

Projekat iz Istrazivanja podataka 1

profesor:Nenad Mitic

astitent:Stefan Kapunac

student: Jelena Mitrovic 357/2020

Sadržaj

1.Uvod	3
2.Analiza atributa	3
3.Klasifikacija	7
3.1. Stabla odlucivanja	7
3.1.2. Slucajna suma	8
3.1.3. Poredjenje modela	9
3.2.1 KNN	10
3.2.2.BaggingClassifier	10
3.2.3 Poredjenje modela	11
3.2.4 PCA - KNN	11
3.3.1. SVM	13
4.Klasterovanje	14
4.1. KMeans	14
4.2. Agglomerative	15
4.2.DBSCAN	15
5.Pravila pridruzivanja	15
5.1.Apriori	16
6. Literatura	18

1.Uvod

Podaci koji se nalaze u bazi podataka sa kojom radim su prikupljeni od učesnika u eksperimentalnim događajima brzih sastanaka od 2002.-2004. godine. Prisutne osobe bi imale četvorominutni "prvi sastanak" sa svakim drugim učesnikom suprotnog pola. Na kraju njihova četiri minuta, učesnici su upitani da li bi želeli da ponovo izadju na sastanak . Od njih je takođe zatraženo da ocenjuju svoj sastanak na osnovu šest atributa: privlačnost, iskrenost, inteligencija, zabava, ambicija i zajednički interesi. Skup podataka takođe uključuje podatke iz upitnika prikupljene od učesnika u različitim tačkama procesa. Ova polja uključuju: demografiju, navike za sastanke, uverenja o tome šta drugi smatraju dragocnim kod partnera i informacije o životnom stilu.

2.Analiza atributa

Baza podataka se sastoji od 8378 instanci i 123 atributa.

Atributi baze podataka:

gender: Pol osobe	attractive: Ocena osobe za samu sebe - atraktivnost
age: Broj godina osobe	sincere: Ocena osobe za samu sebe - iskrenost
age_o:Broj godina partnera	intelligence: Ocena osobe za samu sebe -inteligencija
d_age: Razlika u godinama	funny:Ocena osobe za samu sebe - duhovitost
race: Rasa osobe	ambition: Ocena osobe za samu sebe- ambicija
race_o:Rasa partnera	attractive_partner:Ocena osobe za partnera-atraktivnost
samerace: Da li je vazno da osobe imaju istu rasu(1/0)	sincere_partner:Ocena osobe za partnera-iskrenost
importance_same_race: Koliko je vazno da je partner iste rase	intelligence_partner: Ocena osobe za partnera-inteligencija
importance_same_religion:Koliko je vazno da je partner iste religije	funny_partner: Ocena osobe za partnera-duhovitost
field: Obrazovanje osobe	ambition_partner: Ocena osobe za patnera-ambicija
pref_o_attractive:Koliko je partneru vazna privlacnost	shared_interests_partner: Ocena osobe za partnera-zaj.inter.
pref_o_sinsere: Koliko je partneru vazna iskrenost	sports: Interesovanje osobe za sport
pref_o_intelligence: Koliko je partneru vazna inteligencija	tv sports: Interesovanje osobe za tv sport
pref_o_funny: Koliko je partneru vazno da je osoba duhovita	exercise: Interesovanje osobe za vezbanje
pref_o_ambitious:Koliko je partneru vazna ambicija	dining: Interesovanje osobe za izlazak na veceru
pref_o_shared_interests: Koliko je partneru vazno da imaju zaj. inter.	museum: Interesovanje osobe za obilazak muzeja
attractive_o:Ocena osobe od strane partnera posle sastanka o atraktivnosti	art:Interesovanje osobe za umetnost
sincere_o: Ocena osobe od strane partnera posle sastanka o iskrenosti	hiking:Interesovanje osobe za planinarenje
intelligence_o:Ocena osobe od strane partnera posle sastanka o inteligenciji	gaming: Interesovanje osobe za igranje igrice
funny_o:Ocena osobe od strane partnera posle sastanka o duhovitosti	clubbing: Interesovanje osobe za nocne izlaske
ambitious_o: Ocena osobe od strane partnera posle sastanka o ambicioznost	reading:Interesovanje osobe za citanje

shared_interests_o: Ocena osobe od strane partnera posle sastanka o zaj.inter. tv: Interesovanje osobe za gledanje televizije

attractiveImportant: Kakvog partnera trazite - atraktivnost theater: Interesovanje osobe za odlazak u pozoriste

sincereImportant: Kakvog partnera trazite - iskrenost movies: Interesovanje osobe za gledanje filmova

intelligenceImportant: Kakvog partnera trazite - inteligencija concerts: Interesovanje osobe za koncerte

funnyImportant: Kakvog partnera trazite - duhovitost music: Interesovanje osobe za muziku

ambitionImportant: Kakvog partnera trazite - ambicija shopping: Interesovanje osobe za shopping

shared_interests_important: Kakvog partnera trazite -zaj.inter. yoga: Interesovanje osobe za jogu

interests_correlate: Korelacija interesovanja osobe i partnera

expected_happy_with_sd_people: Koliko ocekujete da cete biti sretni sa osobama koje upoznate na brzom sastanku

expected_num_interested_in_me: Od 20 osoba koje cete upoznati, koliko ocekujete da ce biti zainteresovano za zabavljanje sa vama

