

Klasterovanje “the OpenFlights Airports Database extended” skupa podataka

Seminarski rad u okviru kursa Istraživanje podataka 1
Matematički fakultet

David Popov
Avgust 2019.

Sadržaj

1 Opis skupa podataka.....	2
2 Korišćeni alati.....	3
3 Preprocesiranje i vizuelizacija podataka.....	3
4 Klasterovanje.....	9
4.1 K-means.....	9
4.2 Kohonen.....	21

1 Opis skupa podataka

The OpenFlights Airports Database extended je skup koj sadrži podatke o areodromima, lukama i železničkim stanicama, tako da svaki slog sadrži geografsku širinu i dužinu, nadmorsku visinu, vremensku zonu, ime, grad i državu gde se nalazi. Podaci se dobijaju u CSV formatu.

Detaljan opis skupa podataka:

Ime atributa	Tip podataka	Opis
Airport id	Integer	Jedinstven identifikacion broj.
Name	String	Naziv areodroma, železničke stanice ili luke.
City	String	Ime grada gde se nalazi areodrom, železničke stanica ili luka.
Country	String	Ime države gde se nalazi areodrom, železničke stanica ili luka.
IATA	String	Kod koj se sastoji od tri slova. Ukoliko nema vrednost postavljeno je na null.
ICAO	String	Kod koj se sastoji od četiri slova. Ukoliko nema vrednost postavljeno je na null.
Latitude	Float	Geografska širina. Negativna vrednost predstavlja jug, a pozitivna sever.
Longitude	Float	Geografska dužina. Negativna vrednost predstavlja zapad, a pozitivna istok.
Altitude	Float	Nadmorska visina u fitima.
Timezone	Float	Vremenska zona koja je prikazana kao razlika sati u odnosu na vreme u Griniču.
DST (Daylight Savings Time)	Char	Predstavlja da li se koristi zimski ili letnji način računanja vremena. Vrednosti: E (Evropa), A (US/Kanada), S (Južna Amerika), O (Australia), Z (Novi Zeland), N (Nije obrađeno) ili U (Nepoznato).
Tz database	String	Prikazivanje vremenske zone u „tz“ formatu.
Type	String	Kog je tipa. Vrednost „airport“ je za areodrome, „station“ za železničke stanice, „port“ za luke i „unknown“ ukoliko je nije poznato.
Source	String	Odakle potiču podaci

2 Korišćeni alati

Za predprocesiranje je korišćen jezik Python sa njegovim bibliotekama Pandas, Timezonefinder, Pytz, Datetime. Za obradu podataka (klasterovanje) i vizuelizaciju je korišten IBM SPSS.

3 Preprocesiranje i vizuelizacija podataka

U ovom delu ćemo se upoznati detaljnije sa podacima i koracima koje sam preduzeo u preprocesiranju radi dobijanja što boljih rezultata.

Kako bih se što bolje upoznao sa podacima morao sam da izvršim analizu. Kako neki podaci sadrže nepostojeće vrednosti, ti podaci nisu korišćeni u daljoj analizi tj. izbačeni su tokom preprocesiranja.

Kolone koje izbacujem iz analize su one koje ili imaju jedinstvenu vrednost za svaki slog ili ima toliko različitih vrednosti da je približno jednako jedinstvenoj vrednosti za svaku kolonu.

Airport id se prvo izbacuje jer je jedinstven za svaki slog.

Polje City se izbacuje jer ima previše različitih gradova, tj. jedan grad može da ima svega nekoliko aerodroma, stanica i luka, pa mi taj udeo u 12000 slogova ništa ne znači.

Što se tiče polja Country kako i tu ima previše različitih vrednosti, a nisam hteo da skroz izbacim to polje, napravio sam novo polje Country New Values koje ima vrednosti 19 zemalja sa najvećim brojem aerodroma, stanica i luka, dok sam svim ostalim slogovima dodelio vrednost Others.

Iz polja DST sam uklonio nepostojeće vrednosti.

Polja IATA i ICAO se izbacuju jer ima previše različitih vrednosti za slogove pa mi u daljoj analizi ne znače.

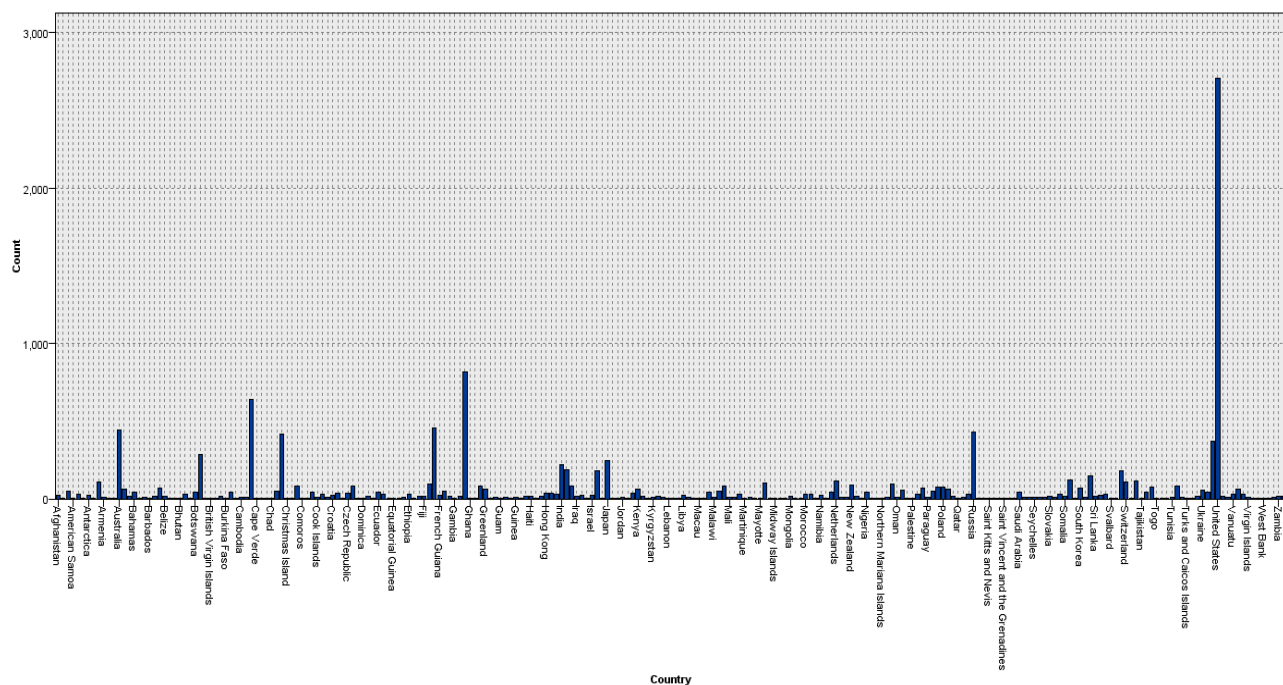
Polje Name se izbacuje jer svaki aerodrom, luka i železnička stanica ima različito ime pa mi u daljoj analizi ne znače.

