

Kombinácia automatický a vizuálnych metód dolovania dát

Tomáš Poláček

školiťel': Mgr. Matej Novotný

Zadanie a cieľ práce

- Kombinácia automatických a vizuálnych metód dolovania dát
- Opísať klady a zápory oboch prístupov, navrhnúť spôsob ako prístupy skombinovať čo najlepšie a ilustrovať túto kombináciu na zvolenej implementácii

Úvod do dolovania dát

- Veľké množstvo nazhromaždených dát, ktoré treba analyzovať a spracovať
- Netriviálne získavanie implicitných, pred tým neznámych, a potencionálne užitočných informácií (Petr Berka, 2003)
- Metódy dolovania: automatické, vizuálne a kombinované

Automatické metódy

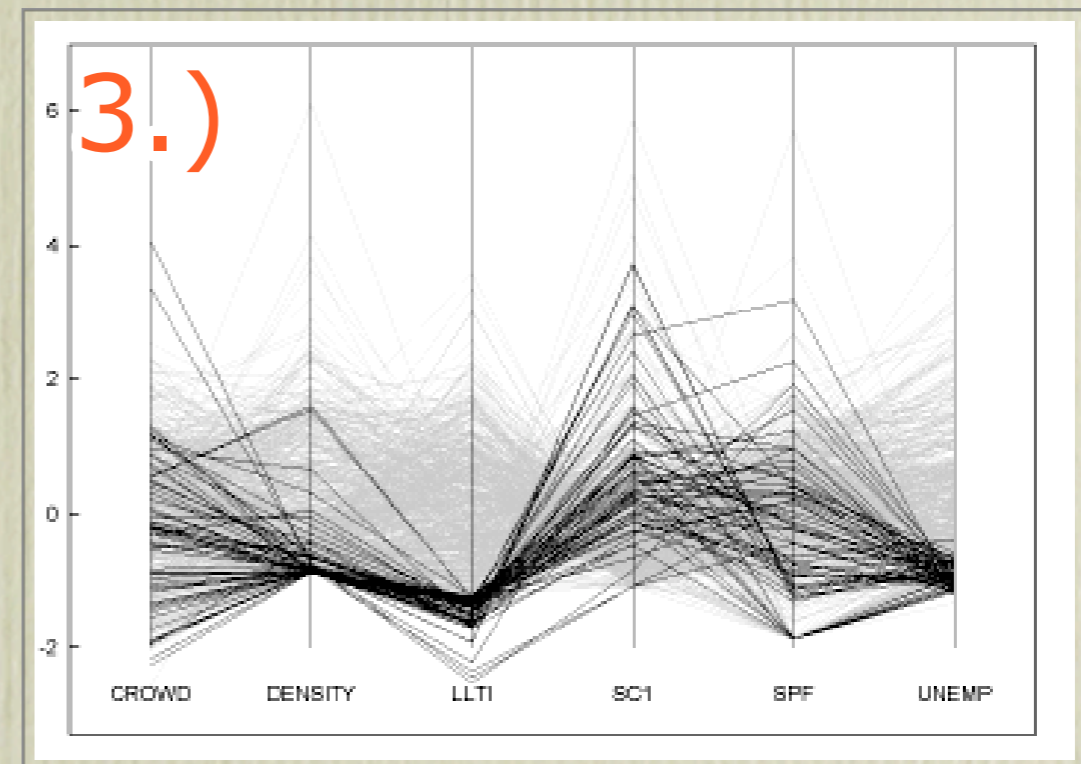
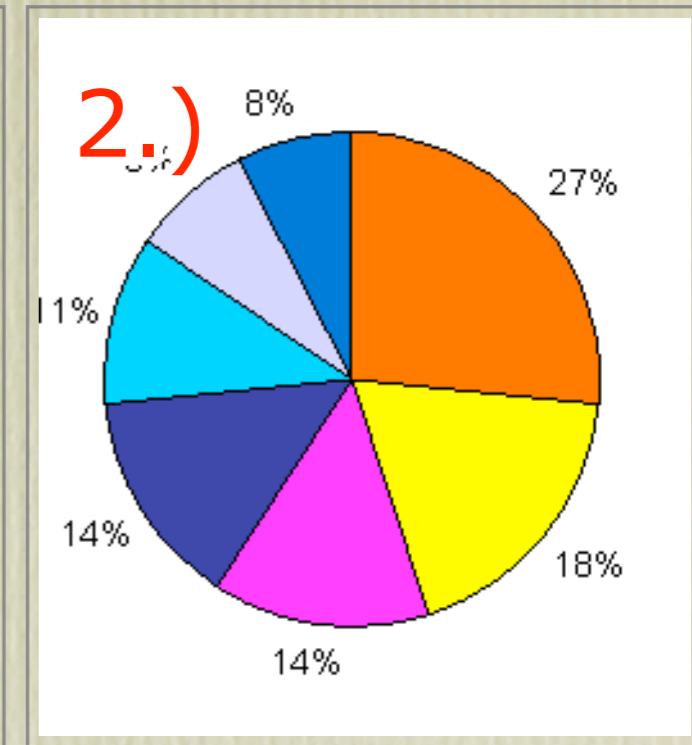
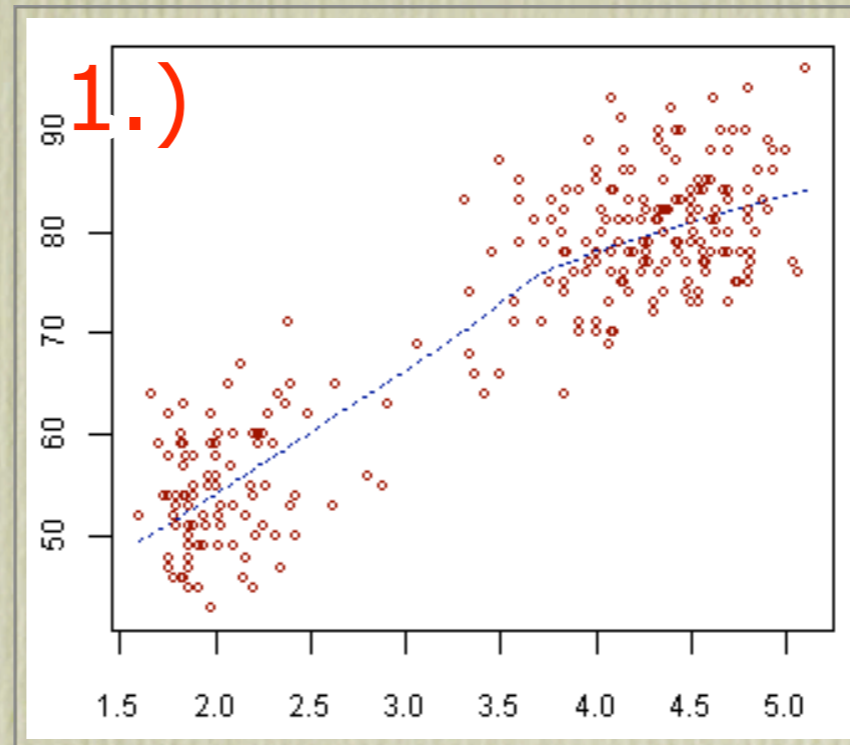
- Patria sem matematické metódy, štatistické metódy a metódy umelej inteligencie
- Väčšinou nie je možné ovplyvniť ich priebeh
- Výsledky sú t'azšie interpretovateľné analytikom a t'azšie prezentovateľné inej osobe
- Napríklad: K-Means, PCA, lineárna regresia ...