expected_num_matches: Koliko ocekujete pozitivnih odgovora za sledeci sastanak

like: Da li vam se svidja partner

guess_prob_liked: Kolika je verovatnoca da se dopadnete vashem partneru

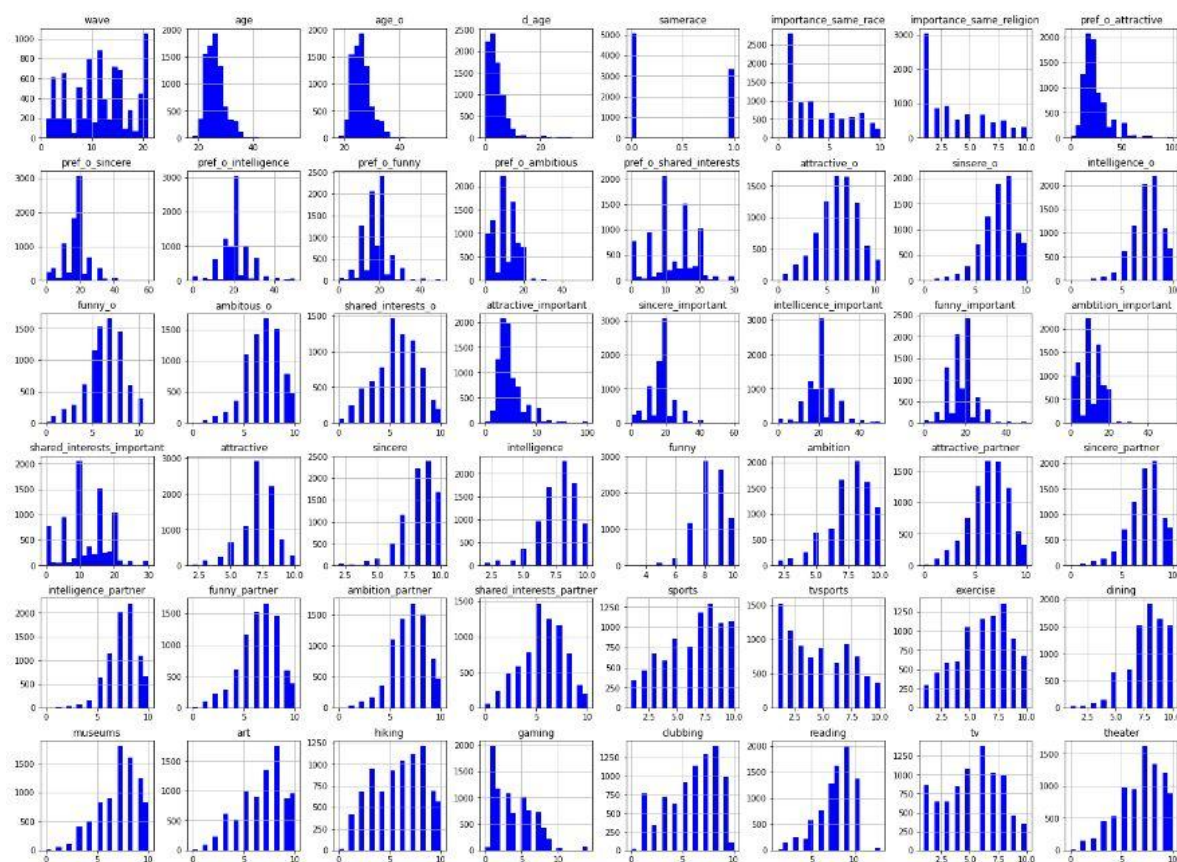
met: Da li ste sreli partera pre ovog sastanka

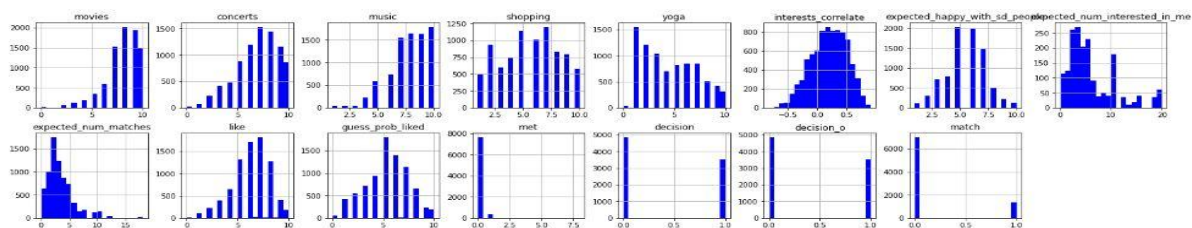
decision: Odluka osobe posle sastanka

decision_o: Odluka partnera posle sastanka

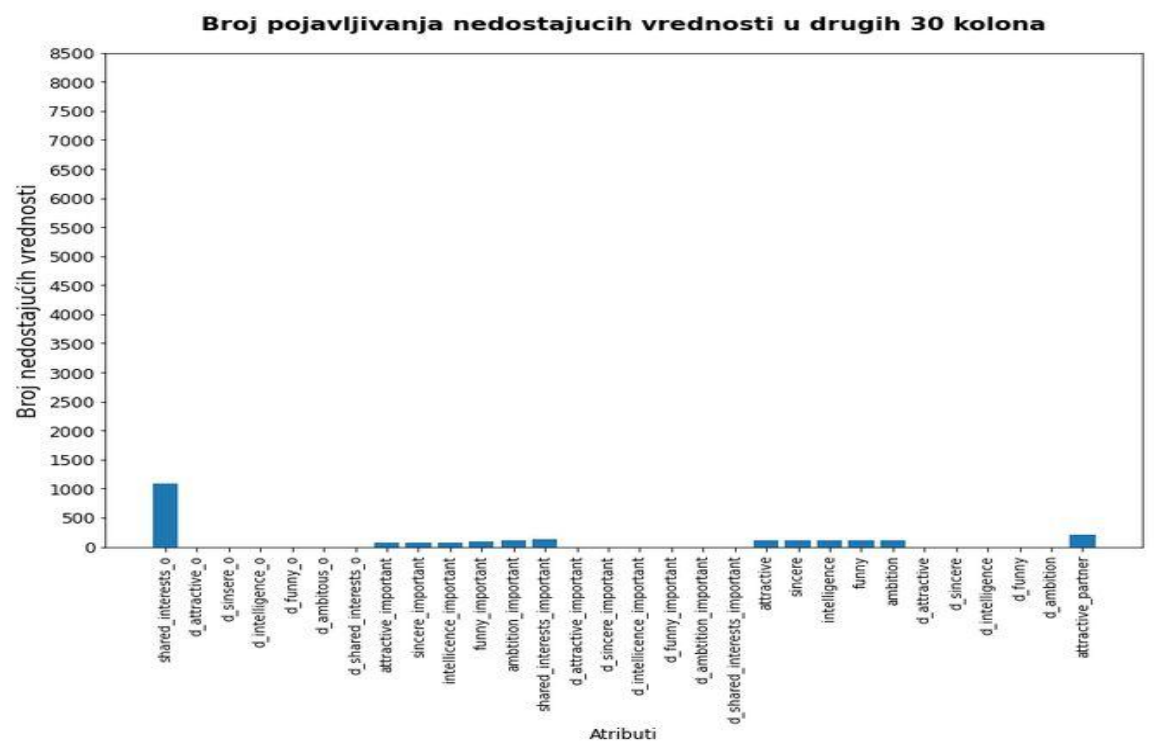
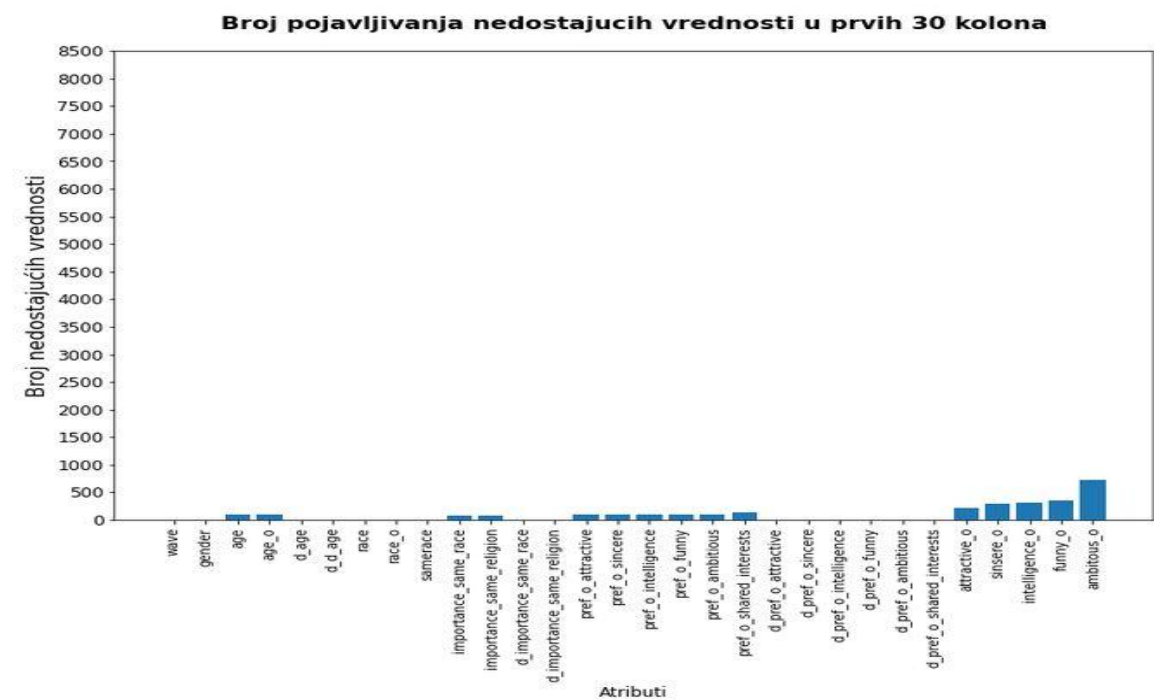
match: Da li ce se osobe opet sresti

