

Z2 c dokumentacia

Nikolas Knapik

November 2024

1 Úvod

Táto dokumentácia analyzuje implementáciu zhľukovania v 2D priestore. Hlavným cieľom je analyzovať priestor a rozdeliť ho na k zhľukov pomocou dvoch variánt aglomeratívneho zhľukovania:

- Použitie centroidu ako stredového bodu
- Použitie medoidu ako stredového bodu

Úspešnosť klastrovania sa hodnotí na základe priemernej vzdialenosti bodov od stredu zhľuku. Úspešný zhľuk nesmie mať priemernú vzdialenosť väčšiu ako 500

2 Popis Algoritmu

2.1 Generovanie bodov

Pre prvých 20 bodov sa náhodne generujú X a Y súradnice v intervale $[-5000, 5000]$. Pomocou `set()` sa zabezpečia unikátne súradnice v prípade duplicit.

Ďalších 20 000 bodov sa generuje od náhodne vybraného už existujúceho bodu od ktorého sa generujú nové súradnice s offsetom 100

2.2 Algomeratívne zhľukovanie

1. Inicializácia:

Každý z 20 020 vygenerovaných bodov je na začiatku zaevídovaný ako samostatný zhľuk. Vzďialenosti medzi všetkými dvojícami zhľukov sú uložené v 2D matici vzdialeností, ktorej diagonála je 0, ktorá predstavuje vzdialenosť zhľuku od samého seba

2. Iterácia zhľukovania:

- Nájde sa najbližšia dvojica zhľukov (pomocou minimálnej hodnoty v matici vzdialeností)
- Zhľuky sa zlúčia a aktualizuje sa ich stred:
 - **Centroid:** Priemerná súradnica bodov v novom zhľuku.
 - **Medoid:** Bod v zhľuku s najmenšou celkovou vzdialenosťou k ostatným bodom.
- Matica vzdialeností sa aktualizuje len pre nový zhľuk.

3. Ukončenie:

- Algoritmus sa zastaví, keď priemerná vzdialenosť bodov od stredu pre všetky zhľuky klesne pod 500

2.3 Vyhodnocovanie úspešnosti

Úspešnosť klastrovania sa hodnotí ako percento zhľukov, ktoré spĺňajú kritérium priemernej vzdialenosti bodov od stredu menšej ako 500.

3 Reprezentácia údajov

- **Body:** Uložené ako zoznam dvojíc (X,Y).
- **Matica vzdialeností:** 2D numpy matica s hodnotami vzdialeností medzi všetkými dvojícami bodov.
- **Zhľuky:** Slovník, kde kľúč je ID zhľuku a hodnota je zoznam bodov.

- **Stredy zhlučkov:**

- Centroid: Priemerná súradnica bodov v zhlučku.
- Medoid: Bod s najmenšou celkovou vzdialenosťou k ostatným bodom.

4 Optimalizácie

Počas optimalizácie som si zvýšil offset na 300 aby som lepšie videl vytvorené kódy. Funkčnosť kódu som vo väčšine prípadoch testoval na 1000 až 5000 bodoch. Importol som knižnicu time aby som vedel testovať čas, za ktorý program zbežne. Pri stále dlhých časoch aj pri menších počtoch bodov som import a implementoval niektoré funkcie z numpy na urýchlenie procesov hlavne pri vytvorení a spravovaní matice a počítanie vzdialeností. Toto mi výrazne urýchlilo bežanie programu. Takisto som rýchlosť testoval bez vizualizácie, ktorá mi už pri malých počtoch zaberala aj do dvoch minút.

4.1 Porovnanie Centroidu a Medoidu

V mojom kóde sa prepínanie používania centroidu a medoidu prepína pomocou use centroids = True alebo = False Každý spôsob má na výpočet svoju vlastnú funkciu, ktorá sa volá na základe use centroids. Keďže pracujem s veľkým množstvom dát, výhodnejšie je pre mňa použitie centroidu keďže výpočet medoidu je náročnejší na čas. Výber medoidu pre výpočet podobných dát môže byť efektívne kvôli faktu, že medoid je konkrétny bod z daného klusteru.