Iz polja Source sam uklonio nepostojeće vrednosti.

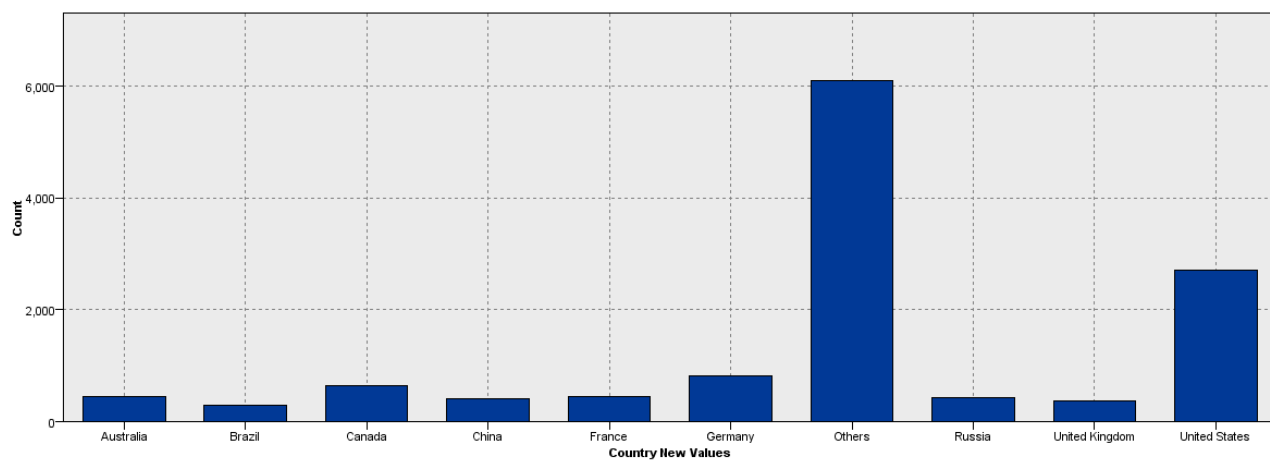
Kako je polje Timezone sadržalo nepostojeće vrednosti, nisam želeo i da ih izbacim već sam u Python-u iskoristio biblioteku timezonefinder gde sam pomoću geografske širine i dužine dobio koja su vremenska zona i samo zamenio sa nepostojećim vrednostima.

Iz polja Type sam uklonio nepostojeće vrednosti.

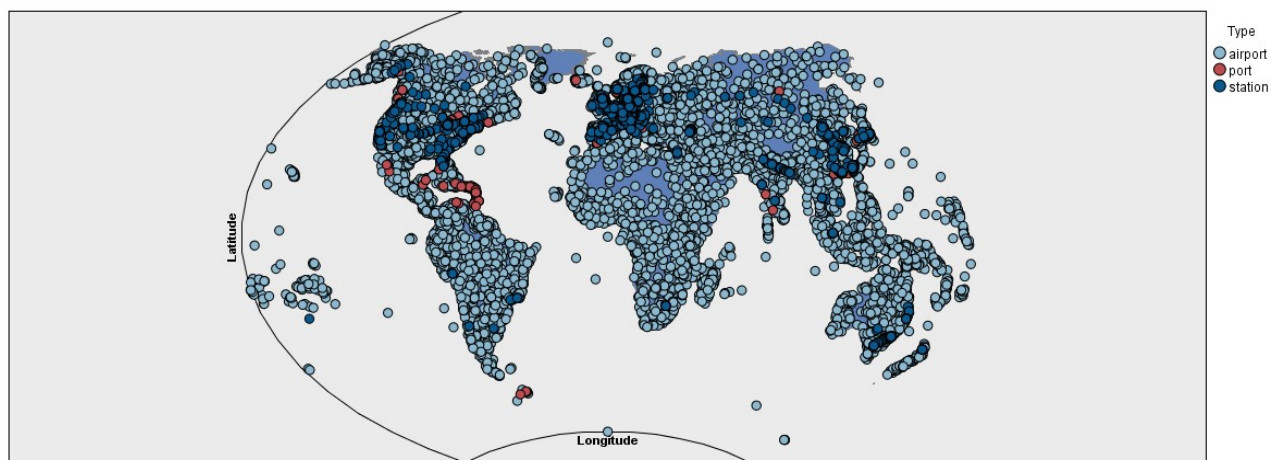
Polje Tz database sam takođe izbacio jer ima previše različitih vrednosti pa mi onda ništa ne znači u daljoj analizi.



