

UNIVERZITET U BEOGRADU  
ELEKTROTEHNIČKI FAKULTET



**PREPOZNAVANJE EMOCIJA NA OSNOVU  
KLASIFIKACIJE OBELEŽJA FIZIOLOŠKIH SIGNALA**

Master rad

Mentor:

Doc. dr Predrag Tadić

Kandidat:

Boris Milićević 3035/2017

Beograd, Septembar 2018.

*Zahvaljujem se svom mentoru profesoru Predragu Tadiću na korisnim sugestijama, profesorici Milici Janković i laboratoriji za biomedicinsko inženjerstvo.*

*Ovo istraživanje je delimično podržano od strane Inovacionog centra Elektrotehničkog fakulteta, Univerziteta u Beogradu, Inovacionog fonda Srbije (broj ID50053) i Ministarstva Obrazovanja, Nauke i Tehnološkog Razvoja Srbije, Beograd, Srbija, (broj OS175016)*  
*Autori se zahvaljuju prof. Vanji Ković i Nikoli Milosavljeviću, Psihološki fakultet, Univerziteta u Beogradu za pomoć u studiji.*

# Sadržaj

<b>1. UVOD .....</b>	<b>4</b>
<b>2. PRIKUPLJANJE PODATAKA .....</b>	<b>6</b>
2.1 EKSPERIMENTALNA PROCEDURA .....	6
2.2 OPIS AKVIZICIJE I PROCESIRANJA SIGNALA .....	7
2.3. OPIS DOBIJANJA SKUPA PODATAKA .....	9
<b>3. METODE .....</b>	<b>12</b>
3.1 METODE RELABELIRANJA PODATAKA .....	12
3.1.1 <i>Polu-supervizirano klasterovanje</i> .....	12
3.2 METODE ANALIZE OBELEŽJA.....	13
3.2.1 <i>PCA analiza</i> .....	13
3.2.2 <i>LDA analiza</i> .....	14
3.2.3 <i>Procena korelisanosti</i> .....	15
3.2.4 <i>Procena informativnosti</i> .....	15
3.3 METODE KLASIFIKACIJE.....	15
3.3.1 <i>Logistička regresija – softmax metoda</i> .....	16
3.3.2 <i>SVM metoda</i> .....	17
3.3.3 <i>Random Forest metoda</i> .....	19
<b>4. REZULTATI.....</b>	<b>21</b>
4.1 RELABELIRANJE PODATAKA.....	21
4.2 ANALIZA OBELEŽJA .....	24
4.2.1 <i>PCA analiza</i> .....	24
4.2.2 <i>Korelisanost obeležja</i> .....	25
4.2.3 <i>Informativnost obeležja</i> .....	26
4.2.4 <i>Rangiranje pomoću slučajnih šuma</i> .....	28
4.3 PERSONALIZOVANA KLASIFIKACIJA .....	29
4.3.1 <i>Podela skupa podataka</i> .....	30
4.3.2 <i>Rezultati primene softmax metode</i> .....	31
4.3.3 <i>Rezultati primene SVM metode</i> .....	33
4.3.4 <i>Rezultati primene Random Forest metode</i> .....	35
4.4 INTERPERSONALNA KLASIFIKACIJA.....	36
4.4.1 <i>Normalizacija podataka</i> .....	36
4.4.2 <i>Rezultati primene softmax metode</i> .....	38
4.4.3 <i>Rezultati primene SVM metode</i> .....	40
4.4.4 <i>Rezultati primene Random Forest metode</i> .....	41
<b>5. ZAKLJUČAK .....</b>	<b>43</b>
<b>6. SPISAK REFERENCI .....</b>	<b>45</b>
<b>7. SPISAK SKRAĆENICA .....</b>	<b>47</b>

<b>8.SPISAK SLIKA .....</b>	<b>48</b>
<b>9.SPISAK TABELA .....</b>	<b>50</b>

# 1. Uvod

Detekcija emocija postaje sve bitnije polje u interakciji čoveka i mašine. Precizno i pouzdano prepoznavanje emocija će omogućiti mašinama da ova interakcija postane prirodnija te da adekvatno reaguju na ljudske potrebe. Sistem za prepoznavanje emocija bi pronašao svoju primenu u različitim oblastima kao što je medicina, rehabilitacija, marketing, pri čemu bi mogla postati deo naših pametnih telefona, kućnih aparata, automobila i mnogih drugih uređaja sa kojima svakodnevno interagujemo.

Ovom tezom je opisan razvoj jednog inteligentnog sistema sposobnog da prepozna tri emotivna stanja koja u velikoj meri utiču na donošenje odluka kod ljudi: strah (panika i anksioznost), fokusiranost (stanje koncentracije) i stanje opuštenosti.

Skup podataka za obučavanje ovog sistema je dobijen ekstrahovanjem obeležja iz elektrofizioloških signala čija je akvizicija vršena u toku eksperimenta. Eksperiment je realizovan u okviru virtuelne realnosti (VR) sa zadatkom da se u ispitaniku pobude ciljane emocije. Unutar VR ispitanik je uključen u bogat trodimenzionalni (3D) svet čiji pažljivo odabran sadržaj dovodi ispitanika u željeno stanje. U toku samog eksperimenta ispitaniku je vršena akvizicija aktivnosti srčanog mišića (elektrokardiografija - EKG), aktivnost respiratornog sistema kao i moždana aktivnost (elektroencefalografija - EEG).

Oblast detekcije emocija je vrlo atraktivna oblast za stručnjake iz mnogih naučnih sfera. Veliki broj objavljenih studija i rezultata istraživanja upravo su rezultat velike popularnosti ove oblasti. Načini na koje autori pristupaju datoj temi su jako raznovrsni, a ta se raznovrsnost prvenstveno ogleda u metodama indukcije ljudskih emocija te informacijama na osnovu kojih se kreiraju inteligentni prediktivni sistemi.

Uobičajen pristup jeste detekcija emocija na osnovu audio-vizuelnog sadržaja. Kamera beleži facijalne ekspresije ispitanika dok je ispitanik u određenom afektivnom stanju. Ovakve studije tipično imaju visoke stope uspešnosti prepoznavanja emocija. Tim istraživača Husam Salih i Lalit Kulkarni su napravili model klasifikatora *Support Vector Machine* (SVM) koji na osnovu facijalnih ekspresija vrši predikciju šest emocija sa tačnošću od 95% [1]. Tim istraživača iz Poljske je pomoću neuralne mreže sa samo jednom skrivenim slojem postigao uspešnost od 96% u klasifikaciji sedam različitih emocija koristeći bazu fotografija na kojima ispitanici oponašaju tipične facijalne ekspresije određenih emocionalnih stanja [2]. Tim istraživača iz Rumunije, predvođen Aleksandrom Bandrabur, postigli su tačnost od 92% na 7 emocija koristeći kombinaciju metoda višeslojnih neuralnih mreža i SVM modela [3]. Iako ovakav pristup pruža visok stepen uspešnosti glavni nedostatak jeste oslanjanje na facijalne ekspresije koje nisu nužno pravi pokazatelj stvarnog emotivnog stanja ispitanika. Facijalne ekspresije podležu čovekovoj volji te je sposobnost ljudi u određenim situacijama prikriti svoje stvarno emotivno stanje.

Drugi pristup jeste da se detekcija emocija vrši na osnovu elektrofizioloških signala. Mnogi autori su prepoznali nedostatke obeležja nad kojima čovek ima voljnu kontrolu. Istraživači J. Kim i E. Andre koriste elektrofiziološke signale kao što su EKG, električnu aktivnost mišića (elektromiogram – EMG), električnu provodljivost kože i aktivnost respiratornog sistema [4]. Za pobuđivanje četiri različita emotivna stanja koristili su muzički sadržaj. Za slučaj personalizovane klasifikacije emocija postignuta je tačnost od 95% dok je u slučaju interpersonalne klasifikacije postignuta tačnost jednaka 70%. Iako su elektrofiziološki signali dobri pokazatelji stvarnog emotivnog stanja javlja se problem prikupljanja kvalitetne baze podataka usled poteškoća u izazivanju ciljanih emocija kod ljudi. Premda su autori ovog rada pokazali određene rezultate

njihova se metodologija može dovesti u pitanje. Naime, poznato je da muzika može izazvati emocije ali je ta sposobnost nekonzistentna kod različitih ispitanika. Dok je kod nekih ljudi lako izazvati emocije na ovaj način drugi mogu ostati neosetljivi na ovakav sadržaj. Slično, tim istraživača iz Koreje, K.Kim i S. Bang koristili su kombinaciju audio vizuelnog sadržaja [5]. Ispitanik je puštan video sa muzikom a uz to profesionalni glumac čita kratke emotivne priče. Ispitanici su bili deca u starosti od pet do osam godina te je ovaj pristup bio efektivan. Ipak za očekivati je da bi uspešnost izazivanja emocija kod odraslih bila znatno manja. Da bi se izbegli nedostaci svih navedenih pristupa, za potrebe prepoznavanja emocija, u ovom radu se koriste podaci dobijeni kroz eksperiment koji je vršen u virtuelnoj realnosti.

U drugom poglavlju teze dat je opis metoda za prikupljanje podataka u vidu opisa postavke ekperimenta, kratkog opisa obrade elektrofizioloških signala i opisa ekstrahovanih obeležja iz istih. U trećem poglavlju se daje opis korištenih metoda i algoritama u sprovedenoj studiji. U četvrtom se poglavlju opisuju postignuti rezultati teze, na prvom mestu to je relabeliranje dostupnog skupa podataka, potom analiza obeležja i konačno personalna i interpersonalna klasifikacija kao postupak projektovanja modela predikcije ljudskih emocija.

## 2. Prikupljanje podataka

U toku eksperimenta ispitanici su izloženi različitim VR sadržajima koji su odabrani da izazovu ciljano raspoloženje. Tokom trajanja eksperimenta vrše se merenje elektrofizioloških signala ispitanika. Elektrofiziološki signali nisu pod voljnom kontrolom čoveka te kao takvi predstavljaju prave indikatore njegovog stvarnog emotivnog stanja.

Navedeni eksperiment je realizovan u zajedničkoj studiji sa kolegom Daliborom Veljkovićem, te se detaljniji opis eksperimenta, ekstrakcije obeležja i labeliranja podataka može pronaći u njegovom master radu pod nazivom: "Multimodalna akvizicija i izdvajanje obeležja fizioloških signala za prepoznavanje emocija." U okviru pomenutog rada se može pronaći detaljan opis etapa eksperimenta, prikazanih VR sadržaja kao i korišćene merne opreme obrade. Kao predmet ovog rada je bila obrada elektrofizioloških signala, ekstrakcija obeležja te labeliranje podataka. U ovom poglavlju ćemo dati kratak pregled ovih metoda.

U potpoglavlju 2.1 je opisana postavka eksperimenta. U potpoglavlju 2.2 su navedeni elektrofiziološki signali, korišćeni za ekstrakciju obeležja. U potpoglavlju 2.3 dat je pregled dostupnog skupa podataka u vidu ekstrahovanih obeležja i njima pridruženih labela.

### 2.1 Eksperimentalna procedura

Na početku od ispitanika se zahteva da potpiše svoj lični pristanak za participaciju. Prvo se ispitaniku postavljaju sledeći uređaji:

1. Kapa za akviziciju EEG signala
2. Pojas za merenje srčane aktivnosti i disanja.
3. Naočare za virtuelnu realnost.

Na slici 2.1.1 prikazana je postavka eksperimenta sa ispitanikom.



Slika 2.1.1 Postavka eksperimenta sa ispitanikom

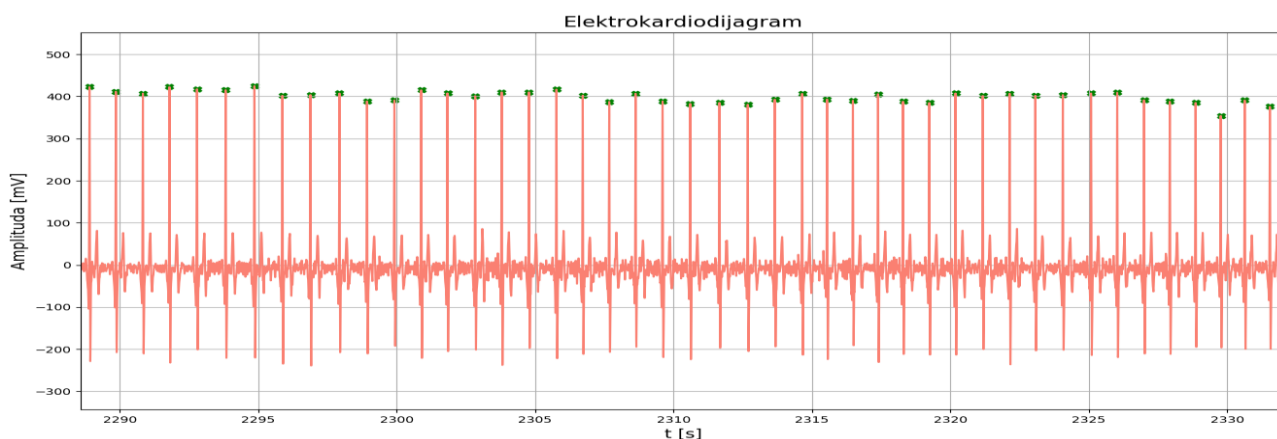
Svaki ispitanik je prošao VR obuku od oko 5 minuta. Eksperiment počinje sa fazom opuštanja. Sledeći sadržaj kojem je ispitanik izložen jeste roler koster. Cilj je da se u ispitaniku izazove strah. Vožnja traje svega 4 minute, nakon čega ispitanik pristupa evaluaciji sopstenog doživljaja kroz interaktivnu aplikaciju. Potom se ispitanik dovodi u prijatno okruženje koje za cilj ima opuštanje i smirivanje njegovih elektrofizioloških signala. U nastavku, ispitanik se vodi u VR streljanu. Cilj je

da se ispitanik skoncentriše na zadatak pred sobom. Ispitanik igra tri runde od po tri minuta. Kao i posle roler kostera ovu etapu slede evaluacija ispitanika i faza opuštanja. Konačno, poslednji sadržaj je interaktivna horor scena koja u proseku traje 6 minuta. Kao i prethodni sadržaji horor scena je praćena samoocenjivanjem ispitanika te poslednjom fazom opuštanja. Trajanje eksperimenta u proseku iznosi šezdeset minuta.

## 2.2 Opis akvizicije i procesiranja signala

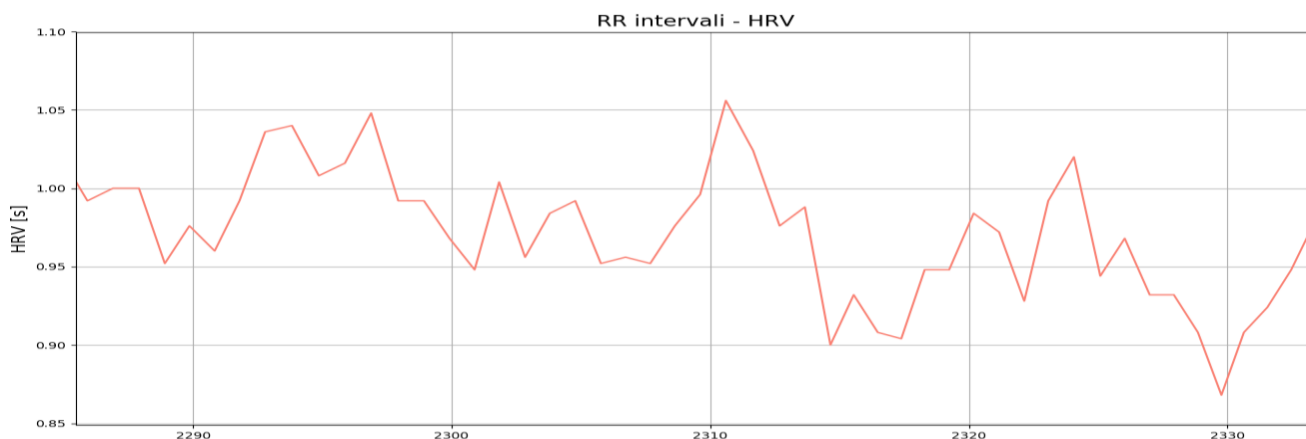
Snimani su elektrofiziološki signali disanja, elektrokardiogram (EKG) i elektroencefalogram (EEG). Signali disanja i EKG su snimani uz pomoć Smartex pojasa koji je opremljen senzorima (Smartex, Piza, Italija) [6]. EEG signal sniman je pomoću sistema u vidu kape koji se sastoji iz EEG pojačavača (Smarting, mBrainTrain, Srbija) koji je povezan na EEG kapu sa 24 elektrode (EasyCap, Nemačka). [7].

Elektrokardiografija je klinička dijagnostička tehnika namenjena praćenju pravilnosti rada srčanog mišića. Akvizicija EKG signala se obavlja učestanošću od 250 Hz. Potom se vrši procesiranje EKG signala u vidu filtracije i ekstrakcije pikova, poznatih kao R pikovi. Na slici 2.2.1 prikazana je filtriran EKG signal sa označenim lokacijama R pikova.



Slika 2.2.1 Filtrirani EKG signal

Konačno, signal RR intervala dobija se kao period između svaka dva sukcesivna R pika. HRV (*Heart Rate Variability*) signal predstavlja vremensku seriju RR intervala. Nastaje tako što se svakom odbirku signala RR intervala dodeli vremenska labela pozicije R pika. HRV signal je prikazan na slici 2.2.2 za isti vremenski interval kao i slika 2.2.1.

