

Klasterovanje Olimpijskih igara

Seminarski rad u okviru kursa
Istraživanje podataka 1
Matematički fakultet

Mladen Canović
mladen.canovic@gmail.com

3. septembar 2019

Sažetak

Olimpijske igre predstavljaju najveći sportski događaj na svetu, koji okuplja sportiste, koji se nadmeću u raznim disciplinama, prethodne 123 godine. Sa druge strane, razvoj tehnologije doveo je do naglog razvoja jedne neolimpijske discipline - istraživanja podataka. Ovaj rad ima za cilj da analizira, obradi i iznese određene zaključke o podacima sa Olimpijskih igara korišćenjem metoda istraživanja podataka. Tehnikama klasterovanja pokušava da odgovori na pitanja o razvoju Olimpijskih igara, broju učesnika, osvajačima medalja, šansama učesnika na osnovu određenih parametara i mnoga druga.

Sadržaj

1	Uvod	2
2	Skup tehnologija i alata	2
3	Analiza i pretprocesiranje podataka	2
3.1	Analiza podataka	2
3.2	Dodatna analiza i vizuelizacija podataka	8
3.3	Izvođenje novih atributa	12
4	Klasterovanje	14
4.1	K-sredine	15
4.2	Hijerarhijsko klasterovanje	20
4.3	Klasterovanje učesnika	23
5	Zaključak	25
	Literatura	25

1 Uvod

Amerikanac Džejms Konoli je pobedom u olimpijskoj disciplini tro-skok, 6. aprila 1896. postao prvi olimpijski šampion nakon više od 1500 godina. Tačnije, od Olimpijskih igara 393. godine nove ere, nakon čega je rimski car Teodosije I zabranio održavanje Olimpijskih igara. Ponovno pokretanje Olimpijskih igara bila je ideja Barona Pjera de Kubertena. Prve moderne Olimpijske igre održane su u Atini, 1896. godine i na njima je učestvovalo 245 muških takmičara iz 14 država [2, 1]. Od 1896. godine do danas, održano je 23 zimskih i 28 letnjih olimpijskih igara u 43 različita grada.

U poglavlju 2 će ukratko biti pobrojan skup tehnologija i alata, koji će biti korišćen u okviru ovog seminarskog rada. U poglavlju 3 biće akcenat na podacima, biće izvršena analiza i pretprocesiranje. Takođe, biće izvršena vizuelizacija podataka i delimično ispitivanje međusobne zavisnosti atributa. U trenucima kada je, po mišljenju autora rada, neophodno dati kratak istorijski osvrt, to će biti učinjeno. U poglavlju 4 će biti izvršeno klasterovanje podataka i analiza dobijenih rezultata. U poglavlju 5 će biti dat osvrt na rad i predloženi dalji koraci u istraživanju.

2 Skup tehnologija i alata

Širok je spektar programskih jezika i programskih alata, koji olakšava analizu i obradu podataka kao i primenjivanje metoda istraživanja podataka nad njima. Skup tehnologija i alata koji će biti korišćen u ovom seminarskom radu je sledeći:

- Programski jezik *Python*;
- Biblioteka *pandas* za manipulaciju podacima;
- Biblioteka *numpy* za manipulaciju matematičkim strukturama;
- Biblioteka *sklearn*, koja olakšava korišćenje raznih metoda istraživanja podataka;
- Biblioteka *matplotlib* za vizuelizaciju podataka i odnosa među njima.

3 Analiza i pretprocesiranje podataka

Osnovu svakog projekta koji se bavi istraživanjem podataka čini dobar početni skup podataka. Podatke je potrebno analizirati, izvršiti pretprocesiranje i, na kraju, koristiti ih za pronalaženje određenih zakonitosti i donošenje konkretnih zaključaka [4].

Podaci, koji će biti korišćeni u nastavku rada sadrže informacije o učesnicima olimpijskih igara počev od 1896. godine i o nacionalnim olimpijskim komitetima (NOK) (eng. *National Olympic Committee (NOC)*), koje su ti učesnici predstavljali¹.

3.1 Analiza podataka

Podaci su razdvojeni u dve *CSV* (*Comma-separated values*) datoteke. U jednoj datoteci se nalaze podaci o prethodnim i trenutnim imenima nacionalnih komiteta. Svaki red ove tabele čine 3 kolone - skraćenica

¹Podaci su preuzeti sa sledećeg linka: <https://www.kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results>.

komiteta, ime države (regiona) koju zastupa komitet, dodatna informacija ukoliko je komitet sa istom skraćenicom na različitim olimpijskim igrama predstavljao državu (region) različitog imena zbog izvesnih geopolitičkih promena. Prikazaćemo nekoliko redova, koji će pomoći da se stekne jasna slika o podacima iz ovog skupa podataka, u listingu 1.

```
1 noc_regions = pd.read_csv('../data/noc_regions.csv')
2 print(noc_regions.iloc[[26, 95, 147, 168, 227]])
3
4 Out>>
5      NOC      region      notes
6 26  BOH    Czech Republic    Bohemia
7 95  IRL      Ireland         NaN
8 147 NFL      Canada      Newfoundland
9 168 ROT         NaN    Refugee Olympic Team
10 227 YUG      Serbia      Yugoslavia
```

Listing 1: Podaci iz skupa podataka o regionima

Proverićemo da li, osim olimpijskog tima izbeglica, postoji još redova sa nedostajućim vrednostima atributa *region* u listingu 2. U nastavku rada ćemo, prilikom spajanja sa drugom datotekom u listingu 4, razrešiti sve nedostajuće vrednosti.

```
1 noc_regions = pd.read_csv('../data/noc_regions.csv')
2 print("Redovi sa nedostajucim vrednostima za atribut 'region':")
3 print(noc_regions[noc_regions['region'].isna()])
4
5 Out>>
6 Redovi sa nedostajucim vrednostima za atribut 'region':
7      NOC region      notes
8 93  IOA     NaN    Individual Olympic Athletes
9 168 ROT     NaN    Refugee Olympic Team
10 208 TUV     NaN         Tuvalu
11 213 UNK     NaN         Unknown
```

Listing 2: Pronalaženje nedostajućih vrednosti za atribut *region*

Što se tiče trećeg reda iz rezultata, došlo je do greške i potrebno je upisati ime države Tuvalu u polje *region*. Na osnovu poslednjeg reda iz rezultata vidimo da za neke učesnike nije poznato odakle dolaze. Kod preostala dva reda, u pitanju su timovi koji su formirani kako bi se omogućio nastup učesnicima, koji to nisu mogli da učine pod zastavom postojećeg olimpijskog komiteta.

Timu nezavisnih učesnika (eng. *Individual Olympic Athletes*) su pripadali učesnici čiji je nacionalni olimpijski komitet bio suspendovan iz političkih razloga. Ovi timovi su nastupali pod zastavom olimpijade. Prvi put je formiran ovaj tim 1992. godine i činili su ga takmičari iz Savezne Republike Jugoslavije i Severne Makedonije². Olimpijski tim izbeglica (eng. *Refugee Olympic Team*) je formiran 2016. kada su pojedini takmičari tokom velike izbegličke krize ostali bez mogućnosti da se takmiče pod zastavama svojih olimpijskih komiteta. Ovo je učinjeno i da bi se podigla svest o ovom velikom problemu današnjice.

Zarad lakše analize, vrednosti iz kolone *notes* za ova 3 reda ćemo ubaciti u polje *region*. U poslednjoj koloni ćemo vrednost postaviti na string *None* na mesto svih nedostajućih vrednosti. Razlog za veliki broj

²Na ovom takmičenju je tim nezavisnih učesnika osvojio 3 medalje u streljaštvu. Jasna Šekarić je osvojila srebrnu medalju, a Aranka Binder i Stevan Pletikosić bronzanu.

nedostajućih vrednosti u poslednjoj koloni je to što je većina komiteta, od svog osnivanja, predstavljala entitet, koji se ni na koji način nije menjao.

Druga tabela predstavlja učesnike olimpijskih igara i njihove nastupe na olimpijskim igrama. Sadrži 15 atributa (kolona) i 271.116 redova (listing 5). Svaki red predstavlja učešće jednog sportiste na jednom događaju i njegov uspeh, tj. odličje koje je osvojio, ukoliko je osvojio. Jedna olimpijska disciplina može imati različite događaje, koji se razlikuju po određenim parametrima (dužina staze kod trčanja, visina skakaonice kod skokova u vodu, itd). Učesnik koji svoju državu predstavlja na više događaja, na istoj olimpijadi, nalazi se u više redova, od kojih svaki red odgovara jednom događaju. Atribute ćemo izlistati u listingu 3, a zatim prikazati značenje svakog atributa u tabeli 1.

```
1 print(data.columns.tolist())
2
3 Out>> ['ID', 'Name', 'Sex', 'Age', 'Height', 'Weight',
4         'Team', 'NOC', 'Games', 'Year', 'Season', 'City',
5         'Sport', 'Event', 'Medal']
```

Listing 3: Ispis atributa

Ime atributa	Značenje
ID	ID učesnika
Name	Ime i prezime učesnika
Sex	Pol učesnika
Age	Starost učesnika
Height	Visina učesnika
Weight	Težina učesnika
Team	Tim učesnika ³
NOC	NOK pod čijom zastavom učesnik nastupa
Games	Godina i sezona olimpijade (npr: <i>1984 Winter</i>)
Year	Godina održavanja olimpijade
Season	Sezona održavanja olimpijade (letnja/zimska)
City	Grad u kome se održala olimpijada
Sport	Naziv sporta
Event	Događaj (npr. <i>Trčanje 200m sprint</i>)
Medal	Osvojena medalja (<i>NA</i> - bez medalje)

Tabela 1: Nazivi atributa i njihovo značenje

Implementiraćemo funkciju u listingu 4, koju ćemo pozivati za učitavanje podataka. Ovo će nam olakšati učitavanje podataka tokom analize i testiranja, jer će se ista procedura učitavanja podataka, sa različitim parametrima, koristiti više puta.

```
1 def get_data(fields=None,
2             noc=None,
3             merge_noc=True,
4             printable=False):
5
6     if fields is None:
7         # ukoliko ne izaberemo attribute, podrazumevano
8         # ce izvuci kompletne redove
```