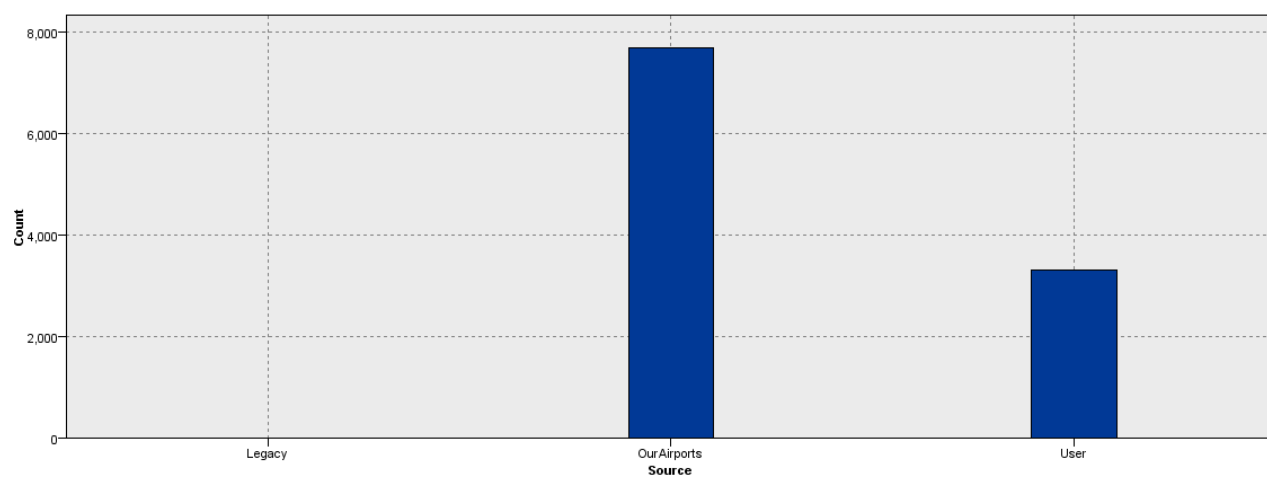
Slika 3.1: Polje Country pre izmene



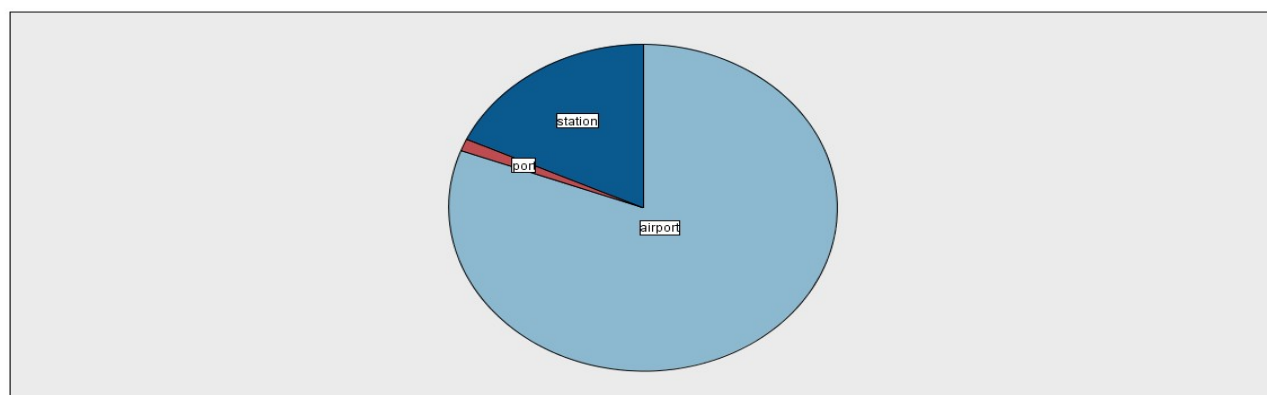
Slika 3.2: Country New Values (Izmenjeno polje Country)



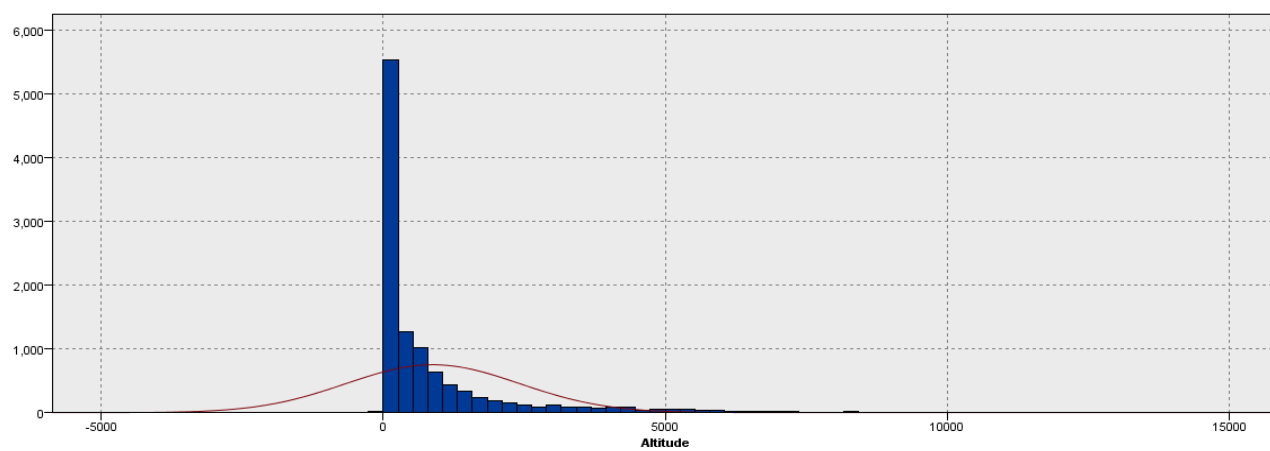
Slika 3.3: Mapa rasporeda areodroma, luka i železničkih stanica



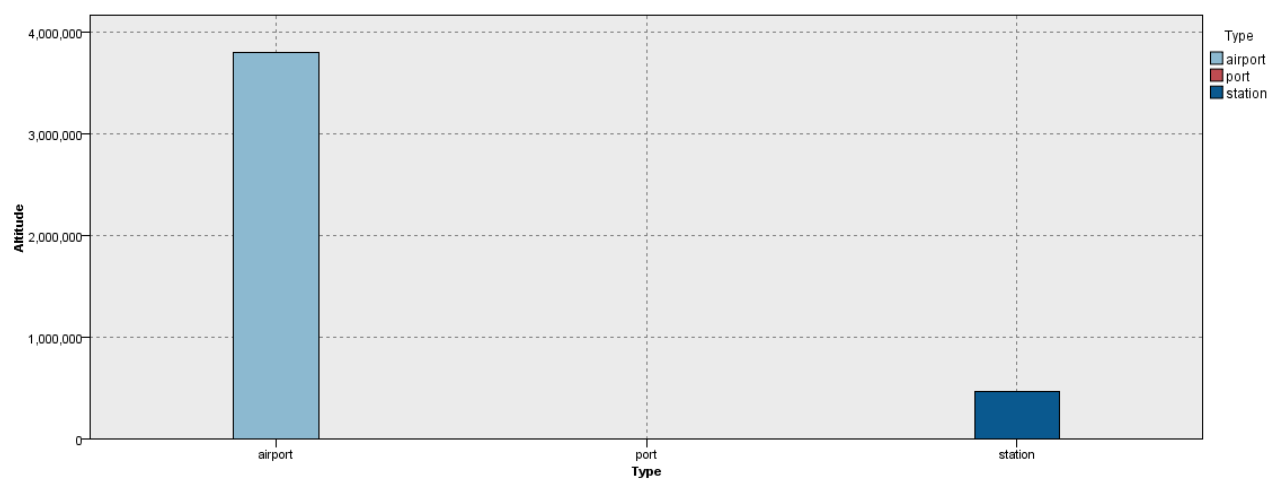
Slika 3.4: Izvor podataka u odnosu na broj slogova