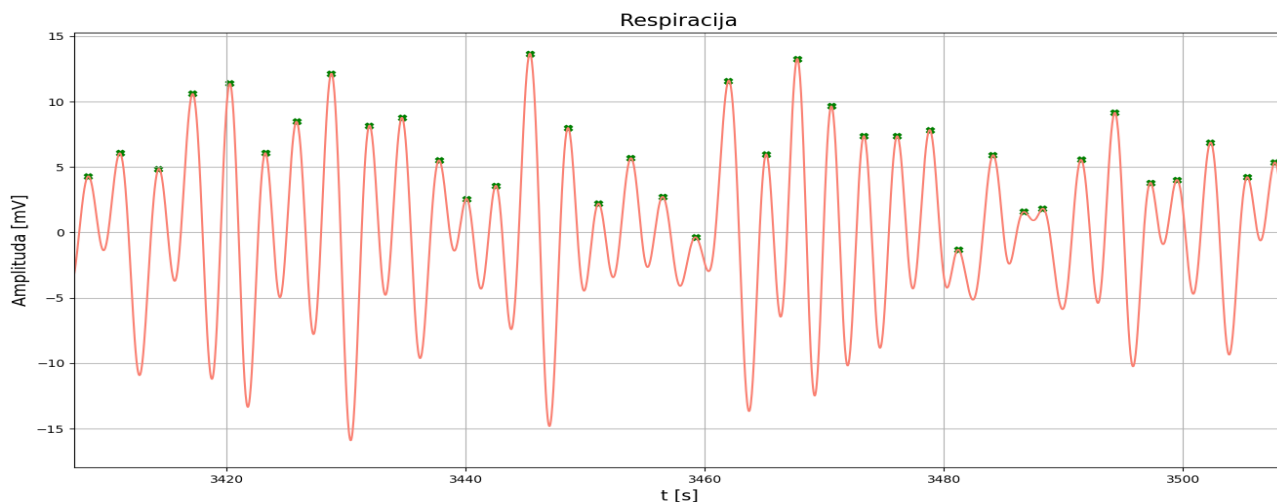


Slika 2.2.2 Signal HRV

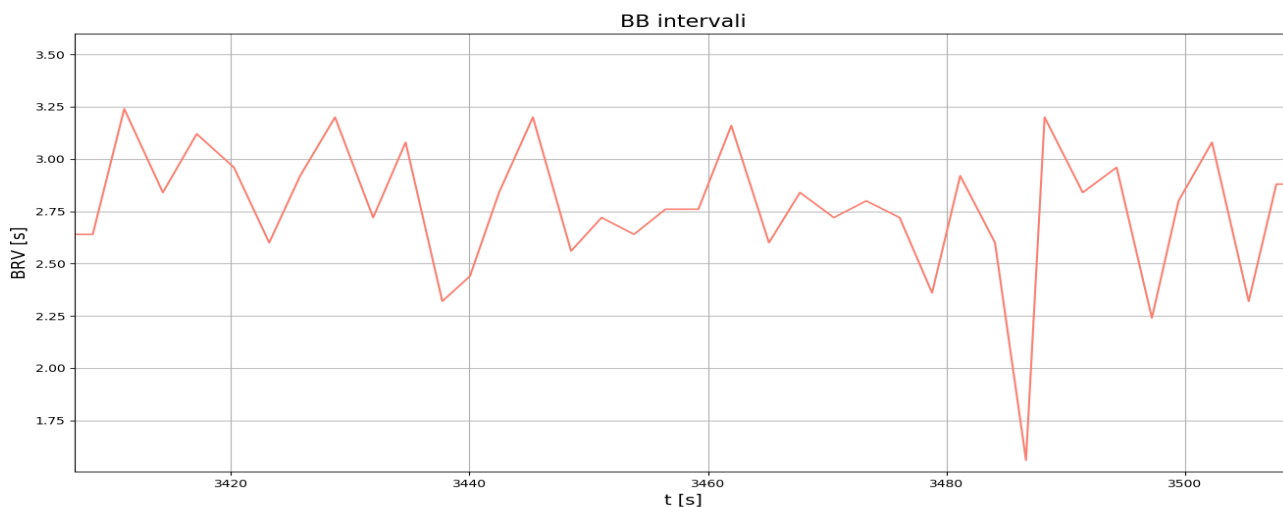


Signal respiracije je odabiran učestanošću od 25 Hz. Usled položaja elastičnog senzora kriva respiracije se direktno odnosi na trenutnu zapreminu grudnog koša. Obrada sirovog signala podrazumeva filtraciju i izračunavanje učestanosti prolaska signala kroz nulu (*Zero Crossing Rate*).

Na osnovu trenutaka u kojima signal ima nultu vrednost izračunavaju se pikovi respiratornog signala. Filtriran signal respiracije sa označenim pikovima je prikazan na slici 2.2.3. Ekvivalentno izračunavanju HRV signala računa se i BRV (*Breathing Rate Variability*) signal na osnovu BB (*Breath to Breath*) intervala, prikazan na slici 2.2.4.

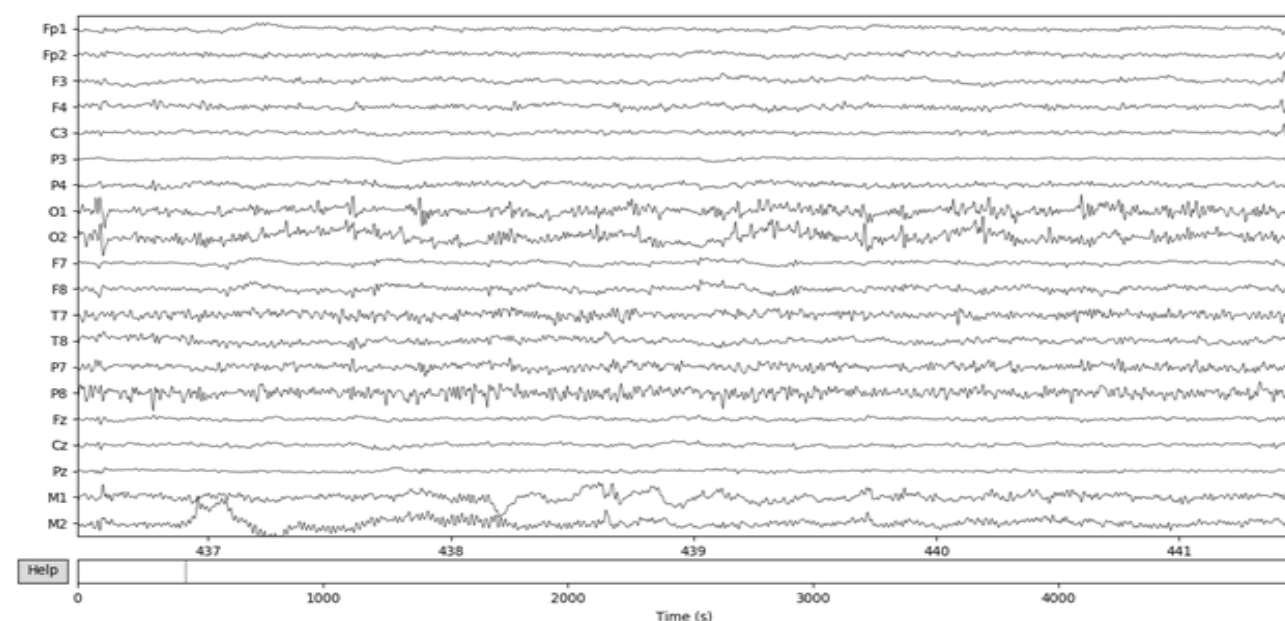


**Slika 2.2.3 Filtriran signal respiracije**



**Slika 2.2.4 BRV (*Breathing Rate Variability*) signal**

Elektroencefalografija (EEG) predstavlja elektrofiziološku metodu snimanja električne aktivnosti mozga. Snimljeni EEG signal je filtriran filtrom propusnikom opsega učestanosti a potom su artefakti koji potiču od očnih pokreta uklonjeni ICA metodom [8]. Na slici 2.2.5 prikazan je segment EEG signala posle filtracije i odstranjivanja očnih atrefakata.



Slika 2.2.5 EEG segment posle filtracije i odstranjivanja očnih atrefakata

## 2.3. Opis dobijanja skupa podataka

Skup podataka se sastoji iz obeležja HRV i BRV signala te njima pridruženih labela koje se odnose na referentno emotivno stanje ispitanika (opuštenost, fokusiranost i strah). Ekstrahovana obeležja dobijanu iz HRV i BRV signala koji su rezultat akvizicije i obrade EKG signala i signala respiracije, što je opisano u prethodnom poglavlju. Izdvajanje obeležja se vršilo nad kratkovremenskim signalom unutar prozora određenog trajanja.

Kod HRV signala se posmatra serija R pikova unutar prozorske funkcije. Ekstrahovana su linearna i nelinearna obeležja iz vremenskog i frekvencijskog domena. Obeležja su navedena u tabeli 2.3.1.

Tabela 2.3.1 Pregled obeležja dobijenih iz HRV signala

Ime	Tip	Opis
<i>meanRR</i>	Vremenski domen	Srednja vrednost RR intervala (procena matematičkog očekivanja)
<i>stdRR</i>	Vremenski domen	Standardna devijacija RR intervala
<i>cvRR</i>	Vremenski domen	Odnos između standardne devijacije i srednje vrednosti RR intervala
<i>pRR<sub>t</sub></i>	Vremenski domen	Udeo RR intervala čija sukcesivna razlika iznosi najmanje $t$ milisekundi
<i>rmssdRR</i>	Vremenski domen	Srednji kvadratni koren sukcesivnih razlika između RR intervala
<i>LF<sub>power</sub></i>	Frekvencijski domen	Snaga HRV signala na opsegu učestanosti od 0.04Hz do 0.15Hz
<i>HF<sub>power</sub></i>	Frekvencijski domen	Snaga HRV signala na opsegu učestanosti od 0.15Hz do 0.4Hz
<i>LFHF<sub>ratio</sub></i>	Frekvencijski domen	Količnik obeležja <i>LF<sub>power</sub></i> i <i>HF<sub>power</sub></i>
<i>SD1</i>	<i>Poincare</i> grafik	Dužina sporedne poluose elipse opisane oko <i>Poincare</i> oblaka

<b><i>SD2</i></b>	<i>Poincare</i> grafik	Dužina primarne poluose elipse opisane oko <i>Poincare</i> oblaka
<b><i>CCM</i></b>	<i>Poincare</i> grafik	Mera kompleksne korelacije <i>Poincare</i> grafika
<b><i>HRA_IL</i></b>	<i>Poincare</i> grafik	Gruzikov indeks simetričnosti <i>Poincare</i> oblaka
<b><i>HRA_mod</i></b>	<i>Poincare</i> grafik	Portin indeks simetričnosti <i>Poincare</i> oblaka
<b><i>ApEn</i></b>	Entropija	Procena uređenosti (entropije) vremenske serije
<b><i>SampEn</i></b>	Entropija	Procena uređenosti (entropije) vremenske serije

Ove varijable se mogu računati na osnovu celokupnog EKG signala, trajanja na primer 24 sata, ili na osnovu kratkih segmenata u trajanju od samo minut do pet minuta.

Obeležja iz vremenskog domena, navedena u tabeli 2.3.1, su najjednostavniji tip obeležja i dele se u dve grupe:

- Obeležja dobijena na osnovu direktnih vrednosti RR intervala (*meanRR*, *stdRR*, *cvRR*)
- Obeležja dobijena na osnovu razlika sukcesivnih RR intervala (*cvRR*, *pRR<sub>t</sub>*)

Obeležje *pRR<sub>t</sub>* je računato za  $t = \{30ms, 50ms, 70ms\}$  [9].

Frekvencijska obeležja su dobijena iz estimirane spektralne snage HRV signala. Za njeno izračunavanje korištena je neparametarska FFT (*Fast Fourier transform*) metoda [10]. Navedena obeležja se dobijaju integraljenjem snage unutar odgovarajućeg frekvencijskog opsega.

*Poincare* grafik predstavlja geometrijski prikaz vremenske serije u vidu zavisnosti dužine svakog RR intervala od njegovog prethodnika [11]. Tačke *Poincare* grafika čine *Poincare* oblak iz koga se važna obeležja ekstrahuju računanjem parametara elipse koja se može opisati oko oblaka. Pored ovih parametara računa se mera kompleksne korelacije kao i simetričnost *Poincare* oblaka [12][13].

Entropija predstavlja meru uređenosti i kompleksnosti. Navedena obeležja predstavljaju estimaciju entropije na kratkim vremenskim serijama [14]. Dakle, za računanje entropije usvojena su dva pristupa, a usled njihove informativnosti oba su zadržana.

Preostala obeležja dobijaju se iz BRV signala. Zbog suštinske sličnosti HRV i BRV signala postupak ekstrahovanja obeležja je identičan. U tabeli 2.3.2 prikazana su obeležja BRV signala.

**Tabela 2.3.2 Pregled obeležja dobijenih iz BRV signala**

Ime	Tip	Opis
<b><i>meanBB</i></b>	Vremenski domen	Srednja vrednost BB intervala (procena matematičkog očekivanja)
<b><i>stdBB</i></b>	Vremenski domen	Standardna devijacija BB intervala
<b><i>cvBB</i></b>	Vremenski domen	Odnos između standardne devijacije i srednje vrednosti BB intervala
<b><i>pBB<sub>t</sub></i></b>	Vremenski domen	Udeo BB intervala čija sukcesivna razlika iznosi najmanje $t$ milisekundi
<b><i>rmssdB</i></b>	Vremenski domen	Srednji kvadratni koren sukcesivnih razlika između BB intervala
<b><i>ULF<sub>power</sub></i></b>	Frekvencijski domen	Snaga signala respiracije na opsegu učestanosti manje od 0.2Hz
<b><i>LF<sub>power</sub></i></b>	Frekvencijski domen	Snaga signala respiracije na opsegu učestanosti od 0.2Hz do 0.325Hz
<b><i>HF<sub>power</sub></i></b>	Frekvencijski	Snaga signala respiracije na opsegu učestanosti od 0.325Hz do 0.8Hz

	domen	
$UHF_{power}$	Frekvencijski domen	Snaga signala respiracije na opsegu učestanosti veće od 0.8Hz
$SD1$	<i>Poincare</i> grafik	Dužina sporedne poluose elipse opisane oko <i>Poincare</i> oblaka
$SD2$	<i>Poincare</i> grafik	Dužina primarne poluose elipse opisane oko <i>Poincare</i> oblaka
$CCM$	<i>Poincare</i> grafik	Mera kompleksne korelacije <i>Poincare</i> grafika
$HRA_{IL}$	<i>Poincare</i> grafik	Gruzikov indeks simetričnosti <i>Poincare</i> oblaka
$HRA_{mod}$	<i>Poincare</i> grafik	Portin indeks simetričnosti <i>Poincare</i> oblaka
$ApEn$	Entropija	Procena uređenosti (entropije) vremenske serije
$SampEn$	Entropija	Procena uređenosti (entropije) vremenske serije

Obeležje  $pBB_t$  je računato za  $t = \{1500ms, 2500ms, 3500ms\}$ .

Potrebno je primetiti da se frekvencijska obeležja ekstrahuju iz neobrađenog signala disanja pomoću FFT metode, a ne iz BRV signala [15].

Sva navedena obeležja imaju bolji kvalitet kada se ekstrahuju na prozoru većeg trajanja. Svakako, zahtevi obeležja signala disanja su strožiji usled sporije dinamike ovog procesa. Izabran je prozor u trajanju od 60 sekundi kako obeležja oba tipa bila zadovoljavajućeg kvaliteta. Odbirci obeležja se dobijaju izračunavanjem na prozoru navedenog trajanja koji se pomera duž vremenske ose HRV i BRV signala za po 2 sekunde,

Na osnovu samoocena ispitanika i posmatranih frekvencijskih karakteristika EEG signala procenjuje se i referentno emotivno stanje ispitanika tokom VR sadržaja. Na ovaj način se svakom odbirku obeležja hronološki pridružuje jedna od labela: opuštenost, fokusiranost ili strah. Dostupne labele se odnose na konkretne sadržaje dok se obeležja formiraju na osnovu celokupnih snimljenih elektrofizioloških signala. Usled toga, dostupni skup podataka sadrži i labelirane i nelabelirane segmente. Važno je naglasiti da se dostupne labele odlikuju pridruženom visokom ili niskom pouzdanošću, koja uzima vrednost od 0 do 10.

## 3. Metode

### 3.1 Metode relabeliranja podataka

U mašinskom učenju je čest problem postojanje skupova podataka koji sadrže nelabelirane podatke. Vrlo često se u realnim studijama labeliranje vrši ručno od strane eksperata iz konkretne oblasti. Pri tome se dešava pojava da određeni podaci ostanu nelabelirani usled nedostatka dokaza na koje bi se eksperti oslonili. Takođe, dok labeliranje određenih podataka predstavlja trivijalan zadatak, postoje i odbirci kod kojih eksperti ne mogu izvršiti labeliranje sa visokom pouzdanošću. Pored ekspertize ovakav postupak iziskuje ulaganje značajnog napora.

Sa druge strane, za formiranje kvalitetnog modela predikcije, potreban je skup podataka što većeg obima. Za postizanje performansi modela takođe je pogubno imati veliki udeo pogrešno labeliranih podataka. Iz ovog razloga pristupa se postupku relabeliranja kako nelabeliranih podataka tako i podataka čije se labele karakterišu niskom pouzdanošću.

#### 3.1.1 Polu-supervizirano klasterovanje

Polu-supervizirano klasterovanje predstavlja automatizovan postupak labeliranja podataka a na osnovu raspoloživih podataka koji imaju pouzdane labele. U ovoj tezi se vrši pomoću algoritma *label spreading* [16]. Ovaj algoritam koristi labelirane podatke i na osnovu njih dodeljuje labelu nelabeliranim podacima.

Neka je  $X$  skup svih odbiraka gde je  $X = \{x_1, \dots, x_l, x_{l+1}, \dots, x_n\}$  dok su raspoložive labelu  $L = \{1, \dots, c\}$ . Neka je labelirano samo prvih  $l$  tačaka skupa  $X$  i neka su njihove labelu  $\{y_1, \dots, y_l\} \in L$ .

Ovaj algoritam uvodi dve pretpostavke o lokalnoj konzistentnosti:

1. Odbirci koji su međusobno blizu u prostoru imaju povećanu verovatnoću da budu istovetno labelirani
2. Odbirci koji pripadaju istoj strukturi (klasteru) imaju povećanu verovatnoću da budu istovetno labelirani

Na osnovu uvedenih pretpostavki, mehanizam ovog algoritma se sastoji u tome da, iterativno, svaki odbirak proširi „informaciju o labeliranju“ koju poseduje na susedne odbirke.

Neka je  $F$  matrica dimenzija  $n \times c$  čiji su svi elementi nenegativni,  $F = \{F_1^T, \dots, F_n^T\}$ . Matrica  $F$  predstavlja svojevrsni rezultat klasifikacije odbiraka iz skupa  $X$  gde je labela odbirka  $x_i$  data kao  $y_i = \operatorname{argmax}_{j \leq c} (F_{ij})$ . Takođe, definišemo matricu  $Y$  dimenzija  $n \times c$  gde je  $Y_{ij} = 1$  ukoliko je odbiraku  $x_i$  dodeljena labela. Matrica  $Y$  sadrži informaciju o postojećim labelama podataka.

Algoritam se sastoji iz sledećih koraka:

1. Prvo se formira matrica afiniteta  $W$  dimenzija  $n \times n$  čiji se elementi računaju kao:

$$w_{ij} = \begin{cases} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}, & i \neq j \\ 0, & i = j \end{cases} \quad (1)$$

Intuitivno značenje prvog koraka iterativne procedure jeste formiranje grafa afiniteta između odbiraka. Čvorovi ovog grafa jesu sami odbirci dok su težine grana jednake elementima matrice  $W$ .

2. U drugom koraku vrši se simetrična normalizacija matrice  $W$ . Ovim postupkom dobija se matrica  $S$ :

$$S = D^{-\frac{1}{2}} W D^{-\frac{1}{2}} \quad (2)$$

gde je  $D$  dijagonalna matrica čiji je elemenat na poziciji  $(i,i)$  jednak sumi svih elemenata u vrsti  $i$  matrice  $F$ .

3. Treći korak predstavlja iterativno izračunavanje matrice  $F$  do konvergencije:

$$F(t+1) = \alpha S F(t) + (1-\alpha) Y \quad (3)$$

U svakoj iteraciji ovog koraka, referenca svakog odbirka se dobija na osnovu informacije koju nose ostali odbirci (prvi sabirak) i na osnovu inicijalne reference za dati odbirak (drugi sabirak). Parametar  $\alpha$  ima vrednosti u intervalu  $[0-1]$  i kontroliše način formiranja reference za posmatrani odbirak. Ukoliko je vrednost ovog parametra bliska jedinici veći uticaj na formiranje reference posmatranog odbirka ima informacija koju nose preostali dobirci (prvi sabirak). U suprotnom referenca se dominantno formira na osnovu postojeće labele ukoliko je ona data.