Vizuálne metódy

- Hľadanie skrytej informácie v dátach pomocou zraku
- Využitie trénovaného ľudského oka a jeho schopnosti prijímať obrovské množstvo informácií
- Vo veľkej miere sa využívajú skúsenosti analytika
- Vyžadujú viac interakcie od používateľa ako automatické metódy

Vizuálne metódy - príklady

- 1.) Scatter plot
- 2.) Koláčový graf
- 3.) Paralelné súradnice



Kombinované metódy

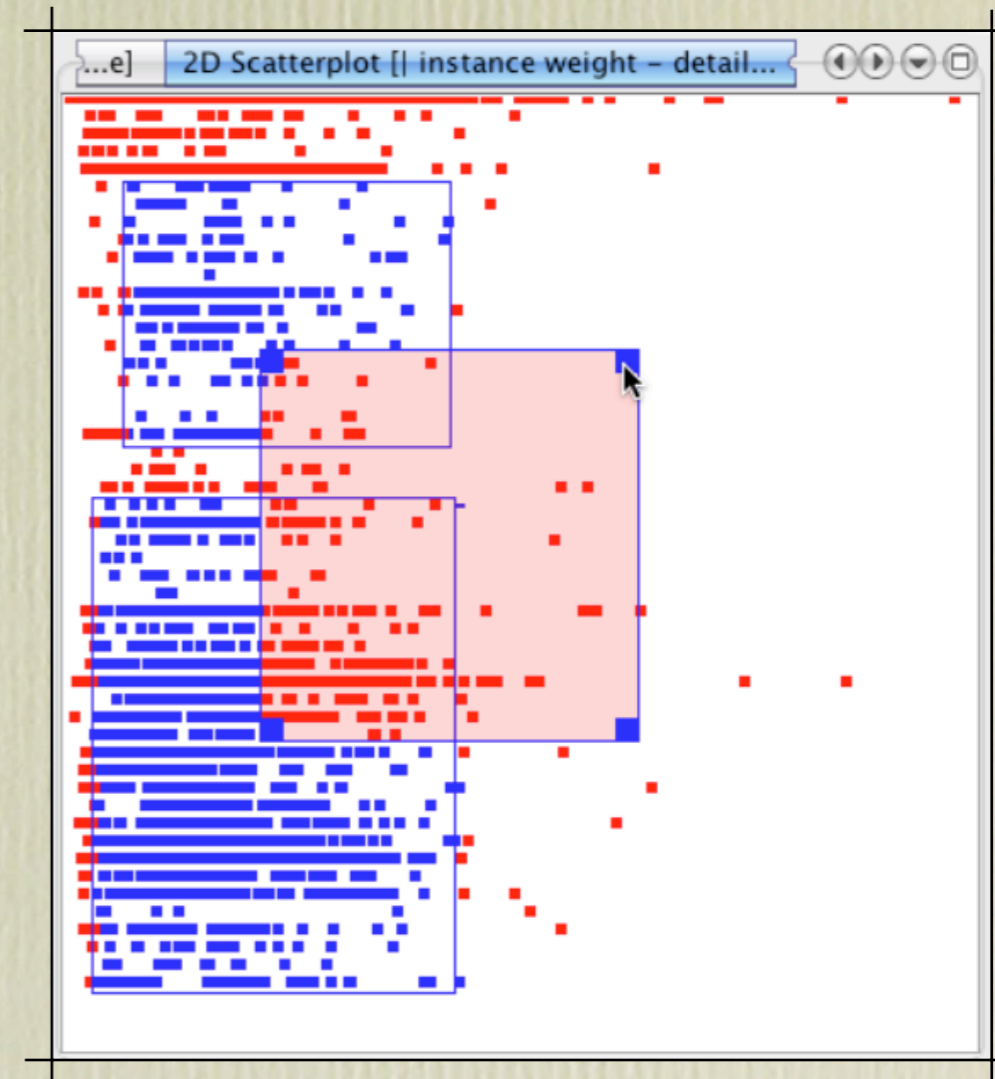
- Rôzne spôsoby kombinácie:
 - **vyvážená kombinácia** - obe metódy zastúpené rovnako, vzájomne sa dopĺňajú
 - **nevyvážená kombinácia** - jedna metóda je použitá ako základ a druhá ju dopĺňuje
- Očakávajú sa pozitívne výsledky

