

Slovenská technická univerzita
Fakulta informatiky a informačných technológií
Ilkovičova 3, 842 16 Bratislava 4

Umelá Inteligencia

Zadanie 4

Klastrovanie

Matsveyeva Lada-Ivanna

ZS 2021/22

Obsah:

Zadanie.....	3
Opis algoritmov.....	4
Reprezentácia údajov	4
Testovanie a porovnanie výsledkov	5
Príklady vizualizácie	6
Zhrnutie a zhodnotenie riešenia	9

Zadanie

Máme 2D priestor, ktorý má rozmery X a Y, v intervaloch od -5000 do +5000. Tento 2D priestor vyplňte 20 bodmi, pričom každý bod má náhodne zvolenú polohu pomocou súradníc X a Y. Každý bod má unikátne súradnice (t.j. nemalo by byť viacero bodov na presne tom istom mieste).

Po vygenerovaní 20 náhodných bodov vygenerujte ďalších 20000 bodov, avšak tieto body nebudú generované úplne náhodne, ale nasledovným spôsobom:

1. Náhodne vyberte jeden zo všetkých doteraz vytvorených bodov v 2D priestore.
Ak je bod príliš blízko okraju, tak zredukujete príslušný interval v nasledujúcich dvoch krokoch.
2. Vygenerujte náhodné číslo X_{offset} v intervale od -100 do +100
3. Vygenerujte náhodné číslo Y_{offset} v intervale od -100 do +100
4. Pridajte nový bod do 2D priestoru, ktorý bude mať súradnice ako náhodne vybraný bod v kroku 1, pričom tieto súradnice budú posunuté o X_{offset} a Y_{offset}

Vašou úlohou je naprogramovať zhukovač pre 2D priestor, ktorý zanalyzuje 2D priestor so všetkými jeho bodmi a rozdelí tento priestor na k zhukov (klastrov). Implementujte rôzne verzie zhukovača, konkrétne týmito algoritmami:

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

Opis algoritmov:

Generujem náhodne body, podľa algoritmu, ktorý je uvedený v zadaní.

k-means

Najprv vyberiem k náhodných bodov, ktoré budú stredmi zhlukov. Ďalej pre každý bod určím najbližší zhluk. Po prvotnom ohodnotení vyrátam stred ešte raz.

Stred centroid – priemerná vzdialenosť všetkých x,y súradníc bodov.

Stred medoid – reálny bod, ktorý má najmenšiu euklidovskú vzdialenosť ku všetkým bodom.

V cykle určujem najbližší zhluk a aktualizujem stred pokiaľ sú zmeny v zhluchoch.

Aglomeratívne zhľukovanie

Každý vygenerovaný bod je zhluk. Vytvorím maticu susednosti, nájdem zhľuky, ktoré sú najbližšie ku sebe a spojím ich do jedného zhľuku. Aktualizujem maticu susednosti pre nový zhluk, ktorý vložím na koniec matice a zase nájdem najbližšie ku sebe zhľuky. Tak sa cyklím kým nedosiahnem potrebný počet zhľukov.

Divízne zhľukovanie

Divízne zhľukovanie je opačne ako aglomeratívne. Všetky vygenerované body sú 1 zhluk. Tento zhluk pomocou k-means so stredom centroid rozdelím na 2 zhľuky. Pôvodne centroidy, oproti tomu ako bolo v k-means, vyberiem nie je náhodne, ale podľa najväčšej vzdialenosti dvoch bodov. Potom kým nebude dosiahnutý požadovaný počet zhľukov v cykle vyberiem zhluk ktorý v sebe má najväčší priemer vzdialenosti a rozdelím na 2 zhľuky.

Reprezentácia údajov:

Vygenerované body mám uložené v dátovom type list. Jedna položka – jeden bod a zhľňa v sebe x -ovú, y -ovú súradnicu a zhluk ku ktorému patrí. Centroidy/medoidy mám v samotných list-och.

Pre maticu susednosti používam maticu z knižnice numpy, je lepšia pre reprezentáciu matice, ako vnorené list-y a taktiež má v sebe funkcie na nájdenie min, max prvkov.

Na výpočet vzdialenosti medzi dvoma bodmi/zhľukmi používam euklidovskú vzdialenosť z knižnice math, funkcia `dist()`.