4. Neka je  $F^*$  finalna vrednost matrice  $F$  (posle konvergencije). Labela za svaki odbirak  $x_i$  je jednaka  $y_i = \operatorname{argmax}_{j \leq c} (F_{ij}^*)$

## 3.2 Metode analize obeležja

Pre nego što se obeležja iskoriste za obučavanje i testiranje prediktivnih modela potrebno steći bolji uvid u njihove karakteristike. Važno je proceniti informativnost obeležja u odnosu na pridružene labele kao i međusobnu zavisnost i redundantnost. Ovakva analiza se sprovodi u cilju dobijanja određene intuicije koja pomaže pri odabiru modela predikcije, tumačenju postignutih performansi kao i pri odabiru hiper-parametara modela. Šta više, ovaj pristup nudi još dve mogućnosti. Prva se ogleda u detekciji i odstranjivanju loših obeležja dok druga nudi mogućnost redukcije dimenzija u vidu transformacije prostora obeležja u nižedimenzioni prostor.

### 3.2.1 PCA analiza

Čest je slučaj da se najveći deo informacije koju nose obeležja nalazi u nižedimenzionom prostoru (u odnosu na dimenzionalnost obeležja) usled redundantnosti informacije unutar obeležja. PCA (*Principal component analysis*) je algoritam koji vrši transformaciju prostora takvu da se od raspoloživih, potencijalno korelisanih obeležja dobije novi skup obeležja koja su međusobno linearno nekorelisana (principijalne komponente) [17].

Ova transformacija se vrši tako da se prva principijelna komponenta prostire duž pravca u prostoru u kome je varijansa podataka najveća te samim tim ova komponenta sadrži najveću količinu informacije iz originalnih obeležja (količina informacije se vezuje za varijansu). Sledeća principijelna komponenta zauzima pravac ortogonalan na pravac prve principijelne komponente sa maksimalnom varijansom. Pravac treće komponente je pravac u prostoru duž koga je maksimalna varijansa podataka a da je pri tome ortogonalan na pravac prve dve komponente itd. (ista pravilnost važi za sve naredne principijelne komponente).

Dakle, PCA izdvaja principijelne komponente takve da količina informacije koju imaju opada od prve ka poslednjoj. Redukcija dimenzija se može izvršiti odbacivanjem principijelnih komponenti koje sadrže najmanji deo informacije originalnih obeležja. Sada će biti opisan način funkcionisanja PCA algoritma.

Pre primene samog algoritma važno je naglasiti da prvo sva obeležja treba dovesti na nultu srednju vrednost i jediničnu varijansu. Ovo se postiže tako što se od svakog odbirka nekog obeležja oduzme procenjena srednja vrednost tog obeležja pa se zatim svaki odbirak obeležja podeli sa procenjenom varijansom tog obeležja. Cilj ovog postupka je da vrednosti svih obeležja postanu međusobno uporedive kako bi se adekvatno mogli odrediti pravci prostora sa najvećom varijansom.

Pronalaženje prve principijelne komponente se svodi na traženje pravca u prostoru gde je varijansa podataka najveća. Ovo se postiže traženjem jediničnog vektora  $\vec{u}$  takvog da je zbir projekcija (kvadrata) na njega maksimalna tj. maksimizujemo:

$$\frac{1}{m} \sum_{i=1}^m (x_i^T u)^2 = \frac{1}{m} \sum_{i=1}^m u^T x_i x_i^T u = u^T \Sigma u$$

Gde  $\Sigma$  predstavlja kovarijacionu matricu vektora obeležja  $x$ . Vektor  $\vec{u}$  se određuje maksimizacijom kriterijumske funkcije:

$$u^T \Sigma u - p(u^T u - 1)$$

Izjednačavanjem izvoda po  $u$  sa nulom dobija se jednakost:

$$\Sigma u = pu$$

tj. da vektor  $\vec{u}$  predstavlja principijelni sopstveni vektor kovarijacione matrice dok je prva principijelna komponenta projekcija podataka na vektor  $\vec{u}$ . Na sličan način se pokazuje da ostale PCA komponente odgovaraju projekcijama podataka na preostale sopstvene vektore matrice  $\Sigma$ .

Važnost svake PCA komponente (tj. količina informacije/varijanse koju komponenta nosi) jednaka je odgovarajućoj sopstvenoj vrednosti sopstvenog vektora koji se koristi za transformaciju pri dobijanju te komponente. Redukcija dimenzija se vrši tako što se za set novih obeležja izabere projekcija originalnih obeležja na prvih  $k$  sopstvenih vektora matrice  $\Sigma$  gde je  $k < n$ .

### 3.2.2 LDA analiza

LDA predstavlja modifikaciju PCA algoritma, opisanog u prethodnom odeljku, gde se originalna obeležja transformacijom prostora prevode u novi skup međusobno nekorelisanih obeležja ali sa ciljem da se maksimizuje separabilnost klasa, tj. da se minimizuje unutarklasno rasejanje a da se maskimizuje međuklasno rasipanje [17].

Matrica unutarklasnog rasejanja računa se kao:

$$S_w = \sum_{i=1}^g N_i \Sigma_i$$

gde je  $g$  ukupan broj klasa a  $N_i$  broj dobiraka svake klase a  $\Sigma_i$  kovarijaciona matrica klase  $i$ .

Matrica međuklasnog rastojanja se računa kao:

$$S_b = \sum_{i=1}^g N_i (M_i - M_0) (M_i - M_0)^T$$

LDA se dobija modifikacijom PCA algoritma koja se ogleda u tome da se pravci prostora na koje se projektuju podaci dobijaju kao sopstvene vrednosti matrice  $S_b S_w^{-1}$  umesto kovarijacione matrice obeležja.

### 3.2.3 Procena korelisanosti

Korelacija se procenjuje Pirsonovim koeficijentom korelacije (*Pearson correlation coefficient* – PCC) [18]. Za dve slučajne varijable  $X$  i  $Y$  ovaj uobičajeni statistički deskriptor je definisan sa:

$$\rho_{X,Y} = \frac{E[(X - E[X])(Y - E[Y])]}{\sigma_X \sigma_Y} \quad (1)$$

Za slučaj vektora odabiraka definiše se kao:

$$r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} \quad (2)$$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i \quad (3)$$

Korelisanost se odnosi na stepen linearne zavisnosti varijabli. Prema tome, ukoliko dve varijable nisu korelisane to ne znači da su one i nezavisne, već samo svedoči da ne postoji međusobna linearna zavisnost.

### 3.2.4 Procena informativnosti

U teoriji verovatnoće *Mutual information* ( $MI$ ) dve slučajne varijable predstavlja meru međusobne zavisnosti te dve varijable [19]. Tačnije,  $MI$  predstavlja, meru u kojoj se informacija koja je sadržana u jednoj slučajnoj varijabli može predstaviti pomoću druge slučajne varijable.

Način funkcionisanja  $MI$  algoritma je zasnovan na određivanju sličnosti između združene raspodele varijabli  $p(X, Y)$  i proizvoda  $p(X)p(Y)$ , te se prema tome  $MI$  izračunava kao:

$$MI(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log_2 \left( \frac{p(x, y)}{p(x)p(y)} \right) \quad (4)$$

Navedena formula se koristi za izračunavanje  $MI$  kod diskretnih varijabli. Postupak je identičan za kontinualne varijable pri čemu se umesto sumiranja vrši integraljenje.

$$MI(X, Y) = \int_{y \in Y} \int_{x \in X} p(x, y) \log_2 \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy \quad (5)$$

Iz formula 4 i 5 lako je zaključiti da  $MI$  ima nultu vrednost u slučaju kada su slučajne varijable međusobno nezavisne jer je u tom slučaju  $p(x, y) = p(x)p(y)$  iz čega sledi:

$$\log_2 \left( \frac{p(x, y)}{p(x)p(y)} \right) = \log_2(1) = 0 \quad (6)$$

Dakle,  $MI$  ima nultu vrednost u slučaju nezavisnih slučajnih varijabli a vrednost  $MI$  je utoliko veća ukoliko je međusobna zavisnost između slučajnih varijabli veća.

## 3.3 Metode klasifikacije

Metode klasifikacije na osnovu labeliranih podataka formiraju modele predikcije koji opisuju zavisnost izlaznih varijabli (labela) u odnosu na ulazne varijable (obeležja). Ove metode podrazumevaju supervizijsko učenje ovih zavisnosti, tj. formiranje modela uz prisustvo reference. Jednom formiran model enkapsulira konkretno ekspertsko znanje na osnovu koga kasnije vrši



predikciju nad nelabiranim podacima. U ovom poglavlju opisana je metoda logističke regresije, kao najjednostavnije metode za klasifikaciju, dok su potom opisane znatno naprednije metode kao što su metoda nosećih vektora (SVM) i metoda slučajnih šuma (*Random Forest*).

### 3.3.1 Logistička regresija – softmax metoda

Logistička regresija predstavlja linearni model za binarnu klasifikaciju. *Softmax* metoda predstavlja generalizaciju ovog algoritma i služi za multinomijalnu linearnu klasifikaciju [17]. *Softmax* metoda vrši klasifikaciju tako što odredi verovatnoću sa kojom određeni odbirak pripada svakoj od klasa te taj odbirak dodeli onoj klasi za koju ima najveću verovatnoću pripadanja.

Neka je skup raspoloživih labela  $y \in \{1, 2, \dots, k\}$ . Da bi se parametrizovala klasifikacija nad  $k$  mogućih ishoda uvodi se  $k$  parametara  $\phi_1, \phi_2, \dots, \phi_k$  gde svaki od parametara označava verovatnoću pripadnosti odbiraka određenoj klasi. Kako odbirak mora pripadati nekoj od klasa jasno je da među navedenim parametrima postoji redundantnost i da važi:

$$\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i \quad (1)$$

te se multinomijalna raspodela karakteriše parametrima  $\phi_1, \phi_2, \dots, \phi_{k-1}$ .

Prvi korak u objašnjenju principa rada *softmax* regresije jeste dokaz da multinomijalna raspodela pripada grupi eksponencijalnih raspodela [17].

Vektorska reprezentacija pripadnosti semplota klasama vrši se pomoću vektora  $T(y)$ . Navedeni vektor se odnosi na odbirak  $i$  i njegova dimenzionalnost je jednaka  $k-1$ . Ovaj vektor ima vrednost 1 na poziciji čiji indeks odgovara labeli klase kojoj odbirak pripada i vrednost 0 na svim ostalim pozicijama. Element  $i$  vektora  $T(y)$  označava se kao  $(T(y))_i$ . Takođe za nastavak izvođenja uvodi se sledeća notacija  $1\{true\} = 1, 1\{false\} = 0$ .

Dalje važi:

$$p(y; \phi) = \phi_1^{1\{y==1\}} \phi_1^{1\{y==2\}} \dots \phi_1^{1\{1-\sum_{i=1}^{k-1} \phi_i\}} = \phi_1^{(T(y))_1} \phi_1^{(T(y))_2} \dots \phi_1^{1-\sum_{i=1}^{k-1} (T(y))_i} \quad (2)$$

$p(y; \phi)$  se dalje može predstaviti kao:

$$p(y; \phi) = e^{(T(y))_1 \log(\phi_1) + (T(y))_2 \log(\phi_2) + \dots + (1 - \sum_{i=1}^{k-1} (T(y))_i) \log(\phi_k)} \quad (3)$$

finalno, dobija se:

$$p(y; \phi) = e^{(T(y))_1 \log\left(\frac{\phi_1}{\phi_k}\right) + (T(y))_2 \log\left(\frac{\phi_2}{\phi_k}\right) + \dots + \log(\phi_k)} = b e^{\eta^T T(y) - a(\eta)} \quad (4)$$

iz formule (4) vidi se da  $p(y; \phi)$  može da se predstavi u vidu raspodele iz familije eksponencijalnih raspodela pri čemu je  $b = 1$ ,  $a(\eta) = -\log(\phi_k)$  i

$$\eta = \begin{bmatrix} \log\left(\frac{\phi_1}{\phi_k}\right) \\ \vdots \\ \log\left(\frac{\phi_{k-1}}{\phi_k}\right) \end{bmatrix} \quad (5)$$

Iz formula (5) i (1) sledi:

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \quad (6)$$

mapiranje  $\eta_i$  u  $\phi_1$  naziva se *softmax* funkcijom.

S obzirom da je *softmax* regresija linearan model i veza između  $\eta_i$  i odbiraka mora biti linearna  $\eta_i = \theta_i^T x$ . Dakle podešavanje modela se sastoji u podešavanju  $k-1$  parametara  $\theta_1, \theta_2, \dots, \theta_{k-1}$ , takođe važno je napomenuti da se uzima  $\theta_k = 0$ . Na osnovu dosadašnjeg izvođenja sledi:

$$p(y = i | x; \theta) = \phi_i = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}} \quad (7)$$

Recimo da se obučavajući skup sastoji od  $m$  labeliranih odbiraka  $(x^{(i)}, y^{(i)})$ ,  $i = 1 \dots m$ . Parametri modela se podešavaju tako što se maskimizira funkcija verodostojnosti obučavajućeg skupa.

$$l(\theta) = \sum_{i=1}^m \log \left( p(y^{(i)} | x^{(i)}; \theta) \right) = \sum_{i=1}^m \log \prod_{l=1}^k \left( \frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1_{\{y^{(i)}=l\}}} \quad (8)$$

maksimizacija ove funkcije vrši se metodom gradijentnog uspona.

Nakon podešavanja parametara izlaz klasifikatora dat je kao  $\phi_1, \phi_2, \dots, \phi_k$  za odbirak koji se klasifikuje.

### 3.3.2 SVM metoda

Support vector machine (SVM) predstavlja jedan od najmoćnijih i najkorišćenijih modela za klasifikaciju [20].

Neka je dat skup podataka:

$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3) \dots (x_N, y_N)\} \in \mathbb{R}^m$$

gde je  $m$  broj obeležja.

Neka je cilj odrediti linearnu funkciju  $f(x) = \langle \omega, x \rangle + b$ ,  $\omega \in \mathbb{R}^m, b \in \mathbb{R}$ , gde  $\langle \cdot, \cdot \rangle$  predstavlja operator unutrašnjeg proizvoda a funkcija zadovoljava nejednakost  $\forall i, |y_i - f(x_i)| \leq \varepsilon$ . Ovaj problem se može zapisati kao konveksan problem:

$$\min \left\{ \frac{1}{2} \|\omega\|^2 \right\}$$

$$\text{ograničenja: } \begin{cases} y_i - f(x_i) \leq \varepsilon \\ f(x_i) - y_i \leq \varepsilon \end{cases}$$

Uz pretpostavku da je ovakva funkcija postoji. Ponekad ovo nije slučaj te želimo da omogućimo da greška za određene odbirke bude veća od  $\varepsilon$ . Uvode se varijable  $\xi_i$  i  $\xi_i^*$  da bi se izborili sa nedostižnim ograničenjima na sledeći način:

$$\min \left\{ \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^N \xi_i + \xi_i^* \right\}$$

$$\text{ograničenja: } \begin{cases} y_i - f(x_i) \leq \varepsilon + \xi_i \\ f(x_i) - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases}$$

Konstanta  $C > 0$  predstavlja konstantu regularizacije te kao takva predstavlja jednu od najbitnijih hiper-parametara ovog modela. Ona određuje u kojoj meri su greške dopustive. Isuviše male vrednosti ove konstante dopustiće velike količine grešaka. Prevelika konstanta  $C$  rezultuje sa strogim problemom koji ne dopušta greške.

Formira se Lagranžijanova funkcija:

$$L = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^N \xi_i + \xi_i^* - \sum_{i=1}^N \eta_i \xi_i + \eta_i^* \xi_i^* - \sum_{i=1}^N \alpha_i^* (\varepsilon + \xi_i^* - y_i + \langle \omega, x \rangle + b) - \sum_{i=1}^N \alpha_i (\varepsilon + \xi_i - y_i + \langle \omega, x \rangle + b)$$

Gde su  $\eta_i, \eta_i^*, \alpha_i, \alpha_i^* \geq 0$  Lagranžijanovi multiplikatori.

$$\begin{aligned} \partial_b L &= \sum_{i=1}^N \alpha_i^* - \alpha_i = 0 \\ \partial_\omega L &= \omega - \sum_{i=1}^N (\alpha_i^* - \alpha_i) x_i = 0 \\ \partial_{\xi_i} L &= C - \alpha_i - \eta_i, \partial_{\xi_i^*} L = C - \alpha_i^* - \eta_i^* \end{aligned}$$

Odavde se dobija:

$$\begin{aligned} \alpha_i + \eta_i &= C, \alpha_i^* + \eta_i^* = C \\ \omega &= \sum_{i=1}^N (\alpha_i^* - \alpha_i) x_i \text{ i } \alpha_i^*, \alpha_i \in [0, C] \\ f(x) &= \sum_{i=1}^N (\alpha_i^* - \alpha_i) \langle x_i, x \rangle + b \end{aligned}$$

Možemo zaključiti da samo odbirci sa  $\alpha_i^* = C \vee \alpha_i = C$  leže izvan  $\varepsilon$  pojasa, tj. imaju grešku veću od  $\varepsilon$ . Odbirci sa  $\alpha_i^* = 0 \wedge \alpha_i = 0$  se nalaze unutar pomenutog pojasa. Odbirci sa  $\alpha_i^* \in (0, C) \vee \alpha_i \in (0, C)$  se nalaze tačno na udaljenosti  $\varepsilon$  od funkcije te kao takvi predstavljaju noseće vektore modela.

Nelinearnost u algoritam se unosi pomoću kernelskih funkcija. U prethodnim izrazima unutrašnji proizvod se menja kernelskom funkcijom  $k(x_i, x)$ . Uobičajeni oblici kernelskih funkcija su:

$$\begin{aligned} k^{poli}(x_i, x) &= \langle x_i, x \rangle^p \\ k^{tanh}(x_i, x) &= \tanh(v + k \langle x_i, x \rangle) \\ k^{rbf}(x_i, x) &= e^{-\frac{\|x_i - x\|^2}{2a^2}} \end{aligned}$$

Dok je rezultat kernela  $k^{poli}(x_i, x)$  ekvivalentan transformaciji vektora obeležja iz originalnog prostora u prostor stepena  $p$  druga dva kernela transformišu originalni prostor stanja u prostor beskonačne dimenzionalnosti. Parametri kernelskih funkcija takođe predstavljaju hiper-parametre modela. U zavisnosti od prirode funkcije  $f(x)$  SVM može vršiti klasifikaciju ili regresiju.

Za slučaj nebalansiranih klasa postoji mogućnost da se hiper-parametar regularizacije  $C$  iskoristi za rešavanje ovog problema. Definiše se  $C$  za svaku klasu zasebno. Uobičajeno je da se definiše jedinstveno  $C$  a posebno koeficijenti balansiranja koji, za pojedinačne klase, množe  $C$  te na taj način definišu hiper-parametar regularizacije za konkretnu klasu. Ovo omogućava da je algoritam strožiji po pitanju klasa sa manje odbiraka tako što se datoj klasi dodeli veća težina i obratno.

### 3.3.3 Random Forest metoda

*Random Forest* je supervizirani algoritam učenja koji predstavlja ansambl stabala odlučivanja obučanih primenom *bagging* algoritma [21]. Stablo odlučivanja predstavlja graf nalik stablu. Na svakom čvoru grafa se donosi odluka, dok su ishodi predstavljeni granama grafa. Odluke nastaju upoređivanjem nekog obeležja u odbirku sa pragom određene vrednosti. Na osnovu da li je obeležje veće ili manje od datom praga skup odbiraka se deli (grana) na sva podskupa. Dakle, svakom odlukom vrši se podela obučavajućeg skupa u cilju povećanja klasne homogenosti podskupova. Čvorovi koji nemaju izlazeće grane se nazivaju listovima. Listovi stabla odluke nose klasne labele. Na taj način čitav skup podataka biva kanalisiran do listova gde im se određuje klasna pripadnost.