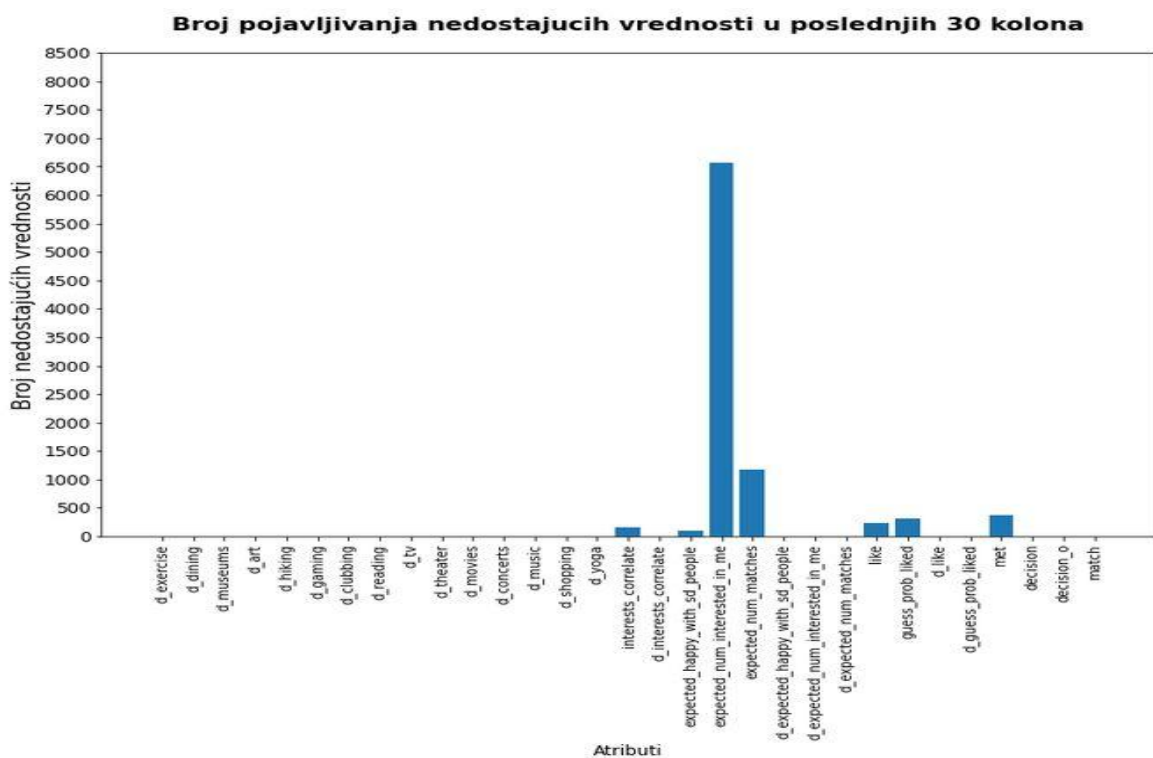
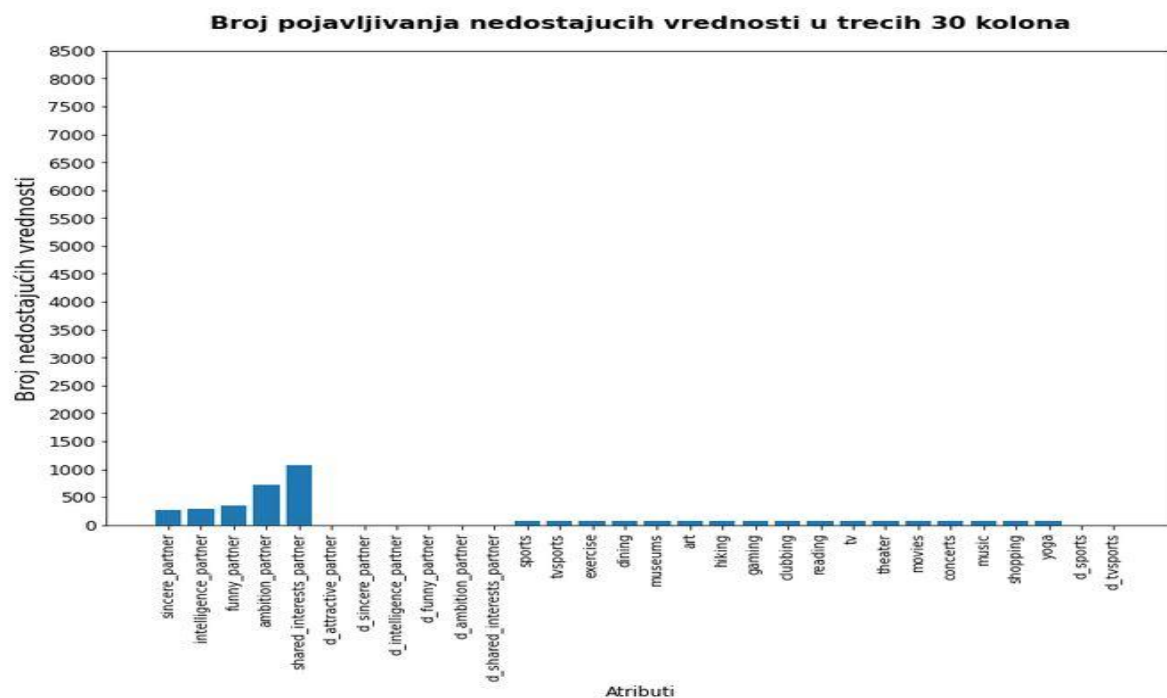
Na slikama ispod prikazan je histogram raspodele podatka nekih od atributa





