$$y_{\mu_1} = \frac{y_A + y_B + y_C}{3} = \frac{0 + 0 + 1}{3} = \frac{1}{3}$$

$$x_{\mu_2} = \frac{x_E + x_D}{2} = \frac{3 + 3}{2} = 3$$

$$y_{\mu_2} = \frac{y_E + y_D}{2} = \frac{1 + 0}{2} = \frac{1}{2}$$

$$\Rightarrow \mu_1(0, \frac{1}{3}), \mu_2(3, \frac{1}{2})$$

II Analiza geometrica:



$$Cluster(\mu_1) = \{A, B, C\}$$

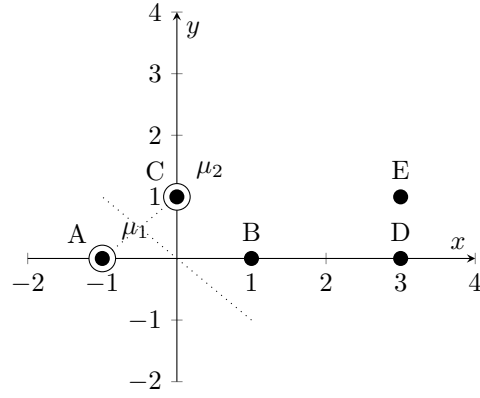
$$Cluster(\mu_2) = \{E, D\}$$

Clusterele nu s-au schimbat  $\Rightarrow$  centroizii raman la fel  $\Rightarrow$  algoritmul a conver in 2 iteratii.

## 2. K-means

- $\mu_2 = C$

I Analiza geometrica:



$$Cluster(\mu_1) = \{A\}$$

$$Cluster(\mu_2) = \{B, C, E, D\}$$

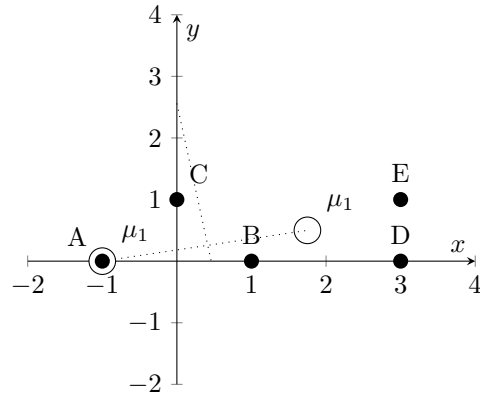
$$\mu_1 = A \Rightarrow \mu_1(-1, 0)$$

$$x_{\mu_2} = \frac{0 + 1 + 3 + 3}{4} = \frac{7}{4}$$

$$y_{\mu_2} = \frac{1 + 0 + 0 + 1}{4} = \frac{1}{2}$$

$$\Rightarrow \mu_1(-1, 0), \mu_2(\frac{7}{4}, \frac{1}{2})$$

II Analiza geometrica:



$$Cluster(\mu_1) = \{A, C\}$$

$$Cluster(\mu_2) = \{B, E, D\}$$

$$x_{\mu_1} = \frac{-1 + 0}{2} = -\frac{1}{2}$$

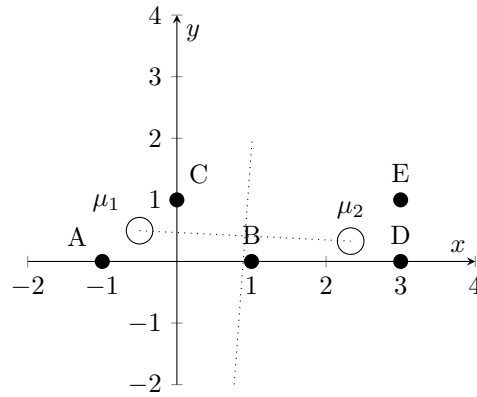
$$y_{\mu_1} = \frac{0 + 1}{2} = \frac{1}{2}$$

$$x_{\mu_2} = \frac{1 + 3 + 3}{3} = \frac{7}{3}$$

$$y_{\mu_2} = \frac{0 + 0 + 1}{3} = \frac{1}{3}$$

$$\Rightarrow \mu_1(-\frac{1}{2}, \frac{1}{2}), \mu_2(\frac{7}{3}, \frac{1}{3})$$

III Analiza geometrica:



$$Cluster(\mu_1) = \{A, C\}$$

$$Cluster(\mu_2) = \{B, E, D\}$$

Clusterelor nu s-au modificat  $\Rightarrow$  Algoritmul a convergat in 3 iteratii

- $\mu_2 = B$

C este la egala distanta de A si B  $\Rightarrow$  aleg aleator, dar consistent.