Bagging algoritam, poznat još kao *bootstrap* agregacija, predstavlja način odabiranja obučavajućeg skupa. Originalni skup podataka se odabira uniformno sa ponavljanjem. Obučavajući skup se formira ovim algoritmom za svako stablo unutar slučaje šume. Ukoliko je odabrani skup iste veličine kao i originalni, pokazuje se da on sadrži oko 63% originalnih odbiraka, dok je ostatak izostavljen i naziva se *out-of-bag* (oob) skup. Ovo je interno implementirani mehanizam za sprečavanje preobučavanja modela.

Skup izostavljenih odbiraka (oob skup) služi za validaciju pri obučavanju stabala. Obučavanje svakog stabla se zaustavlja kada se postigne nabolji rezultat na oob skupu.

Algoritam slučajne šume prima informaciju na ulaz u obliku odbirka. Odbirak se propušta kroz svako stablo odlučivanja pri čemu svako formira odluku o njegovoj klasnoj pripadnosti. Konačna odluka se donosi većinskim glasom.

Slučajnost slučajne šume se ogleda u tome što se svako stablo obučava sa slučajnim podskupom odbiraka i sa slučajnim podskupom obeležja tih odbiraka. Ovo rezultuje većim diverzitetom između stabala i kvalitetnijom konačnom odlukom modela.

RFC u implementaciji *scikit-learn* biblioteke [22] formira heuristiku važnosti svakog obeležja i predstavlja njegovu korisnost pri deljenju skupa podataka na dva skupa veće homogenosti, što se vrši u čvoru odluke. Ova heuristika je u potpoglavlju 6.1.4 korištena za rangiranje obeležja, i generalno može poslužiti za odbacivanje manje bitnih obeležja i na taj način metod za redukciju dimenzionalnosti.

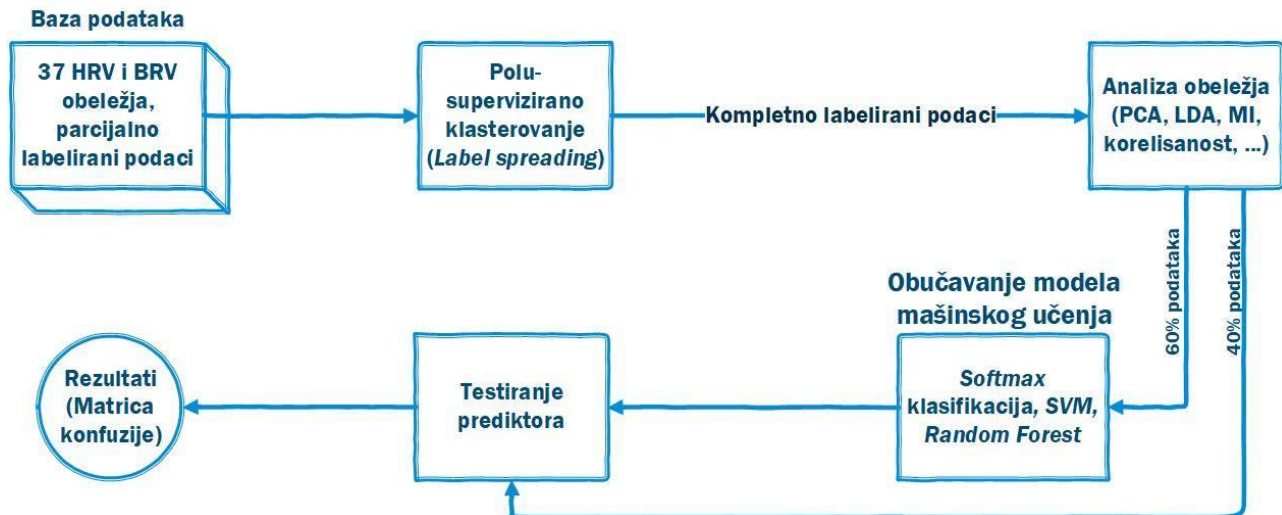
Hiper-parametri slučajne šume su:

- Broj stabala odlučivanja-*n\_estimators*. Uobičajeno više stabala poboljšava performanse šume, čineći predikciju stabilnijom. Više stabala usporava obučavanje i komputaciju predikcije.
- Broj obeležja (*max\_features*) predstavlja broj slučajno odabranih obeležja koja se uzimaju za obučavanje stabla odluke. Može se dati i kao procenat ukupnog broja obeležja. Manji broj obeležja povećava različitost između stabala.
- Maksimalna dubina stabla-*max\_depth*. Dubina stabla predstavlja broj čvorova odluke koji vodi od korenog čvora do lista. Stablo se obučava dok se dostigne željenu homogenost podskupova. Ukoliko se ona ne dostigne obučavanje se obustavlja na maksimalnoj dubini. Veća vrednost ovog parametra učiniće svako stablo sklonim preobučavanju.

- Minimalni broj odbiraka u listu (*min\_samples\_leaf*) predstavlja hiper-parametar koji se retko podešava a može imati značajne posledice po performanse stabla. Naime, ukoliko se dozvoli da unutar lista bude jedan jedini odbirak stablo će biti u mogućnosti da formira potprostor oko jedinstvenog odbirka te da se na taj način preobuči. Veće vrednosti učiniće svako pojedinačno stablo robusnijim. Vrednost ovog parametra je svakako u vezi sa obimom dostupnih podataka.
- Minimalni broj odbiraka za razdvajanje u čvoru-*min\_samples\_split*. Odnosi se na to koliko se jednom odlukom može odvojiti odvojiti odbiraka od ukupnog skupa. Slično prethodnom hiper-parametru, ukoliko je vrednost mala, u čvoru se može desiti razdvajanje svega nekoliko odbiraka od celokupnog skupa podataka, što može rezultovati preobučavanjem.

## 4. Rezultati

Rezultati ove teze su dati u vidu relabeliranja podataka, analize i rangiranja obeležja te projektovanju modela za personalizovanu i interpersonalnu klasifikaciju. Postupak koji vodi do ovih rezultata je prikazan na slici 4.1 u vidu dijagrama toka.



Slika 4.1 Dijagram toka opšte procedure

U potpoglavlju 4.1 opisana je primena algoritma polu-supervizijskog klasterovanja u cilju relabeliranja podataka te na taj način obogaćivanja originalnog skupa podataka. U potpoglavlju 4.2 prikazani su rezultati analize i procene kvaliteta obeležja u cilju sticanja uvida u dobijeni skup podataka. U potpoglavlju 4.3 dati su rezultati projektavanja modela za predikciju emocija kod pojedinačnih ispitanika, dok je u potpoglavlju 4.4 vršeno projektovanje modela koji vrši generalnu predikciju ljudskih emocija.

### 4.1 Relabeliranje podataka

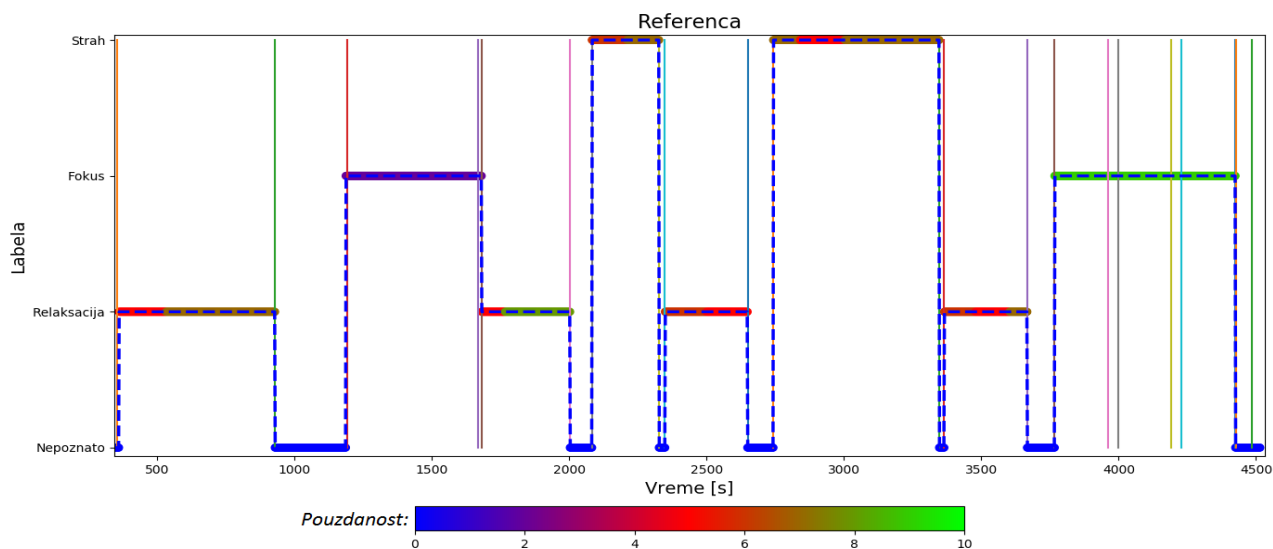
Kao što je rečeno u potpoglavlju 2.3 na raspolaganju imamo podatke koji nisu upotpunosti labelirani. Tokom procesa labeliranja prepoznati su delovi sadržaja koji su labelirani sa visokom pouzdanošću. Takođe, pronađene su i sekvence kojima je dodeljena referenca ali koja je usled nepotpunih informacija bila niske verodostojnosti. Ostatak snimljenih elektrofizioloških signala pripada delovima eksperimenta čija se referenca nije mogla proceniti.

Podaci sa labelama niske pouzdanosti nose rizik sa sobom. Naime, velika je verovatnoća da, iako im je dodeljena jedna labela, pojedini njihovi delovi pripadaju nekom drugom emotivnom stanju. Vrlo se često radi o podacima koji se nalaze na prelazu emotivnih stanja. Tipičan primer toga jeste pojava opuštenosti kod evaluacije. Ispitanik koji je već dobro upoznat sa procedurom somoocenjivanja može joj sa lakoćom pristupiti. Na taj način njegovi elektrofiziološki signali mogu da svedoče o većoj opuštenosti nego stepenu koncentracije. Unos pogrešno labeliranih podataka može značajno degradirati performanse modela predikcije.

Pojava nelabeliranih podataka direktno smanjuje obim dostupnih odbiraka. Dakle, delovi elektrofizioloških signala koji međusobno sinhronizovani, procesirani i iz kojih su ekstrahovana obeležja se prosto odbacuju. Zasad projektovanja kvalitetnijeg i robusnijeg modela predikcije u interesu nam je da raspolažemo sa što većim skupom podataka.

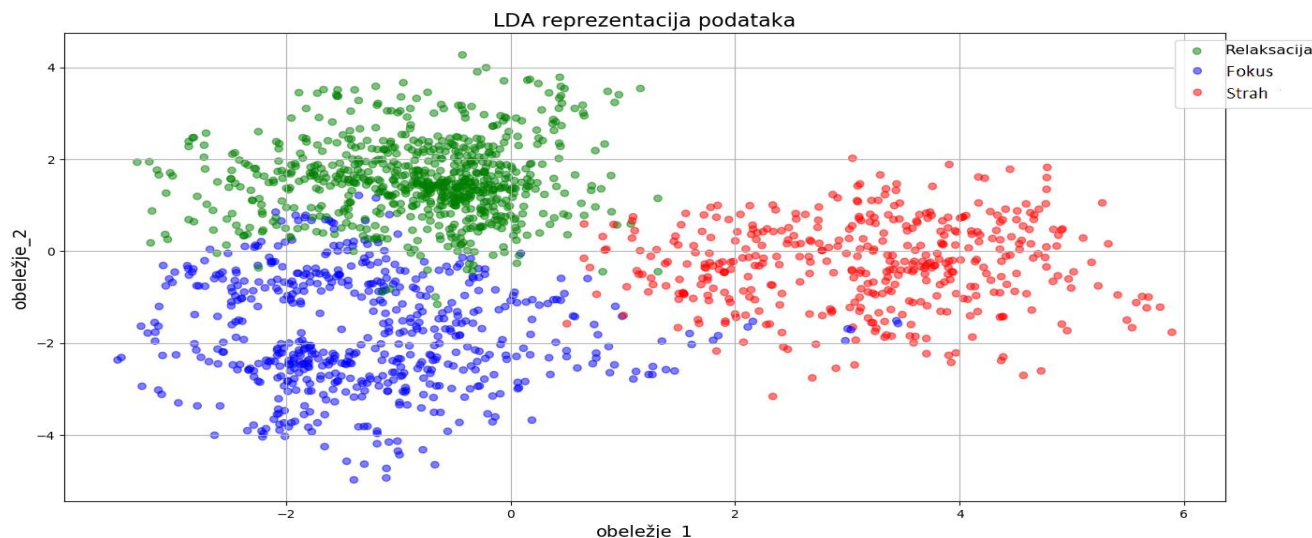
Kako bi se rešio problem pogrešno labeliranih i odbačenih podataka pristupa se postupku klasterizacije. Reč je o polu-supervizijskoj klasterizaciji gde se na osnovu podataka visoke pouzdanosti vrši dodeljivanje labele podacima niske pouzdanosti i nelabeliranim podacima. Polu-supervizirano klasterovanje vršeno je pomoću algoritma *label spreading* koji je opisan u odeljku 3.1.1, a realizovan u okviru *python* biblioteke *scikit-learn* [22].

Na slici 4.1.1 prikazana je ručno formirana referenca za jednog ispitanika. Na grafiku vidimo sve tipove labela uključujući i sekvence sa nepoznatom labelom. Vertikalne linije na istom grafiku označavaju etape eksperimenta. Njihova legenda je izostavljena radi preglednosti dodeljnih labela. Takođe, vidimo da je pouzdanost raznovrsna. Postoje čak labele za pouzdanošću nižom od 5.



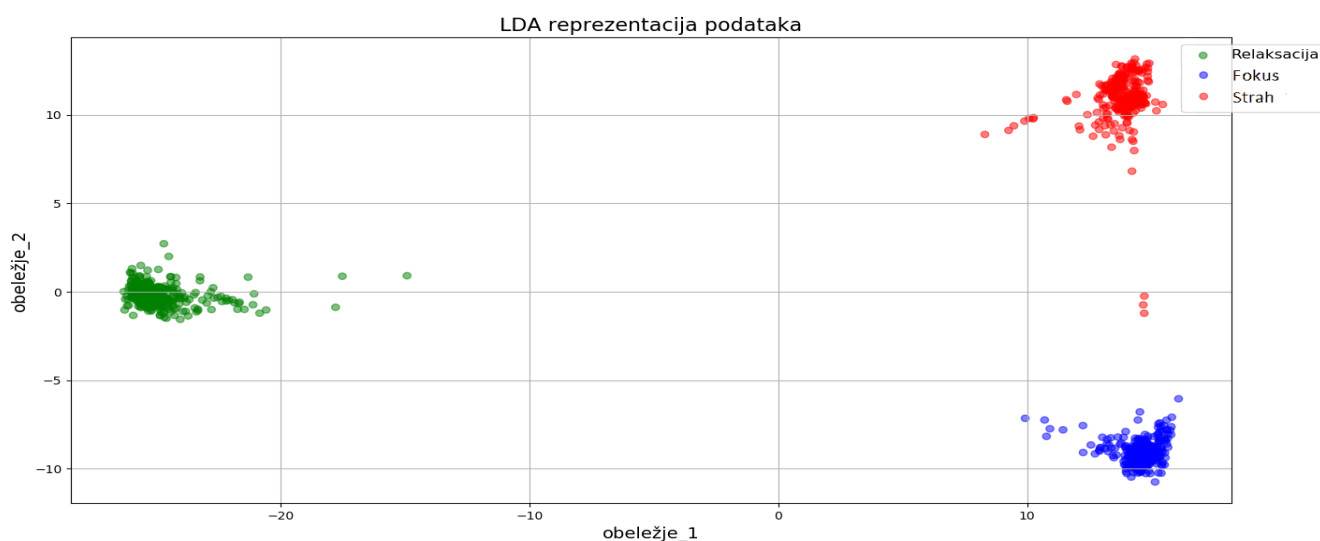
Slika 4.1.1 Ručno formirana referenca

Koristićemo Fišerovu LDA metodu, opisanu u odeljku 3.2.2, za prikaz podataka u dvodimenzionalnom prostoru. Na slici 4.1.2 prikazani su podaci u prostoru LDA obeležja. Ovaj prikaz je generisan pomoću svih labeliranih podataka uključujući one sa niskom pouzdanošću. Sa slike vidimo da postoji separabilnost u nekoj meri ali da se klase relaksacije i fokusa dosta mešaju. Ovo je posledica već diskutovane pojave da se određeni odbirci proglašavaju fokusom niske verodostojnosti a da zapravo predstavljaju odbirke opuštenosti.



Slika 4.1.2 LDA prikaz podataka labeliranih ručnom metodom bez izuzimanja odbiraka niske pouzdanosti

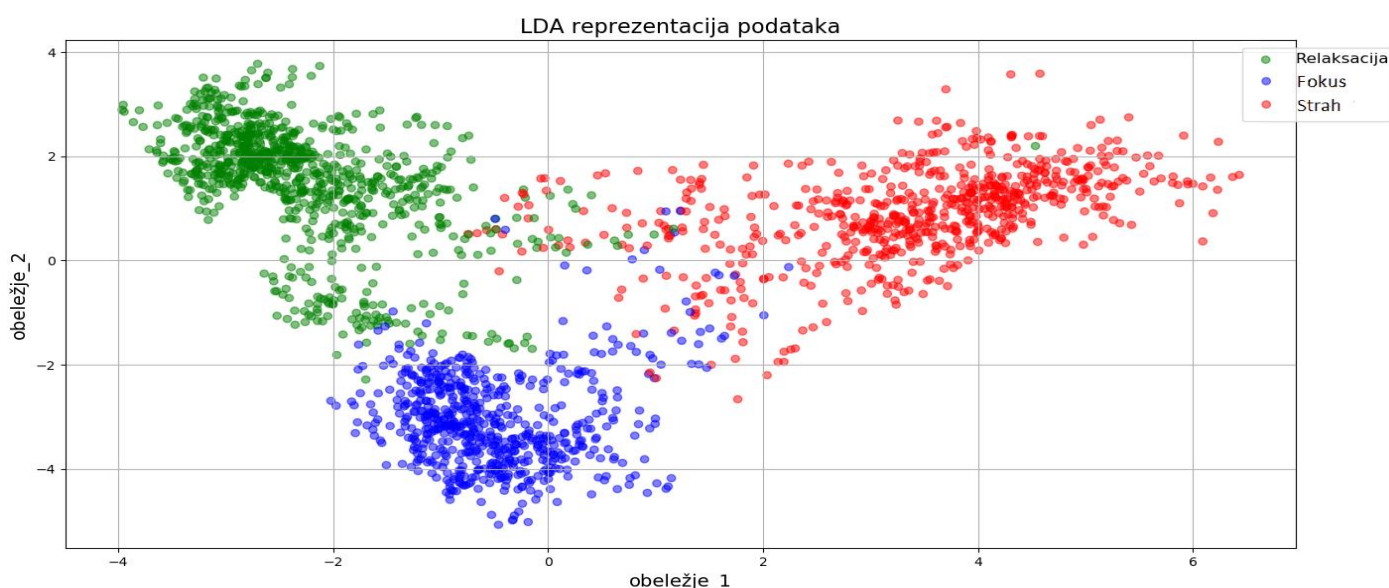
Kada se izvrši LDA prikaz samo podataka visoke pouzdanosti (veće ili jednake od 8) dobijaju se izrazito separabilni klasteri. Pomenuti prikaz dat je na slici 4.1.3. Vidimo da kvalitetno labelirani podaci predstavljaju prave predstavnike svojih klasa.