Baza podatka sa kojom radim, ima veliki broj nedostajucih vrednosti. Procenat instanci koje sadrze nedostajuce vrednosti je priblizno 87.5 %. Graficki prikaz nedostajucih vrednosti:

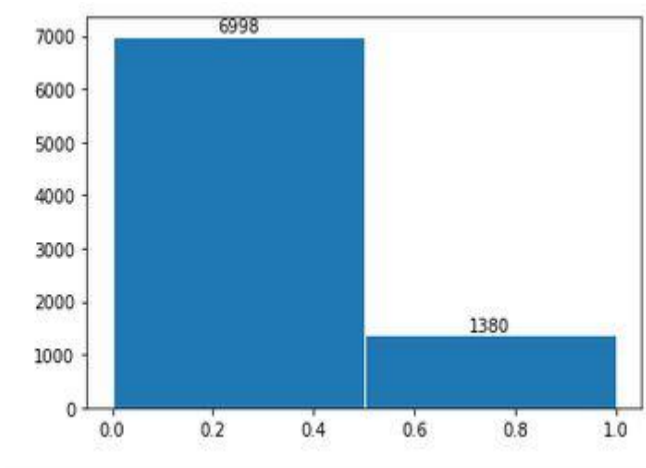




U koraku pretprocesiranja sam sve nedostajuce vrednosti zamenila medijanom za tu kolonu, a sve kategoricke (imenske) attribute sam prevela u numericke koriscenjem LabelEncoder i na taj nacin sam pripremila podatke za klasifikaciju.

3. Klasifikacija

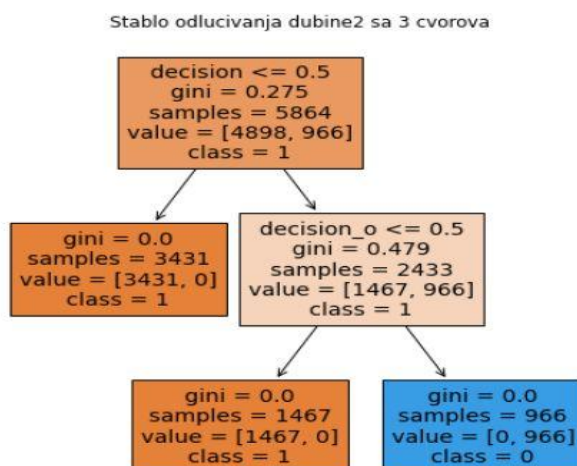
Zadatak klasifikacije u slučaju mog skupa podataka, jeste predviđanje ciljne promenljive ('match'), na osnovu ostalih atributa. Dakle, zadatak klasifikacije je predviđanje da li će se osobe ponovo vidjeti. Ciljna promenljiva 'match' ima 2 vrednosti (0 i 1), pa zaključujem da se radi o binarnoj klasifikaciji. Na slici ispod je prikazana balansiranost između klasa. Na osnovu slike zaključujemo da postoji nebalansiranost između klasa, međutim, mislim da i u jednoj i u drugoj klasi ima dovoljno instanci, pa stoga neću koristiti neke od metoda za smanjenje razlike u dimenziji klasa, međutim, zbog ovoga neću koristiti tačnost kao kriterijum provere modela, nego f1-meru.



Obrađeni algoritmi klasifikacije: Stabla odlučivanja, Slučajna suma, K najbližih suseda, Bagging Classifier i SVM.

3.1. Stabla odlučivanja

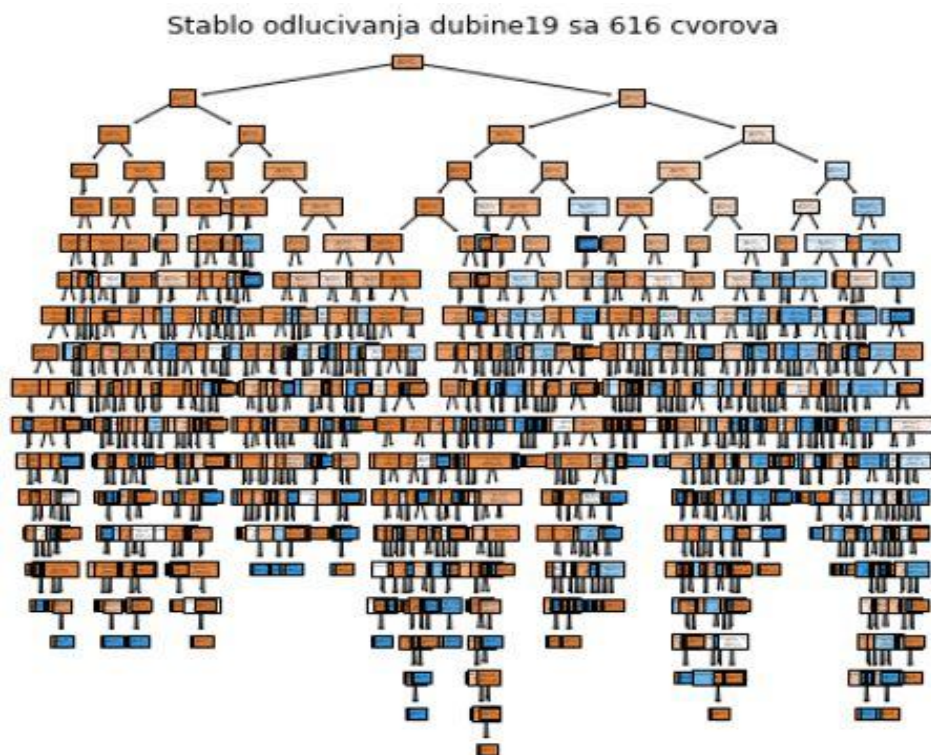
Primenom `DecisionTreeClassifier()` sa podrazumevanim parametrima, dobila sam sledeće stablo odlučivanja:



Na slici primecujemo, da se prilikom formiranja stabla koriste samo 2 atributa : 'decision' i 'decision_o'. Uticaji ovih atributa na formiranje stabla su redom: 0.27812303 i 0.72187697, a uticaj svih ostalih jednak je nula. Koriscenjem ovog modela dobijen je f1-score od 1 i na trening i na test skupu, a odgovarajuće matrice konfuzije za trening i test skup redom su:

1	0	1	0
1	4898	0	1
0	0	966	0
			2100
			0
			414

Posto sam videla da na formiranje stabla uticu samo 2 atributa, htela sam da vidim kako bi izgledalo stablo odlucivanja kada bih izbacila te attribute.. Dobila sam stablo dubine 19 sa 616 cvorova.



Ovaj model ne daje dobre rezultete, zato sto se previse prilagodio trening podacima i izgubio moc generalizacije. Matrice konfuzije na trening i test skupu:

1	0	1	0
1	5248	0	1
0	0	1035	0
			1506
			244
			201
			143

Model je komplikovaniji nego sto bi trebalo da bude.

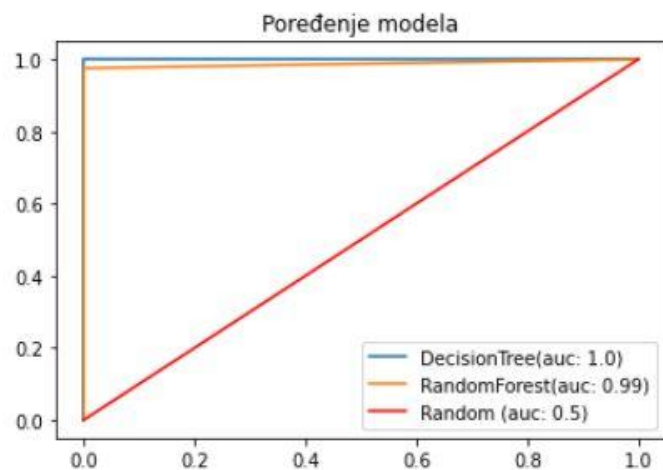
3.1.2. Slucajna suma

Pozivanjem GridSearchCV() za model RandomForestClassifier, dobila sam da su optimalni parametri ovog modela 'max_depth'=15, 'min_samples_split'= 2, 'n_estimators': 300. Proverom matrice konfuzije, vidimo da model ne pravi greske na trening skupu, ali pravi 10 gresaka na test skupu.

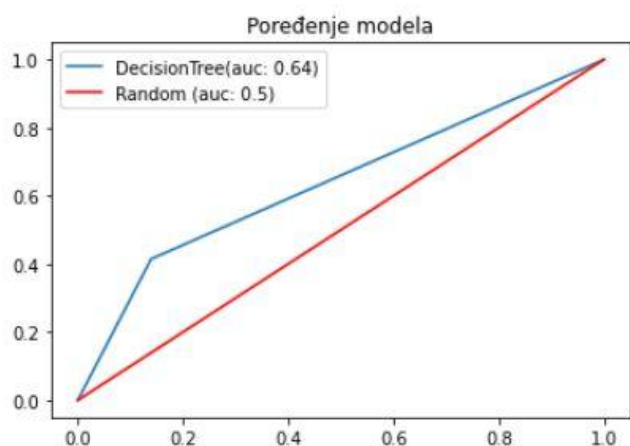
1	0	1	0
1	4898	0	1
0	0	966	0
			2100
			0
			404

3.1.3. Poređenje modela

Modele poredimo koristeći ROC krivu. Sa slike ispod primećujemo da su navedeni algoritmi dali jako dobre rezultate, zato što je površina ispod krive blizu jedinice.



Poredim model koji ne sadrži dva glava atributa sa random modelom. Ovaj model je malo bolji od random modela, ali je losiji u odnosu na gore navedene.



3.2.1 KNN

Pre primene KNN algoritma moramo da normalizujemo podatke. Za normalizaciju sam koristila StandardScaler. Na normalizovane podatke sam primenila KNeighborsClassifier() sa podrazumevanim parametrima. Model pravi jako veliki broj gresaka gde instance klase 1, klasifikuje kao instance klase 0. Odgovarajuće matrice konfuzije na trening i test skupu su:

	1	0
1	4780	118
0	516	450

	1	0
1	2024	76
0	307	107

Model je znatno lošiji u odnosu na stabla odlučivanja, pa sam htela da proverim da li postoje elementi van granica koji možda kvare tačnost ovog modela. Za proveru elemenata van granica koristila sam IQR algoritam i videla da postoji veliki broj elemenata koji se tretiraju kao outlieri (u jednoj koloni je procenat takvih elemenata je čak 49%). Zbog velikog broja outlieri nisam mogla da obrisem te podatke, zato što bih izgubila mnogo podataka, pa sam odlučila da outlieri zamenim medijanom i na taj način smanjim njihov broj. Ovim postupkom nisam u potpunosti izgubila elemente van granica, ali se procenat istih znatno smanjio. Primenom KNeighborsClassifier() sa podrazumevanim parametrima nad ovim podacima dobila sam matrice konfuzije koje su skoro identične. Zatim, probala sam da podesim hiper parametre. Zbog velike nebalansiranosti, za score za koristila f1 meru. Primenom GridSearchCV algoritma, dobila sam da su optimalni parametri KNeighborsClassifier() algoritma: 'n_neighbors': 10, 'p': 2, 'weights': 'distance', a najveća vrednost f1 scorea je približno 0.33. Podesavanjem hiper parametara, dobili smo model koji se prilagodio trening podacima.