Ciele práce

- Rozšírenie a vylepšenie K-Means algoritmu pomocou kombinácie s vizuálnou metódou Scatterplot
- Zachytiť analytický proces, ktorý vzniká pri segmentácii dát, vizualizovať ho a umožniť jeho opätovné použitie
- Implementovať navrhnuté spôsoby kombinácie do aplikácie
- Ukázať či boli implementované vylepšenia úspešné a poskytujú lepšie výsledky

Použité analytické nástroje a metódy

- K-Means algoritmus
- Obdĺžniková selekcia
- Scatterplot vizualizácia
- Stromové a grafové vizualizácie



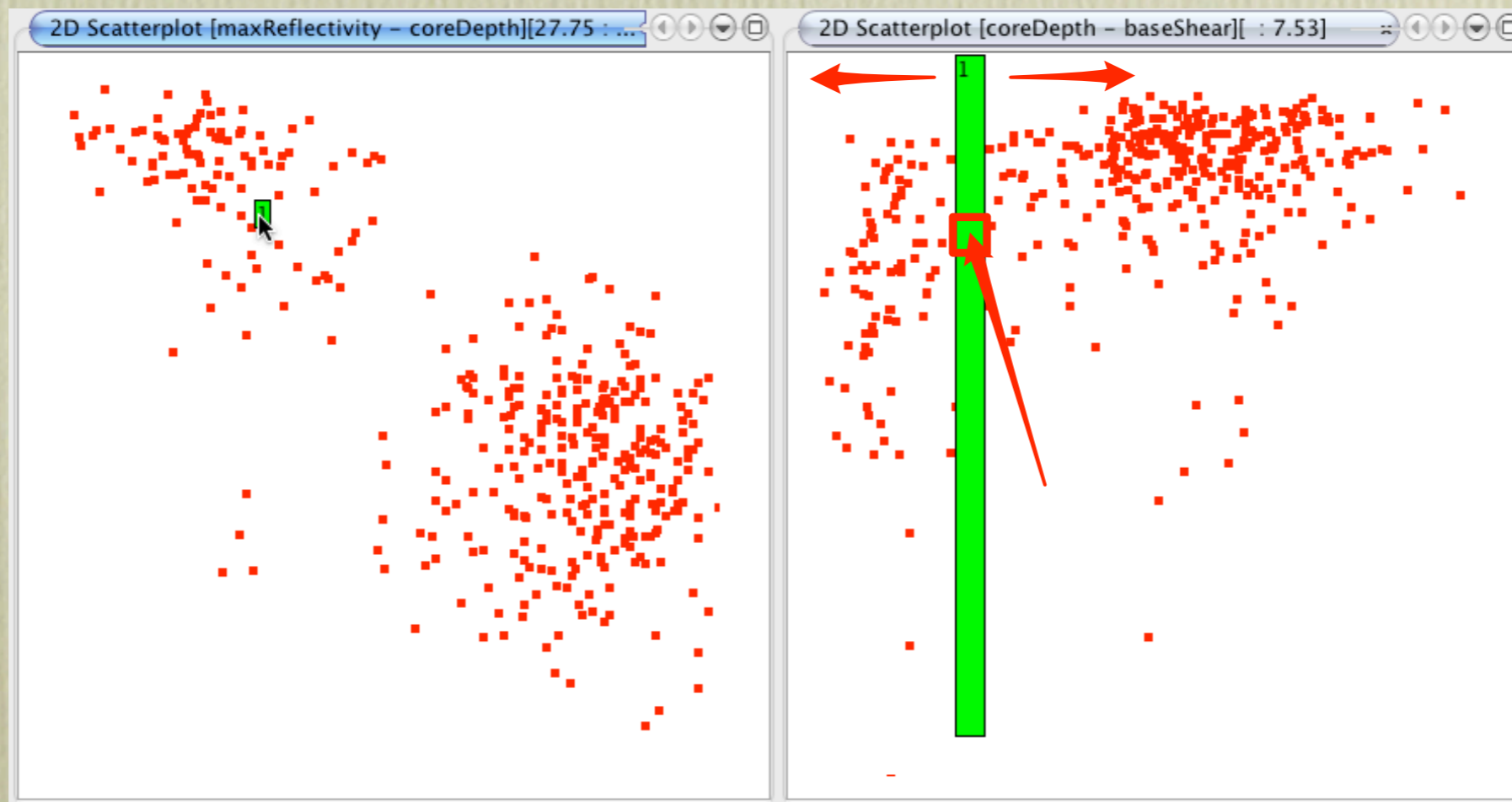
obdĺžniková selekcia

Vylepšenie K-Means

- Zadávanie počiatočných stredov zhlukov pre algoritmus
- Upravenie výpočtu algoritmu tak, aby počítal aj s oddeľujúcimi nadrovinami