³Jedan NOK može kod nekih sportova imati više timova na jednom događaju

```

9     athlete_events =
10         pd.read_csv('../data/athlete_events.csv',
11                     error_bad_lines=False)
12     else:
13         if noc is not None:
14             # ukoliko je izabran NOK po kome zelimo
15             # filtrirati potrebno je dodati i taj
16             # atribut u spisak atributa
17             fields.append('NOC')
18
19         athlete_events =
20             pd.read_csv('../data/athlete_events.csv',
21                         error_bad_lines=False,
22                         usecols=fields)
23
24     df = pd.DataFrame(athlete_events)
25
26     if noc is not None:
27         nationality_filter = athlete_events['NOC'] == noc
28         df = df[nationality_filter]
29
30     if merge_noc is True:
31         noc_regions =
32             pd.read_csv('../data/noc_regions.csv')
33         df = pd.merge(df, noc_regions,
34                     on='NOC', how='left')
35
36         df['region'].fillna(noc_regions['notes'],
37                             inplace=True)
38         df['notes'].fillna('None', inplace=True)
39         df['region'].fillna('None', inplace=True)
40
41     # ispisivanje dobijenih podataka ukoliko je trazeno
42     if printable is True:
43         pd.set_option('display.max_rows', None)
44         pd.set_option('expand_frame_repr', False)
45         print(df)
46
47     return df

```

Listing 4: Funkcija za učitavanje podataka

Primetimo da smo u funkciji izvršili spajanje ove dve tabele po atributu koji označava nacionalni olimpijski komitet, kako bismo redu pridružili atribut, koji označava naziv države (regiona, organizacije) za koju učesnik nastupa. Učit ćemo podatke korišćenjem implementirane funkcije, a zatim broj atributa i redova novoformirane tabele izračunati u listingu 5.

```

1 data = get_data(printable=False)
2 print(data.shape)
3
4 Out>> (271116, 17)

```

Listing 5: Učitavanje podataka i ispis dimenzije

U listingu 6 prikazimo broj takmičara za svaki od olimpijskih komiteta, koji su na početku imali nedostajuće vrednosti u polju *region*.

```

1 print("IOA count: ", data[data['NOC'] == 'IOA'].shape[0])
2 print("ROT count: ", data[data['NOC'] == 'ROT'].shape[0])
3 print("TUV count: ", data[data['NOC'] == 'TUV'].shape[0])
4 print("UNK count: ", data[data['NOC'] == 'UNK'].shape[0])

```

```

5
6 Out>>
7 IOA count: 76
8 ROT count: 10
9 TUV count: 6
10 UNK count: 2

```

Listing 6: Broj takmičara za komitete sa početnim nedostajućim vrednostima polja *region*

Pogledajmo osnovne statističke informacije učitanih podataka, koristeći funkciju *describe* iz biblioteke *pandas*. Podatke ćemo prikazati u listingu 7 uz zaokruživanje na 2 decimale.

```

1 print(data.describe().round(2))
2
3 Out>>
4
5 count      271116.00  261642.00  210945.00  208241.00  271116.00
6 mean        68248.95    25.56    175.34    70.70    1978.38
7 std         39022.29     6.39    10.52    14.35    29.88
8 min           1.00    10.00    127.00    25.00    1896.00
9 25%         34643.00    21.00    168.00    60.00    1960.00
10 50%         68205.00    24.00    175.00    70.00    1988.00
11 75%        102097.25    28.00    183.00    79.00    2002.00
12 max        135571.00    97.00    226.00    214.00    2016.00

```

Listing 7: Statističke informacije učitanih podataka

Dakle, prosek godina je 25.56, prosečna visina 175.34m, a težina 70.70kg itd. Možemo primetiti dve zanimljive ekstremne vrednosti - najstariji učesnik u istoriji olimpijskih igara imao je 97 godina, a najmlađi učesnik 10. Takođe, iz ovih podataka možemo zaključiti koliko je redova sa nedostajućim vrednostima po svakoj koloni od navedenih. Na primer, nedostaje nam informacija o starosti 9474 (271116 – 261642) učesnika. Broj redova iz cele tabele, koji sadrže nedostajuće vrednosti u bilo kojoj koloni, dobijamo u listingu 8.

```

1 null_rows = data[data.isna().any(axis=1)]
2 print(null_rows.shape[0])
3
4 Out>> 240935

```

Listing 8: Broj redova sa nedostajućim vrednostima

Nakon što smo pronašli broj redova sa nedostajućim vrednostima potrebno je utvrditi u kojim kolonama se nedostajuće vrednosti javljaju:

```

1 print(data.columns[data.isna().any()].tolist())
2
3 Out>> ['Age', 'Height', 'Weight', 'Medal']

```

Listing 9: Kolone koje sadrže nedostajuće vrednosti

Odokativnom analizom možemo doći do zaključka da najveći broj nedostajućih vrednosti dolazi iz činjenice da je vrednost u koloni koja predstavlja osvojenu medalju, učesnicima koji nisu osvojili nijedno odličje na događaju, označena sa *NA*. U polje koje označava osvojenu medalju, učesnicima koji nisu osvojili nijedno odličje na navedenom događaju, izmeničemo u string *None* (listing 10).

```

1 data['Medal'].fillna('None', inplace=True)

```

Listing 10: Ažuriranje polja koja označavaju osvojenu medalju

Kao što smo primetili, najstariji učesnik olimpijade imao je 97 godina. Proverimo ispravnost ovog podatka, pritom smanjujući starosnu granicu na 85 godina (listing 11).

```
1 age_filter = data['Age'] > 85
2 data = data[age_filter]
3 print(data[['Name', 'Age', 'Games', 'Sport']])
4
5 Out>>
6
7      Name      Age      Games      Sport
8 60861 Thomas Cowperthwait Eakins 88.0 1932 Summer Art Competitions
9 60862 Thomas Cowperthwait Eakins 88.0 1932 Summer Art Competitions
10 60863 Thomas Cowperthwait Eakins 88.0 1932 Summer Art Competitions
11 98118 Winslow Homer      96.0 1932 Summer Art Competitions
12 257054 John Quincy Adams Ward 97.0 1928 Summer Art Competitions
```

Listing 11: Učesnici stariji od 85 godina

Svi učesnici stariji od 85 godina su učestvovali na takmičenjima iz oblasti umetnosti⁴. Možemo pretpostaviti da su i u opštem slučaju učesnici koji su se takmičili u ovoj disciplini, među najstarijim učesnicima. Dodatnim smanjivanjem granice i analizom dobijenih podataka, potvrđujemo ovu pretpostavku. Discipline koje su imale najstarije učesnike, pored umetnosti, su streljaštvo, streličarstvo, jahanje, jedrenje (ovo će biti pokazano u radu). Broj učesnika iz oblasti umetnosti ćemo izračunati i prikazati u listingu 12.

```
1 sport_filter = data['Sport'] == 'Art Competitions'
2 data = data[sport_filter].drop_duplicates(
3     subset=['ID'], keep='first')
4 num_of_artists = data.shape
5 print("Number of art competitors: ", num_of_artists[0])
6
7 Out>> Number of art competitors: 1814
```

Listing 12: Broj učesnika u umetničkim takmičenjima

Sličnim postupkom možemo pronaći i podatke o najmlađem učesniku u istoriji olimpijskih igara. U pitanju je Dimitrios Loundras, koji je kao najmlađi učesnik u istoriji došao i do bronzane medalje (listing 13).

```
1 age_filter = data['Age'] == 10
2 data = data[age_filter]
3 print(data[['Name', 'Age', 'Games', 'Sport', 'Medal']])
4
5 Out>>
6
7      Name      Age      Games      Sport      Medal
8 142882 Dimitrios Loundras 10.0 1896 Summer Gymnastics Bronze
```

Listing 13: Najmlađi učesnik u istoriji

Kao što smo videli pomoću funkcije iz listinga 8, u početnom skupu podataka ima 240.935 redova sa nedostajućim vrednostima. Ove vrednosti ne smemo odbaciti, već ih je neophodno ispitati i obraditi na odgovarajući način. Nakon što smo izmenili vrednosti atributa *Medal* sa *NA* u string *None*, preostao je ukupno 64.951 red sa nedostajućim vrednostima, što možemo proveriti korišćenjem istog postupka iz listinga 8. Redove kod kojih atribut *Age* ima nedostajuće vrednosti, to polje ćemo zameniti srednjom vrednosti starosti svih takmičara za koje imamo ovu informaciju.