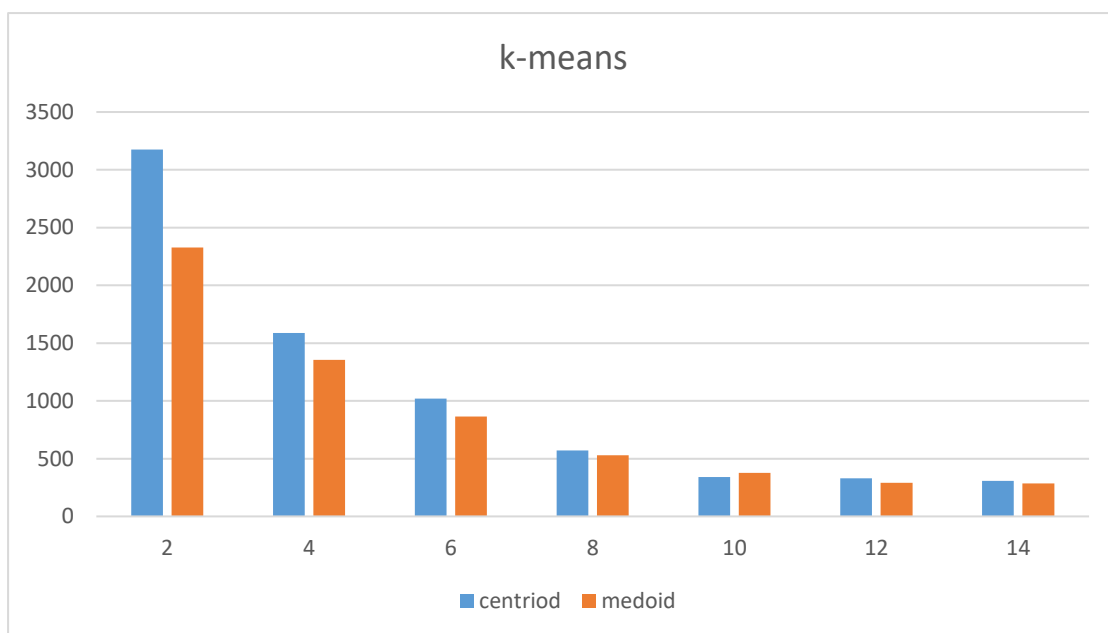
Pre vizualizáciu používam knižnicu matplotlib. Rôzne zhľuky reprezentujem rôznymi farbami.

Testovanie a porovnanie výsledkov

Na zhodnotenie výsledkov algoritmov porovnávam vzdialenosť bodov od stredu zhlukou pre rôzny počet zhlukou.

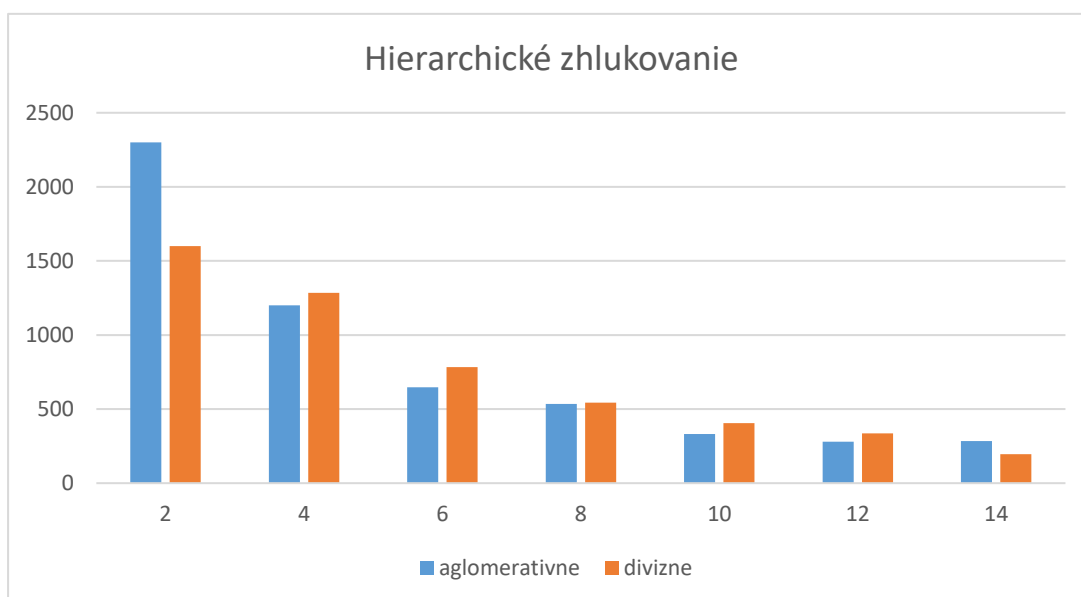
Podľa výsledkov je vidieť, že keď zvyšujem počet zhlukou, tak sa znižuje priemerná vzdialenosť v zhluchoch. Pre každý algoritmus po 8 zhlukov priemerná vzdialenosť sa neklesá. Podľa tohto je optimálny počet zhlukov 8 až 10.