	1	0
1	4898	0
0	0	966

	1	0
1	2044	56
0	331	83

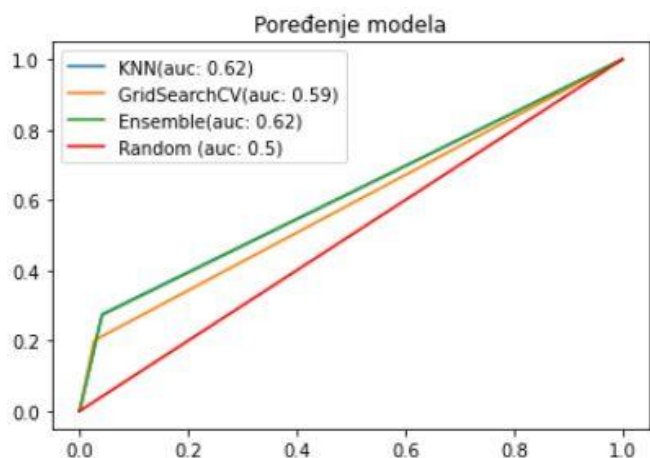
3.2.2. BaggingClassifier

Primenom ove ansambl metode dobila sam približno iste rezultate kao i primenom KNN algoritma. Matrice konfuzije na trening i test skupu:

	1	0
1	4812	86
0	537	429

	1	0
1	2014	86
0	301	113

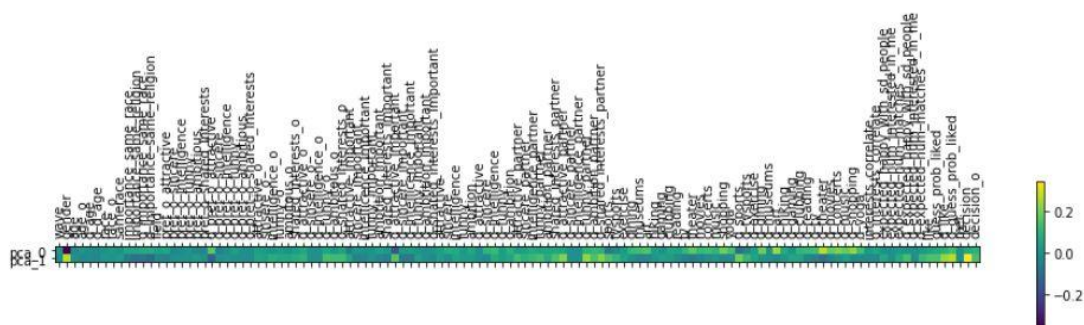
3.2.3 Poredjenje modela



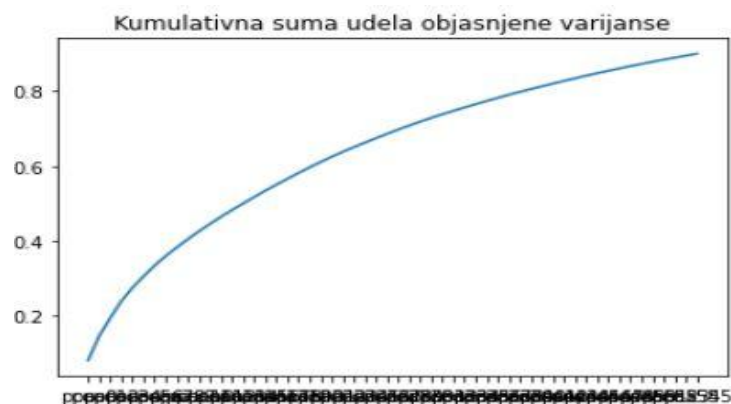
Modeli su znatno losiji u odnosu na stabla odlucivanja. Zbog velike razlike izmedju broja instanci u klasama, mislim da K najblizih suseda nije najbolji izbor za klasifikaciju mog skupa podataka. Upravo zbog toga sto imamo mnogo manje istanci klase 1, vecina suseda ce pripadati klasi 0 i zbog toga model jako veliki broj FN gresaka.

3.2.4 PCA - KNN

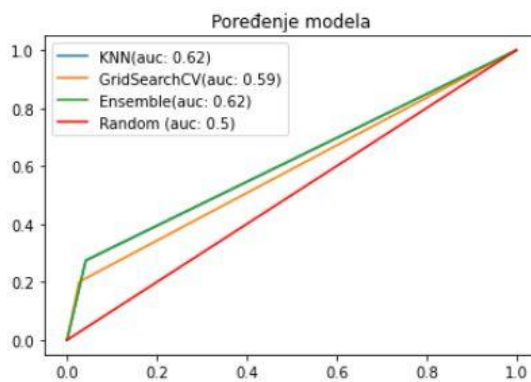
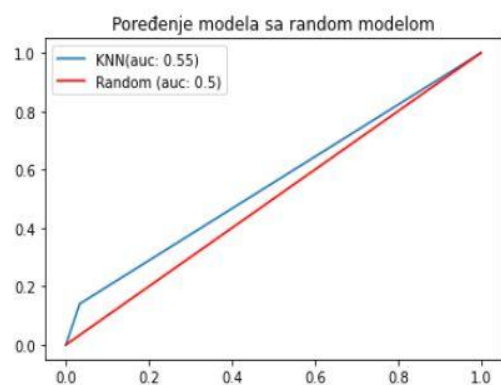
Moja baza podataka se sastoji od 121. atributa, sto je poprilično puno. Htela sam da proverim da li primenom PCA algoritma mogu da smanjim broj atributa. Prvo sam probala da broj atributa smanjim na dva. Pokusala sam da vizuelno prikazem koliko je koji atribut uticao na pravljenje novih atributa, medjutim zbog velikog broja atributa vizualizacija najjasnija.



Sa slike vidimo da smo sa 2 glavne komponente uspjeli da objasnimo samo 15% ukupne varijanse, pa zaključujem da nije moguće smanjenje broja atributa na dva. Zatim sam probala da smanjim broj atributa, ali da udeo objasnjenje varijanse bude bar 90%. Uz ovaj uslov broj glavnih atributa je porastao na 56. Na slici ispod prikazana je kumulativna suma objasnjene varijanse:

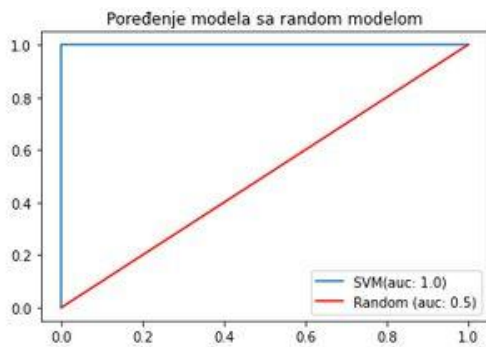


Na transformirane podatke primenila sam KNN algoritam. Na levoj slici ispod prikazano je porednje ovog modela sa random modelom. Dobijeni model je losiji u odnosu na model koji sam dobila primenom KNN algoritma na originalnim podacima (slika desno), ali ne puno losiji.