⁴Iako je malo poznato, takmičenja iz oblasti umetnosti su bila deo olimpijskih igara od 1912. do 1948. godine. U takmičenje su ulazila dela inspirisana sportom, a grane umetnosti koje su bile zastupljene bile su arhitektura, slikarstvo, književnost, vajarstvo i muzika. Ova olimpijska disciplina je ukinuta, jer je smatrano da su svi umetnici koji su se prijavljivali, profesionalci i da ne bi trebalo da učestvuju na olimpijadi. Osnovna ideja olimpijskih igara je bila da se nadmeću amateri. Sredinom dvadesetog veka, modernizacija društva i samog sporta, kao i političke promene, uticale su na promenu stava međunarodnog olimpijskog komiteta (eng. *International Olympic Committee (IOC)*) i počelo se sa omogućavanjem učešća profesionalcima [3, 2]

Za attribute *Height* i *Weight* ćemo koristiti modu svih redova za svaki od atributa. Postupak i rezultat njegovog izvršavanja dat je u listingu 14. Možemo primetiti da smo se uspešno rešili svih nedostajućih vrednosti iz polaznog skupa podataka.

```

1 def fill_na_attr_rows(data):
2     data['Medal'].fillna('None', inplace=True)
3
4     age_mean = int(data['Age'].mean())
5     data['Age'] = data['Age'].fillna(age_mean)
6     data['Age'].isna().sum()
7
8     height_mode = data['Height'].mode()[0]
9     data['Height'] = data['Height'].fillna(height_mode)
10    data['Height'].isna().sum()
11
12    weight_mode = data['Weight'].mode()[0]
13    data['Weight'] = data['Weight'].fillna(weight_mode)
14    data['Weight'].isna().sum()
15
16    return data
17
18
19 data = fill_na_attr_rows(data)
20 # proveravamo da li postoji jos redova sa nedostajucim
21 # vrednostima
22 print(data[data.isna().any(axis=1)].shape[0])
23
24 Out>> 0

```

Listing 14: Izmena nedostajućih vrednosti srednjim

Nakon što smo otklonili sve nedostajuće vrednosti iz početnog skupa podataka, novi skup podataka ćemo sačuvati u novu datoteku. Postupak je prikazan u listingu 15.

```

1 data.to_csv(r'../data/full.csv')

```

Listing 15: Čuvanje ažuriranih podataka u novi dokument

3.2 Dodatna analiza i vizualizacija podataka

Kako bismo još bolje razumeli podatke, dodatno ćemo ih analizirati i vizualizovaćemo ih korišćenjem biblioteke *matplotlib*. Učitamo podatke iz novoformirane datoteke i prikazaćemo neke zavisnosti njihovih atributa. Prvo što ćemo prikazati je broj učesnika na svakoj od letnjih olimpijada kroz istoriju. Isto ćemo učiniti odvojeno za muške, odnosno ženske učesnike. U listingu 16 je prikazan postupak izračunavanja za sve takmičare. Dodavanjem filtera za pol, možemo izračunati broj ženskih, odnosno muških takmičara na svakoj olimpijadi (zbog sličnosti sa prethodnim, u ova dva slučaja neće biti navođen listing već samo plot). Na slici 1 prikazan je broj svih učesnika na svakoj od olimpijada kroz istoriju. Na slici 2 prikazan je broj muških, a na slici 3 broj ženskih učesnika kroz istoriju.

```

1 # učitavanje podataka iz nove datoteke
2 data = pd.read_csv('../data/full.csv')
3
4 # izdvajanje redova sa letnjih olimpijada sa cuvanjem
5 # jednog reda za svakog takmicara
6 filtered_data = data[(data['Season'] == 'Summer')].
7     drop_duplicates(['Year', 'ID'])
8 years = filtered_data[['Year']]
9 # broj pojavljivanja svake godine odrzavanja i

```

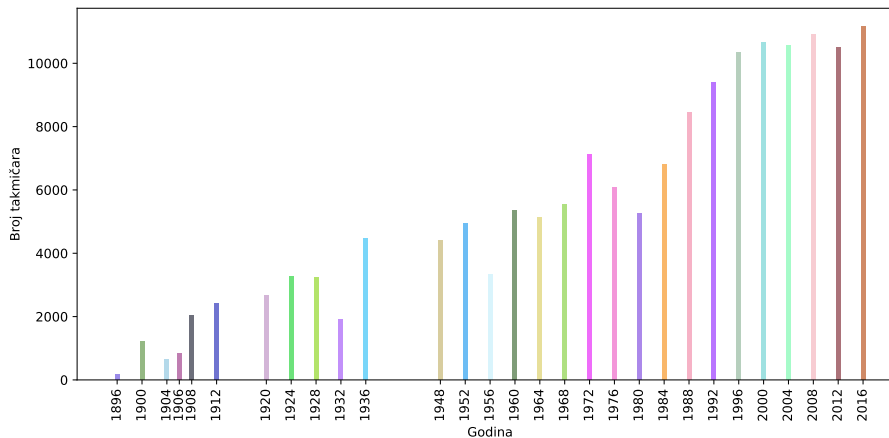


```

10 # sortiranje po godinama
11 contestants_count = pd.Series([x for x in years['Year']])
12     .value_counts()
13 contestants_count = contestants_count.sort_index()
14
15 # kreiranje DataFrame-a od prethodno izracunatih
16 # vrednosti
17 contestants_count_df =
18     pd.DataFrame(contestants_count,
19         columns=['Count']).reset_index()
20
21 # godine na x-osu, broj ucesnika na y-osu
22 x = contestants_count_df['index']
23 y = contestants_count_df['Count']
24
25 f = plt.figure(figsize=(10, 5))
26 plt.plot(x, y, linestyle='None')
27 # napravimo listu odgovarajuće dužine sa nasumično
28 # napravljenim bojama u heksadekadnom formatu
29 number_of_colors = contestants_count_df.count(axis=0)[0]
30 colors = ["#" + ''.join([random.choice('0123456789ABCDEF')
31     for j in range(6)])
32     for i in range(number_of_colors)]
33
34 plt.bar(x, y, alpha=0.6, align='center', color=colors,
35     linewidth=10)
36 plt.xlabel("Godina")
37 plt.ylabel("Broj takmicara")
38 plt.xticks(x)
39 plt.tick_params(axis='x', which='major', rotation=90)
40 plt.show()

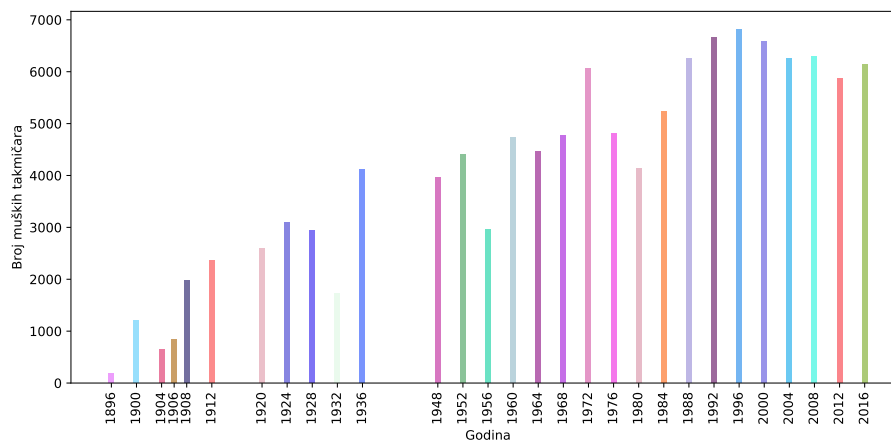
```

Listing 16: Učitavanje podataka i računanje broja učesnika za svaku olimpijadu

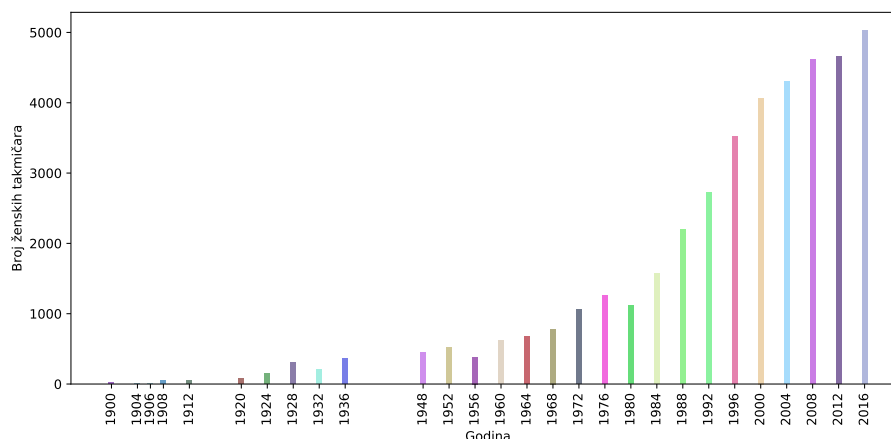


Slika 1: Broj učesnika kroz istoriju

Iz prikazanih plotova vidimo da je broj učesnika na Letnjim olimpijskim igrama u skoro neprestanom porastu, sa izuzetkom nekoliko olimpijada kada je broj učesnika opadao. Možemo videti i da na prvoj olimpijadi nije bilo ženskih takmičara. Decenijama nakon toga, broj ženskih takmičara je bio neuporedivo manji od kolega suprotnog pola. Tek krajem



Slika 2: Broj muških učesnika kroz istoriju



Slika 3: Broj ženskih učesnika kroz istoriju

sedamdesetih i početkom osamdesetih dolazi do povećanog rasta učesnica i popularizacije ženskog sporta na Olimpijskim igrama.

Na plotovima još vidimo da nekoliko olimpijada "nedostaje". U pitanju su godine u kojima su Olimpijske igre bile otkazivane zbog Prvog svetskog rata, a kasnije i Drugog svetskog rata (1916, 1940, 1944). Poslednja letnja olimpijada pre Drugog svetskog rata je održana u Trećem rajhu (Nacistička Nemačka). Još jedna godina koja ispada iz obrasca održavanja Olimpijskih igara na svake četiri godine je 1906. godina. Tada su održane prve i jedine Olimpijske međugre (eng. *1906 Intercalated Games*), u Atini. Medalje i uspehe ostvarene na ovoj olimpijadi Međunarodni olimpijski komitet ne ubraja u zvanične [2]. Iz ovog razloga, ovu godinu ćemo obrisati iz našeg skupa podataka, postupkom prikazanim u listingu 17.

```
1 data = data[data['Year'] != 1906]
```