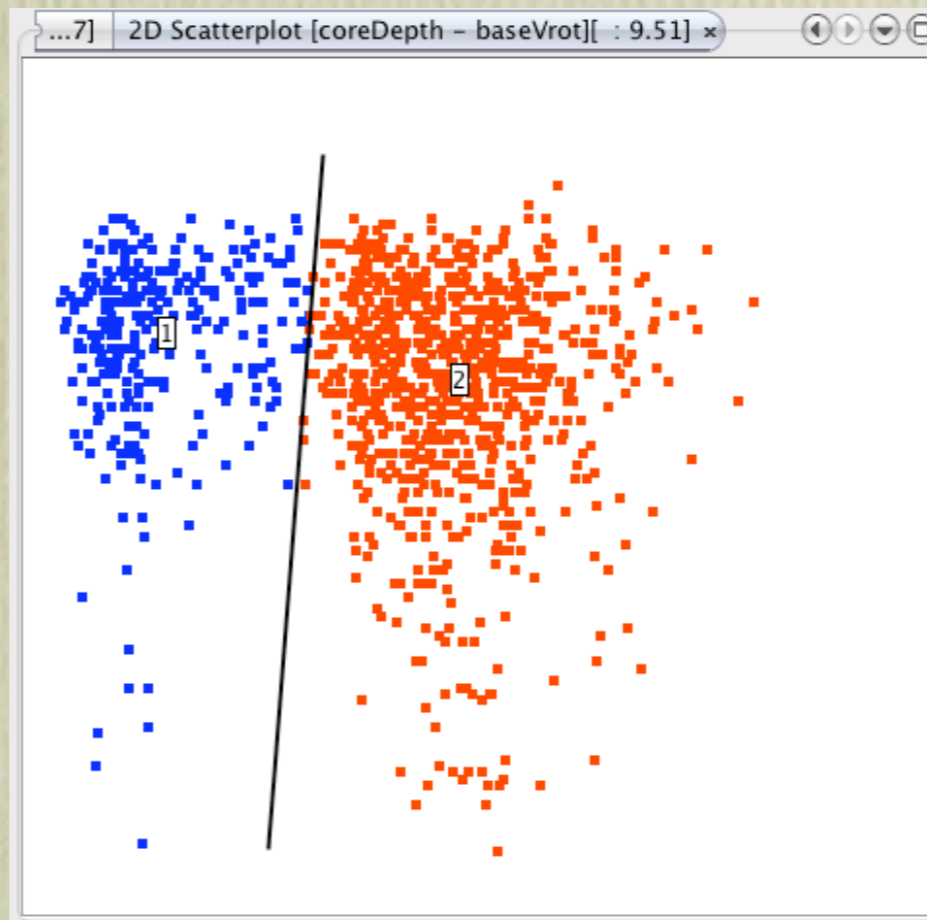
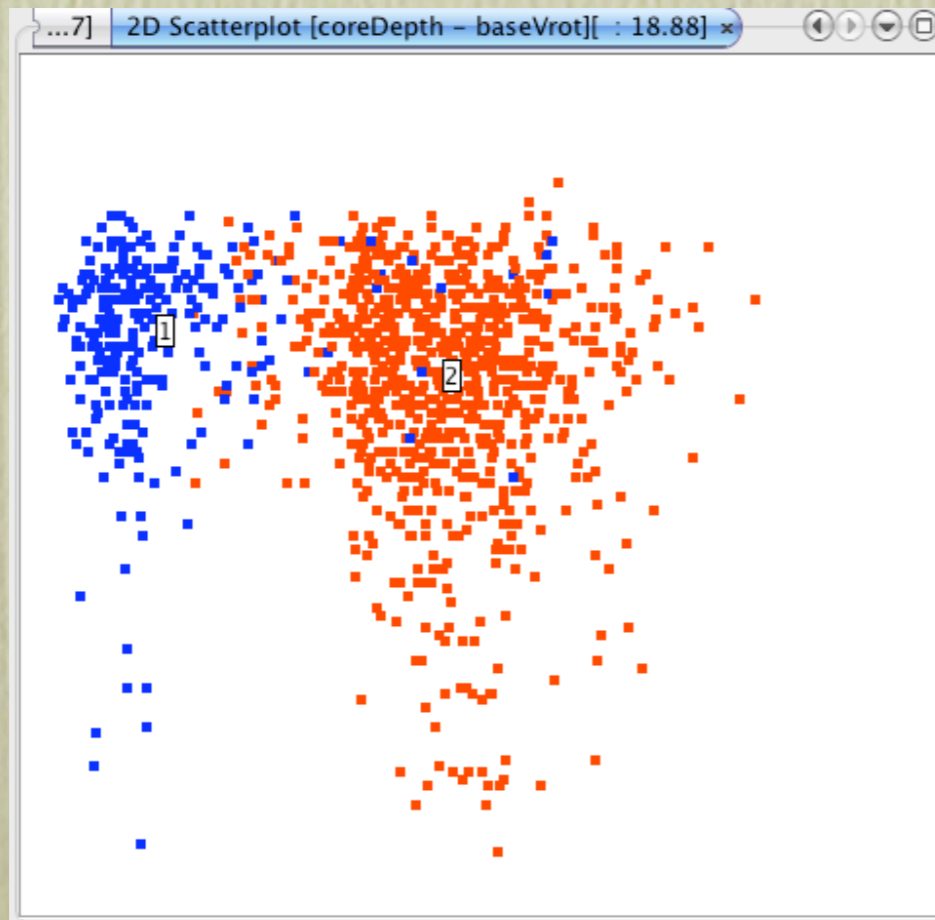
Zadávanie stredov zhlukov