**Slika 4.1.3 LDA prikaz podataka labeliranih ručnom metodom sa izuzimanjem odbiraka niske pouzdanosti**

Polu-supervizijsko klasterovanje je izvršeno na osnovu pouzdanih podataka (koji su prikazani na slici 4.1.3). Algoritam jednako tretira i nelabelirane podatke i podatke sa labelama niske pouzdanosti. Ovakvi odbirci bivaju klasterovani na osnovu njihovih obeležja, bez apriorne informacije o njihovoj klasnoj pripadnosti.

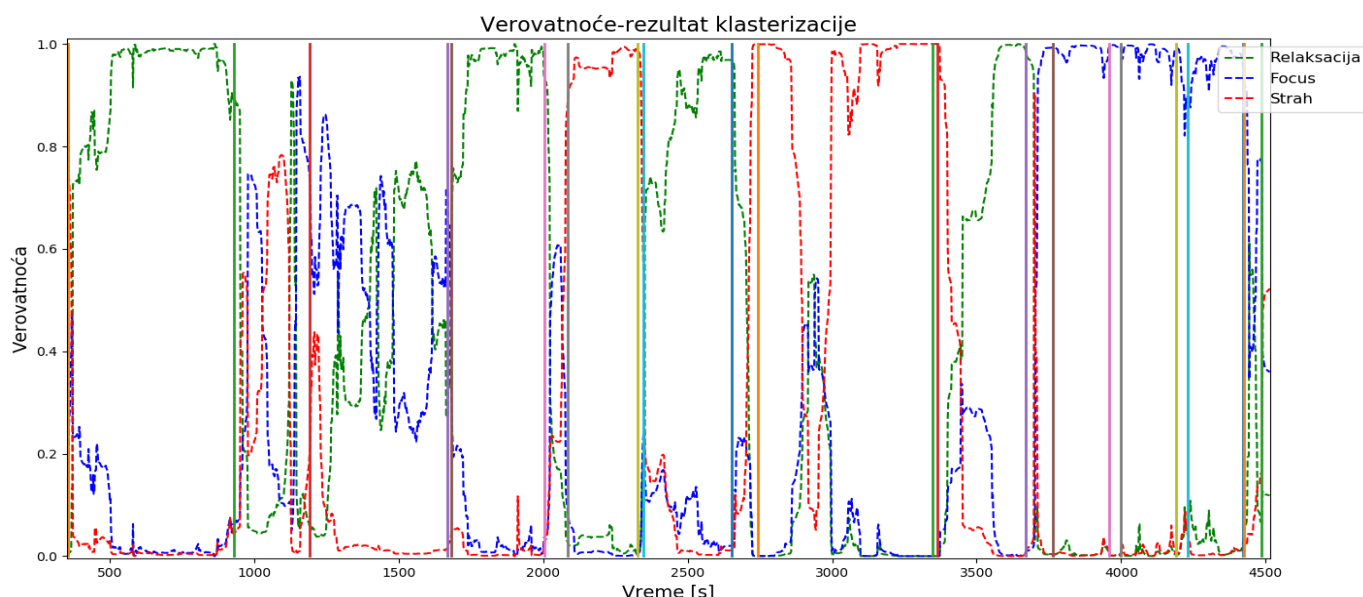
Algoritam za svaki odbirak vraća verovatnoće pripadanja svakoj od ponuđenih klasa, dakle tri realna broja koji u zbiru daju jedan. Konačna labela uzima se kao klasa kojoj odbirak pripada sa najvećom verovatnoćom. Iste verovatnoće sada preklapaju vrednost pouzdanosti. Naime, pouzdanost neke labele je veća ukoliko je verovatnoća pripadnosti toj klasi veća. Dakle, kao nova vrednost pouzdanosti usvaja se kao najveća verovatnoća pripadnosti. Na slici 4.1.4 dat je LDA prikaz podataka nakon klasterizacije.



**Slika 4.1.4 LDA prikaz podataka labeliranih polu-supervizijskom klasterizacijom**



Dok sa jedne strane algoritam vrši labeliranje nedovoljno pouzdanih podataka sa druge strane algoritam ostavlja mogućnost reklasterizacije i polaznih podataka, tj. odbiraka visoke pouzdanosti. Ova mogućnost je većinom simbolična, jer su polazni podaci i isprva bili separabilni. Svakako, ovaj mehanizam nam vraća verovatnoću tj. novu vrednost varijable pouzdanosti i za ovaj set podataka, što nam je od interesa. Nove labele date u obliku verovatnoća klasne pripadnosti prikazane su na slici 4.1.5.



Sa slike 4.1.5 vidimo da algoritam pokazuje znakove neodlučnosti u vremenskom opsegu od 1000s

**Slika 4.1.5 Rezultat klasterizacije u vidu verovatnoća klasne pripadnosti**

do 1600s. Ukoliko konsultujemo sliku 4.1.1 zapazićemo da je ovaj opseg u startu bio nepouzdan. Takođe vidimo da je pomenuti opseg ispraćen maksimalnim verovatnoćama manjim od 0.8, što se u novoj metrici pouzdanosti može smatrati nekvalitetnim. Ovo je primer podataka koji ni nakon klasterizacije ne postaju reprezentativni predstavnici neke klase, tj. klasterizacija ne rešava njihov problem.

Ipak, kako je na pomenutom opsegu reč o odbircima koji imaju po dve verovatnoće uporedivih vrednosti, može se zaključiti da se ovi odbirci nalaze na prelazu datih klasa, te bliže određuju prelazni pojas. Ovakvi odbirci se ne odbacuju jer mogu doprineti projektovanju robusnijeg klasifikatora. Drugim rečima, ovi odbirci nisu autlajeri (*outlier* – opservacija udaljena od drugih pripadnika svoje klase), već prosto odbirci koji imaju karakteristike obe klase.

## 4.2 Analiza obeležja

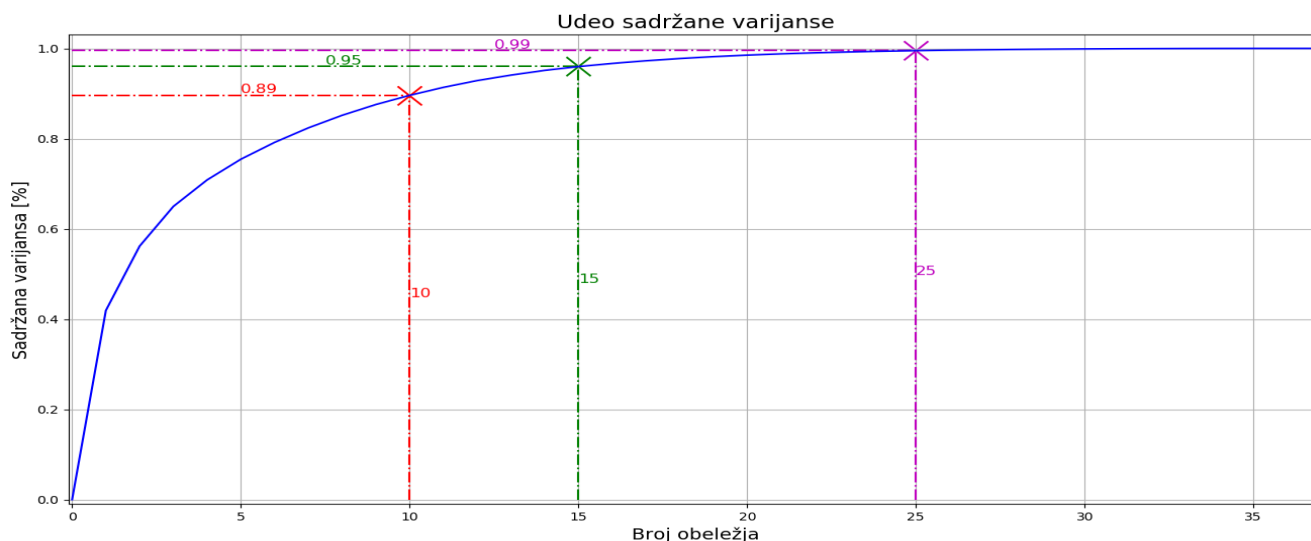
U ovom potpoglavlju dajemo uvid u karakteristike ekstrahovanih obeležja, u smislu međusobne korelisanosti i informativnosti. Takođe vršimo procenu njihove svrsishodnosti za potrebe klasifikacije. Konkretno, izračunaćemo koeficijent korelacije, *mutual information* i PCA kao indikatore navedenih osobina.

### 4.2.1 PCA analiza

PCA analiza je definisana i objašnjena u okviru odeljka 3.1.1. Kao što je rečeno PCA analiza projektuje podatke na međusobno ortogonalne pravce koji maksimiziraju njihovu unutrašnju

varijansu. Naš cilj jeste da procenimo dominantnost ovih pravaca. Naime, ukoliko je ukupna varijansa podataka većinski sadržana u prvih nekoliko komponenti PCA analize postoji opravdanje da se izvrši značajna redukcija dimenzija. Na primer, da već prva komponenta poseduje varijansu koja iznosi više od 90% ukupne varijanse svih ortogonalnih komponenti mogli bismo zaključiti da je ukupna informativnost obeležja mala u odnosu na njihov broj.

Izvršena je PCA analiza nad 37 obeležja opisanih u potpoglavlju 2.3. Na slici 4.2.1. prikazan je udeo sadržane varijanse u odnosu na broj PCA komponenti. Na slici je posebno označeno koji procenat varijanse biva obuhvaćen ukoliko se zadrži 10, 15 ili 25 obeležja.



Slika 4.2.1 Udeo sadržane varijanse u odnosu na broj PCA komponenti

Rezultat je da prvih 15 PCA komponenti sadrži 95% ukupne varijanse obeležja. Ovo je ohrabrujući rezultat koji nagoveštava da obeležja nisu isuviše redudantna. Kako raspolažemo sa oko 20.000 odbiraka broj od 37 obeležja nije veliki, tj. nije potrebno vršiti redukciju dimenzija.

#### 4.2.2 Korelisanost obeležja

Od interesa je proceniti u kojoj meri su konkretna obeležja međusobno zavisna. Kako su za obeležja izabrana uobičajena obeležja vremenske serije, ova analiza je bitna sa stanovišta budućeg rada sa EKG signalom i signalom respiracije. Drugim rečima, bliži uvid u standardna obeležja ovih elektrofizioloških signala predstavlja jedan od rezultata naše teze. Biće određen koeficijent međusobne korelacije između obeležja opisan u odeljku 3.2.3.

Radi preglednosti prikazaćemo samo grupe obeležja među kojima postoje značajni stepeni korelacije. Prva grupa prikazana je u tabeli 4.2.1. Treba napomenuti da obeležja ekstrahovana iz EKG signala imaju prefiks *ecg*, dok obeležja signala respiracije imaju prefiks *rsp*.

Tabela 4.2.1 Koeficijent korelacije prve grupe obeležja

<i>r</i>	<i>ecg_stdRR</i>	<i>ecg_cvRR</i>	<i>ecg_SD1</i>	<i>ecg_SD2</i>	<i>ecg_SD1SD2</i>	<i>ecg_CCM</i>
<i>ecg_stdRR</i>	1	0.96	0.68	0.99	0.92	0.75
<i>ecg_cvRR</i>	0.96	1	0.62	0.96	0.86	0.69
<i>ecg_SD1</i>	0.68	0.62	1	0.68	0.88	0.89

<i>ecg_SD2</i>	0.99	0.96	0.68	1	0.92	0.75
<i>ecg_SD1SD2</i>	0.92	0.86	0.88	0.92	1	0.92
<i>ecg_CCM</i>	0.75	0.69	0.89	0.75	0.92	1

Usled prirode funkcije za procenu korelacije, koja je data izrazom (2) u odeljku 3.2.3, tabela 4.2.1 je simetrična matrica. Na početku zapažamo visok nivo linearne zavisnosti između parova obeležja *ecg\_stdRR* i *ecg\_cvRR* te *ecg\_SD2* i *ecg\_SD1SD2*. Parovi ovih obeležja su među sobom direktno proporcionalna prema tome ovo nije argument za ili protiv ovih obeležja [9].

Takođe, zapaža se i visok stepen korelacije *ecg\_stdRR* i *ecg\_SD2*. Ova dva obeležja su direktno povezana nelinearnom vezom, takvom da će povećanjem jedne rasti i druga i obratno [11]. U istoj formuli učestvuje i *rmssdRR* obeležje koje je direktno proporcionalno obeležju *ecg\_SD1*. S obzirom da *ecg\_SD1* nije izrazito korelisano sa *ecg\_stdRR* nameće se zaključak da *ecg\_stdRR* i *ecg\_SD2* nisu redundantni [11].

Druga grupa bi bila obeležja tipa  $pRR_t$  zajedno sa *ecg\_SD1* obeležjem. U tabeli 4.2.2 prikazi su koeficijenti korelacije ove grupe.

**Tabela 4.2.2 Koeficijent korelacije druge grupe obeležja**

<b>r</b>	<i>ecg_pNN<sub>30</sub></i>	<i>ecg_pNN<sub>50</sub></i>	<i>ecg_pNN<sub>70</sub></i>	<i>ecg_SD1</i>
<i>ecg_pNN<sub>30</sub></i>	1	0.81	0.65	0.81
<i>ecg_pNN<sub>50</sub></i>	0.81	1	0.82	0.87
<i>ecg_pNN<sub>70</sub></i>	0.65	0.82	1	0.85
<i>ecg_SD1</i>	0.81	0.87	0.85	1

Primećuje se kako je *ecg\_pNN<sub>30</sub>* više korelisano sa *ecg\_pNN<sub>50</sub>* nego sa *ecg\_pNN<sub>70</sub>*. Ovo je sasvim logično imajući u vidu način njihovog izračunavanja [9]. Obeležje *ecg\_pNN<sub>50</sub>* je gotovo podjednako korelisano sa druga dva *pNN* obeležja što je odgovara činjenici da se zasniva na vremenskim razlikama (50 ms) koje se nalaze sa sredini između 30 ms i 70 ms. Interesantno je da *ecg\_SD1* poseduje značajnu dozu korelisanosti sa svim *pNN* obeležjima.

Posebno, kao treću grupu imamo obeležja koja estimiraju entropiju vremenske serije (*ApEn* i *SampEn*) koja beleže visok nivo korelacije sa koeficijentom  $r = 0.93$ . Ovo je posve očekivano jer oba obeležja su mere iste karakteristike. Da je korelacija među njima mala značilo bi da jedna od njih ne estimira uredenost na odgovarajući način. Kako je opisano u literaturi, obeležje *SampEn* je pouzdanije i sposobnije da radi sa manjim brojem odbiraka vremenske serije od obeležja *ApEn*. Prema tome, u spornoj situaciji veće poverenje bi trebalo ukazati obeležju *SampEn* [14].

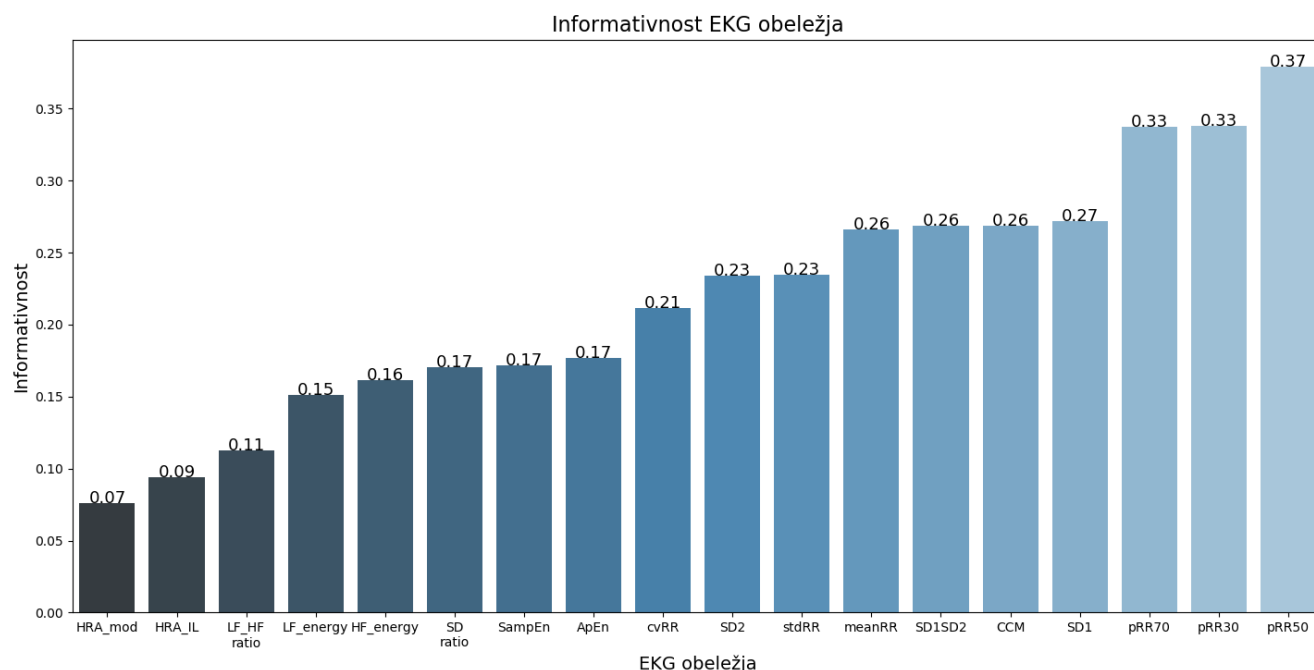
Slični rezultati se dobijaju se i za obeležja signala respiracije. Javlja se visoka korelisanost unutar skupova obeležja koja su ekvivalentna tri navedene grupe obeležja EKG signala. Prema tome, mogu se izvesti slični zaključci. Takođe, bitno je napomenuti da nije bilo zapažene korelisanosti između obeležja EKG signala i signala respiracije, što je očekivani rezultat.

### 4.2.3 Informativnost obeležja

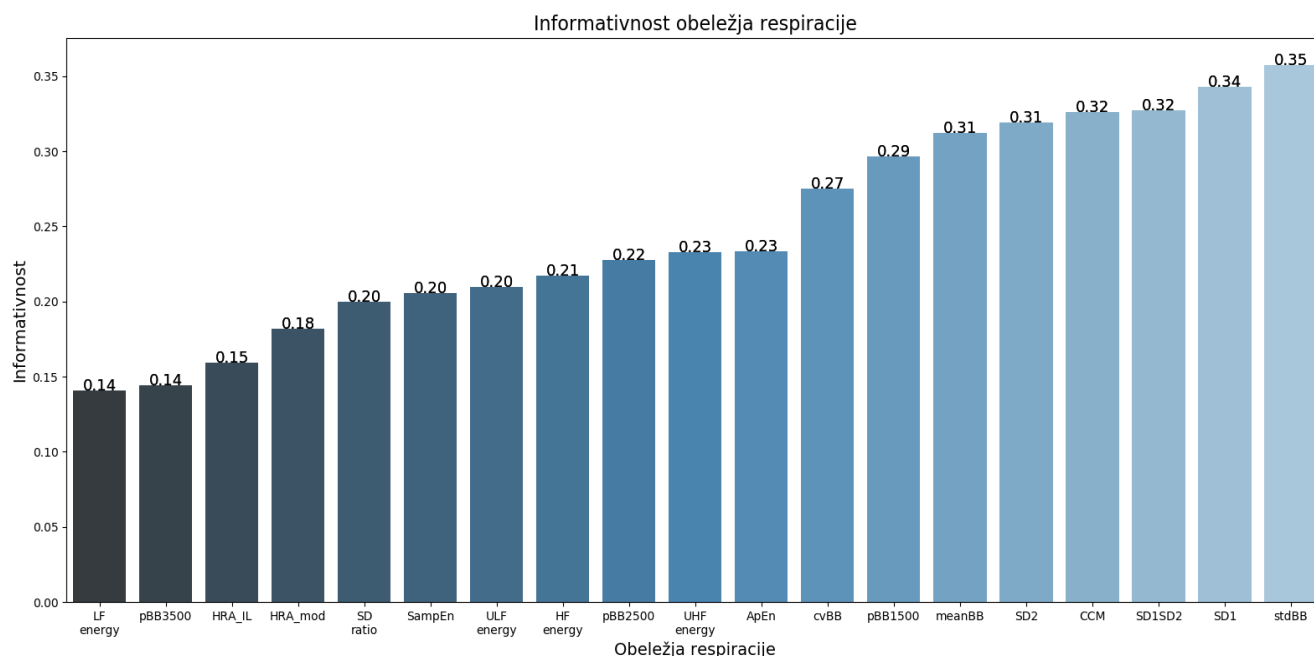
Cilj je odrediti stepen informativnosti obeležja u smislu klasne pripadnosti. U kojoj meri određeno obeležje nosi informaciju o labeli koju dati odbirak nosi. Ovo predstavlja vid rangiranja obeležja

jer je za očekivati da obeležja visoke informativnosti budu značajnija pri projektovanju modela predikcije.