Listing 17: Brisanje rezultata iz nezvanične godine

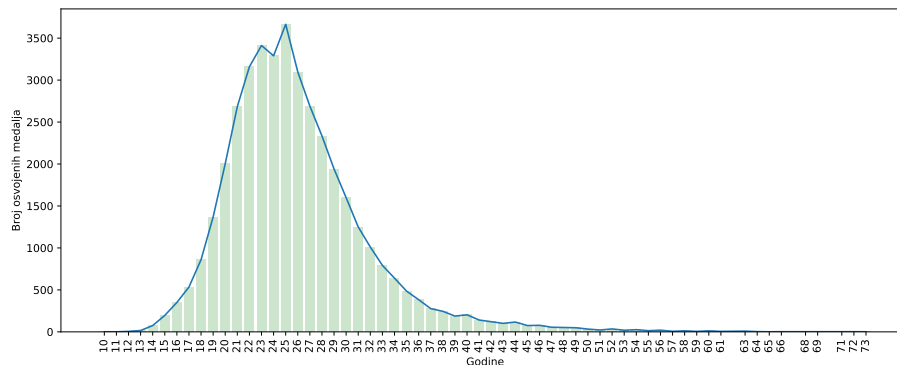
Sledeća stvar koju ćemo ispitati je broj osvojenih medalja za sve uzraste takmičara da bismo ispitati zavisnost starosti učesnika i osvajanja medalje. Postupak je prikazan u listingu 18, a plot na slici 4.

```

1 # izdvajanje redova sa ucesnicima, koji su osvojili neku
2 # medalju i izdvajanje samo kolona 'Age' i 'Medal'
3 medals = data[(data['Medal'] != 'None')]
4 medals = medals[['Age', 'Medal']]
5 medal_count = pd.Series([x for x in medals['Age']]).
6     value_counts()
7 medal_count = medal_count.sort_index()
8
9 medal_df = pd.DataFrame(medal_count, columns=['Count']).
10     reset_index()
11
12 x = medal_df['index']
13 y = medal_df['Count']
14 plt.figure(figsize=(12, 5))
15 plt.plot(x, y)
16 plt.bar(x, y, alpha=0.2, align='center', color='green',
17     linewidth=10)
18 plt.xlabel("Godine")
19 plt.ylabel("Broj osvojenih medalja")
20 plt.xticks(x)
21 plt.tick_params(axis='x', which='major', rotation=90)
22 plt.show()

```

Listing 18: Učitavanje podataka i računanje medalja za sve uzraste



Slika 4: Broj osvojenih medalja po uzrastima takmičara

Vidimo da su medalje osvajali učesnici starosti od 10 do 73 godine, sa izuzetkom uzrasta od 62, 67 i 70 godina. Možemo primetiti i da su medalje najčešće osvajali učesnici starosti između 23 i 25 godina. Takođe, vidimo da je nezanemarljiv broj takmičara starijih od 50 godina osvojio medalju. Izračunajmo broj medalja za svaki sport za koji je to slučaj, na sličan način kao u prethodnoj analizi (listing 19) i prikazimo plot na slici 5.

```

1 older_medalists_sports = data[(data['Medal'] != 'None') &
2     (data['Age'] > 50)]
3 older_medalists_sports = older_medalists_sports[['Sport',
4     ]]
5 sports_medal_count = pd.Series([x for x in
6     older_medalists_sports['Sport']]).value_counts()

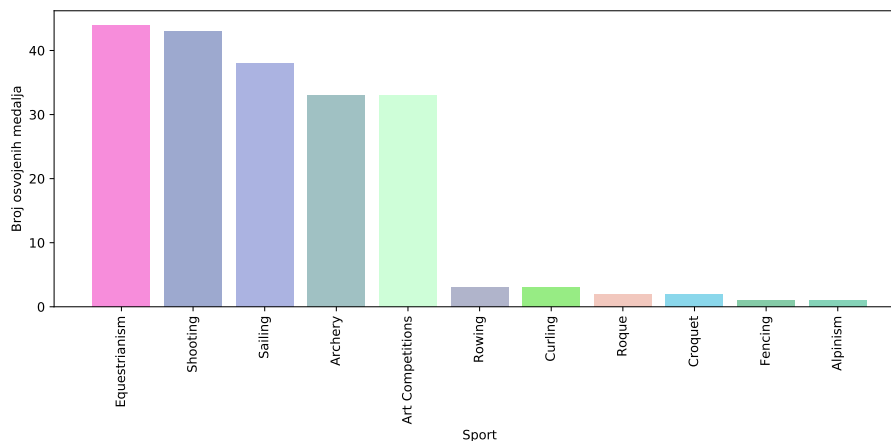
```

```

5 sports_medal_df = pd.DataFrame(sports_medal_count,
6                                columns=['Count']).reset_index()
7
8 x = sports_medal_df['index']
9 y = sports_medal_df['Count']
10
11 f = plt.figure(figsize=(10, 5))
12 plt.plot(x, y, linestyle='None')
13
14 number_of_colors = sports_medal_df.count(axis=0)[0]
15 colors = ["#" + ''.join([random.choice('0123456789ABCDEF')
16                             for j in range(6)])
17           for i in range(number_of_colors)]
18
19 plt.bar(x, y, alpha=0.5, align='center', color=colors,
20         linewidth=10)
21 plt.xlabel("Sport")
22 plt.ylabel("Broj osvojenih medalja")
23 plt.xticks(x)
24 plt.tick_params(axis='x', which='major', rotation=90)
25 plt.show()

```

Listing 19: Računanje broja medalja po sportu za takmičare starije od 50 godina



Slika 5: Broj osvojenih medalja za svaki sport za koji je takmičar stariji od 50 godina osvojio neko od odličja

Sa plotu na slici 5 vidimo da se izdvaja pet sportova u kojima su takmičari stariji od 50 godina posebno uspešni. U pitanju su, po zastupljenosti medalja: jahanje, streljaštvo, jedriličarstvo, streličarstvo, umetnost.

3.3 Izvođenje novih atributa

Podatke, koje smo do sada obradili, ćemo iskoristiti za kreiranje nove tabele u kojoj će se nalaziti atributi, koji su izvedeni iz postojećih podataka. Novoformiranu tabelu ćemo kasnije koristiti prilikom klasterovanja.

Postupak formiranja nove tabele će biti prikazan u listingu 20. Nakon učitavanja podataka, u tabelu sačinjenu od početnih podataka, dodaćemo novo polje, koje će u slučaju da je medalja osvojena imati vrednost 1, a

inače 0. Nazvaćemo ga *MedalWon*. Grupišimo redove po atributu *region* i izračunajmo ukupan broj medalja na svim olimpijadama tih zemalja. Ovu kolonu ćemo nazvati *MedalsCount*. Dalje, izračunaćemo ukupan broj učesnika na svim olimpijadama svake države na sličan način. Ovaj atribut ćemo nazvati *ParticipantsCount*. Sledeći atribut koji ćemo izračunati ćemo nazvati *GamesCount* i on će sadržati broj igara na kojima je država učestvovala. Atribut *PartPerGamesScaled* će sadržati broj svih učesnika podeljen sa brojem igara na kojima je država učestvovala, skaliran na 1 (1 odgovara državi koja je imala najviše učesnika). Atribut *SuccessRate* određuje procenat uspešnosti države, skaliran na 1 (1 odgovara državi koja ima najveći procenat uspešnosti). Izračunaćemo ga tako što ćemo broj svih osvojenih medalja podeliti sa brojem učesnika.

```

1 # učitavanje podataka i dodavanje indikatorske kolone -
2 # da li je medalja osvojena ili ne
3 data = pd.read_csv('../data/full.csv')
4 data['MedalWon'] = np.where(
5     data.loc[:, 'Medal'] == 'None',
6     0, 1)
7
8 # izracunavanje broja osvojenih medalja za svaku drzavu
9 # ucesnicu grupisanjem po atributu 'region' i
10 # sumiranjem prethodno napravljene kolone
11 country_medals_won_athletes =
12     data.groupby(['region'])['MedalWon'].
13     agg('sum').reset_index()
14 country_medals_won_athletes =
15     country_medals_won_athletes.sort_values(
16         'MedalWon', ascending=False)
17 country_medals_won_athletes.rename(
18     columns={'MedalWon': 'MedalsCount'},
19     inplace=True)
20
21 # izracunavanje ukupnog broja predstavnika na svim
22 # olimpijadama za svaku drzavu
23 athletes_count = data.drop_duplicates(['ID', 'Year'])
24 athletes_count = athletes_count.groupby(
25     ['region'])['ID'].count().reset_index()
26 athletes_count.rename(
27     columns={'ID': 'ParticipantsCount'}, inplace=True)
28
29 # izracunavanje ukupnog broja olimpijada na kojima
30 # je svaka drzava ucestvovala
31 games_count = data.drop_duplicates(['region', 'Games'])
32 games_count = games_count.groupby(['region'])['Year'].
33     count().reset_index()
34 games_count.rename(columns={'Year': 'GamesCount'},
35     inplace=True)
36
37 # spajanje prethodno izracunatih vrednosti u DataFrame,
38 # po atributu 'region'
39 country_data = pd.DataFrame(country_medals_won_athletes)
40 country_data = pd.merge(country_data,
41     athletes_count,
42     on='region')
43 country_data = pd.merge(country_data,
44     games_count,
45     on='region')
46
47 # racunanje broja ucesnika po igrima i efikasnosti
48 # drzave, i njihovo skaliranje na 1

```

```

48 country_data['PartPerGamesScaled'] =
49     country_data['ParticipantsCount'] /
50     country_data['GamesCount']
51 country_data['SuccessRate'] =
52     country_data['MedalsCount'] /
53     country_data['ParticipantsCount']
54
55 x = country_data[['PartPerGamesScaled', 'SuccessRate']].
56     values
57 min_max_scaler = preprocessing.MinMaxScaler()
58 x_scaled = min_max_scaler.fit_transform(x)
59 country_data[['PartPerGamesScaled', 'SuccessRate']] =
60     x_scaled
61
62 # ispisivanje prvih 10 redova novoformirane tabele
63 # MC      - MedalsCount
64 # PC      - ParticipantsCount
65 # GC      - GamesCount
66 # PPGS    - PartPerGamesScaled
67 # SR      - SuccessRate
68 print(country_data.head(10))
69
70 Out>>
71      region      MC      PC      GC      PPGS      SR
72 0      USA    5637   12851    50    1.000000    0.900179
73 1    Russia    3947    8100    36    0.874540    1.000000
74 2   Germany    3756   10653    46    0.900349    0.723555
75 3       UK    2068    8489    51    0.645134    0.499933
76 4   France    1777    8280    51    0.629077    0.440428
77 5     Italy    1637    7226    51    0.548101    0.464910
78 6   Sweden    1536    5414    50    0.417209    0.582225
79 7   Canada    1352    6464    49    0.509828    0.429233
80 8 Australia    1349    5534    48    0.444682    0.500254
81 9   Hungary    1135    4109    49    0.321515    0.566862

```

Listing 20: Kreiranje nove tabele od postojećih atributa

Korišćenjem postupka prikazanom u listingu 21 ćemo napraviti novu *CSV* datoteku, koju ćemo koristiti prilikom klasterovanja.

```

1 data.to_csv(r'../data/country_stats.csv')

```

Listing 21: Čuvanje novoformirane tabele u *CSV* dokument

4 Klasterovanje

U prethodnom poglavlju smo se rešili nedostajućih vrednosti, dodatno se upoznali sa podacima i vizualizovali ih, zarad sticanja još bolje slike o njihovim osobinama. Takođe, u poglavlju 3.3 smo podatke izdvojili u novu tabelu, koju ćemo u nastavku koristiti. U ovom poglavlju ćemo se baviti klasterovanjem (klaster analizom). Cilj klasterovanja je pronalaženje grupa objekata sličnih karakteristika, odnosno pronalaženje grupa objekata, koje se razlikuju u odnosu na sve druge grupe objekata [4].