Medzi k-means časovo rýchlejší je so stredom ako centroid, lebo len spočíta všetky súradnice a nájde z nich priemer. Vtedy keď pre medoid je potrebné prejsť každý bod, a nájsť pre neho vzdialenosť ku každému bodu v zhluke. Ale podľa efektívnosti je lepší algoritmus so stredom ako medoid, lebo je to už reálny bod v zhluke.



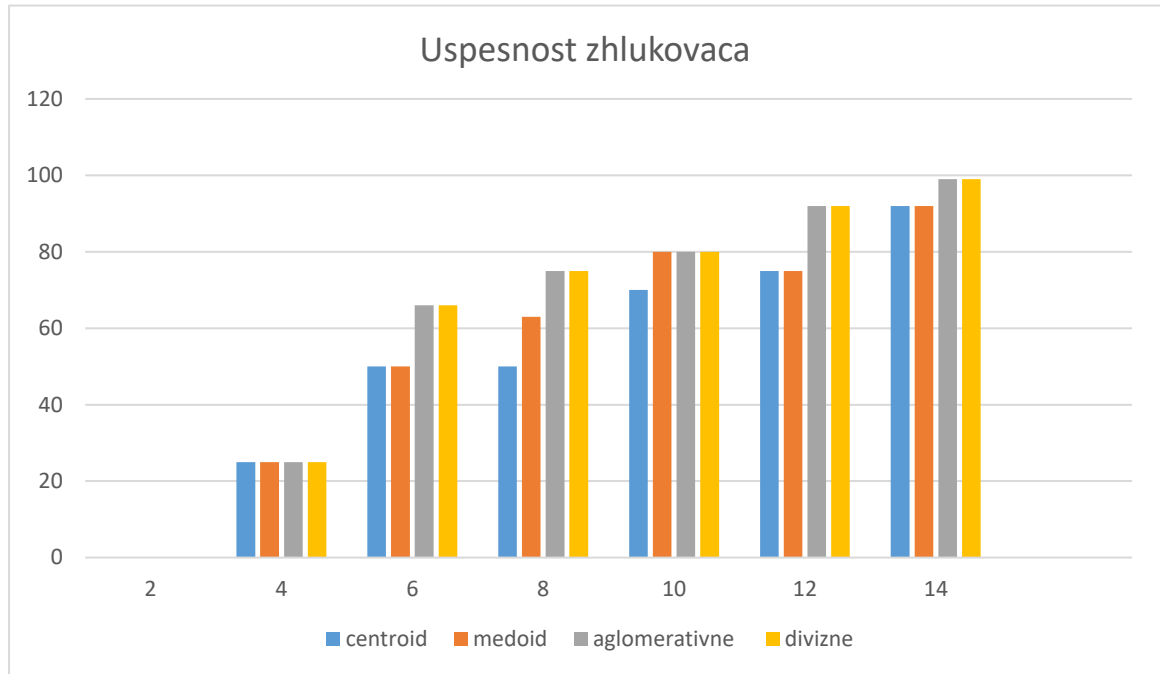
Medzi aglomeratívnym a divíznym je rýchlejší divízny ale majú približne rovnaké výsledky primernej vzdialenosti bodov od stredu. Aglomeratívne kvôli tomu že začína z toho, že každý bod je zhluk, tak matica susednosti pre počet bodov 20 000, je veľkosťou 20000*20000, čo je celkovo pamäťové náročné, a kvôli tomu je aj pomalšia.

V aglomeratívnom zhlukovaní môže nastať situácia, keď v jednom zhluke bude len pár bodov, vtedy keď pri divíznom toto nemôže nastať (len ak počet zhlukov sa takmer rovná počtu bodov).