3.3.1. SVM

Poslednji algoritam klasifikacije koji sam obradila jeste SVM. Ideja je pronaci hiperravan koja razdvaja podatke iz 2 klase. Ako su podaci linearno razdvojeni koristimo marginu, odnosno imacemo uslove da se hiperravan nalazi na bezbednom rastojanju od podataka i da su podaci iz jedne ravni sa jedne strane hiperrvani, a drugi sa druge. Ukoliko podaci nisu linearno razdvojeni koristimo meku marginu, odnosno dopustamo da imamo neke pogresne instance. Primenom GridSearchCV algoritma za model SVM dobila sam da su optimalmi parametri SVM modela $\{C=0.001, \text{'kernel'}='linear'\}$. GridSearch algoritmu nisam naglasila da kao score koristi f1-meru, medjutim, mislim da nije problem sto je koristio accuracy, zato sto pogledom na ROC krivu i AUC vidimo da je model dobar i sa ovom merom.

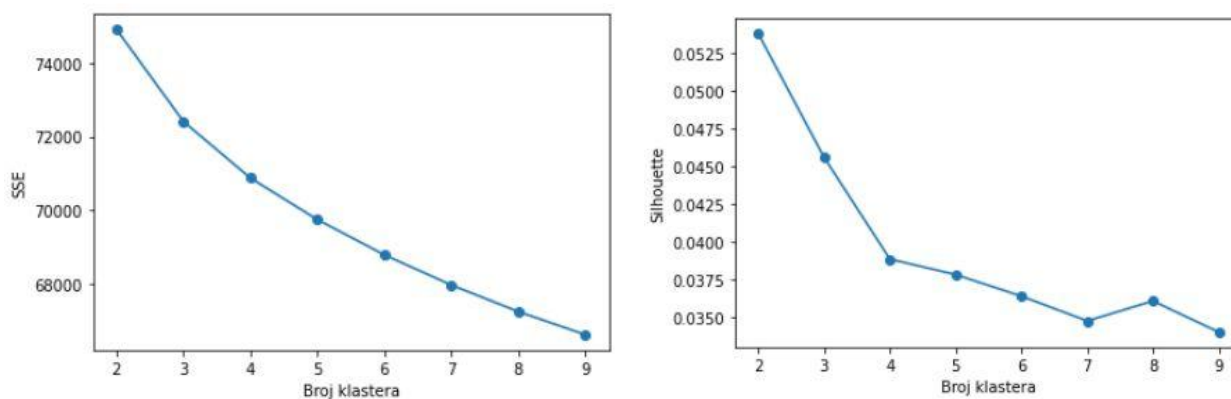


4. Klasterovanje

Za razliku od klasifikacije, gde smo imali ciljnu promenljivu koju je trebalo predvideti, klasterovanje spada u nenadledano učenje, odnosno, bez ciljne promenljive potrebno je naći neke sličnosti između podataka i grupisati podatke. Obradjeni algoritmi klasterovanja: K sredina, hijerarhijsko klasterovanje i DBSCAN. Primenom ovih algoritama zaključila sam da moj skup podataka nije pogodan za klasterovanje.

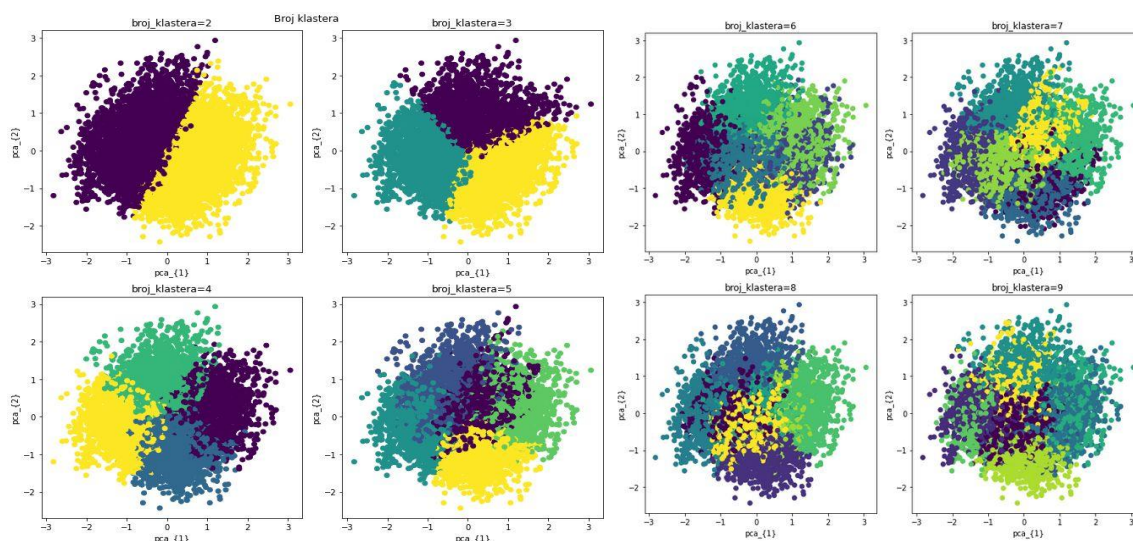
4.1. KMeans

KMeans predstavlja iterativni algoritam koji deli podatke u K klastera. Svaka tačka pripada tačno jednom klasteru. Ovaj algoritam zasniva se na reprezentativnim predstavnicima. Prvi korak prilikom primene ovog algoritma jeste određivanje optimalnog broja klastera. Na slici ispod prikazana je promena SSE tokom povećanja broja klastera. Promena SSE nije najbolja mera za izbor optimalnog broja klastera, zbog toga što se povećanjem broja klastera SSE smanjuje, pa zato moramo da posmatramo i promenu Silhouette. Na slici ispod prikazana je ta promena.



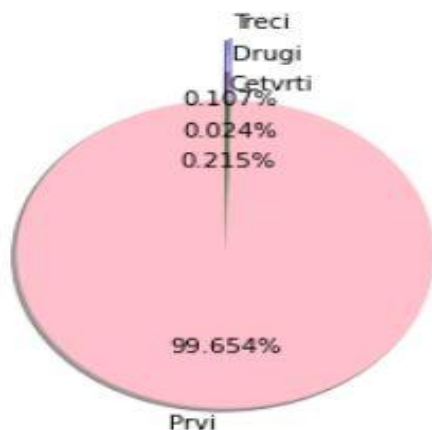
Primenom pravila lakta zaključujemo da je optimalan broj klastera jednak 4 (ili 7?). Silhouette skor za $k=4$ je približno jednak 0.0375, što je jako loše.

Vec smo videli da primenom PCA algoritma i smanjenjem broja komponenti na 2 možemo da objasnimo samo 15% ukupne varijanse podatka, međutim, uradila sam smanjenje broja komponenti na 2 samo da bih probala da vizualizujem klasterove. Na slikama ispod prikazana je promena sa povećanjem broja klastera:



4.2. Agglomerative

Primenom algoritma Agglomerative sa parametrima {n_clusters=4, linkage='average', compute_distances=True}, dobila sam da najveći broj instanci pripada prvom klasteru, čak 8349, drugom klasteru pripada samo 2 instance, trecem 9, a cetvrtom 18.