Postoji više tehnika klasterovanja. U ovom radu koristićemo sledeće tehnike:

- **K-sredine** - Predstavlja tehniku, koja se zasniva na centroidama i particionisanju elemenata (tačaka), tako da te tačke pripadaju klasteru čijoj centroidi su najbliži. Slovo *K* u imenu K-sredine predstavlja broj klastera, odnosno broj centroida, na koliko želimo razdvojiti

tačke. Centroide predstavljaju centralne tačke klastera, koje su na početku izabrane nasumično. U svakom koraku primene algoritma za određivanje K klastera, za svaku tačku se određuje kom klasteru pripada na osnovu udaljenosti od centroide. Zatim se računa nova centroida za svaki klaster na osnovu vrednosti tačaka, koje pripadaju ažuriranom klasteru. Sa primenom se staje nakon koraka u kojem nema promena centroida. Bliskost dveju tačaka se računa korišćenjem Euklidskog rastojanja, kosinusnog rastojanja ili drugih mera bliskosti [4].

- **Hijerarhijsko klasterovanje** - Postoje dva glavna tipa hijerarhijskog klasterovanja. Kod prvog, početni broj klastera u ovoj metodi odgovara broju tačaka za koje se klasterovanje vrši. U svakom koraku se dva najbliža klastera spajaju, sve dok ne nastane jedan, sveobuhvatni klaster. Ovaj tip, koji ćemo i koristiti u ovom radu, se naziva klasterovanje spajanjem (eng. *agglomerative clustering*). Drugi pristup polazi od sveobuhvatnog klastera i u svakom koraku se ovaj klaster deli dok ne ostanu samo pojedinačne tačke kao klasteri ili dok ne ostanu K klastera. Ovaj tip se naziva klasterovanje razdvajanjem (eng. *divisive*). Podaci se vizuelizuju korišćenjem dendrograma⁵ [4].

4.1 K-sredine

Prvo ćemo primeniti tehniku K-sredina. Da bismo odredili optimalan broj klastera koristimo *elbow* metod. Rezultat primene ovog metoda je kriva (eng. *elbow curve*) sa koje možemo pročitati optimalan broj klastera. Tu vrednost čitamo na y -osi na mestu gde se kriva prelama i prestaje sa rastom. Ta stagnacija znači da nam veći broj klastera ne bi značajno doprineo u raspoređivanju podataka u klastere [4].

Pokušaćemo da klasterujemo podatke korišćenjem nekoliko kolona. Prvo ćemo to učiniti korišćenjem novoformiranih atributa *MedalsCount* i *ParticipantsCount*. Za početak, korišćenjem gore navedenog *elbow* metoda pronađimo optimalan broj klastera. Postupak određivanja vrednosti, koje će formirati *elbow* krivu prikazan je u listingu 22. Rezultat izvršavanja ovog listinga prikazan je na slici 6

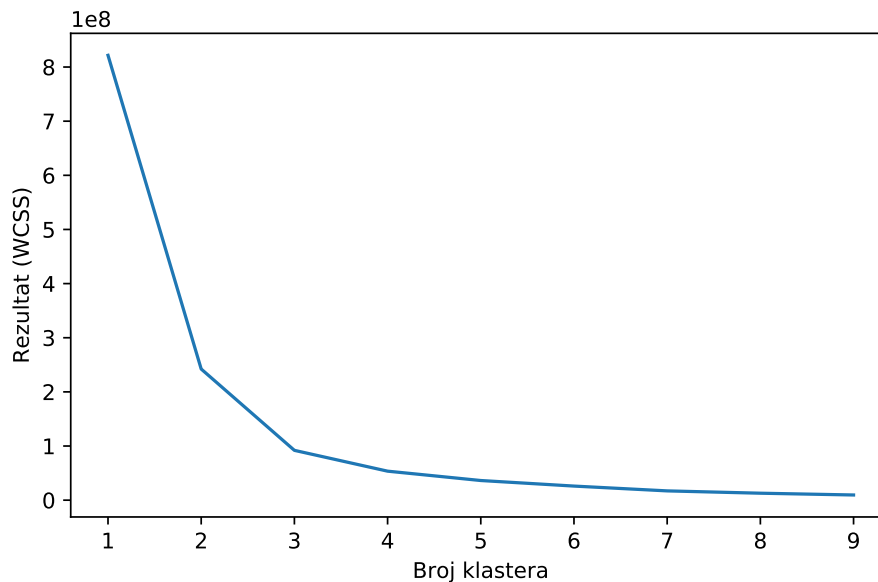
```

1 # biramo podskup podataka, koji zelimo iskoristiti za
2 # klasterovanje
3 cluster_data = country_data[['ParticipantsCount',
4                               'MedalsCount']]
5
6 # biramo raspon klastera za koji zelimo proveravati
7 # kvalitet i koristimo metod K-sredina za svaku
8 # od opcija za K
9 ec_cluster_options = range(1, 10)
10 kmeans_res = [KMeans(n_clusters=i, init='k-means++')
11               for i in ec_cluster_options]
12 # wcss - within-cluster sum of square
13 wcss = [kmeans_res[i].fit(cluster_data).inertia_
14         for i in range(len(kmeans_res))]
15
16 f_ec = plt.figure(figsize=(6, 4))
17 plt.plot(range(1, 10), wcss)
18 plt.xlabel('Broj klastera')
19 plt.ylabel('Rezultat (WCSS)')
```

⁵Dendrogram predstavlja dijagram, koji pokazuje hijerarhijski odnos između objekata.

```
plt.show()
```

Listing 22: Postupak određivanja *elbow* krive



Slika 6: *Elbow* kriva za određivanje broja klastera za atribut *MedalsCount* i *ParticipantsCount*

Sa plota *elbow* krive vidimo da vrednost počinje da minimalno opada nakon broja 5. To je broj koji ćemo izabrati za broj klastera, koji ćemo koristiti za klasterovanje. U listingu 23 prikazan je postupak, kojim će se dobiti plot na kojem možemo videti dobijene klastere. Plot je prikazan na slici 7. Tačke obojene u jarko crveno, narandžasto, plavo, braon i zeleno predstavljaju centroeide.

```
1 # biramo podskup podataka za klasterovanje
2 cluster_data = country_data[['ParticipantsCount',
3                               'MedalsCount']]
4
5 kmeans = KMeans(n_clusters=5, init='k-means++')
6 y_kmeans = kmeans.fit_predict(cluster_data)
7
8 centroids = kmeans.cluster_centers_
9
10 # izracunate klastere dodajemo u tabelu kao atribut
11 # 'Cluster'
12 cluster_data['Cluster'] = y_kmeans
13 # srednje vrednosti za oba atributa za svaki klaster
14 kmeans_mean_cluster =
15     pd.DataFrame(round(cluster_data.groupby('Cluster').
16                        mean(),
17                        1))
18
19 f_clusters = plt.figure(figsize=(6, 4))
20 plt.xlabel('Broj ucesnika')
21 plt.ylabel('Broj medalja')
22 plt.scatter(cluster_data['ParticipantsCount'],
```

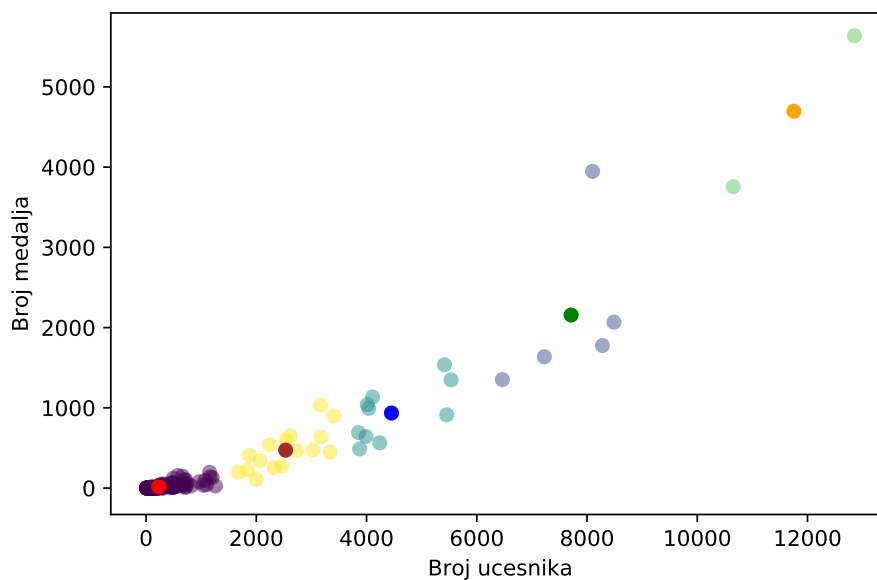


```

22     cluster_data['MedalsCount'], c=kmeans.labels_
      .astype(float), alpha=0.5)
23
24 plt.scatter(centroids[:, 0], centroids[:, 1],
25             c=['red', 'green', 'blue',
26               'orange', 'brown'])
27 plt.show()
28
29 print(kmeans_mean_cluster)
30
31 Out>>
32 Cluster  ParticipantsCount  MedalsCount  CentroidColor
33 0          233.0           15.5          red
34 1       7711.8          2156.2        green
35 2       4452.2           935.5          blue
36 3      11752.0          4696.5        orange
37 4       2532.2           473.1          brown