Računa se  $MI$  svakog obeležja sa referencom na način opisan u odeljku 3.2.4. Informativnost EKG obeležja poredana u rastućem trendu prikazana je na slici 4.2.2, informativnost obeležja respiracije prikazana je na slici 4.2.3.



**Slika 4.2.2 Informativnost EKG obeležja**



**Slika 4.2.3 Informativnost obeležja respiracije**

$MI$  predstavlja samo način za rangiranje obeležja. Sa slike 4.2.2 vidimo da najveću informativnosti nose obeležja tipa  $pNN$ . Kako su obeležja iz ove grupe međusobno korelisana najverovatnije je reč

o redundantnoj informaciji. Vidi se da i druge grupe međusobno korelisanih obeležja, navedenih u odeljku 4.2.2, imaju slične nivoe informativnosti.

Sa slike 4.2.3 vidimo da i obeležja respiracije imaju uporedive nivoe informativnosti sa EKG obeležjima. Prema tome, postoji opravdanost ekstrakcije kompleksnijih obeležja iz vremenske serije BRV signala.

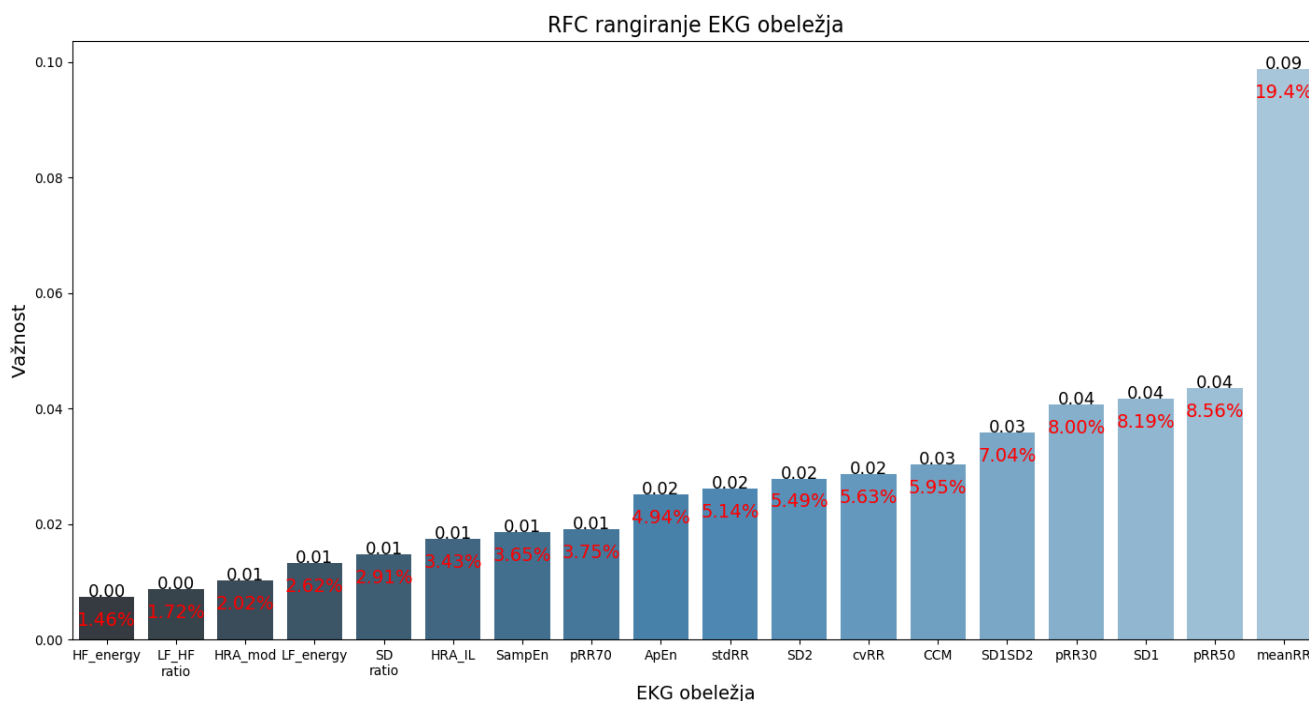
Obeležja *HRA\_IL* i *HRA\_mod* su konzistentno loše rangirana kod obeležja oba tipa elektrofiziološkog signala. Takođe, za EKG signal slabo su rangirani frekvencijska obeležja. Pretpostavka je da je razlog ovome nedovoljan broj odbiraka HRV signala na vremenskom prozoru za adekvatnu estimaciju spektralne gustine snage signala.

#### 4.2.4 Rangiranje pomoću slučajnih šuma

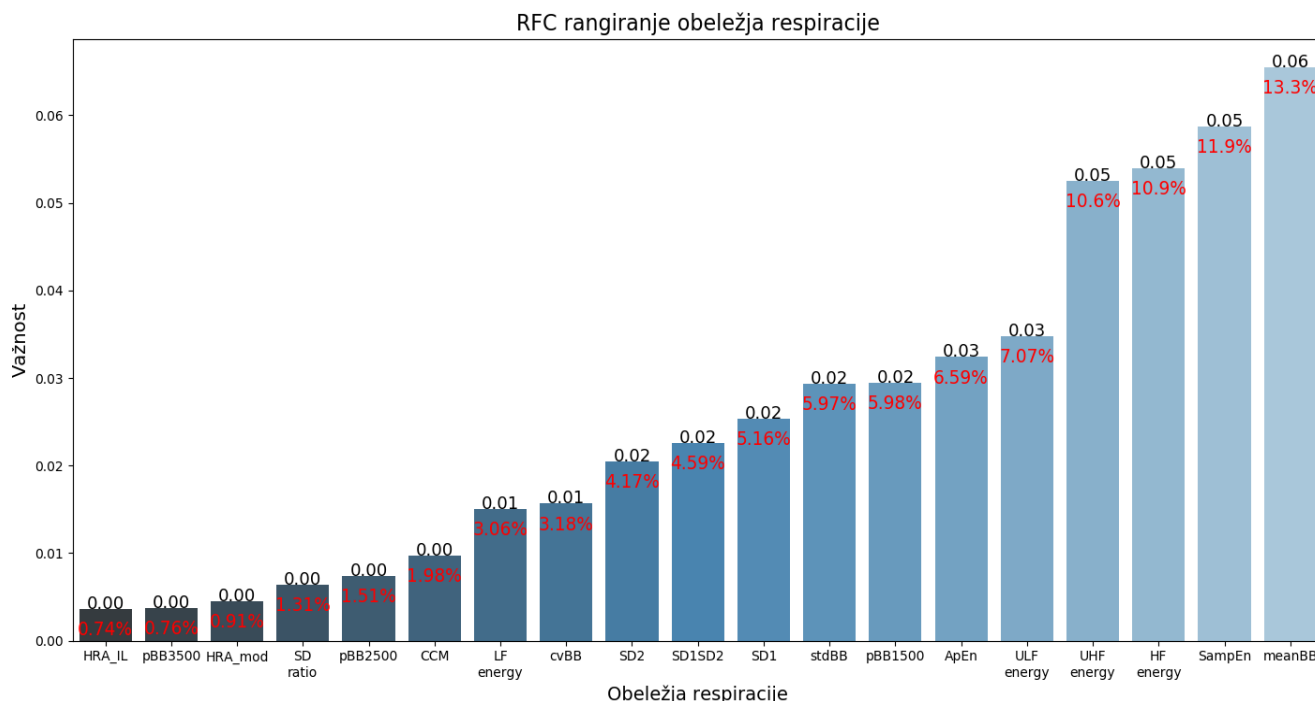
Slučajne šume (*Random Forest*) predstavljaju tip modela koji može da služi kako za regresiju (estimaciju kontinualne funkcije) tako i za klasifikaciju (RFC – *Random Forest Classifier*). Ovaj model ima tu osobinu da nakon obučavanja poseduje informaciju o tome koliko je koje obeležje bilo korisno u smislu klasifikacije. Ova informacija se može koristiti za rangiranje obeležja. Konkretni algoritam je objašnjen u odeljku 3.3.3.

Informaciju o važnosti obeležja koju algoritam vraća ipak treba uzeti sa rezervom. Naime, kao povratnu informaciju algoritma dobijamo korisnost obeležja za obučavanje RFC modela a ne uopšteni značaj i važnost obeležja. Na primer ukoliko je neko obeležje po oceni RFC modela veoma bitno to ne znači nužno da je isto obeležje podjednako bitno i za klasifikator nosećih vektora (SVM – *Support Vector Machine*).

Na slici 4.2.4 prikazana je važnost EKG obeležja, procenjena od strane RFC modela. Slično, na slici 4.2.5 prikazana je važnost obeležja respiracije.



Slika 4.2.4 Važnost EKG obeležja prema RFC modelu



**Slika 4.2.5 Važnost obeležja respiracije prema RFC modelu**

Procenti na slikama 4.2.4 i 4.2.5 predstavljaju udeo važnosti u odnosu na podskupove EKG obeležja i obeležja respiracije respektivno. Sa slike 4.2.4 vidimo da obeležje srednje vrednosti RR intervala ima dominantan značaj i doprinos treniraju slučajne šume. Ovo obeležje je zapravo direktno proporcionalno estimaciji HR (*Heart Rate*) signala tj predstavlja srčani ritam ispitanika. Prema tome, sasvim je opravdano da ovo obeležje sadrži gotovo petinu ukupnog doprinosa EKG obeležja. Ukoliko ovo uporedimo sa slikom 4.2.2 vidimo da se rang obeležja *meanRR* značajno povećao. Obeležja *SD1* i *pRR<sub>50</sub>* su konzistentno ostali pri vrhu, dok je obeležje *pRR<sub>70</sub>* značajno opalo po važnosti.

Za slučaj obeležja respiracije sa slike 4.2.5 vidimo da se u odnosu na sliku 4.2.3 poredak veoma promenio. Slično *meanRR* obeležje *meanBB* se značajno više rangira. Međutim najbolje rangiranih pet obeležja respiracije sa slike 4.2.3 se ne pojavljuju u prvih pet obeležja rangiranih RFC modelom. Takođe, sa iste slike se zapaža i veći značaj *SampEn* obeležja u odnosu na *ApEn* obeležje, što je očekivani rezultat usled nedostatka odabiraka vremenske serije nad kojoj se računaju [14]. Ovakav zdravorazumski poredak nije prisutan na slici 4.2.3 što ide u prilog verodostojnosti RFC rangiranja.

### 4.3 Personalizovana klasifikacija

Cilj personalizovane klasifikacije jeste projektovanje klasifikatora emotivnog stanja konkretnog ispitanika. Dakle, obučavanje se vrši na osnovu pojedinačnog eksperimenta nekog ispitanika sa motivacijom da se dobije model za predikciju budućih emocija i raspoloženja istog ispitanika.

Prvo se vrši podela trajanja eksperimenta na obučavajući, validirajući i testirajući skup, što je opisano u odeljku 4.3.1. Potom se vrši sukcesivno obučavanje sve kompleksnijih klasifikatora, počevši od jednostavne *softmax* metode (odeljak 4.3.2), preko SVM modela sa različitim kernelnim funkcijama (odeljak 4.3.3) do RFC klasifikatora (odeljak 4.3.4). Obučavanje za svakog ispitanika traje sve dok se ne postignu zadovoljavajući rezultati na testirajućem skupu.

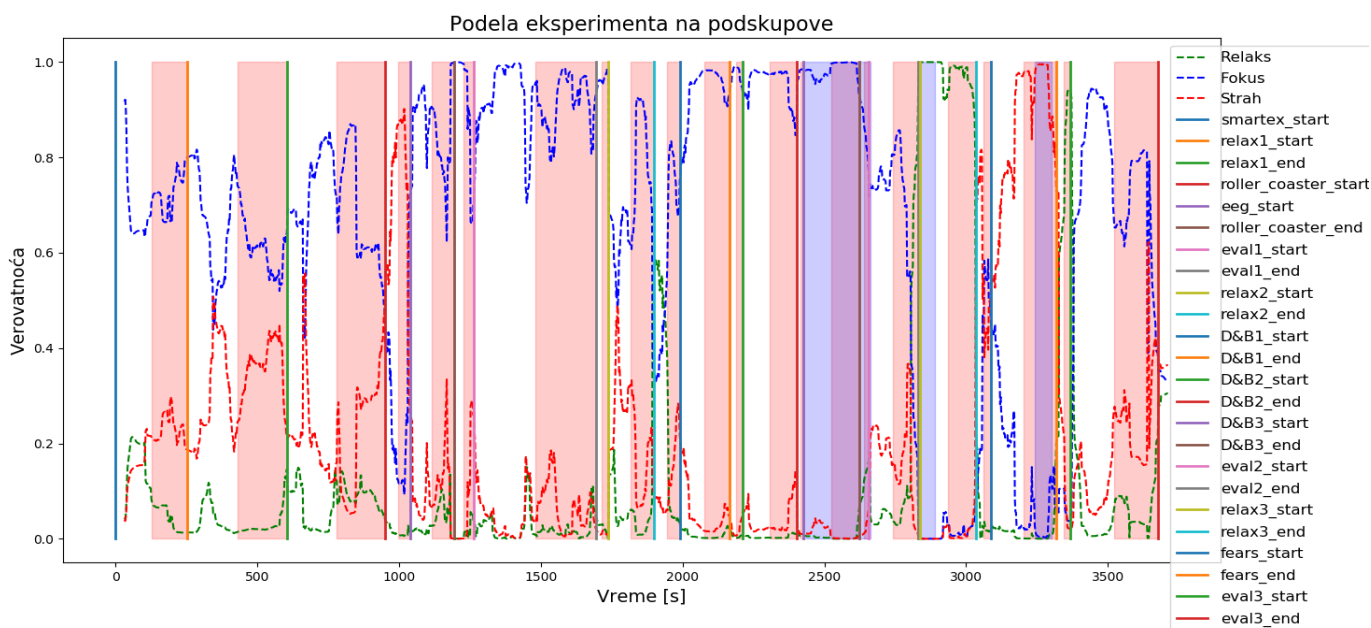
### 4.3.1 Podela skupa podataka

Uobičajeno celokupni skup podataka se deli na obučavajući, validirajući i testirajući podskup na slučajan način. Treba napomenuti da testirajući skup služi isključivo za finalnu estimaciju uspešnosti modela predikcije te se kao takav ne koristi niti u jednoj fazi obučavanja ili validacije. Validirajući skup je odabran da bude slične veličine kao i obučavajući.

Ovo omogućava i 2-fold krosvalidaciju u cilju podešavanja hiper-parametara klasifikatora, a ne samo validaciju radi ranog zaustavljanja procesa obučavanja. K-fold krosvalidacija predstavlja proces u kom se određeni skup podataka deli na  $K$  delova.  $K$  puta se vrši obučavanje nekog klasifikatora sa konkretnim vrednostima parametara, koji se još nazivaju i hiper-parametri. Obučavanje se vrši na  $K-1$  podskupova a testiranje samo na jednom. Pri svakom obučavanju menja se testirajući podskup, a rezultat evaluacije se pamti. Konačan uspeh takve konfiguracije hiper-parametara se dobija kao usrednjen uspeh iz  $K$  evaluacija. Ovaj postupak služi za automatski odabir konfiguracije hiper-parametara, tako što se niz predefinisanih konfiguracija hiper-parametara testira jedan za drugim.

Skup podataka sa kojim raspolazemo dobijen je ekstrahovanjem obeležja pomoću prozora trajanja 60s i sa pomerajem od 2s, kako je opisano u potpoglavlju 2.3. Prema tome dva uzastopna odbirka imaju preklapanje od 58s tj. 96%. Na osnovu toga, da se zaključiti da ovi odbirki imaju veoma sličan informativni sadržaj. Usled hronološke povezanosti postoji visok stepen korelacije između odbiraka. Kada bi se podskupovi birali na slučajan način, svaki od njih bi sadržao sličan informativni sadržaj. Tada evaluacija modela predikcijom na testirajućem skupu ne bi imala smisla jer bismo dobili identične rezultate kao na obučavajućem skupu.

Ukoliko bismo povećali pomeraj prozora, u cilju smanjenja pomenute korelacije odbiraka, došlo bi do značajne redukcije obima podataka. Pomeraj od 4s koji ne bi rešio pomenutu pojavu bi prepolovio naše podatke. Kako se ljudske emocije kontinualno menjaju, ova pojava bi postojala čak i za pomeraj od 60s. Opisana podela vršena je ručno, bez prethodnog proređivanja odbiraka. Na slici 4.3.1 prikazan je primer ove podele.



**Slika 4.3.1 Podela eksperimenta na obučavajući, validirajući i testirajući skup: crvena šrafura predstavlja validirajući skup, plava i ljubičasta šrafura predstavlja testirajući skup, nešrafirani delovi eksperimenta predstavljaju obučavajući skup**



Svaka podetapa eksperimenta podeljena je automatski na dva dela jednakih dužina. Prvi delovi ulaze u obučavajući skup, dok drugi automatski ulaze u validirajući skup podataka. Potom se ručno bira reprezentativan testirajući skup. U interesu nam je da u testirajućem skupu budu prisutni svi tipovi labela, kao i da te labele budu visoke pouzdanosti (visoka verovatnoća pripadnosti). Na slici 4.3.1 crvena šrafura predstavlja validirajući skup, plava i ljubičasta šrafura predstavlja testirajući skup, nešrafirani delovi eksperimenta predstavljaju obučavajući skup.

Ovakva konfiguracija omogućava veće razmake između podskupova. Izražena korelacija na prelaznim delovima nije presudna za objektivnost evaluacije na testirajućem skupu. Treba imati u vidu da iako je podela izvršena na ovaj način, konkretni skupovi se mogu grupisati na proizvoljan način. Na primer moguće je testirajući skup pridružiti obučavajućem a evaluaciju vršiti na daleko brojnijem validirajućem skupu. Ovo je posebno korisno kada u projektovanju klasifikatora nema potrebe za automatskim podešavanjem hiper-parametara pomoću *2-fold* krosvalidacije, već se hiper-parametri ručno podešavaju. Takav je slučaj u svim primerima personalne klasifikacije koja je predmet ovog potpoglavlja.

