

Slovenská Technická Univerzita v Bratislave
Fakulta informatiky a informačných technológií

Umelá inteligencia

Zadanie 4b – Zhlukovanie

Riešený problém

Zadanie sa týka riešenia problému zhukovania bodov v 2D priestore. Cieľom je vytvoriť zhukovače, ktoré efektívne rozdelia body v 2D priestore do k zhukov.

Veľkosť 2D priestoru predstavujú X a Y súradnice v rozmedzí od -5000 do $+5000$, pričom do priestoru sa umiestni 20020 vygenerovaných bodov. Súradnice prvých 20 bodov sú vygenerované náhodne. Ďalších 20000 bodov sa potom generuje v blízkosti už vygenerovaných bodov, z ktorých je vždy vybratý iba jeden, do ktorého blízkosti sa nový bod vygeneruje.

Po vygenerovaní všetkých 20020 bodov sa spustí prehľadávanie 2D priestoru príslušným zhukovačom, kedy sa body v priestore rozdelia na požadovaný počet zhukov. Za úspešný zhukovač sa považuje taký, ktorého priradené body nemajú vzdialenosť od stredu zhuku väčšiu ako 500, pričom stred môže byť reprezentovaný ťažiskom zhuku (**centroid**), alebo bodom, ktorý sa nachádza najbližšie ku ťažisku zhuku (**medoid**).

Opis riešenia

Úlohou je naprogramovať zhukovač, ktorý prehľadá všetky vygenerované body v 2D priestore a následne rozdelí tento priestor na k zhukov, pričom je potrebné implementovať verzie zhukovača s použitím nasledujúcich algoritmov:

- **k-means** zhukovanie, kde stred je **centroid**
- **k-means** zhukovanie, kde stred je **medoid**
- **aglomeratívne** zhukovanie, kde stred je **centroid**
- **divízne** zhukovanie, kde stred je **centroid**

Priebeh algoritmov

Na začiatku sa vygeneruje do 2D priestoru prvých 20 náhodných bodov. Následne sa vygeneruje ďalších 20000 a to tak, že vždy sa z už doteraz vygenerovaných bodov vyberie náhodne jeden, k jeho X a Y súradnici sa vygeneruje x -ový a y -ový súradnicový offset z intervalu -100 až $+100$, k pôvodným súradniciam sa priráta príslušný offset a výsledkom sú súradnice nového bodu. Pokiaľ vzniknuté súradnice nového bodu presahujú niektorý okraj 2D priestoru, tak sa offset zredukuje, aby sa novovzniknutý bod vmestil do rozmerov 2D priestoru. Po vygenerovaní všetkých 20020 bodov sa spustí príslušný zhukovač.

K-means zhukovanie (centroid/medoid)

Z vygenerovaných 20020 bodov sa náhodne vyberie k bodov, ktoré budú predstavovať jednotlivé zhuky, pričom tento prvý bod reprezentuje začiatkový centroid a medoid zhuku.

Následne sa v cykle prehľadávajú všetky ostatné body v 2D priestore, ktoré ešte nie sú priradené do žiadneho zhuku. Pre každý bod sa zo všetkých zhukov vyberá najviac vyhovujúci zhuk, do ktorého bude prehľadávaný bod následne priradený.

Najviac vyhovujúcim zhukom je taký, ktorého vzdialenosť od jeho centra (centroid/medoid) k prehľadávanému bodu je najmenšia. Po priradení príslušného bodu do najviac vyhovujúceho zhuku sa vyberie ďalší bod a cyklus sa opakuje.

Po prehľadaní všetkých bodov sa pre každý zhuk vypočíta jeho nové centrum (centroid/medoid), všetky jeho priradené body sa odstránia a znovu sa spustí celý cyklus prehľadávania a priradzovania bodov. To sa opakuje až do vtedy, pokiaľ sa novovzniknuté centrum nelíši od starého, alebo sa líši len minimálne. Vtedy cyklus končí a jednotlivé zhuky sú finálne.

Aglomeratívne zhukovanie (centroid)

Každý z vygenerovaných 20020 bodov na začiatku predstavuje samostatný zhuk, pričom tento prvý bod reprezentuje začiatkový centroid a medoid zhuku.