- Zadávanie stredov klikaním do grafu
- Zadávanie súradníc stredu v ďalších rozmeroch
- Odhadnutie počtu a pozície stredov zhlukov
- Zahrnutie do výpočtu len to čo nás zaujíma



Oddel'ujúce nadroviny

- Zadávajú sa kreslením
- Dva významy:
 - “Odrezanie” outliers dát a pevné oddelenie zhlukov

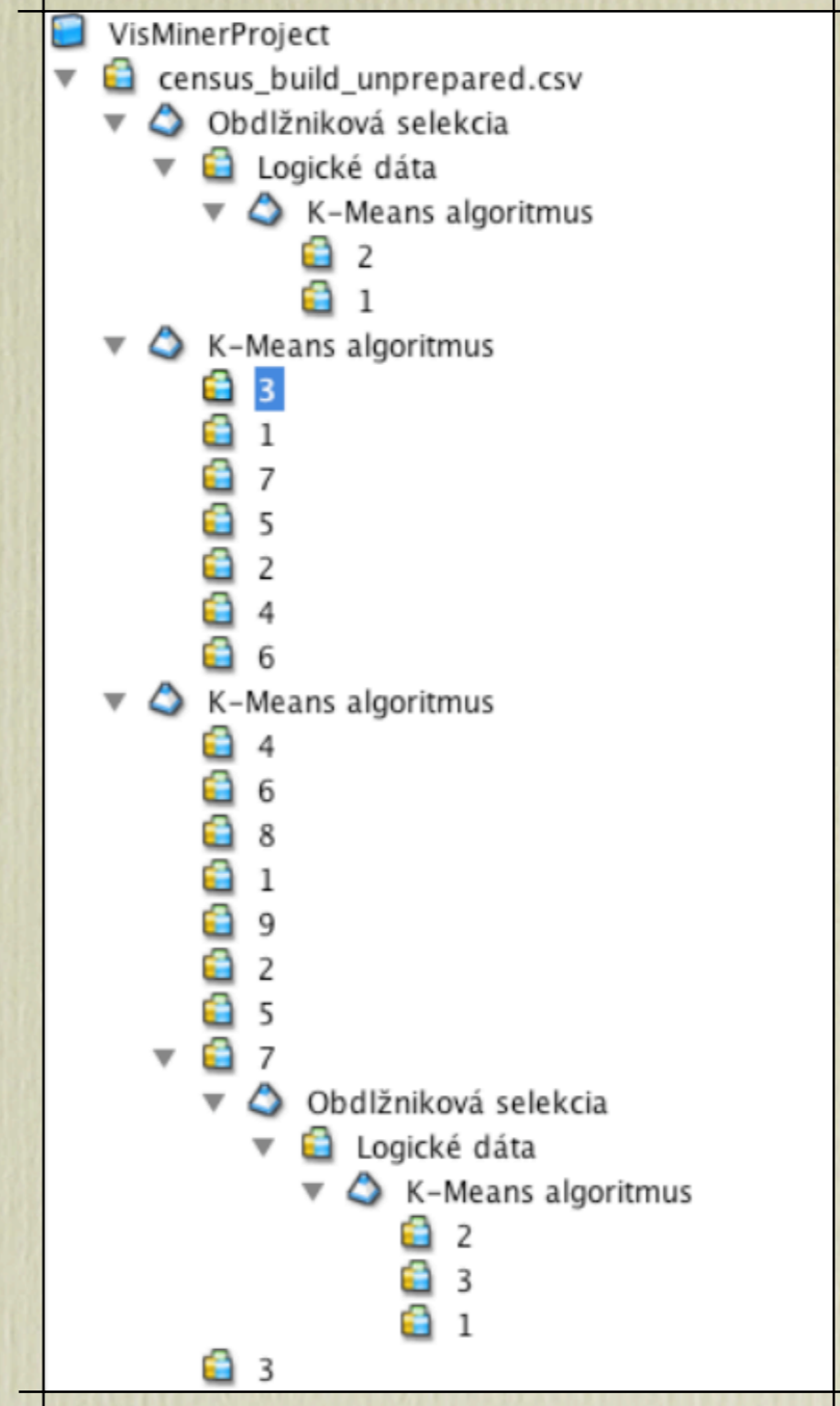


Analytické sedenie

- Záznam uvažovania analytika nad dátami
- Postupnosť krokov práce analytika s dátami
- Vysokoúrovňové programovanie
- Štruktúra sedenia:
 - stromový charakter
 - striedanie dvoch typov vrcholov: dátové vrcholy a filtrovacie vrcholy