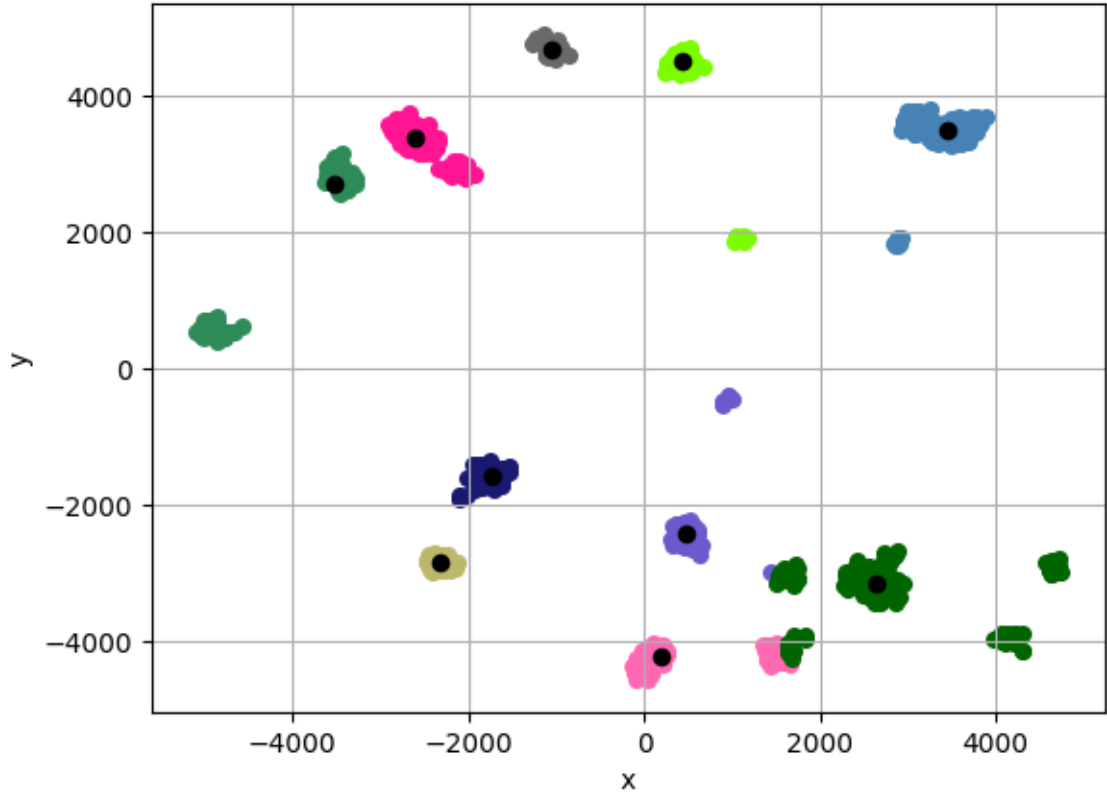
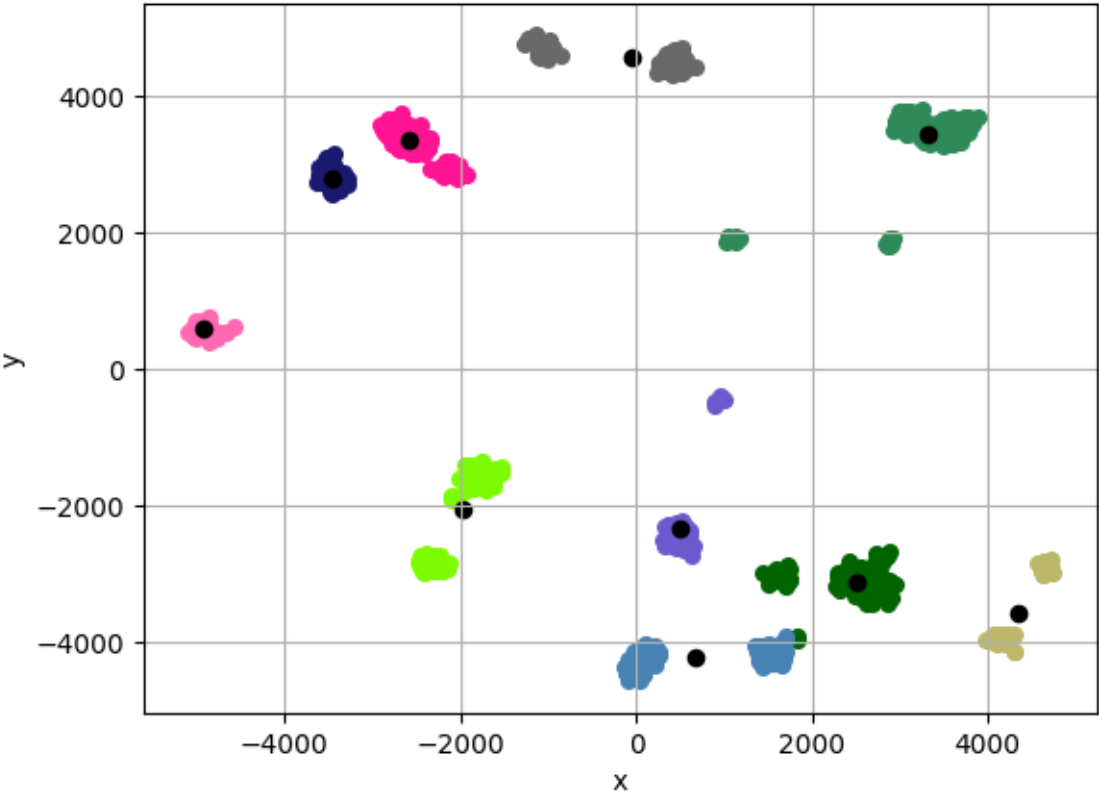
Porovnanie úspešnosti zhukovačov.