Následne sa v cykle hľadajú také 2 zhuky, ktoré sú od seba najmenej vzdialené. Po ich nájdení sa zlúčia do jedného spoločného zhuku a to tak, že body z jedného zhuku sa priradia druhému zhuku, vypočíta sa nové spoločného centrum zhuku a prázdny zhuk sa odstráni. Následne sa cyklus zlučovania zhukov opakuje, až pokiaľ neostane iba požadovaných k zhukov. Vtedy cyklus končí a jednotlivé zhuky sú finálne.

Divízne zhukovanie (centroid)

Na začiatku sa body v 2D priestore rozdelia do dvoch zhukov pomocou K-means algoritmu. Následne sa v cykle vždy vyberá z už vzniknutých zhukov ten, ktorý má najhoršie ohodnotenie – najhoršia priemerná vzdialenosť od centra (centroid). Ten sa potom rozdelí na ďalšie dva zhuky pomocou K-means algoritmu. Cyklus sa opakuje do vtedy, pokiaľ nevzniklo požadovaných k zhukov.

Reprezentácia údajov

Reprezentácia priestoru – trieda **Space**

```
public class Space {  
    public ArrayList<Point> points;  
    public boolean[][] spaceMap;  
    public int spaceMapOffset;  
    public Random random;
```

Reprezentácia bodu – trieda **Point**

```
public class Point {  
    public int posX;  
    public int posY;
```

Reprezentácia zhľuku – trieda **Cluster**

```
public class Cluster {  
    public Point centroid;  
    public Point medoid;  
    public ArrayList<Point> points;  
    public int sumX;  
    public int sumY;
```

Reprezentácia kMeans zhľukovania – trieda **KMeans**

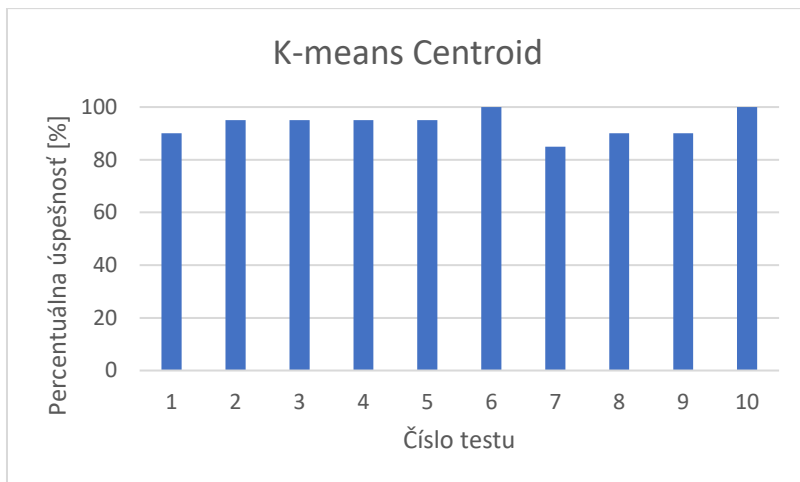
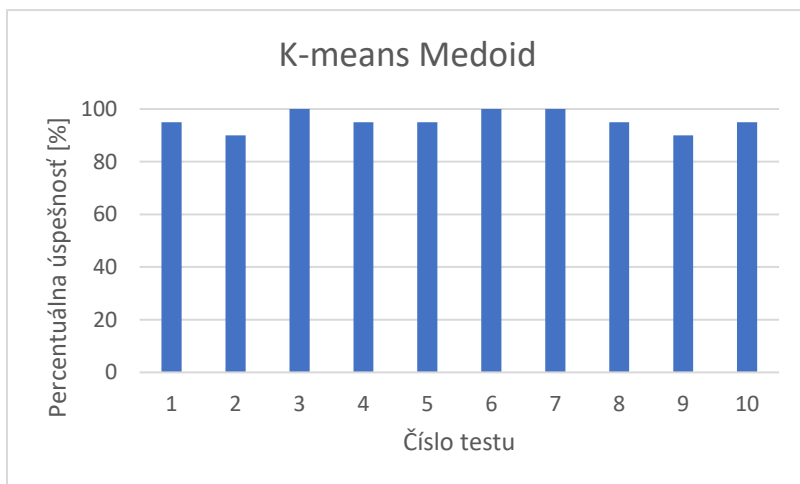
```
public class KMeans {  
    public int k;  
    public ArrayList<Point> points;  
    public ArrayList<Cluster> clusters;  
    public Random random = new Random();  
    public final double maxDistance = Integer.MAX_VALUE;
```