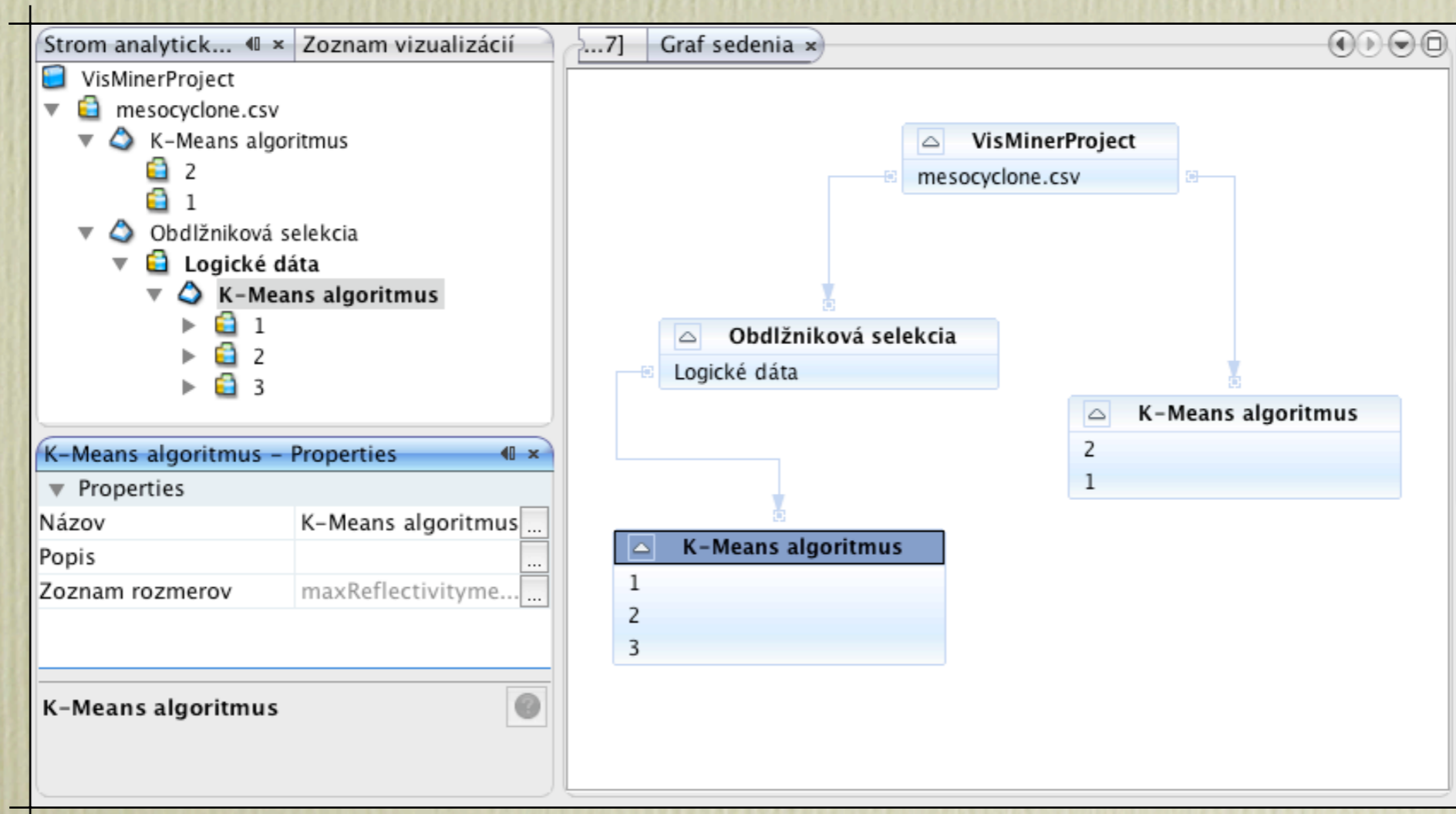
Analytické sedenie - štruktúra

- Dátové vrcholy:
 - Fyzické dáta
 - Logické dáta
- Filtrovacie vrcholy:
 - obdĺžniková selekcia
 - K-Means algoritmus



Vizualizácia sedenia

- Vizualizácia pomocou stromu a grafu
- Prezentovanie postupu ako vznikli výsledky (Visual Storytelling)



Interakcia so sedením

- Vrcholy v strome sedenia umožňujú interakciu
 - Modifikácia parametrov vrchola
 - Mazanie vrcholu zo sedenia

Strom analytického sedenia

- 2
 - Obdĺžniková selekcia
 - Logické dáta
 - K-Means algoritmus
 - 3
 - 2
 - 1
 - K-Means algoritmus

3 - Properties

Properties	
Názov	3
Popis	
Farba	[128,117,230]
Počet záznamov v zhľuku	649
Stredná hodnota	age: 52.4453detailed indu...
Stredná kvadratická odchýlka	0.70160997

3

Strom analytického sedenia

- VisMinerProject
 - generated_1.csv
 - K-Means algoritmus
 - 3
 - 1
 - 2
 - Obdĺžniková selekcia
 - Logické dáta
 - K-Means algoritmus
 - 3
 - 2
 - 1
 - K-Means algoritmus
 - 2
 - 1

Obdĺžniková selekcia - Properties

Properties	
Názov	Obdĺžniková sel...
Popis	

Obdĺžniková selekcia

2D Scatterplot [instance weight - detail...]