Evidentno je da je broj odbiraka labeliranih kao strah značajno manji od broja odbiraka druge dve labele. Ovo je svakako očekivano jer je izazvati emociju straha znatno teže od druga dva uobičajnija raspoloženja. Usled ovoga, javlja se problem nebalansiranih klasa, te je ovome, pri projektovanju modela, potrebno posvetiti posebnu pažnju. Model koji ne vodi računa od ovom balansu bi favorizovao najmnogobrojniju klasu.

#### 4.3.2 Rezultati primene softmax metode

Kao najjednostavniji linearni klasifikator *softmax*, metoda je vrlo česta polazna tačka u postupku projektovanja adekvatnog modela predikcije. Ono što krasi ovu metodu jeste njena nesposobnost da se preobuči. Dakle, u smislu pomenosti i varijanse, ova metoda pati od pomenosti, tj. retko kada daje zadovoljavajuće rezultate, ali je izrazito robusna i otporna na pojavu preobučavanja (*overfitting*). Samim tim performanse ovakvog modela predstavljaju minimalni zahtev za performanse drugih moćnijih i kompleksnijih modela klasifikacije.

Za potrebe teze korištena *scikit-learn* implementacija *softmax* regresije [22]. Pre obučavanja klasifikatora prvo je potrebno podesiti njegove parametre. Prvi parametar je regularizacija  $c$  i predstavlja toleranciju algoritma na greške odbiraka, mala vrednost ovog parametra, pomoću ovog hiper-parametra algoritam postaje imun na autlajere. Drugi parametar jesu težine klasa, prilikom obučavanja veći značaj se pridaje odbircima klase čija je težina veća i suprotno. Podešavanje težina klasa potrebno je u slučaju nebalansiranih klasa.

U tabeli 4.3.1 data su podešavanja *softmax* klasifikatora za personalnu klasifikaciju kod 9 različitih ispitanika. Konstanta regularizacije je označena sa  $c$ , preostale kolone sadrže težine klasa.

**Tabela 4.3.1 Hiper-parametri softmax klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika**

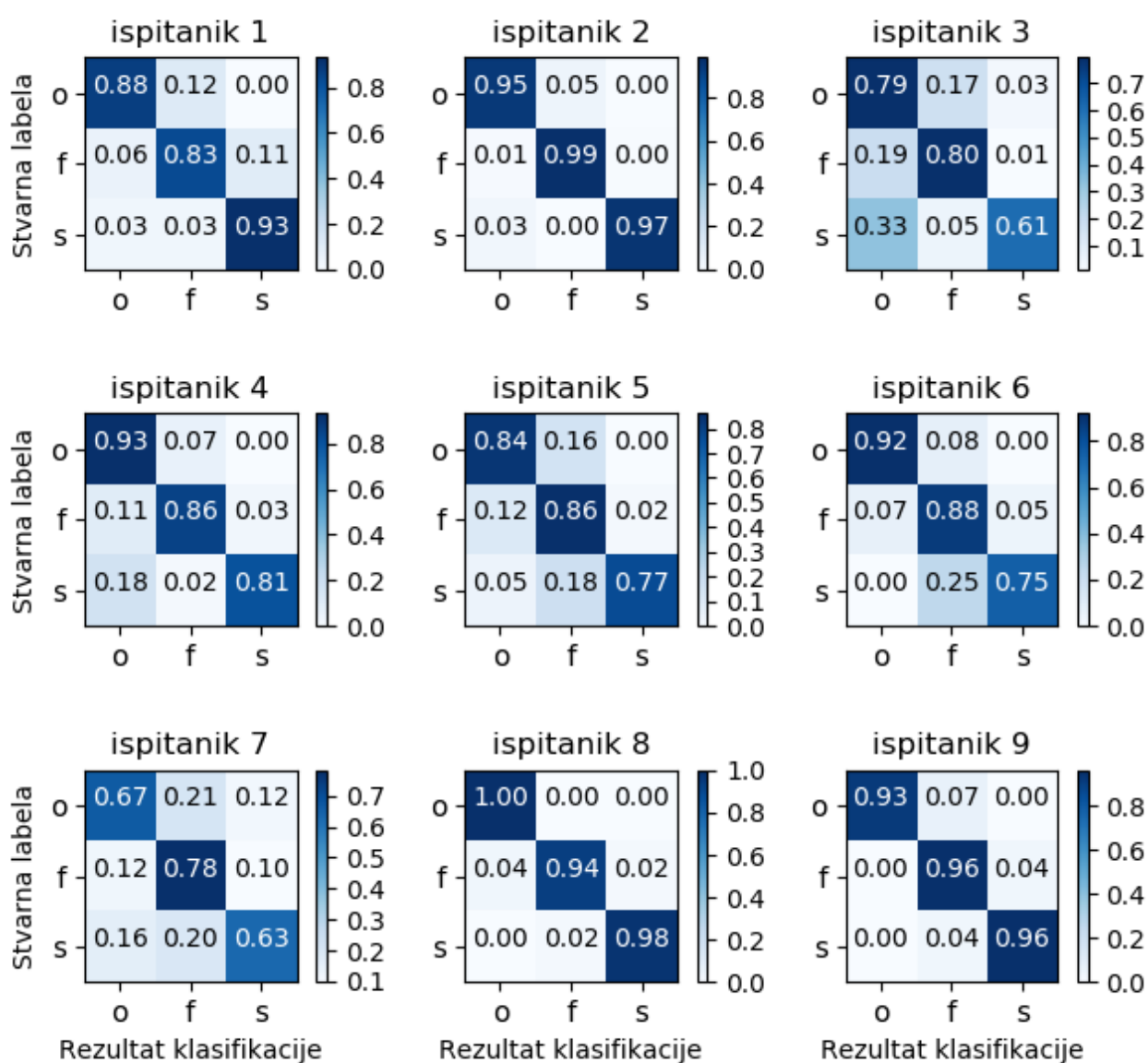
	$c$	opuštenost	fokusiranost	strah
ispitanik 1	1	2.5	0.5	2
ispitanik 2	1	0.85	0.5	2
ispitanik 3	1	1	0.5	6
ispitanik 4	1	1	0.7	5
ispitanik 5	1	1	1	1



ispitanik 6	1	1	0.5	4
ispitanik 7	1	1	0.5	3
ispitanik 8	1	1	0.5	3
ispitanik 9	1	1	0.5	3

Kao što je napomenuto u ranijim poglavljima, broj odbiraka u klasama je nebalansiran pri čemu su najbrojniji odbirci klase fokusiranosti a najmalobrojniji odbirci klase straha. U tabeli 4.3.1 može se videti da je za kvalitetnu klasifikaciju bilo potrebno dati veću težinu malobrojnijim klasama.

Rezultati klasifikacije na testirajućem skupu za 9 ispitanika prikazani su na slici 4.3.2 u vidu matrica konfuzije. Na slici se vidi da se kod nekih ispitanika klasifikacija može izvršiti sa visokom stopom uspešnosti čak i pomoću linearnog klasifikatora (ispitanici 2, 8 i 9) dok je za uspešnu klasifikaciju kod drugih ispitanika potrebno upotrebiti složeniji klasifikator.



Slika 4.3.2 Rezultati klasifikacije za *softmax* regresiju na testirajućem skupu, značenje oznaka na slici je: o- opuštenost, f-fokusiranost, s-strah

Za ispitanika broj 2 linearni *softmax* klasifikator je odabran kao finalni model predikcije. To znači da se pristupati projektovanju kompleksnijih modela za ovog ispitanika.

#### 4.3.3 Rezultati primene SVM metode

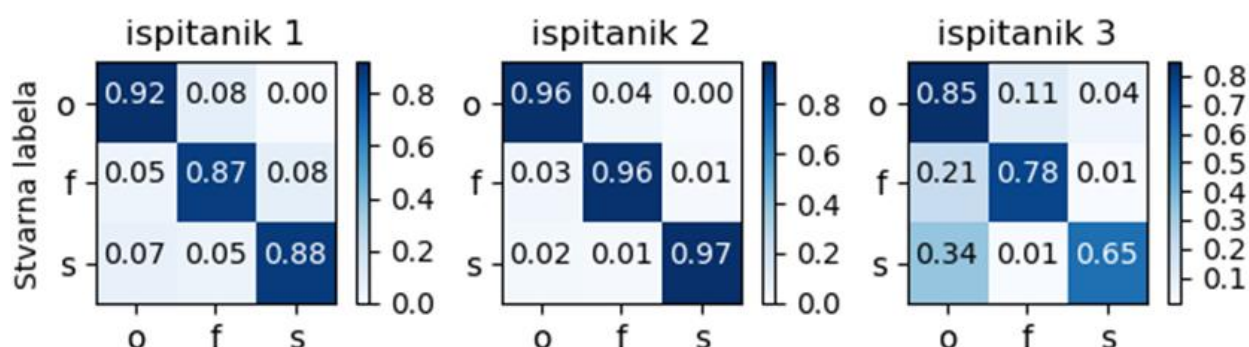
Metod nosećih vektora je vrlo čest i opšte prihvaćen model za predikciju. U zavisnosti od odabrane kernelske funkcije SVM može rezultovati linearnim, kvadratnim, tj. uopšte polinomijalnim te nelinearnim klasifikatorom. Detaljan opis algoritma je dat u odeljku 3.2.1.

Posle *softmax* regresije pokušana je personalizovana klasifikacija linearnim SVM klasifikatorom. I za ovaj klasifikator korištena je *scikit-learn* implementacija [22]. Parametri funkcije su isti kao kod *softmax* regresije i opisani su u odeljku 3.2.1. U tabeli 4.3.2 prikazani su usvojeni hiper-parametri linearnog SVM klasifikatora za 9 različitih ispitanika (u pitanju su isti ispitanici kao u poglavlju o *softmax* regresiji).

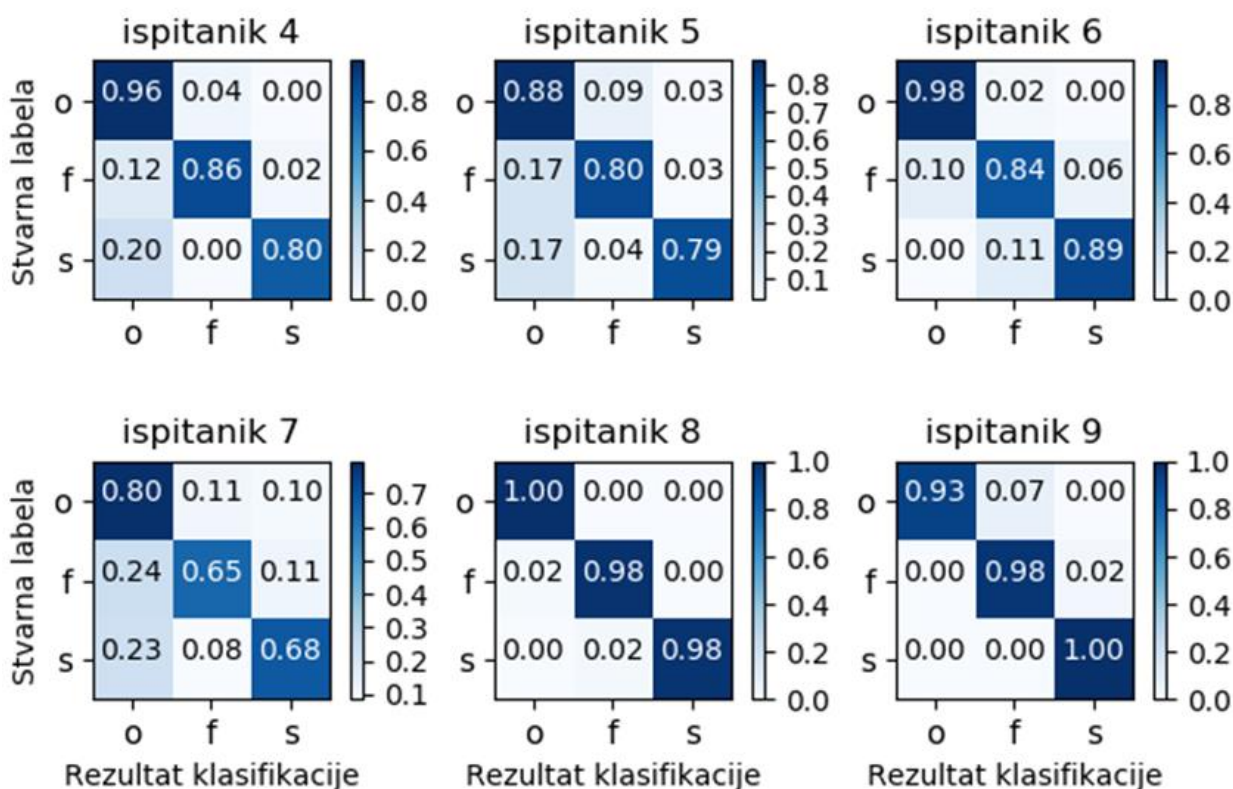
**Tabela 4.3.2 Parametri linearnog SVM klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika**

	<i>C</i>	opuštenost	fokusiranost	strah
ispitanik 1	13	2	1	3
ispitanik 2	10	2	1	3
ispitanik 3	10	2	0.5	3
ispitanik 4	2	1.6	1.2	7
ispitanik 5	1.2	2	0.5	3
ispitanik 6	5	0.7	0.1	4
ispitanik 7	1	3	0.8	4
ispitanik 8	1	1	0.8	3
ispitanik 9	25	5	0.45	2

Rezultati klasifikacije na testirajućim skupovima skupu za 9 ispitanika prikazani su na slikama 4.4.3 i 4.3.4 u vidu matrica konfuzije. Kao što se može videti na slici linearni SVM pokazuje bolje performanse od linearnog *softmax* klasifikatora.



**Slika 4.3.3 Rezultati klasifikacije za linearni SVM klasifikator na testirajućem skupu za ispitanike 1-3, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah**



Slika 4.3.4 Rezultati klasifikacije za linearni SVM klasifikator na testirajućem skupu za ispitanike 4-9, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah

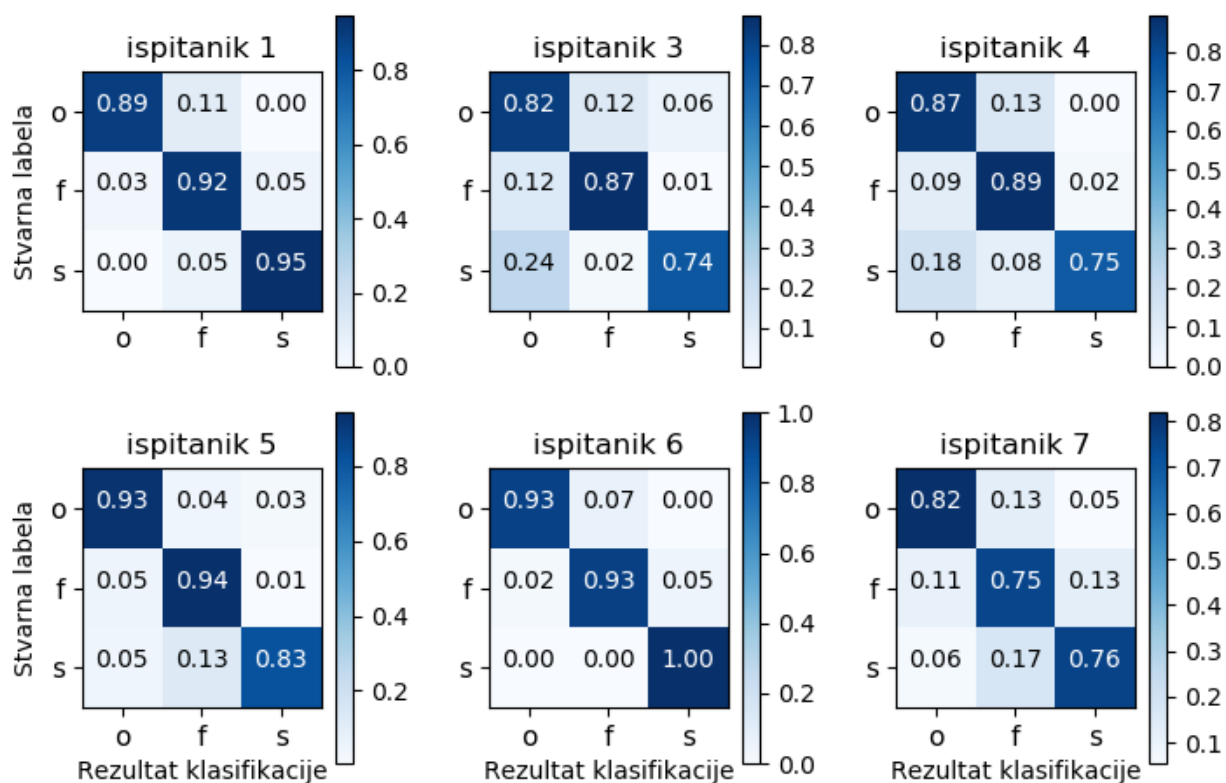
Vidimo da je za ispitanika broj 2 linearni *softmax* model dao bolje rezultate. Ispitanicima broj 8 i 9 je linearni SVM model usvojen kao finalni model predikcije zbog jako dobrih rezultata. Za ostale ispitanika pristupa se projektovanju nelinearnog SVM modela sa kernelskom funkcijom:

$$k^{rbf}(x_i, x) = e^{-\frac{\|x_i - x\|^2}{2a^2}}$$

Lista usvojenih hiper-parametara prikazana na u tabeli 4.3.3. Konfuzione matrice kao procena performansi modela na testirajućem skupu prikazane su na slici 4.3.5:

Tabela 4.3.3 Parametri nelinearnog SVM klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika

	<i>c</i>	<i>a</i>	opuštenost	fokusiranost	strah
ispitanik 1	10	0.015	4	4	4
ispitanik 3	10	0.0033	1.1	0.5	4
ispitanik 4	5	0.025	2	0.1	6
ispitanik 5	10	0.00297	0.45	2	6
ispitanik 6	1.6	0.009	0.45	2	6
ispitanik 7	0.1	0.006	1	1	6



Slika 4.3.5 Rezultati klasifikacije za nelinearni SVM klasifikator na testirajućem skupu, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah

Rezultati na ispitanicima 1, 5 i 6 su sasvim zadovoljavajući te se ovakav model usvaja kao finalni. Za ispitanike 3, 4 i 7, u potrazi za modelom boljih performansi, pristupamo projektovanju *Random Forest* modela.

#### 4.3.4 Rezultati primene Random Forest metode

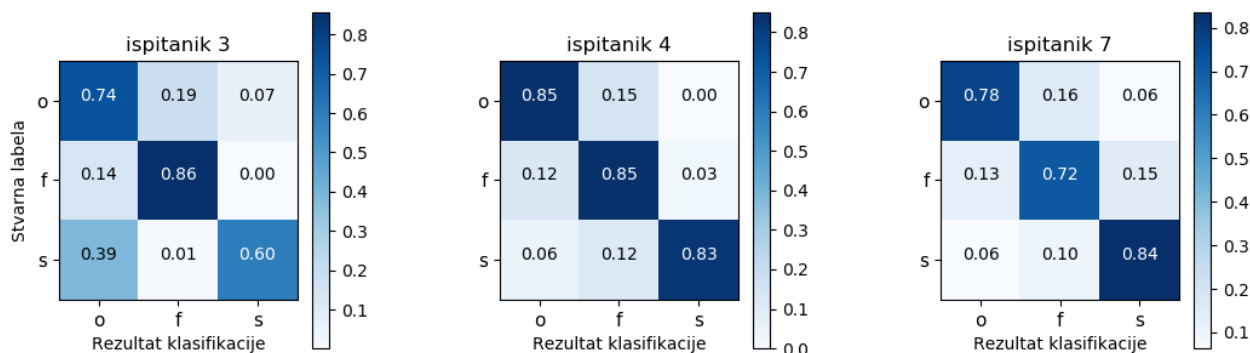
*Random Forest* je fleksibilan, lagan za korišćenje, algoritam mašinskog učenja koji i bez podešavanja hiper-parametara često daje solidne rezultate. RFC ne predstavlja nužno kompleksniji algoritam, tj model veće varijanse od *rbf* SVM modela, ali usled potpuno drugačijeg pristupa često daje bolje rezultate. *Random Forest* algoritam je opisan u odeljku 3.3.3.