Slika 3.5: Zastupljenost areodroma, luka i železničkih stanica



Slika 3.6: Histogram nadmorske visine



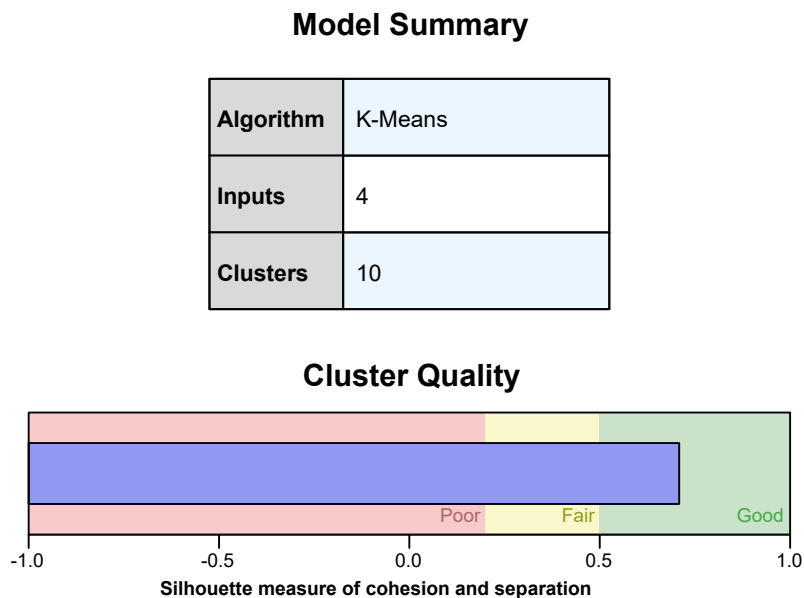
Slika 3.7: Odnos nadmorske visine i tipa (areodrom, luka, železnička stanica)

4 Klasterovanje

Ideja klasterovanja jeste da podelimo naše podatke u grupe (klaster) prema određenim zavisnostima. Prilikom analize koristili su se K-means i Kohonen algoritam

4.1 K-means

Prva ideja je bila da se vidi koj je maksimalan broj atributa koji može da se koristi tako da klasterovanje prikazuje dobre rezultate.

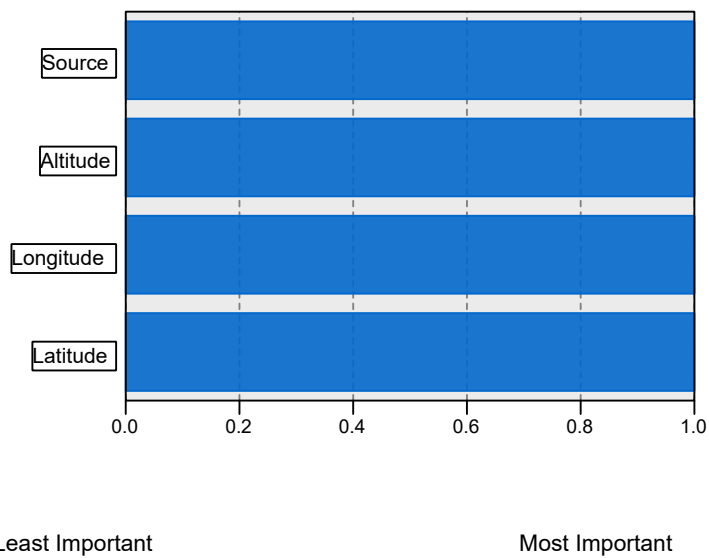


Slika 4.1.1: Koeficijent seneke za prvu ideju

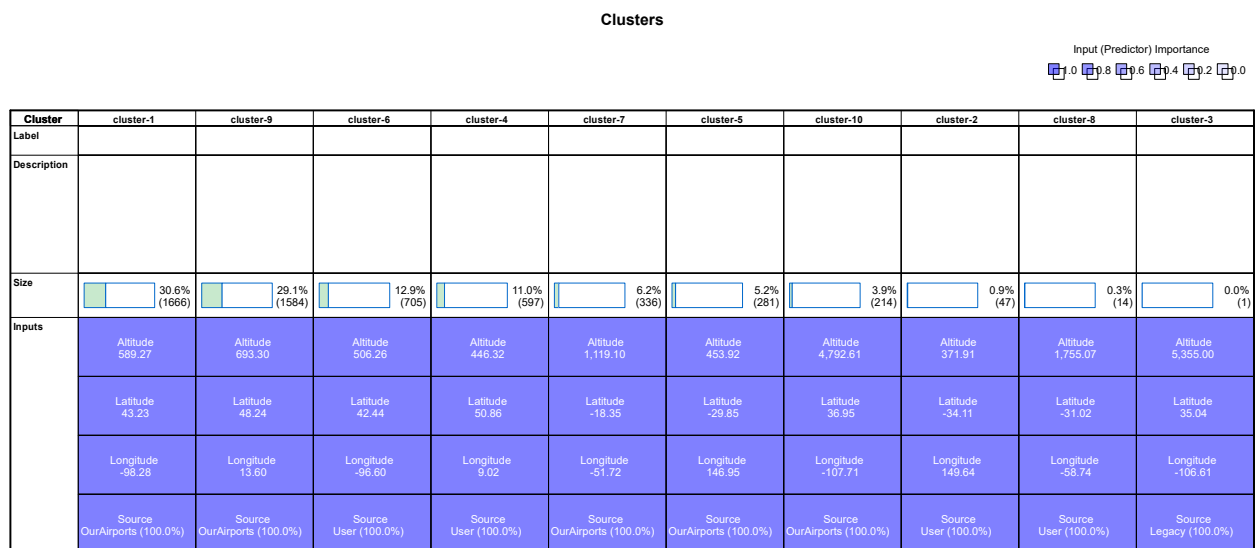
Dobilo se da je najbolje kad je broj klastera 10 i kad se se kao ulazni podaci koriste atributi: Altitude, Latitude, Longitude i Source. Koeficijent senek koj se dobije je 0.7.