Opätovné použitie sedenia

- Opätovné použitie na kompatibilné dáta
- Uloženie sedenia do XML súboru

Výsledky - zhrnutie

- Vylepšenie K-Means algoritmu kombináciou s vizuálnym zadávaním počiatočných stredov zhlukov a oddeľujúcich nadrovín.
- Zachytenie vizualizácie a opätovné použitie analytického sedenia.
- Splnili sa očakávania pozitívneho prínosu kombinácie automatických metód s vizuálnymi

Ukážka aplikácie

- zadávania počiatočných stredov zhlukov
- zadávanie oddeľujúcich nadrovín
- výpočet algoritmu a vizualizácia výsledkov

Strom analytickéh... x

- VisMinerProject
 - mesocyclone.csv
 - K-Means algorit...

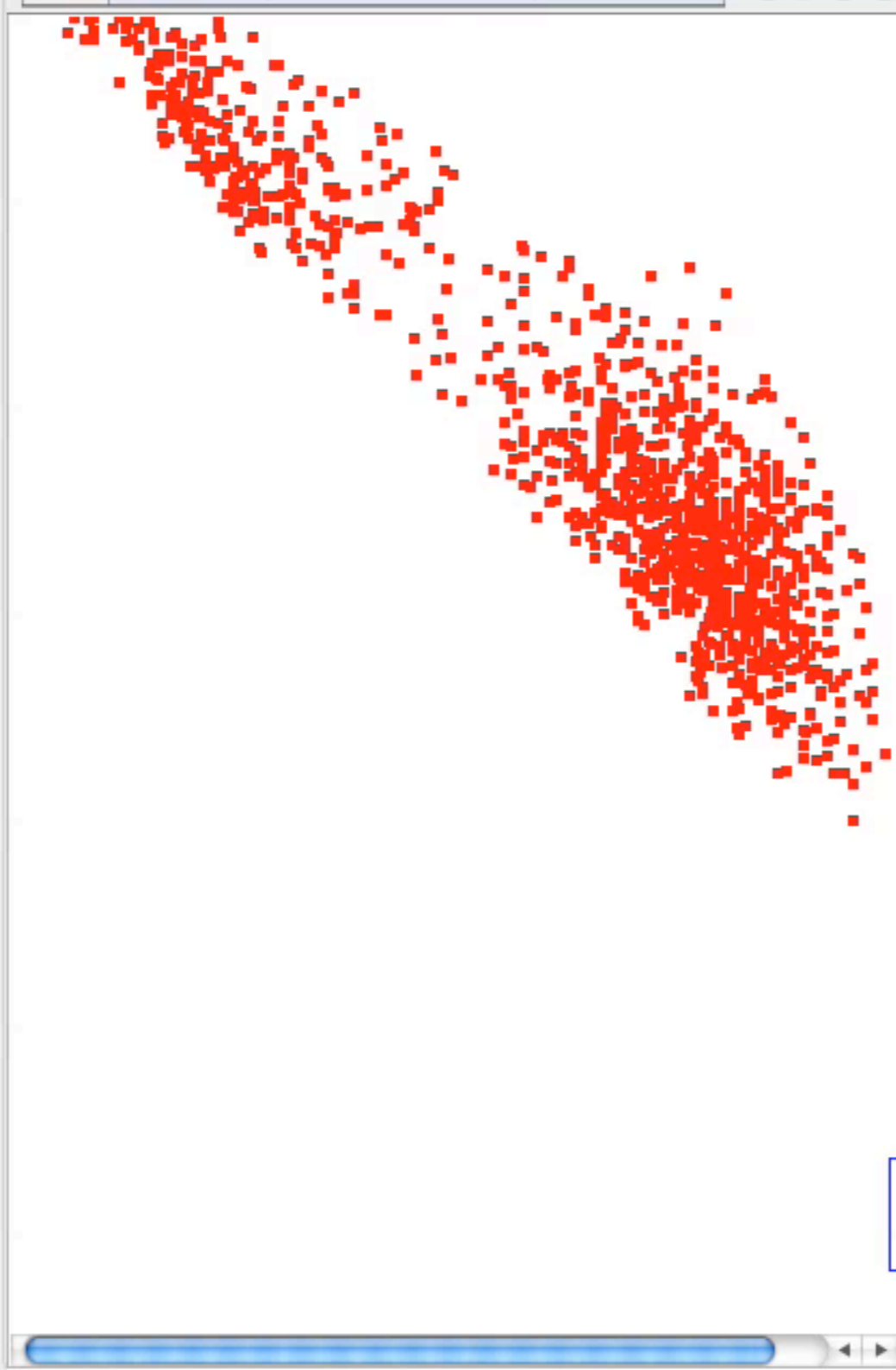
K-Means algoritm... x

▼ Properties

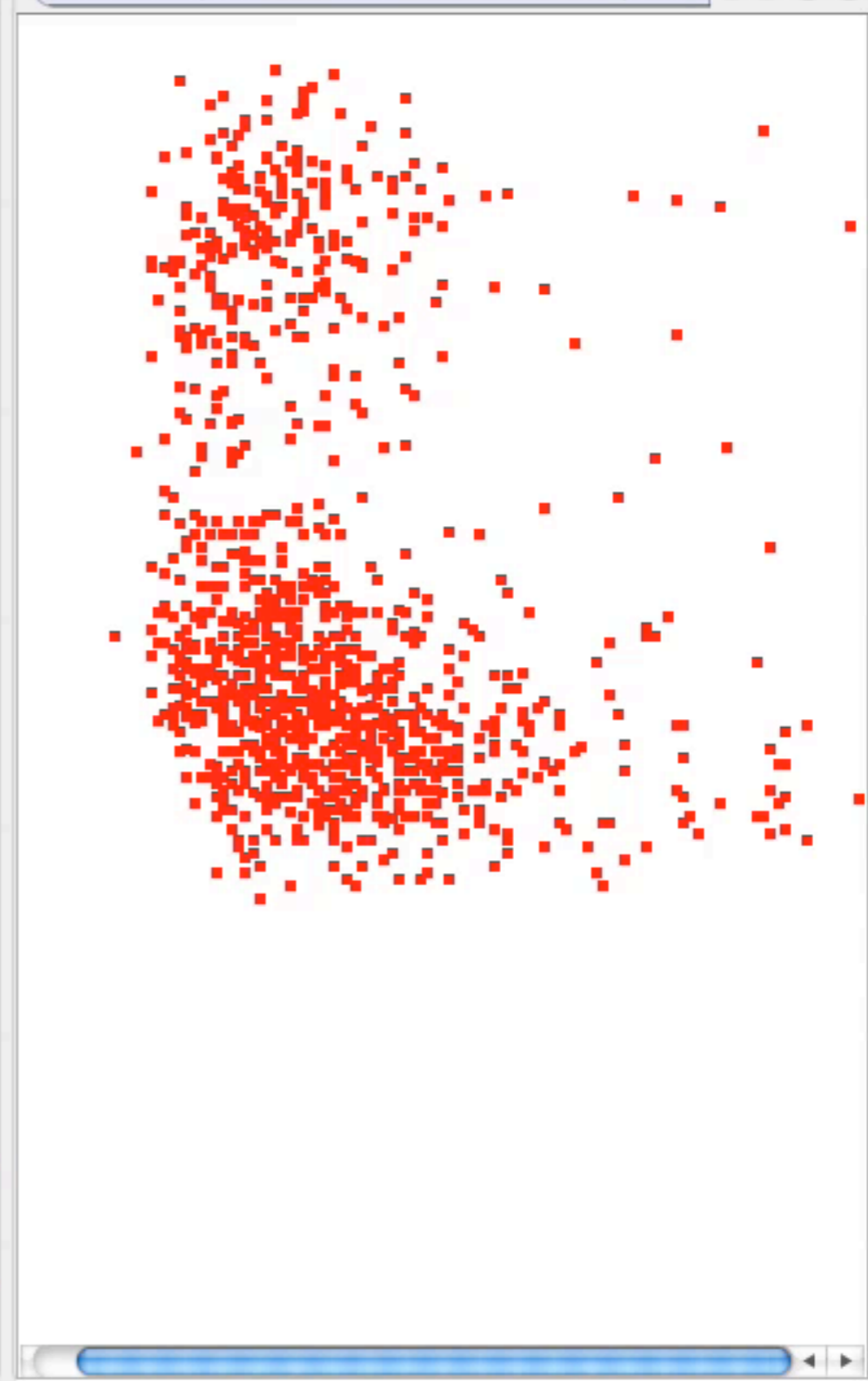
Názov	K-Mean...	...
Popis		...
Počet kôl alg	2	

K-Means algoritmus

...ov 2D Scatterplot [maxReflectivity - meanRe...



2D Scatterplot [baseVrot - maxReflectivity][8...



Selekcia dát

Zoznam dostupných algoritmov

Zoznam vizualizácií

Ďakujem za pozornosť



Otázky oponenta