Vrednost silhouette score za ove parametre jeste 0.076253, što je i dalje jako malo. Proverila sam da li se bolji score dobija ako koristimo single ili complete vezu, ali u oba slucaja dobijamo manji silhouette score.

4.2.DBSCAN

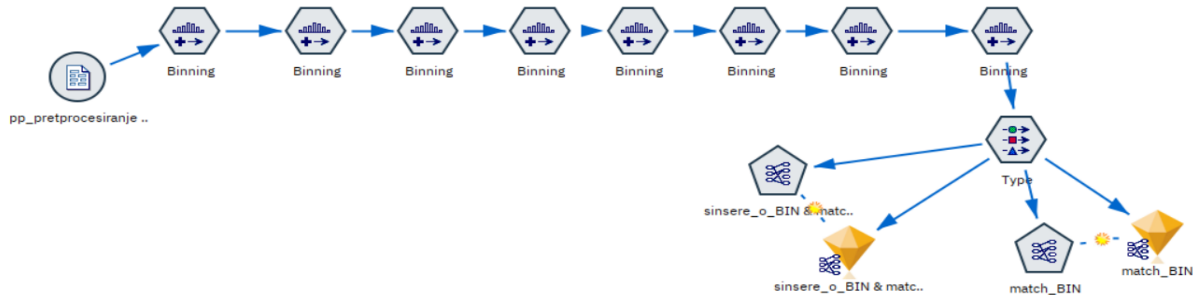
U petlji sam pokusala da odredim koji su dobri parametri ovog modela tako da imamo sto manje suma, odnosno, instanci koje ne pripadaju nijednom klasteru, medjutim, bilo kojom kombinacijom (min_samples_value = range(2,50) eps_values = [0.1, 0.15, 0.2, 0.25, 0.3]) sam dobila da su sve instance šumovi, odnosno ovaj algoritam nije uspeo da pronadje nijedan klaster. Na osnovu svega ovoga zaključujem da moji podaci nisu pogodni za klasterovanje.

5.Pravila pridruzivanja

Skup koji sam koristila za pravila pridruzivanja jeste onaj koji sam dobila pretprocesiranjem podataka na pocetku, odnosno zamenom null vrednosti medijanom i koriscenjem LabelEncodera. Zbog jako velikog broja atributa, iskoristila sam Feature Selection da izdvojim glavne attribute. Izdvojila sam 35 najvaznijih atributa, odnosno 36 sa atributom 'match'. Odlucila sam da cilj pravila pridruzivanja u slucaju mog skupa podataka bude nalazenje skupa osobina koje povlace vrednost atributa 'match = 1', zato sto mi je to imalo najvise smisla.. Obradila sam algoritam Apriori.

5.1. Apriori

Svih 35 atributa koje sam dobila kao najvaznije za predviđanje cilja bila su neprekidna, pa sam iskoristila cvor Binning..Nad ovakvim tipovima mogla sam da primenim Apriori algoritam.



Prvo sam primenila algoritam sa uslovima da se u glavi moze naci samo atribut 'match', a u telu se mogu naci svi ostali atributi. Pravila koja sam dobila nisu preterano zanimljiva, iako imaju visoku ocenu Lift mere, zato sto je dobijeni zakljucak isti kao i onaj koji sam dobila primenom stabla odlucivanja. Ono sto mora da vazi da bi odgovor za ponovni susret bio pozitivan, tj. da je vrednost promenljive 'match'= 1, jeste da atributi 'decision' i 'decision_o' takodje moraju da budu True. Dakle, dobila sam da ukoliko je vrednost flag atributa 'decision' i 'decision_o' true, onda ce se osobe ponovo izadji, bez obzira na vrednost ostalih atributa. Neka od pravila:

Consequent	Antecedent	Support %	Confidence %	Lift
match_BIN	decision_o_BIN	16.472	100.0	6.071
match_BIN	decision_BIN	12.593	100.0	6.071
match_BIN	like_BIN = 3	11.9	100.0	6.071
match_BIN	decision_o_BIN	11.9	100.0	6.071
match_BIN	decision_BIN	11.9	100.0	6.071
match_BIN	attractive_o_BIN = 3	13.046	100.0	6.071
match_BIN	decision_o_BIN	11.017	100.0	6.071
match_BIN	decision_BIN	11.017	100.0	6.071
match_BIN	d_attractive_partner_BIN = 2	10.098	100.0	6.071
match_BIN	decision_o_BIN	10.098	100.0	6.071
match_BIN	decision_BIN	10.098	100.0	6.071
match_BIN	d_ambitious_o_BIN = 2			
match_BIN	d_ambition_partner_BIN = 2			

Probala sam da izbacim attribute ‘decision’ i ‘decision_o’ i dobila sam sledeca pravila.Pravila imaju manju pouzdanost i lift meru od gore navedenih, ali nam daju neke nove informacije u odnosu na stabla odlucivanja.

Consequent	Antecedent	Support %	Confidence %	Lift
match_BIN <div>match_BIN = 3</div>	shared_interests_o_BIN = 3 funny_o_BIN = 3 like_BIN = 3 attractive_o_BIN = 3 attractive_partner_BIN = 3	5.908	65.051	3.949
match_BIN	d_guess_prob_liked_BIN = 3 funny_o_BIN = 3 like_BIN = 3 attractive_o_BIN = 3 funny_partner_BIN = 3	5.312	64.944	3.943
match_BIN	guess_prob_liked_BIN = 3 funny_o_BIN = 3 like_BIN = 3 attractive_o_BIN = 3 funny_partner_BIN = 3	5.3	64.865	3.938
match_BIN	d_guess_prob_liked_BIN = 3 funny_o_BIN = 3 attractive_o_BIN = 3 attractive_partner_BIN = 3 intelligence_o_BIN = 3	5.085	64.319	3.905
match_BIN	guess_prob_liked_BIN = 3 funny_o_BIN = 3 attractive_o_BIN = 3 attractive_partner_BIN = 3 intelligence_o_BIN = 3	5.073	64.235	3.9

6. Literatura

1. Materijali sa predavanja profesora Nenada Mitića
2. Materijali sa vezbi iz kursa Istraživanje podataka:
https://github.com/MATF-istrazivanje-podataka-1/materijali_2022-2023
3. Scikit-Learn Machine Learning in Python: <https://scikit-learn.org/stable/>