Dalje možemo videti da sva 4 atributa maksimalno utiču na rezultate klasterovanja.

Predictor Importance



Slika 4.1.2: Zavisnost atributa od važnosti za predviđanje



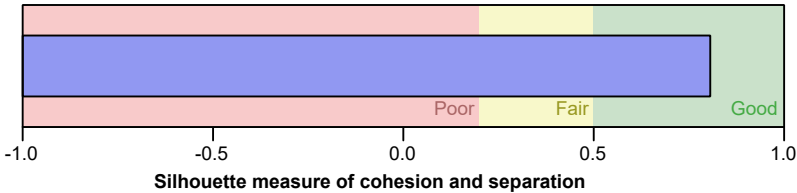
Slika 4.1.3: Rezultati klasterovanja prve ideje

Druga ideja je bila da se izvrši klaster analiza samo za atribute Timezone i Longitude. Najbolji rezultati se dobijaju kad imam 6 klastera, dok se za koeficijent senke dobija vrednost 0.8.

Model Summary

Algorithm	K-Means
Inputs	2
Clusters	6

Cluster Quality



Slika 4.1.4: Koeficijent senke za drugu ideju

Clusters

Input (Predictor) Importance						
<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>						
Cluster	cluster-6	cluster-1	cluster-5	cluster-2	cluster-3	cluster-4
Label						
Description						
Size	<div><div></div><div></div>40.3% (2192)</div>	<div><div></div><div></div>37.5% (2040)</div>	<div><div></div><div></div>16.0% (873)</div>	<div><div></div><div></div>6.1% (331)</div>	<div><div></div><div></div>0.1% (8)</div>	<div><div></div><div></div>0.0% (1)</div>
Inputs	Longitude 12.14	Longitude -78.59	Longitude -128.59	Longitude 146.76	Longitude -89.76	Longitude 174.11
	Timezone 1.15	Timezone -4.97	Timezone -8.03	Timezone 10.07	Timezone 9.97	Timezone -10.00

Slika 4.1.5: Rezultati klasterovanja za drugu ideju

Ono što možemo da zaključimo jeste da su se klasteri grupisali po meridijanima i vremenskim zonama, što nam i jeste očekivan rezultat.

Prvi klaster se prostire kroz Evropu i Afriku.

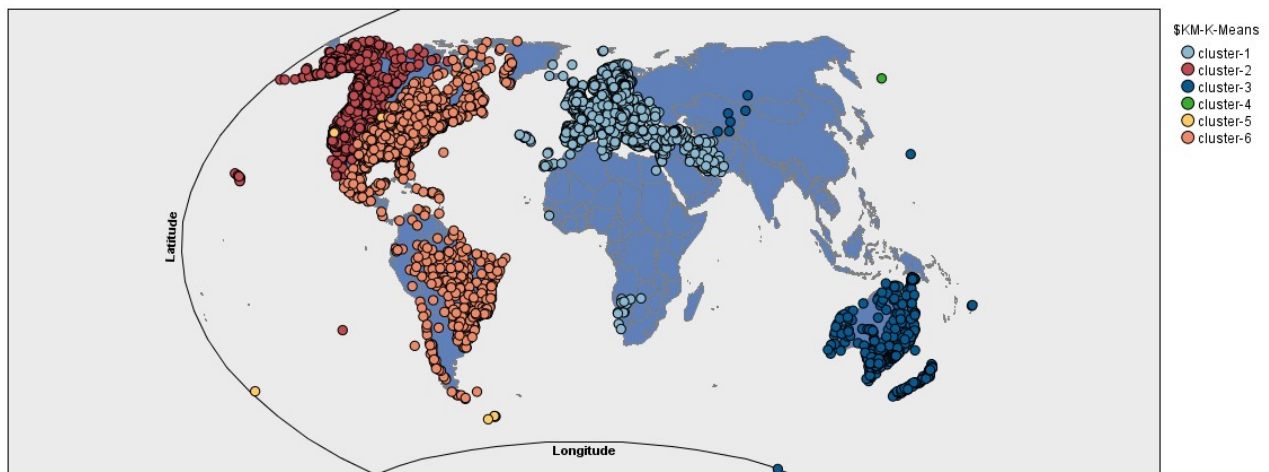
Drugi klaster je zauzeo predeo oko Aljaske i zapadnog dela Severne Amerike.

Treći klaster se prostire kroz Australiju, Okeaniju i Novi Zeland.

Četvrti klaster je zauzeo mali deo kod severne Azije.

Peti klaster je zauzeo delove u tihom okeanu.

Šesti klaster se protire kroz istočni deo Severne Amerike i kroz Južnu Ameriku.



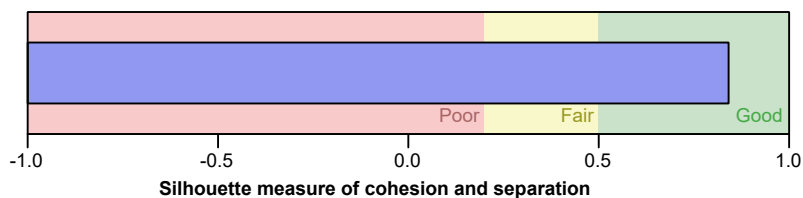
Slika 4.1.6: Mapa klastera za drugu ideju

Treća ideja je bila da se izvrši klaster analiza samo za atribute Type i Altitude. Najbolji rezultati se dobijaju kada se koriti 3 klastera, a za koeficijent senke dobijamo rezultata 0.8.