- $C \in Cluster(\mu_1) \Rightarrow$  convergenta in 2 iteratii (identic incepand cu itearatia 2 a cazului  $\mu_2 = C$ )
- $C \in Cluster(\mu_2) \Rightarrow$  convergenta in 3 iteratii (dupa un rationament analog cu cel dezvoltat mai sus)

### Concluzie:

Folosind k-means++ pentru initializarea centroizilor am favorizat convergenta algoritmului in 2 iteratii, comparativ cu posibilitatea de a se incheia dupa 3 iteratii in cazul  $\mu_2 = C$  sau  $\mu_2 = B$  (cu  $C \in Cluster(\mu_2)$ ).

$$(1) P_{k-means}(\mu_2 \in \{B, C\}) = 0.5$$

$$(2) P_{k-means++}(\mu_2 \in \{B, C\}) = 0.1535$$

(1) + (2)  $\Rightarrow$  Am favorizat alegerea ce duce la convergenta in 2 iteratii

Mai mult, pentru seturi de date mai mari, diferenta intre numarul de iteratii pana la convergenta cu si fara initializare k-means++ poate fi considerabila  $\Rightarrow$  merita alocarea timpului necesar pentru calculele facute de k-means++, chiar daca sunt destul de complexe.

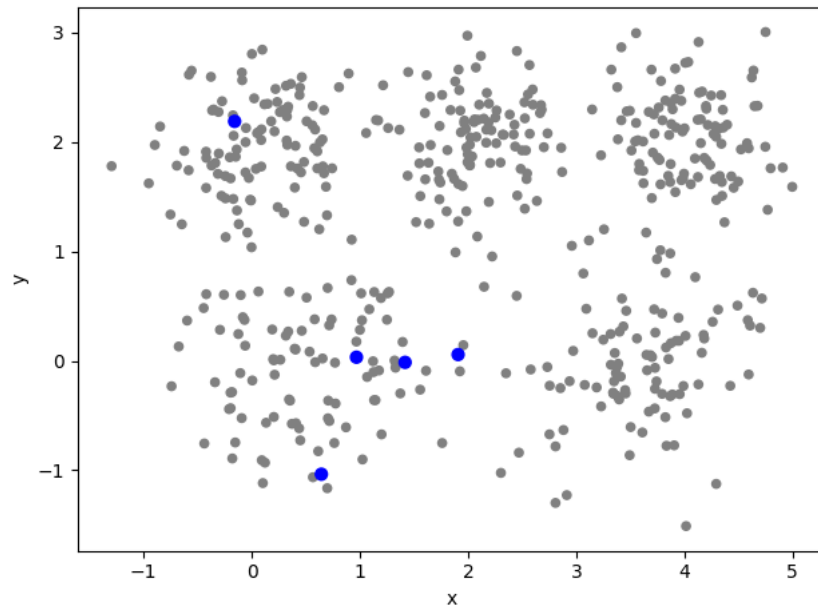
## 3 Studiu

Am rulat de 200 de ori algoritmul k-means ( $k = 5$ ) pe un set de date cu 500 de intrari si am reprezentat grafic in plan punctele din datele de intrare (gri), centroizii alesi la initializarea algoritmului in cadrul primei iteratii a unei rulari a experimentului (albastru) si toate centrele cluster-urilor obtinute ca rezultat in urma fiecarei iteratii (rosu).

Prima data am initializat centroizii in maniera aleatorie uniforma pentru fiecare iteratie si am obtinut urmatoarele rezultate:

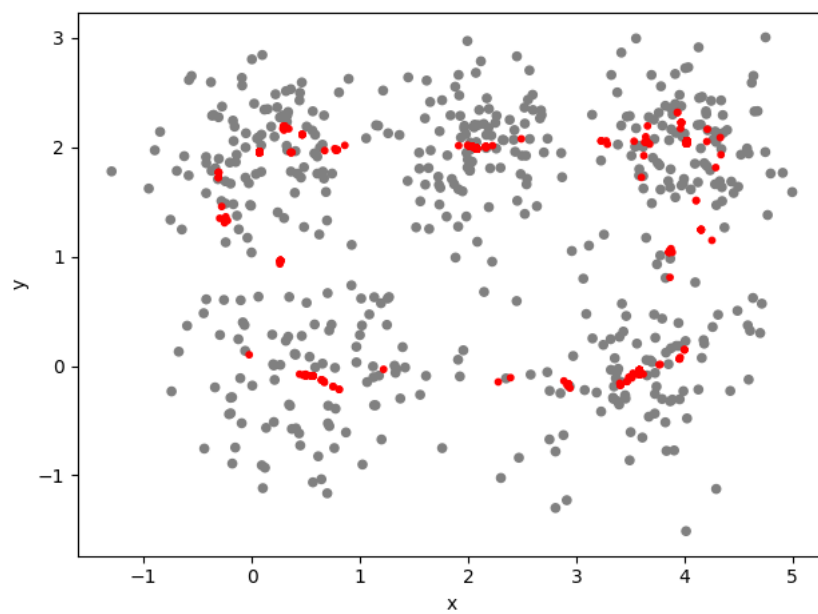
- Centroizii alesi la prima initializare a algoritmului:

Figure 1



- Cele 5 x 200 clusterelor obtinute la finalul experimentului:

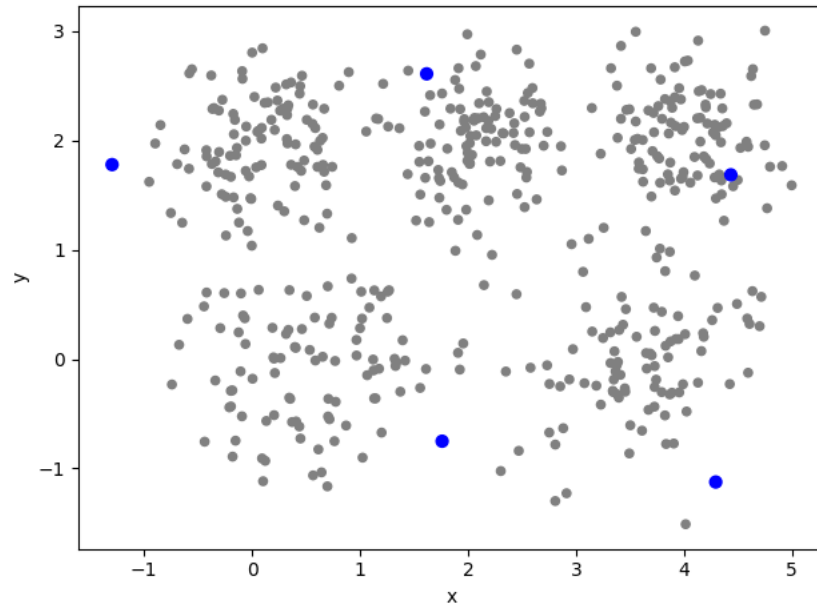
Figure 2



Pentru urmatoarea rulare a experimentului am initializat algoritmul folosind distributia de probabilitate furnizata de k-means++ si am obtinut urmatoarele rezultate:

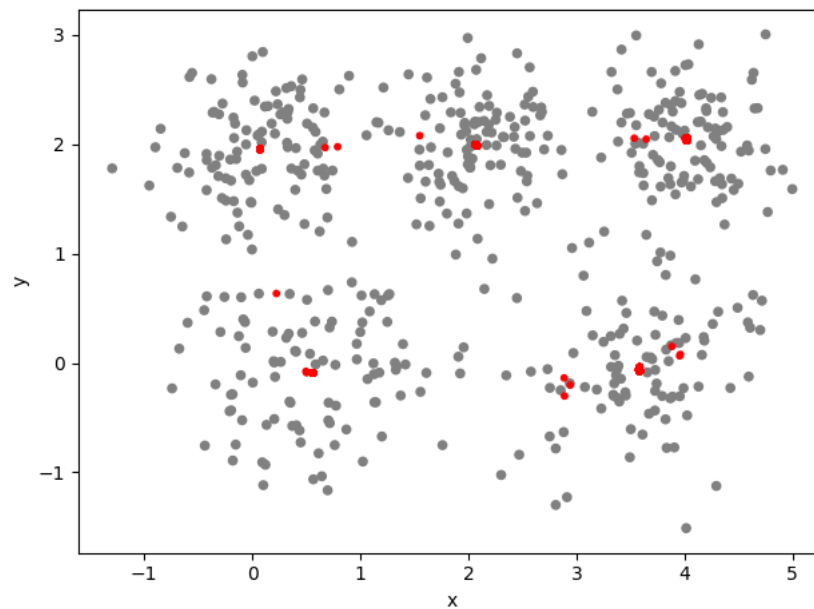
- Centrozii alesi la prima initializare a algoritmului:

Figure 3



- Cele 5 x 200 clusterelor obtinute la finalul experimentului:

Figure 4



De asemenea, am afisat la finalul fiecarei rulari o statistica ce include: minimul obiectivelor obtinute (min\_obj), media acestora (avg\_objective), deviatia standard si numarul mediu de iteratii a algoritmului pana a convera:

Figure 5: Statistics

(a) K-means

```
min_obj: 222.36598776929088
avg_objective: 252.2581628589124
standard deviation: 68.8992271769586
avg_iterations_till_convergence: 9.58
run_time: 65.3121440410614
```

(b) K-means++

```
min_obj: 222.36598776929088
avg_objective: 225.47792603039224
standard deviation: 77.03268890965262
avg_iterations_till_convergence: 6.02
run_time: 70.8408350944519
```

### Informatii extrase:

Prin compararea figurii 1 cu 3 se poate observa ca centroizii alesi cu k-means++ sunt mult mai bine distribuiti in setul de date, cate unul in fiecare din cele 5 grupari in care sunt dispuse punctele. Astfel, sunt din prima plasati in vecinatatea centrelor clusterelor ce vor fi obtinute de algoritmul ceea ce se reflecta atat in numarul scazut de iteratii necesare pana la convergenta ( $6.02 < 9.58$ ), cat si in scaderea mediei obiectivului:  $225.47 < 252.25$ .