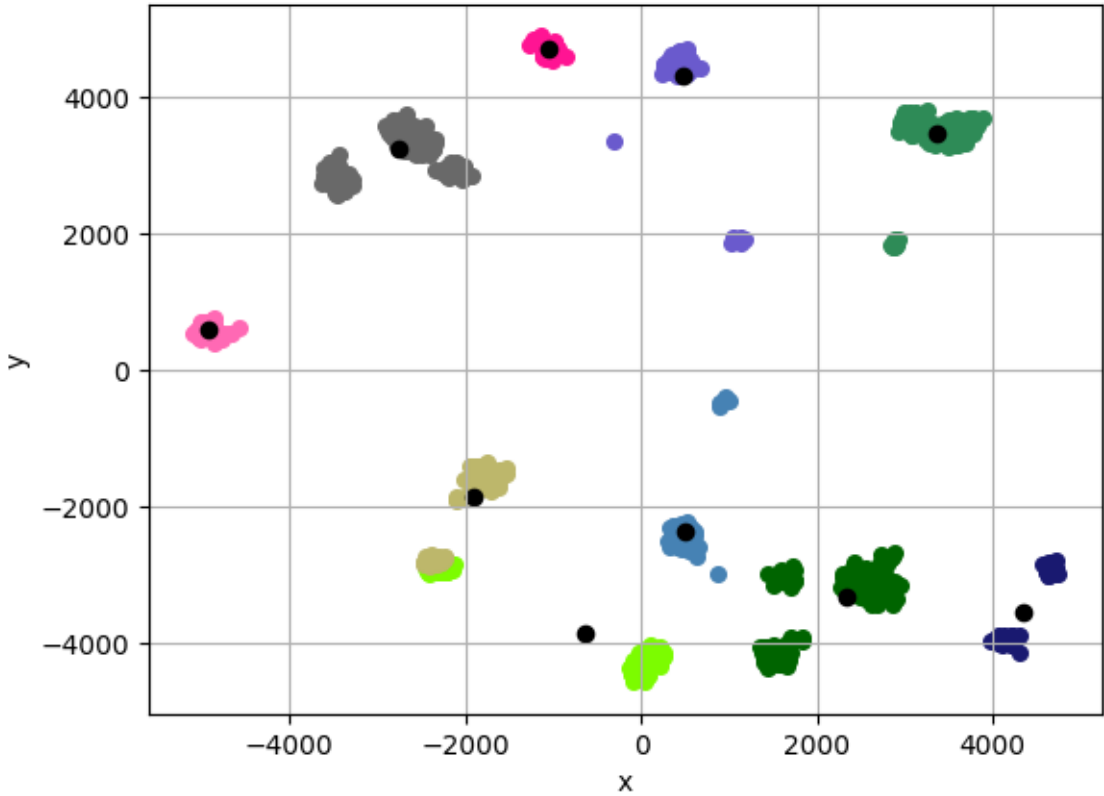
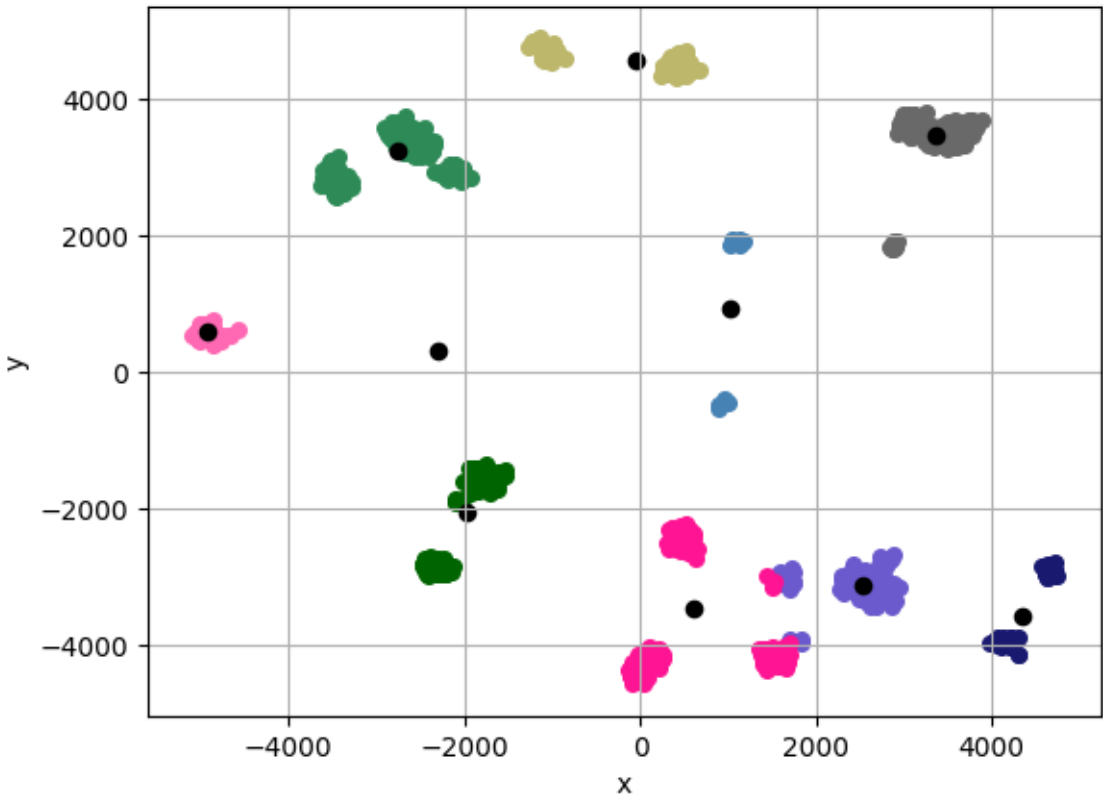
Podľa zadania úspešný zhukovac vtedy keď každý zhuk má priemernú vzdialenosť bodov od stredu nie je viac ako 500. Porovnala som všetky algoritmy spolu, pre rôzny počet zhukov. Pri menšom počte zhukov od 2 po 4 úspešnosť je menšia ako 50%, je to spôsobené tým, že zhukov je malo, a počet bodov a veľkosť plochy sú veľké. Pre počet zhukov 4 a viac úspešnosť je viac ako 50%. Úspešnosť v hierarchyckom zhukovaní je väčšia ako pre k-means.



Príklady vizualizácie:

1. k-means centroid
2. k-means medoid
3. aglomeratívne zhukovanie
4. divízne zhukovanie





Zhrnutie a zhodnotenie riešenia

Pri k-means algoritmu veľkú rolu hrá to, ako sa vyberú stredné body zhlukov, lebo sa vyberajú náhonne. Je to časovo a pamäťové výhodnejšie, ako napríklad pri aglomeratívnom a divíznom algoritmoch.

Pri hierarchických algoritmoch rozdelenie do zhlukov je presnejšie, ale je časovo a pamäťové náročnejšie. Pri týchto algoritmoch je veľmi podstatné vedieť ich optimalizovať.