```

Listing 23: Klasterovanje podataka po atributima *MedalsCount* i *ParticipantsCount*



Slika 7: Dobijeni klasteri za attribute *MedalsCount* i *ParticipantsCount*

Možemo primetiti da su klasteri i elementi unutar njih uglavnom linearno raspoređeni. Odatle vrlo jasno možemo zaključiti da postoji linearna veza između broja učesnika i broja osvojenih medalja. Činjenica je da što je veći broj takmičara, šanse za medalju rastu. Takođe, možemo pretpostaviti da što je jedna država uspešnija i broj njenih takmičara će na budućim olimpijadama biti u porastu. Prikažimo dve najuspešnije države iz klastera iz gornjeg desnog ugla u listingu 24. Vidimo da su najbliže po broju učesnika.

```

1 print(countries_data[
2     ['region', 'ParticipantsCount', 'MedalsCount']
3     ].sort_values(

```

```

4         'ParticipantsCount', ascending=False).
5         head(2))
6
7 Out>>
8      region  ParticipantsCount  MedalsCount
9 0      USA             12851           5637
10 2  Germany             10653           3756

```

Listing 24: Dve najuspešnije države ukoliko klasterujemo po atributima *MedalsCount* i *ParticipantsCount*

U prethodnom primeru zanemarili smo činjenicu da broj takmičara zavisi i od toga na koliko Olimpijskih igara je ta država učestvovala. Sa više učešća i broj takmičara je mogao biti veći, pa samim tim i broj medalja (kao što smo u prethodnom primeru videli, postoji linearna zavisnost ova dva atributa). U narednom primeru ćemo iskoristiti kolonu *PartPerGamesScaled*, koja će nam pomoći da ne favorizujemo države, koje su na olimpijadi učestvovala više puta od ostalih. Kolona *SuccessRate* će biti korišćena kao parametar uspeha učesnika svake od država. Postupak je analogan prethodnom, pa neće biti detaljno navođen. Dobija se *elbow* kriva prikazana na slici 8. Broj klastera će ponovo biti 5. Klasteri su prikazani na slici 9. Srednje vrednosti klastera i koji klaster odgovara kojoj boji prikazano je u listingu 25.

```

1 print(kmeans_mean_cluster)
2
3 Out>>
4 Cluster  PartPerGamesScaled  SuccessRate  CentroidColor
5 0                0.0           0.0             red
6 1                0.3           0.5             green
7 2                0.9           0.9             blue
8 3                0.1           0.2             orange
9 4                0.5           0.4             brown

```

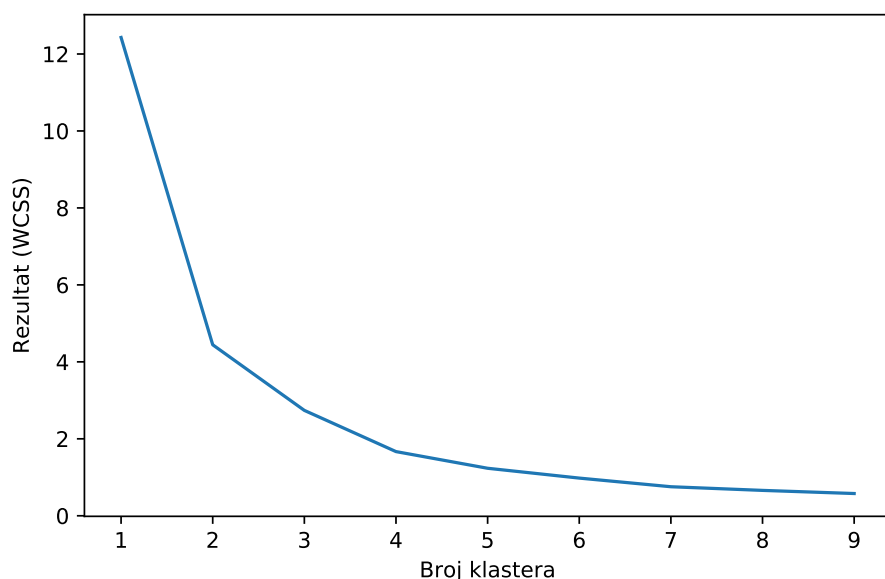
Listing 25: Srednje vrednosti klastera

Plot na slici 9 nam daje bolju predstavu o efikasnosti takmičara jedne države. Klaster u gornjem desnom uglu predstavlja 3 najuspešnije države na olimpijadi - Nemačku, Rusiju i SAD. Najveću efikasnost osvajanja odličja imaju takmičari iz Rusije. Iza njih su SAD, a zatim Nemačka, koja je imala oko 2500 učesnika više, ali i oko 200 medalja manje od Rusije. Žuti klaster sa braon centroidom predstavlja države, koje u proseku imaju slabiju efikasnost, dok svetlo plavi klaster sa zelenom centroidom predstavlja države, koje u proseku imaju dobru efikasnost. Država, koja se izdvaja po efikasnosti u klasteru sa zelenom centroidom je Norveška. Prikažimo najboljih 15 država, ukoliko se posmatra njihova efikasnost osvajanja medalja u listingu 26. Vidimo da se Srbija nalazi na 15. mestu najefikasnijih država. Ovaj uspeh treba pripisati velikim uspesima srpskih sportista u timskim sportovima, te je broj odlikovanih takmičara veći.

```

1 top_sr = country_data.sort_values(
2     ['SuccessRate'],
3     ascending=False).reset_index()
4 print(top_sr[['region', 'SuccessRate']].head(15))
5
6 Out>>
7      region  SuccessRate  MedalsCount  ParticipantsCount
8 0      Russia    1.000000           3947              8100
9 1        USA    0.900179           5637              12851
10 2    Germany    0.723555           3756              10653

```



Slika 8: *Elbow* kriva za određivanje broja klastera za atribut *PartPerGames-Scaled* i *SuccessRate*

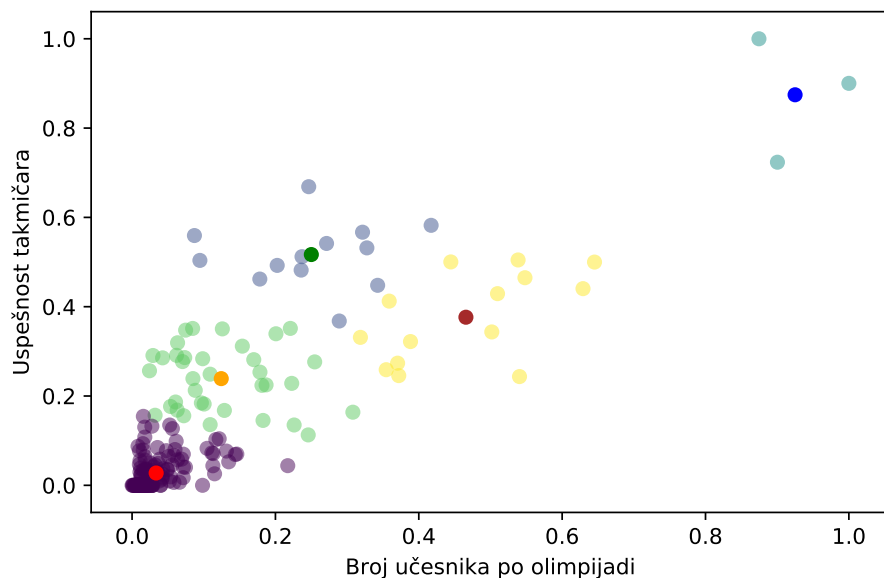
11	3	Norway	0.668743	1033	3170
12	4	Sweden	0.582225	1536	5414
13	5	Hungary	0.566862	1135	4109
14	6	Jamaica	0.559365	157	576
15	7	Finland	0.541634	900	3410
16	8	Netherlands	0.531709	1040	4014
17	9	Romania	0.512068	653	2617
18	10	China	0.504662	993	4038
19	11	Pakistan	0.503682	121	493
20	12	Australia	0.500254	1349	5534
21	13	UK	0.499933	2068	8489
22	14	Serbia	0.492489	539	2246

Listing 26: 15 najefikasnih država na olimpijadi

4.2 Hijerarhijsko klasterovanje

Sada ćemo primeniti metodu hijerarhijskog klasterovanja i uporediti sa rezultatima dobijenim metodom K-sredina. Za određivanje broja klastera, koristićemo dendrogram. Kreirajmo dendrogram za naše podatke, za atribut *ParticipantsCount* i *MedalsCount*. Postupak njegovog kreiranja će biti prikazan u listingu 27, a dendrogram će biti prikazan na slici 10. U pozivu funkcije za kreiranje dendrograma ćemo koristiti argument *truncate_mode='lastp'*, pa ćemo sa $p=20$ videti poslednjih 20 odsecanja. Ovo koristimo da bismo smanjili gustinu dendrograma, odnosno da bismo povećali preglednost. Takođe, koristimo parametar *show_contracted=True* kako bismo na plotu iscrtali tačke, čija visina će nam označavati visinu klastera, koji nisu iscrtani. Izvršićemo odsecanje na visini 3000, koristeći parametar *max_d*.

```
1 from scipy.cluster.hierarchy import dendrogram,
```



Slika 9: Dobijeni klasteri za atribute *PartPerGamesScaled* i *SuccessRate*

```

2                                     linkage
3
4 cluster_data = country_data[['ParticipantsCount',
5                               'MedalsCount']]
6
7 f_dendro = plt.figure(figsize=(15, 10))
8 # izracunamo rastojanja izmedju elemenata i kreiramo
9 # dendrogram
10 Z = linkage(cluster_data, method='average')
11 dendrogram = dendrogram(Z, truncate_mode='lastp',
12                          p=20, leaf_rotation=90,
13                          leaf_font_size=21,
14                          show_contracted=True)
15
16 # vrednost na y-osi na kojoj cemo vrsiti odsecanje
17 max_d = 3000
18
19 plt.axhline(y=max_d, c='k')
20 plt.xlabel('Broj ucesnika', fontsize=21)
21 plt.ylabel('Broj osvojenih medalja', fontsize=21)
22 plt.show()