Model Summary

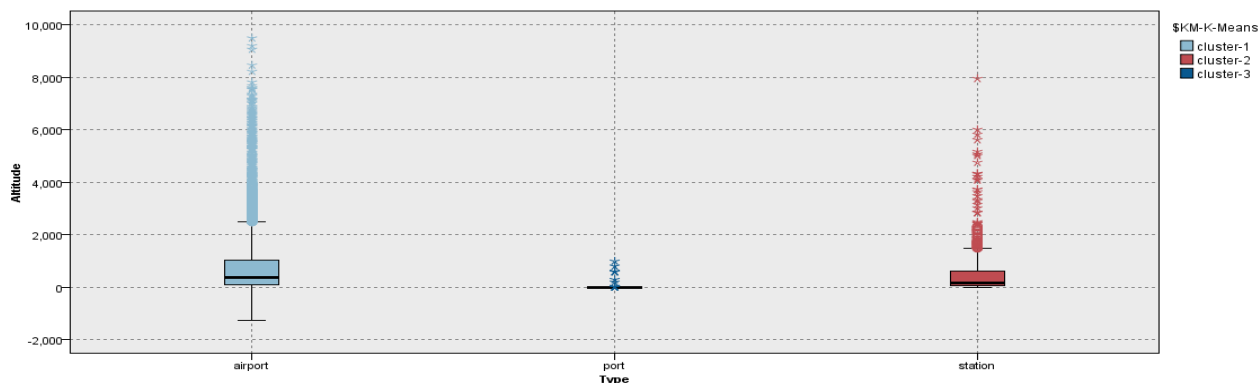
Algorithm	K-Means
Inputs	2
Clusters	3

Cluster Quality



Slika 4.1.7: Koeficijent seneke za treću ideju

Rezultat koj se dobija je da sva tri klastera sadrže različite tipove. Prvi sadrži samo areodrome, gde je prosečna visina 866.54 ft. Drugi sadrži samo železničke stanice, gde je prosečna visina 474.71 ft. Treći sadrži samo luke, gde je prosečna visina 71.30 ft. Rezultati koji se dobijaju su krajnje logični, luke će uvek biti na nadmorskoj visina koja je bliža nadmorskoj visina vode (0 ft), dok areodromi i železničke stanice mogu da se grade i na većim nadmorskim visinama.



Slika 4.1.8: Boxplot klasterovanja za treću ideju

Četvrta ideja je bila da se izvrši klaster analiza samo za attribute Longitude i Latitude. Najbolji rezultati se dobijaju kada se koriti 7 klastera, a za koeficijent senke dobijamo rezultata 0.7.

Model Summary

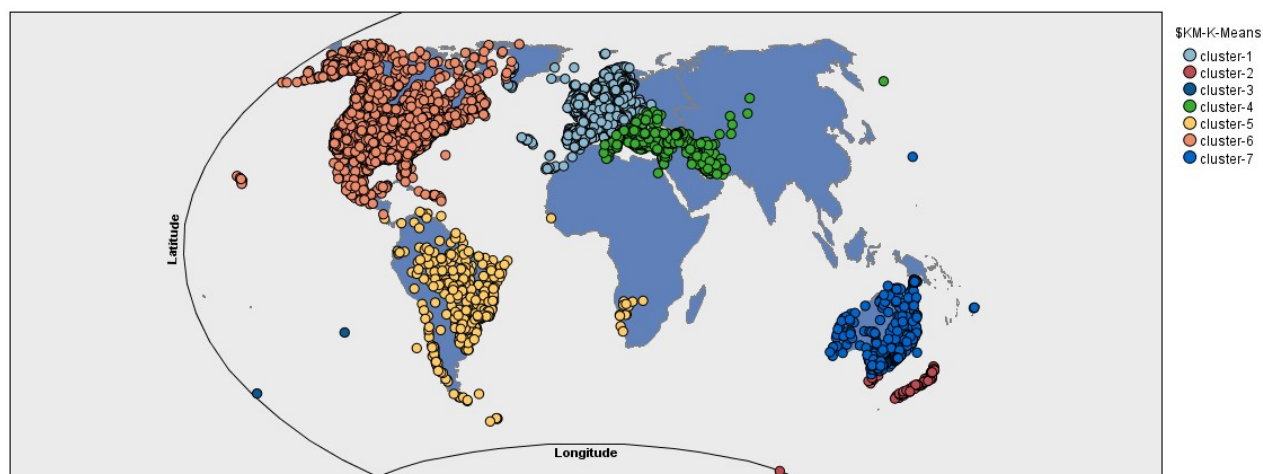
Algorithm	K-Means
Inputs	2
Clusters	7

Cluster Quality



Slika 4.1.9: Koeficijent seneke za četvrtu ideju

Rezultati koj se dobiju su da svaki klaster zauzima po neki deo kontinenta. Prvi zauzima područije Evrope. Drugi zauzima prostor Novog Zelanda. Treći zauzima prostor Austrlije. Četvrti zauzima deo prostora Azije, tj. Malu Aziju i bliski istok. Peti deo Afrike i Južnu Ameriku. Šesti zauzima prostor Severne Amerike. Sedmi obuhvata ostrva u Tihom okeanu.



Slika 4.1.10: Mapa klastera za četvrtu ideju

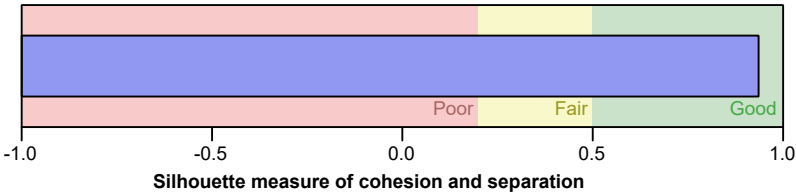
Rezultati su ocekivani s obzirom da smo klasterovali po geografskoj širini i dužini.

Peta ideja je bila da se izvrši klaster analiza samo za atribute Country New Values i Type. Najbolji rezultati se dobijaju kad imam 10 klastera, dok se za koeficijent senke dobija vrednost 0.9.