- **otázka:** Ako by sa zmenila funkčnosť algoritmu, keď by algoritmus K-means počítal zhľady vo všetkých rozmeroch, nielen v tých, ktoré majú zadané stredy?
- **odpoveď:**
 - dopočítanie chýbajúcich súradníc stredov pomocou strednej hodnoty
 - obmedzenie K-means bolo zámerné: zahrnutie do výpočtu len to, čo používateľ chce

- **otázka:** Vysvetlit' význam obdĺžnikovej selekcie, ktorú som v kapitole 5 nevyužil pri dokumentácii vytvoreného postupu.
- **odpoved':**
 - poskladanie selekcie dát z rôznych oblastí
 - vytvorenie logických dát
 - zahrnutie len tých dát do výpočtu, o ktoré ma analytik záujem

- **otázka:** Ako sa vybrali vizualizácie zvolené pre demonštráciu vytvoreného postupu?
- **odpoveď:**
 - so zámerom čo najlepšie demonštrovať funkcie aplikácie

- **otázka:** Porovnanie kvality výsledných zhlukov pri náhodnom zadaní počítačových stredov a pri použití ručne zadaných stredov.
- **odpoveď:**
 - kvalita zhlukov rovnaká
 - počet kôl výpočtu K-means sa výrazne skrátil (pri ukázkových dátach a dvoch zhlukoch viac než o 70%)