```

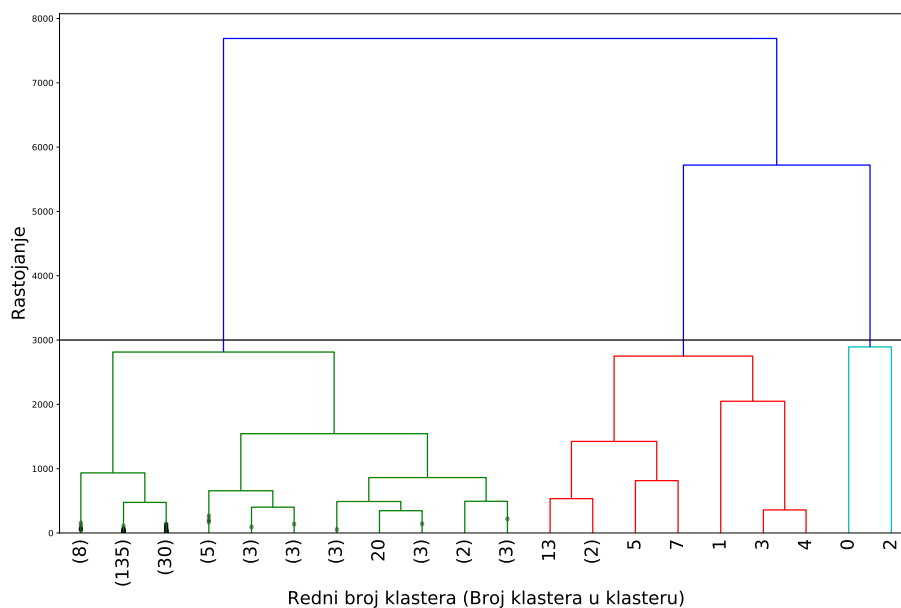
Listing 27: Kreiranje dendrograma za atribute *ParticipantsCount* i *MedalsCount*

Broj na x-osi na slici 10, koji se nalazi u zagradi predstavlja broj elemenata u tom klasteru. Sa ove slike možemo videti da je prava $y=3000$, koju smo odabrali za odsecanje, presekla 3 vertikalne linije. Zbog ovoga će nam broj klastera biti 3. U listingu 28 prikazimo postupak za iscertavanje klastera, a na slici 11 prikazimo te klasterne.

```

1 from scipy.cluster.hierarchy import linkage, fcluster
2
3 cluster_data = country_data[['ParticipantsCount',

```



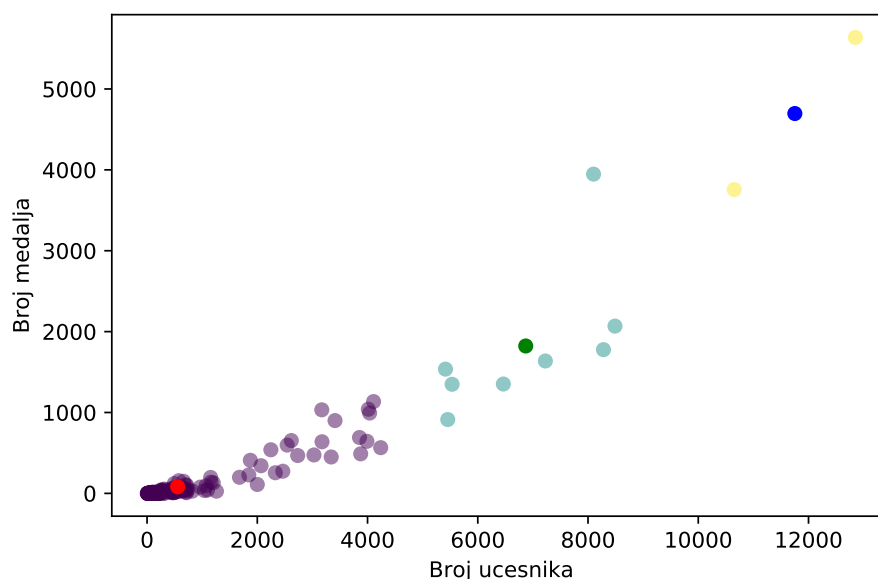
Slika 10: Dendrogram za atribute *ParticipantsCount* i *MedalsCount*

```

4         'MedalsCount']]
5
6 clusters = fcluster(Z, 3000, criterion='distance')
7 # oduzimamo 1 da bi brojanje klastera krenulo od 0
8 # ovo radimo samo da bismo mogli da trazimo
9 # ispravne indekse iz liste sa bojama
10 cluster_data['Cluster'] = clusters - 1
11
12 colors = ['red', 'green', 'blue']
13
14 cluster_means = pd.DataFrame(
15     round(cluster_data.groupby('Cluster').mean(), 1))
16 centroids = np.array(cluster_means)
17 cluster_means['CentroidColor'] = pd.DataFrame(colors)
18
19 plt.figure(figsize=(6, 4))
20 plt.scatter(cluster_data['ParticipantsCount'],
21             cluster_data['MedalsCount'],
22             c=clusters, alpha=0.5)
23 plt.scatter(centroids[:, 0], centroids[:, 1], c=colors)
24 plt.xlabel('Broj ucesnika')
25 plt.ylabel('Broj medalja')
26 plt.show()
27
28 # stampamo srednje vrednosti klastera i boje centroida
29 print(cluster_means)
30
31 Out>>
32 Cluster  ParticipantsCount  MedalsCount  CentroidColor
33 0          555.8           80.7          red
34 1        6870.2          1822.4          green
35 2       11752.0          4696.5          blue

```

Listing 28: Kreiranje klastera za atribute *ParticipantsCount* i *MedalsCount*



Slika 11: Klasteri za atribute *ParticipantsCount* i *MedalsCount*

U odnosu na klasterovanje K-sredina, sa parametrima koje smo izabrali, sada imamo 3 klastera, koji opet, vrlo jasno prikazuju određenu linearnu vezu broja učesnika i broja medalja jedne države. Ponovo se izdvajaju dve najuspešnije države gledajući ove attribute - SAD i Nemačka. Osam država se nalazi u klasteru sa zelenom centroidom, koji predstavlja uspešnije države po ovim parametrima, dok se u klasteru sa crvenom centroidom nalaze manje uspešne države.

Sledeći primer, koji ćemo analizirati je hijerarhijsko klasterovanje podataka na osnovu atributa *PartPerGamesScaled* i *SuccessRate*. Prikažimo dendrogram za ove attribute na slici 12, a klasterne na slici 13. Srednje vrednosti klastera i koji klaster odgovara kojoj boji prikazano je u listingu 29.

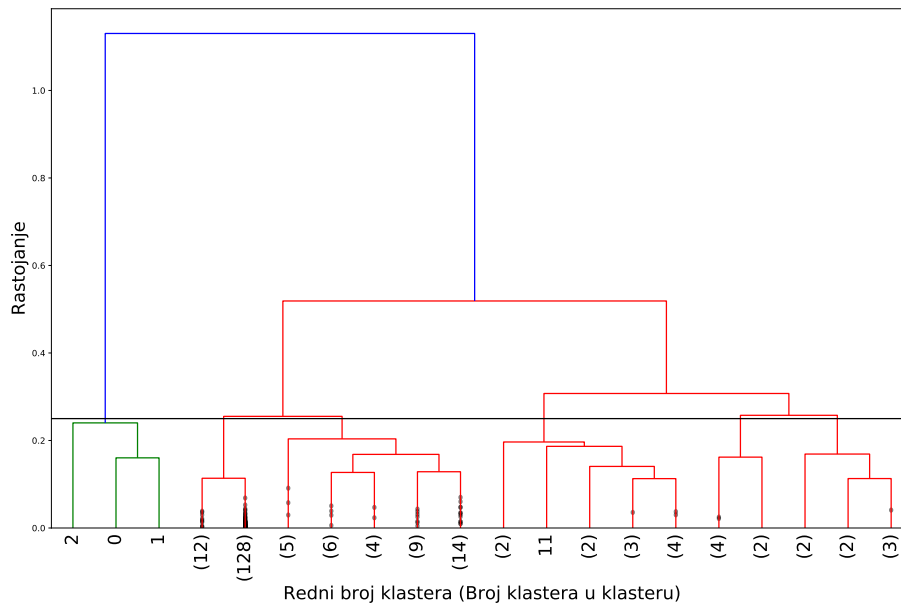
```

1 print(cluster_means)
2
3 Out>>
4 Cluster  PartPerGamesScaled  SuccessRate  CentroidColor
5 0          0.9              0.9             red
6 1          0.0              0.0             green
7 2          0.1              0.2             blue
8 3          0.2              0.5             orange
9 4          0.4              0.3             pink
10 5          0.5              0.5             black

```

Listing 29: Srednje vrednosti klastera za atribute *PartPerGamesScaled* i *SuccessRate*

Sa slike 12 zaključujemo da ćemo ovaj put koristiti 6 klastera. Na slici 13 vidimo 6 klastera, koji nešto bolje razdvajaju države po uspešnosti gledajući dva parametra, koja su korišćena. Ponovo imamo iste 3 najuspešnije države, dok su manje uspešne države jasnije razdvojene, tako da sada imamo klaster sa narandžastom centroidom, koji predstavlja države koje su izuzetno uspešne uprkos manjem broju takmičara. Sa druge strane, u klasteru sa crnom centroidom nalaze se države, koje sa prosečno