Model Summary

Algorithm	K-Means
Inputs	2
Clusters	10

Cluster Quality

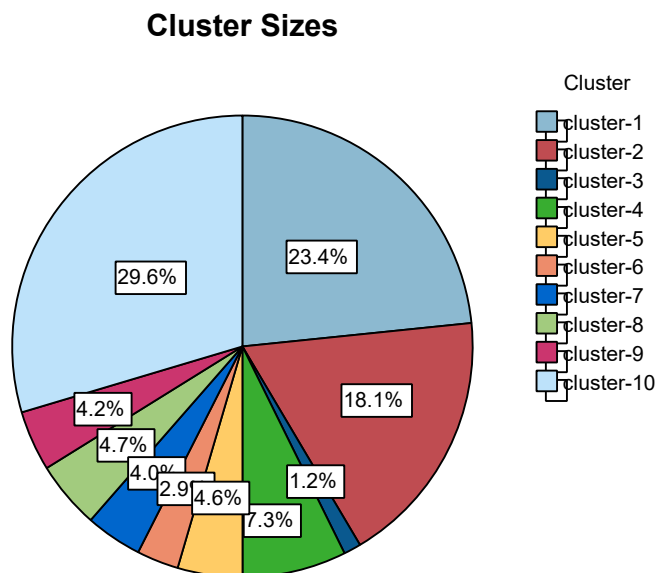


Slika 4.1.11: Koeficijent seneke za petu ideju

Clusters

Input (Predictor) Importance										
<div><div></div>1.0<div></div>0.8<div></div>0.6<div></div>0.4<div></div>0.2<div></div>0.0</div>										
Cluster	cluster-10	cluster-1	cluster-2	cluster-4	cluster-8	cluster-5	cluster-9	cluster-7	cluster-6	cluster-3
Label										
Description										
Size	<div><div></div>29.6% (1612)</div>	<div><div></div>23.4% (1273)</div>	<div><div></div>18.1% (985)</div>	<div><div></div>7.3% (398)</div>	<div><div></div>4.7% (257)</div>	<div><div></div>4.6% (248)</div>	<div><div></div>4.2% (230)</div>	<div><div></div>4.0% (216)</div>	<div><div></div>2.9% (160)</div>	<div><div></div>1.2% (66)</div>
Inputs	Country New Values	Country New Values Others (100.0%)	Country New Values	Country New Values Canada (100.0%)	Country New Values Brazil (100.0%)	Country New Values Germany (100.0%)	Country New Values Australia (100.0%)	Country New Values France (100.0%)	Country New Values	Country New Values Others (45.5%)
	Type airport (100.0%)	Type airport (100.0%)	Type station (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type port (100.0%)

Slika 4.1.12: Rezultati klasterovanja pete ideje

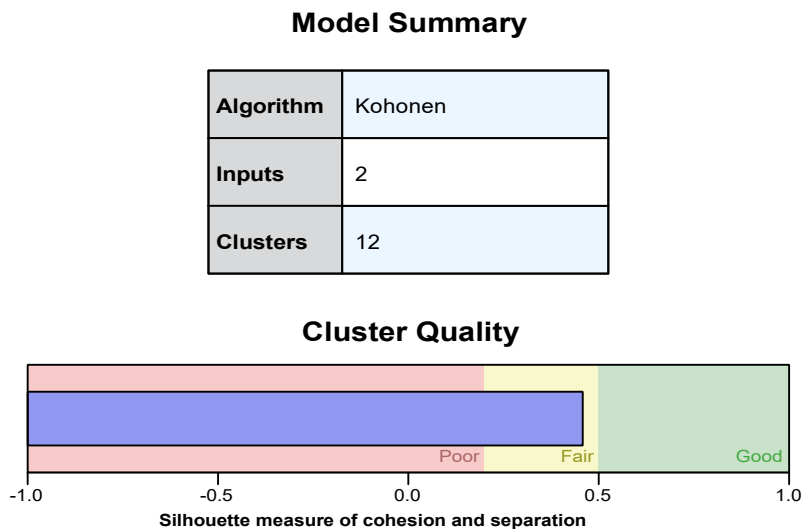


Size of Smallest Cluster	66 (1.2%)
Size of Largest Cluster	1612 (29.6%)
Ratio of Sizes: Largest Cluster to Smallest Cluster	24.42

Slika 4.1.13: Klasteri pete ideje

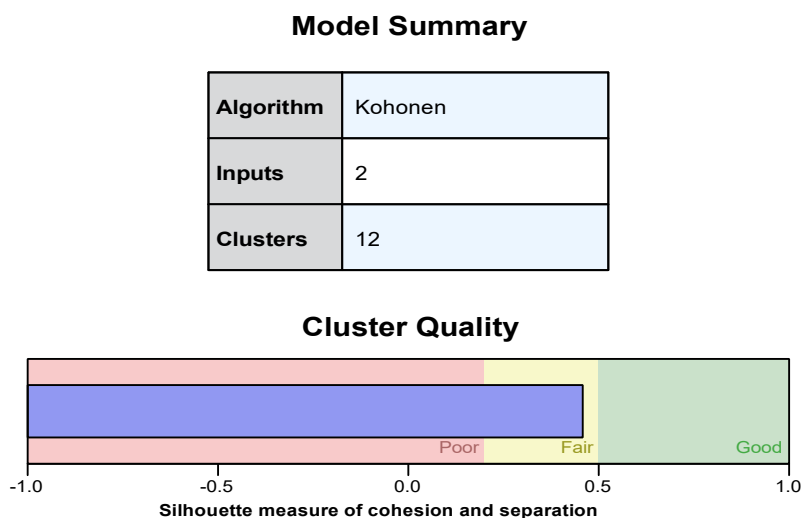
4.2 Kohonen

Prva ideja je bila da se izvrši kohonen algoritam po podrazumevanim vrednostima, ali je klasterovanje dalo lose rezultate, tj. koeficijent senke je imao vrednost 0.4.



Slika 4.2.1: Koeficijent senke za prvu ideju

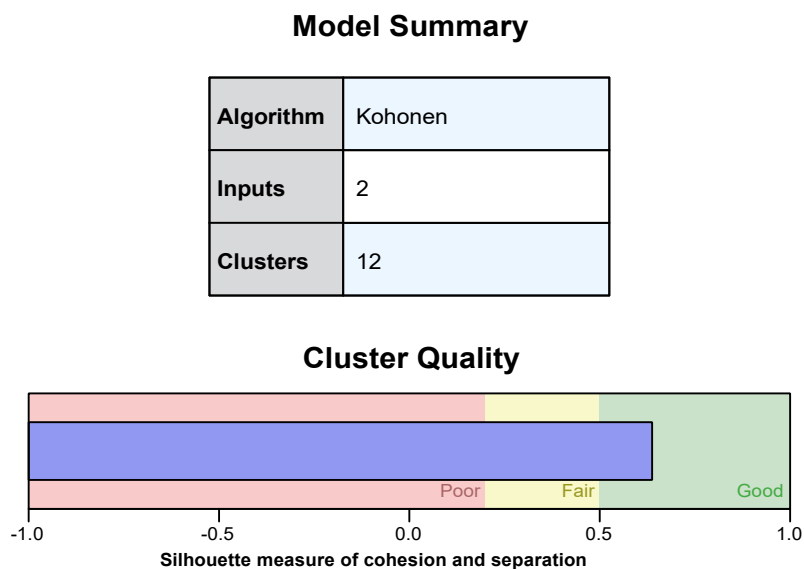
Druga ideja je bila da se izvrši kohonen algoritam nad atributima Longitude i Latitude. Ali i on je dao lose rezultate, vrednost koeficijenta senke je bila nesto manje od 0.4.