5 Vizualizácia

Všetky body sú vykreslené v 2D priestore aj pred aj po zhlučovaní. Po zhlučovaní sú zhlučkovy rozdelené podľa farby a centroid/medoid každého zhlučkovu je označený X

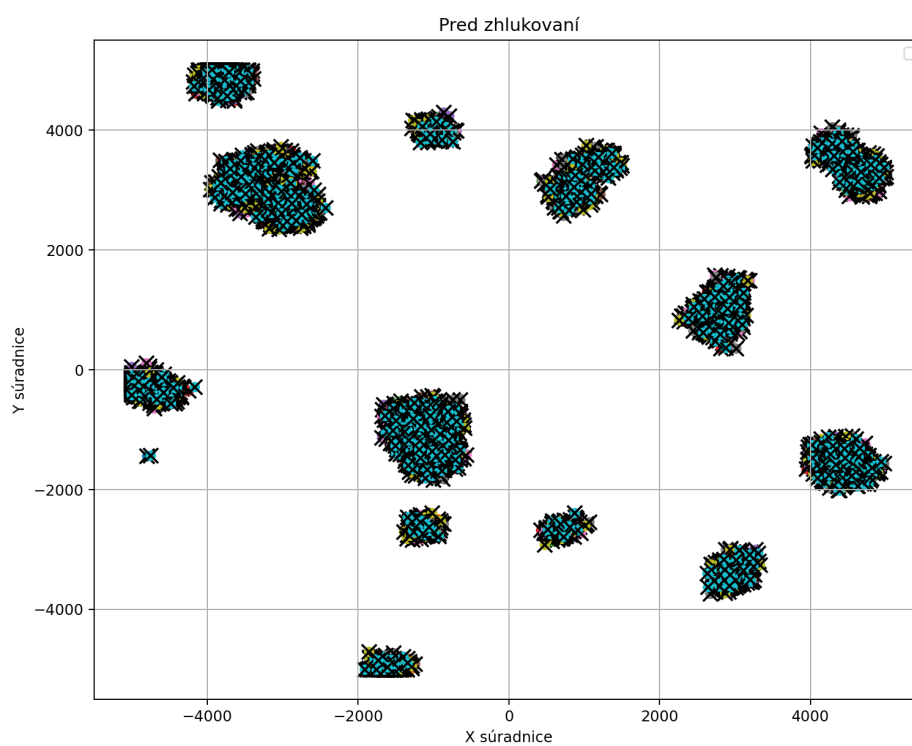


Figure 1: Příklad 20 000 bodov pred zhlukovaním

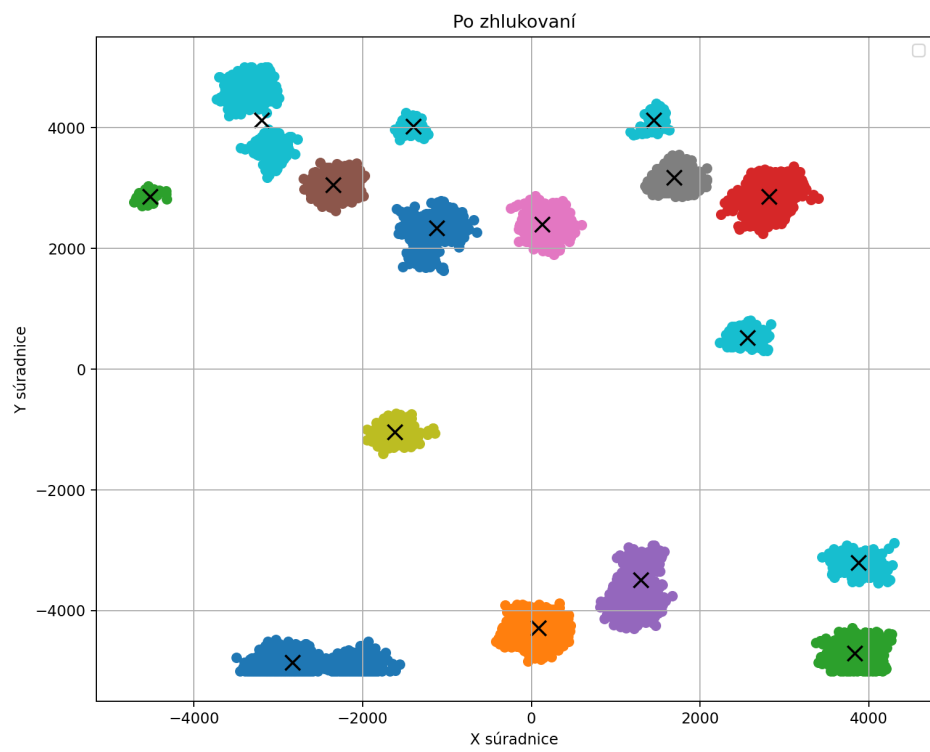


Figure 2: 20 000 bodov. Centroid ako stred

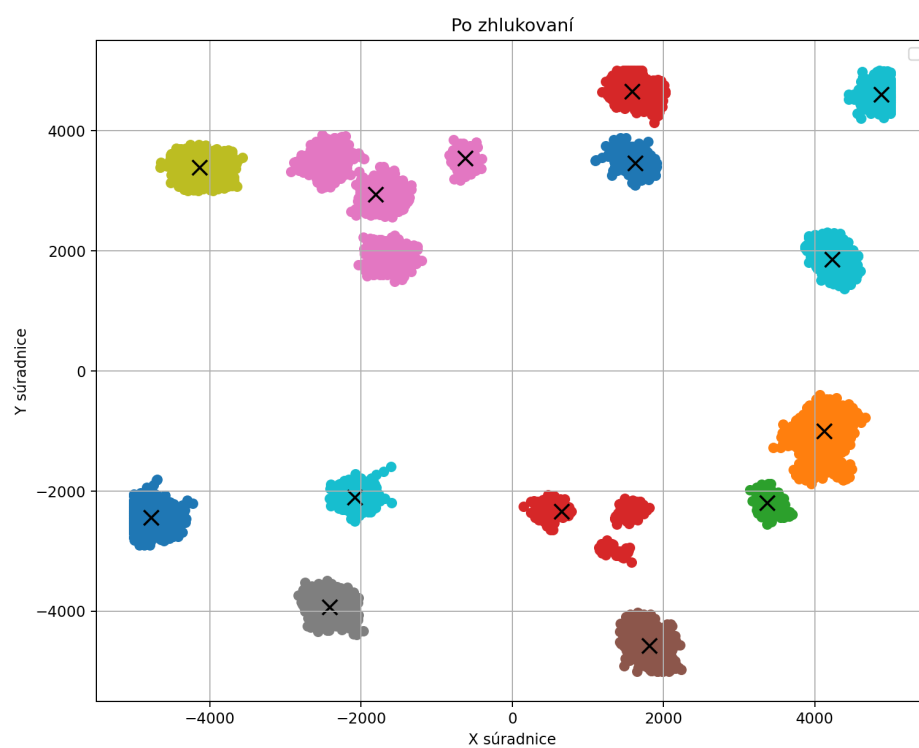


Figure 3: 20 000 bodov. Medoid ako stred

6 Čas behu kódu

Kód mi aj po optimalizácií bežal pomerne dlho. Kedže algoritmus má kubickú zložitosť, pri 20 000 bodoch mi algoritmus trval až cez 3 hodiny. Na základe funkcie `time` som zachytil čas pre centroid 216 minút. Pre medoid mi kód pre 20 000 bodov bežal 360 minút. Taktiež som zistil že vizualizácia zaberá pomerne veľký čas hlavne pri prvom vykreslení. Algoritmus bez vizualizácie už pri 1000 vyhadzoval zrýchlenie skoro o 2 minúty.

7 Záver

Na základe môjho kódu a vykonanej optimalizácie sa centroid preukázal ako efektívnejšie predstavenie stredu zhuku. Úspešnosť je hodnotená na konci programu a je vypočítaná podľa priemernej vzdialenosti bodov od stredu zhuku. Vizualizácia poskytla jasný tvorbu zhukov a potvrdil funkčnosť implementovaných algoritmov. Celkovo táto implementácia ukázala, že aglomeratívne zhukovanie je robustný prístup na anlyzu dát v 2D priestore.