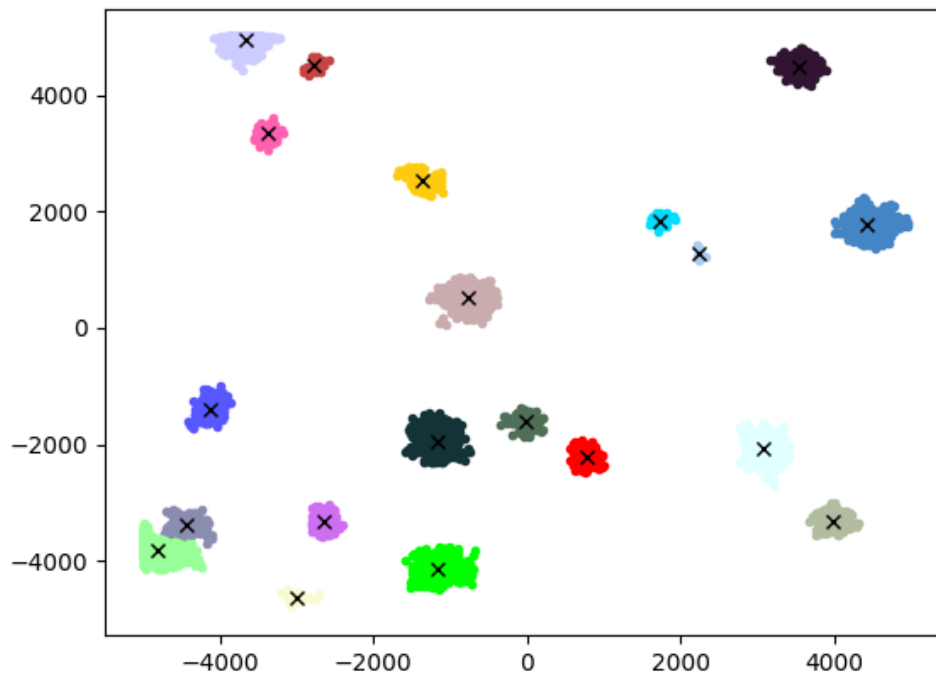
Reprezentácia aglomeratívneho zhľukovania – trieda **Agglomeration**

```
public class Agglomeration {  
    public int k;  
    public ArrayList<Point> points;  
    public ArrayList<Cluster> clusters;  
    public double[][] distanceMap;  
    public final double maxDistance = Integer.MAX_VALUE;  
    public Map<Integer, Cluster> dictionary;  
    public ArrayList<Integer> list;
```

Reprezentácia divízneho zhľukovania – trieda **Division**

```
public class Division {  
    public int k;  
    public ArrayList<Point> points;  
    public ArrayList<Cluster> clusters;
```

Testovanie (pre rovnaké rozmiestnenie – 20020 bodov a 20 klastrov)**K-means Centroid** (Priemerný čas **41ms**, akceptovanie **2/10**)**K-means Medoid** (Priemerný čas **34ms**, akceptovanie **3/10**)**Divízne zhlukovanie** (Priemerný čas **14ms**, akceptovanie **8/10**)

Aglomeratívne zhlukovanie (iba jeden beh, čas **145min**, akceptovanie **1/1**)**Všeobecné informácie**

Program:

- Program bol vyhotovený v jazyku Java (SDK 16)
- Použité bolo IDE IntelliJ IDEA Ultimate 2021.2

Prerekvizity:

- **Python** vo verzii 3+
- Python knižnica **matplotlib**
- Java knižnica **matplotlib4j** (<https://github.com/sh0nk/matplotlib4j>)

Používateľské rozhranie:

- CLI navigácia

```
(0) Exit
(1) K-means - centroid
(2) K-means - medoid
(3) Agglomerative - centroid
(4) Divisive - centroid
Enter your choice: |
```