Slika 4.2.2: Koeficijent senke za drugu ideju

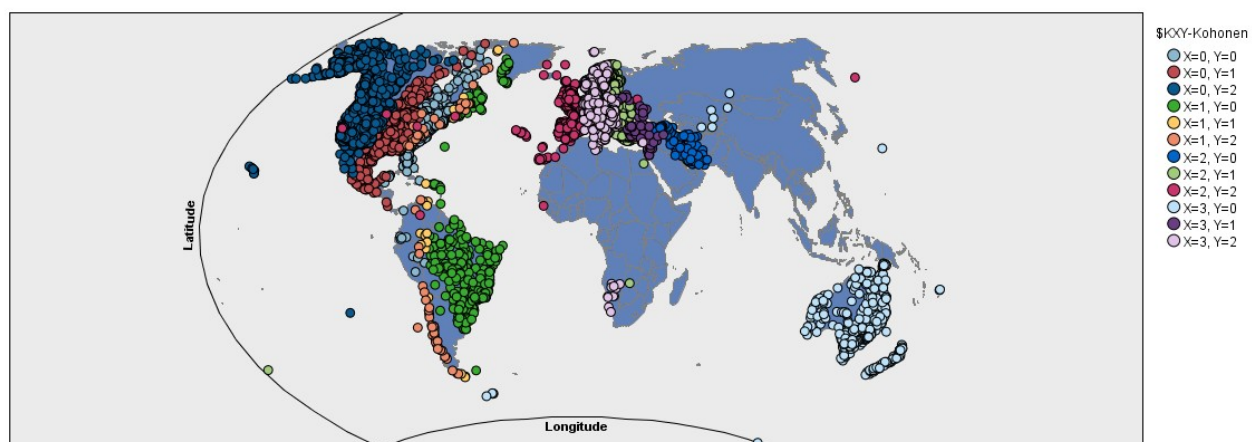
Prethodna klasterovanja nećemo dodatno objašnjavati jer su dala loše rezultate.

Treća ideja je bila da se izvrši kohonen algoritam nad atributima Longitude i Time. Ovaj put klasterovanje je dalo nešto bolje rezultate i koeficijent senke je imao vrednost 0.6.



Slika 4.2.3: Koeficijent seneke za treću ideju

Rezultata koj se dobija sličan je ko i kada se koristi k-means algoritam, s tim da ima više klastera pa su samim tim gušće raspoređeni. Klasteri se raspoređuju po medijanama, što je i za očekivati.



Slika 4.2.4: Mapa klastera za drugu ideju

Četvrta ideja je bila da se izvrši kohonen algoritam nad atributima Altitude i Type. Dobili smo 5 klastera, a vrednost koeficijenta senek je 0.7.

Model Summary

Algorithm	Kohonen
Inputs	2
Clusters	5

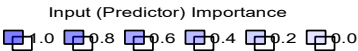
Cluster Quality



Slika 4.2.5: Koeficijent seneke za četvrtu ideju

Rezultata koj se dobija sličan je ko i kada se koristi k-means algoritam, s tim da su klaster koj sadrže luke i železničke stanice ostaju isit, dok klaster koj sadrži areodrome je razbijen na tri, dok u kad smo koristili k-means imali smo samo jedan klaster.

Clusters



Cluster	X=2, Y=2	X=0, Y=2	X=2, Y=1	X=2, Y=0	X=0, Y=0
Label					
Description					
Size	<div><div></div>46.9% (2556)</div>	<div><div></div>18.1% (985)</div>	<div><div></div>18.1% (985)</div>	<div><div></div>15.7% (853)</div>	<div><div></div>1.2% (66)</div>
Inputs	Altitude 163.99	Altitude 474.71	Altitude 876.43	Altitude 2,960.31	Altitude 71.30
	Type airport (100.0%)	Type station (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type port (100.0%)

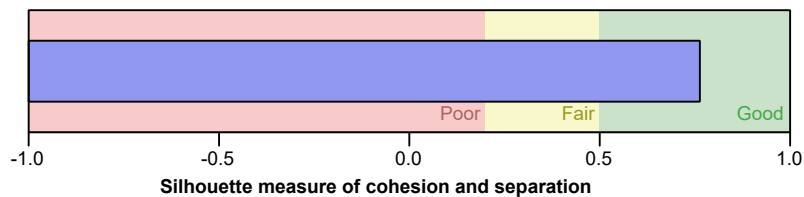
Slika 4.2.6: Rezultati klasterovanja četvrte ideje

Peta ideja je bila da se izvrši kohonen algoritam nad atributima Country New Values i Type. Dobili smo 9 klastera, a vrednost koeficijenta senek je 0.8. Ovde dobijamo najbolje klasterovanje kohonenovim algoritmom. U odnosu na k-means imamo jedan klaster manje. Ostali klasteri su dosta slični sa klasterima dobijenim k-means algoritmom.

Model Summary

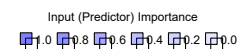
Algorithm	Kohonen
Inputs	2
Clusters	9

Cluster Quality



Slika 4.2.7: Koeficijent seneke za petu ideju

Clusters



Cluster	X=0, Y=2	X=3, Y=2	X=3, Y=0	X=0, Y=0	X=1, Y=0	X=2, Y=2	X=2, Y=0	X=0, Y=1	X=1, Y=1
Label									
Description									
Size	29.6% (1611)	23.7% (1293)	23.4% (1273)	13.5% (735)	4.7% (254)	4.1% (225)	0.6% (30)	0.4% (21)	0.1% (3)
Inputs	Country New Values	Country New Values Canada (30.8%)	Country New Values Others (100.0%)	Country New Values	Country New Values Others (98.4%)	Country New Values France (96.0%)	Country New Values Others (100.0%)	Country New Values	Country New Values France (100.0%)
	Type airport (100.0%)	Type airport (100.0%)	Type airport (100.0%)	Type station (100.0%)	Type station (98.4%)	Type airport (96.4%)	Type port (100.0%)	Type port (100.0%)	Type port (100.0%)

Slika 4.2.8: Rezultati klasterovanja pete ideje