Slika 12: Dendrogram za atribute *PartPerGamesScaled* i *SuccessRate*

istim brojem učesnika ostvaruju u proseku iste rezultate, kao države iz narandžastog klastera. Države u klasteru sa plavom centroidom predstavljaju države koje u proseku imaju pozitivan učinak gledajući ova dva parametra. Roze klaster predstavlja šest država, koje imaju u proseku veliki broj učesnika, ali slab učinak što se tiče osvojenih medalja. Prikažimo ove države u listingu 30.

```

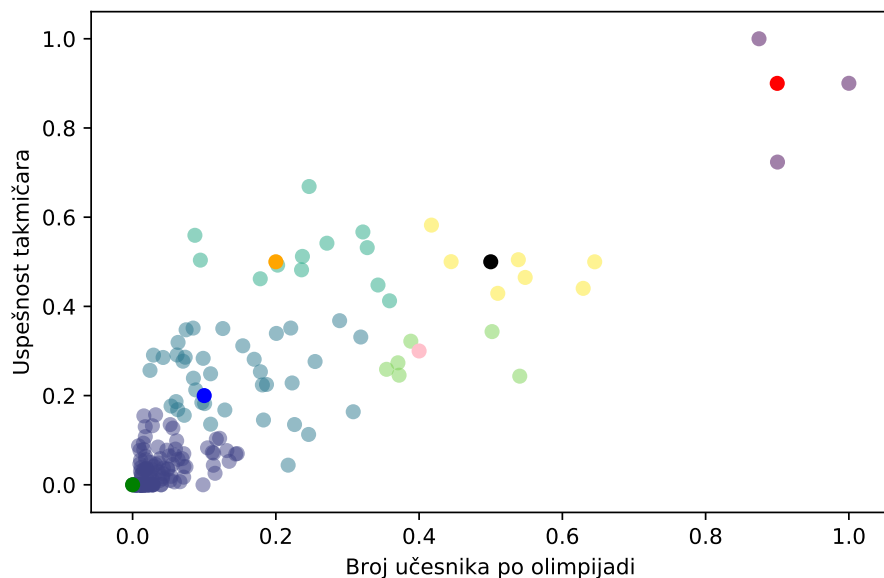
1 least_efficient =
2     country_data[
3         (country_data['PartPerGamesScaled'] > 0.35) &
4         (country_data['SuccessRate'] < 0.4)]
5 print(least_efficient[['region',
6                        'PartPerGamesScaled',
7                        'SuccessRate']])
8
9 Out>>
10      region  PartPerGamesScaled  SuccessRate
11 13    Japan          0.501845      0.343474
12 20    Poland          0.370429      0.273529
13 22    Spain          0.354540      0.258907
14 23    Brazil          0.388423      0.321926
15 31    Ukraine          0.540514      0.243522
16 35    Belarus          0.372032      0.245697

```

Listing 30: Države sa negativnim odnosom učesnika i osvojenih medalja

4.3 Klasterovanje učesnika

Poslednje što ćemo ispitati jeste uticaj visine, težine i godina takmičara na osvajanje medalje. Za početak, izvući ćemo podatke o visinama i težinama za sve takmičare, normalizovati ih, a zatim sabrati normalizovana polja za visinu i težinu. Zbir ova dva atributa će činiti novi atribut



Slika 13: Klasteri za attribute *PartPerGamesScaled* i *SuccessRate*

SizeNorm, koji predstavlja dimenzije takmičara. Takođe, izvući ćemo broj godina takmičara i normalizovati ga. Ovaj podatak će činiti kolonu *AgeNorm*. Postupak za računanje odgovarajućih podataka će biti prikazan u listingu 31. Postupak određivanja i iscrtavanja *elbow* krive i klastera neće biti ispisan, već će biti prikazan samo plot, koji je rezultat klasterovanja. Nakon pozivanja *elbow* metoda odabrano je da će biti pet klastera. Na slici 14 su prikazani klasteri, koji su dobijeni za učesnike bez medalja, a na slici 15 klasteri za učesnike, koji su osvojili neku od medalja.

```

1 data['HeightNorm'] =
2   (data['Height'] - data['Height'].mean()) /
3   np.std(data['Height'], axis=0)
4 data['WeightNorm'] =
5   (data['Weight'] - data['Weight'].mean()) /
6   np.std(data['Weight'], axis=0)
7 data['SizeNorm'] =
8   data['HeightNorm'] + data['WeightNorm']
9 data['AgeNorm'] =
10  (data['Age'] - data['Age'].mean()) / np.std(data['Age',
11  ], axis=0)
12 size_age_medal = data[data['MedalWon'] == 1]
13 size_age_no_medal = data[data['MedalWon'] == 0]
14 size_age_medal = size_age_medal[['SizeNorm', 'AgeNorm']]
15 size_age_no_medal =
16   size_age_no_medal[['SizeNorm', 'AgeNorm']]
17
18 print(size_age_medal.head)
19
20 Out>>
21      SizeNorm  AgeNorm
22 3    0.339734  1.347182
23 37   0.339734  0.710409

```

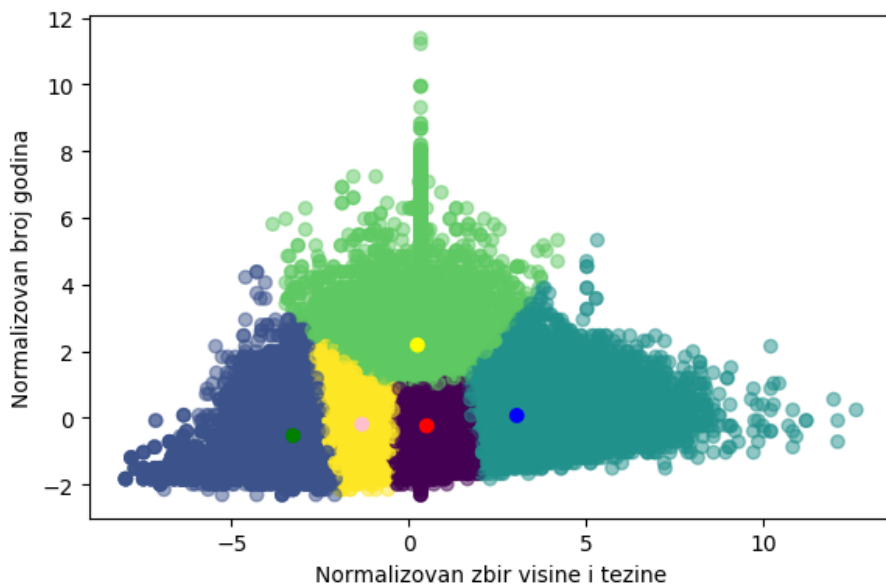


```

24 38 0.339734 0.710409
25 40 1.954303 0.392023
26 41 -0.664815 0.392023

```

Listing 31: Postupak za računanje normalizovanih parametara

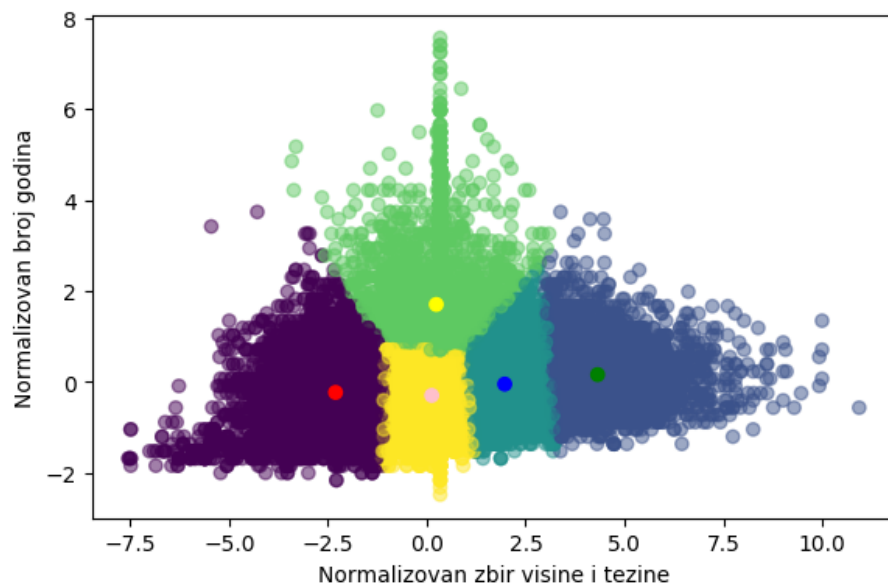


Slika 14: Klasteri za normalizovane attribute *SizeNorm* i *AgeNorm* za učesnike, koji nisu osvojili medalju

Ako pogledamo sliku 14 i sliku 15, možemo zaključiti da se dobijaju skoro potpuno isti klasteri, te da učesnici bilo kog uzrasta, visine ili težine, imaju podjednaku šansu za osvajanje medalja.

5 Zaključak

U ovom radu su primenjeni i prikazani koraci pojedinih osnovnih metoda istraživanja podataka na skupu podataka o takmičarima sa svih modernih Olimpijskih igara. Izvršena je analiza i obrada podataka, a zatim primenjene tehnike klasterovanja i upoređeni su rezultati dobijeni različitim tehnikama. Doneseni su određeni zaključci o olimpijadi i njenim učesnicima: broj učesnika je u porastu; broj ženskih takmičara je u velikom porastu; takmičari svih visina, težina i uzrasta imaju podjednake šanse za uspeh; najstariji osvajač medalje je od najmlađeg stariji 63 godine, itd. Sledeći korak u istraživanju bi bilo dodatno istraživanje postojećih, ali i uvođenje novih podataka, koji bi doneli dodatnu informaciju o državama učesnicama (bruto domaći proizvod, populacija, prosečna primanja, itd) ili o samim takmičarima.



Slika 15: Klasteri za normalizovane attribute *SizeNorm* i *AgeNorm* za učesnike, koji su osvajali medalju

Literatura

- [1] Allen. Guttmann. *The Olympics : a history of the modern games* / Allen Guttmann. University of Illinois Press Urbana [Ill.], 1992.
- [2] William Mellors Henry and Patricia Henry Yeomans. *An approved history of the Olympic games*. Los Angeles, Calif Southern California Committee for the Olympic Games, 1984 ed edition, 1984.
- [3] Richard Stanton. *The Forgotten Olympic Art Competitions*. Victoria: Trafford Publishing, 2001.
- [4] P.N. Tan, M. Steinbach, and V. Kumar. *Introduction to Data Mining*. Always learning. Pearson Addison Wesley, 2006.