RFC algoritam je primenjen na ispitanike 3, 4 i 7 sa hiper-parametrima datim u tabeli 4.3.4. Parametar *random\_state* zapravo predstavlja seme generatora slučajnih brojeva u *scikit-learn* biblioteci. Na slici 4.3.6 prikazane su konfuzione matrice ovih modela na testirajućem skupu.

Tabela 4.3.4 Parametri *random forest* klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika

	Ispitanik 3	ispitanik 4	ispitanik 7
<i>max_depth</i>	7	5	8
<i>n_estimators</i>	80	50	120
<i>max_features</i>	0.25	0.3	0.134
<i>min_samples_split</i>	31	42	34
<i>min_samples_leaf</i>	11	24	12

<i>random_state</i>	0	2	42
opuštenost	0.4	1	0.81
fokusiranost	0.1	0.5	0.18
Strah	3.5	5	3.9



**Slika 4.3.6** Rezultati klasifikacije za *Random Forest* klasifikator na testirajućem skupu, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah

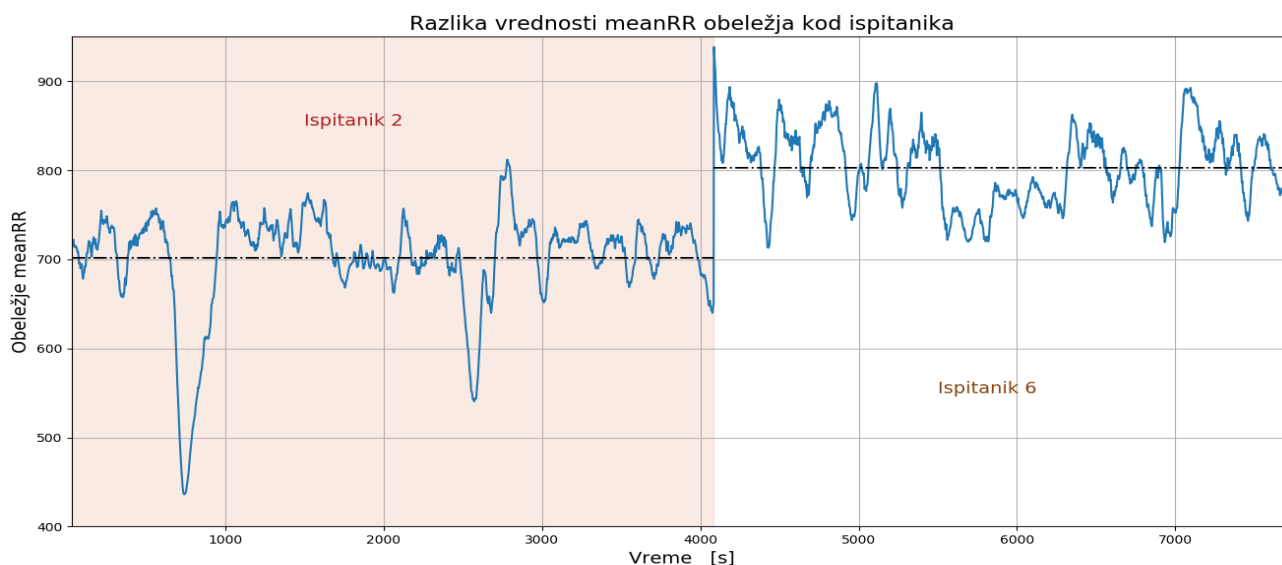
U poređenju sa rezultatima nelinearnog SVM klasifikatora sa slike 4.3.5 vidimo da ispitanik broj 3 ima značajno slabije performanse kod *Random Forest* klasifikatora. Zbog toga se u slučaju ispitanika broj 3 kao finalni model usvaja SVM sa *rbf* kernelskom funkcijom. Ispitanik broj 4, sa *Random Forest* klasifikatorom ima slabije rezultate za prve dve labele ali bolji rezultat za labelu straha. Dakle metoda slučajnih šuma da ujednačenu preciznost za sve tri klase te se iz tog razloga usvaja kao finalni model za ispitanika 4. Ispitanik broj 7 ima bolje performanse sa *rbf* SVM klasifikatorom. Međutim, *Random Forest* klasifikator daje nešto veću preciznost za predikciju straha. U zavisnosti od aplikacije ovo može biti presudan faktor, te ukoliko nam je prioritet prepoznati strah tada bismo usvojili RFC kao finalni model.

## 4.4 Interpersonalna klasifikacija

Cilj interpersonalne klasifikacije jeste projektovanje jedinstvenog modela za predikciju emocija koji bi se mogao primeniti nad svim ispitanicima. Skup podataka se formira spajanjem svih individualnih skupova čije je formiranje opisano u odeljku 4.3.1. U odeljku 4.4.1 vrši se diskusija o normalizaciji podataka među ispitanicima. Potom, u narednim odeljcima, pristupa se projektovanju modela predikcije istim redosledom kao u poglavlju o personalizovanoj klasifikaciji.

### 4.4.1 Normalizacija podataka

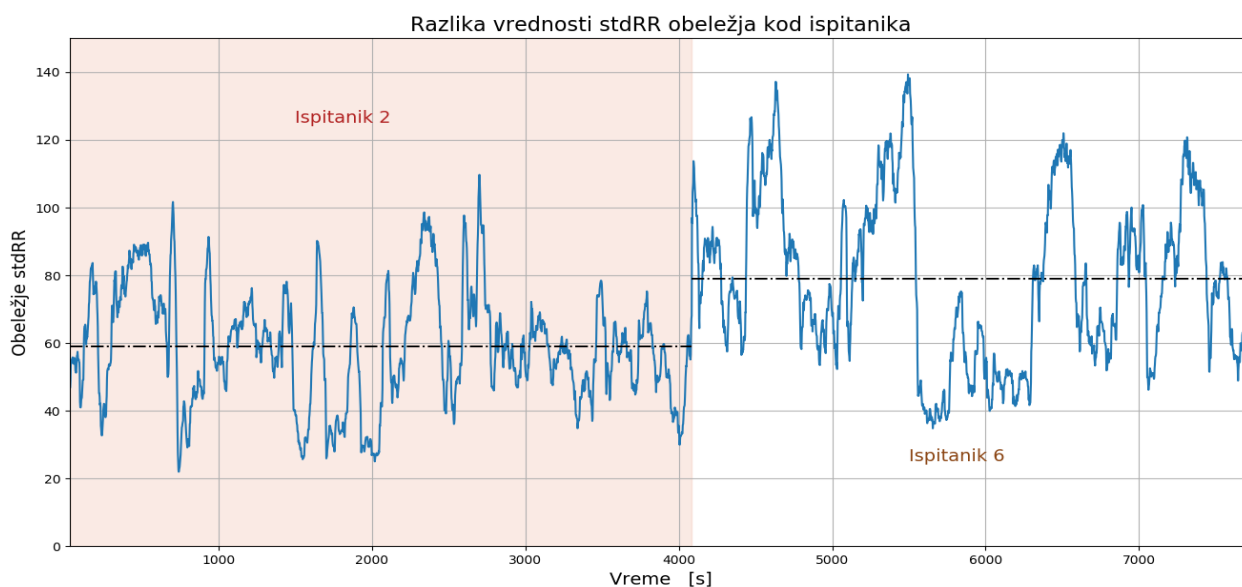
Opštepoznato je da elektrofiziološki signali ljudi nisu nužno uporedivi. Njihove karakteristike zavise od mnogih faktora kao što su zdravstveno stanje, sportske navike, godine, ishrana i slično. U našem ekperimentu učestvovali su ispitanici različitog pola, životnih navika i starosti. Na slici 4.4.1 je demonstrirana pojava različite vrednosti istog obeležja kod različitih ispitanika tokom eksperimenta.



**Slika 4.4.1 Razlika vrednosti *meanRR* obeležja tokom eksperimenta kod dva ispitanika**

Na slici 4.4.1 je na obojenom delu grafika prikazano obeležje *meanRR* za ispitanika 2, dok je na njega nadovezано isto obeležje ispitanika 6. Prikazana obeležja se odnose na trajanje celog eksperimenta tokom kog su ispitanici prolazili kroz različite sadržaje.

Vidimo da obeležja nisu uporediva kako po svojoj srednjoj vrednosti tako i po svojoj dinamici promene. Prema prikupljenim informacijama, ispitanik broj 6 praktikuje sportske aktivnosti znatno češće i u značajnoj meri više od ispitanika broj 2. Poznato je da je sportistima u stanju mirovanja srčani ritam znatno sporiji što slika 4.4.1 i potvrđuje. Kada je u pitanju dinamika srčanog ritma sa slike 4.4.2 vidimo značajnu razliku i u standardnim devijacijama njihovih RR intervala – *stdRR*.



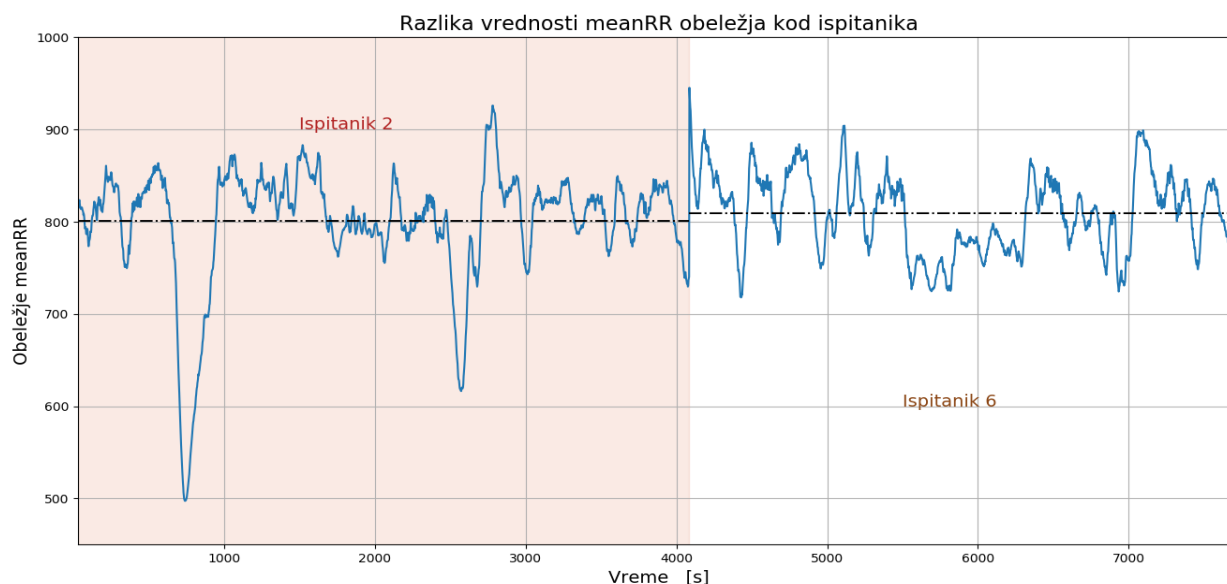
**Slika 4.4.2 Razlika vrednosti *stdRR* obeležja tokom eksperimenta kod dva ispitanika**

Odatle potreba za nekim vidom normalizacije kako bi se postigla međusobno uporediva obeležja kod različitih ispitanika. Primenjeno je svođenje srednje vrednosti HRV signala na ekvivalentni srčani ritam od 75 otkucaja u minuti (1.25 otkucaja u sekundi). HRV signal je pomnožen sa koeficijentom:



$$1.25/\overline{hrv}^*$$

Pri tome  $\overline{hrv}^*$  nije obična srednja vrednost  $hrv$  signala, već je u pitanju srednja vrednost signala kada se ostave odbirci između 20. i 80. percentila amplitude  $hrv$  signala. Na slici 4.4.3 prikazan je rezultat normalizacije na obeležju  $meanRR$ .

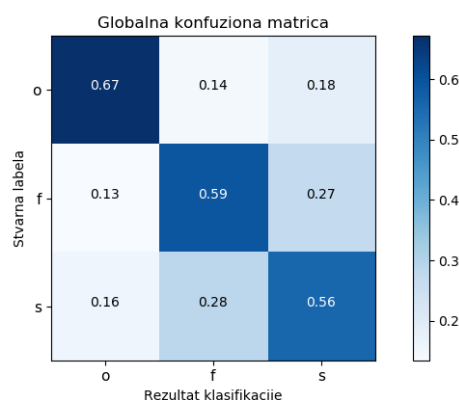


**Slika 4.4.3 Razlika vrednosti  $stdRR$  obeležja tokom eksperimenta kod dva ispitanika**

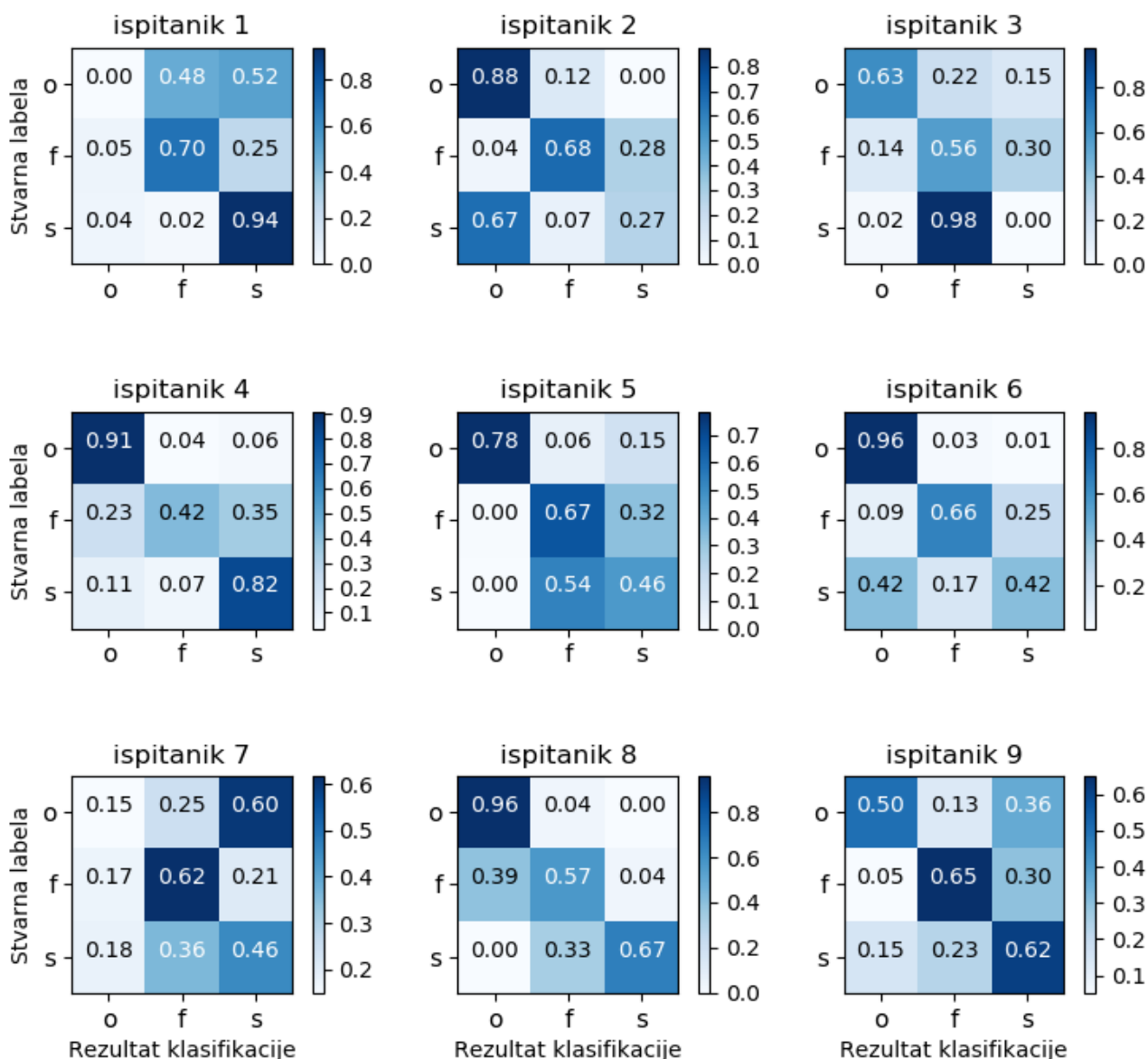
Slična normalizacija sprovedena je i na  $BRV$  signalu svodeći ga na ekvivalentni ritam disanja od 15 respiratornih ciklusa u minuti.

#### 4.4.2 Rezultati primene softmax metode

Primenjuje se linearni algoritam za klasifikaciju, opisan u odeljku 3.3.1. Odabran je hiper-parametar regulacije  $c=10$ , a parametri balansiranja kao 0.6, 0.4, i 3.2 za opuštenost, fokusiranost i strah respektivno. Na slici 4.4.4 prikazan je rezultat klasifikacije testirajućeg skupa dok su na slici 4.4.5 prikazani rezultati projektovanog klasifikatora na individualnim testirajućim skupovima ispitanika.



**Slika 4.4.4 Rezultat interpersonalne klasifikacije softmax linearnog klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah**



**Slika 4.4.5** Rezultat softmax linearnog klasifikatora na pojedinačnim testirajućim skupovima

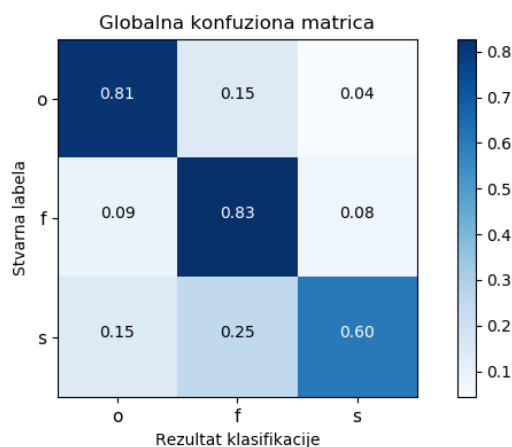
Sa slike 4.4.4 vidimo drastičan pad performansi u poređenju sa rezultatima personalizovane klasifikacije. Ipak, rezultat na globalnom testirajućem skupu daje ujednačen uspeh na sve tri labela. Međutim kada se ovako obučeni model testira na individualnim testirajućim skupovima ispitanika, kako je prikazano na slici 4.4.5, zapažamo nekonzistentne performanse.

Dok je kod ispitanika broj 8 i 9 rezultat uporediv sa globalnim rezultati ostalih ispitanika su nezadovoljavajući. To govori o ograničenoj primeni ovakvog modela, i nagoveštava da se složene pojave u elektrofiziološkim signalima, koji su odraz ljudskih emocija, ne mogu tako lako kategorizovati.



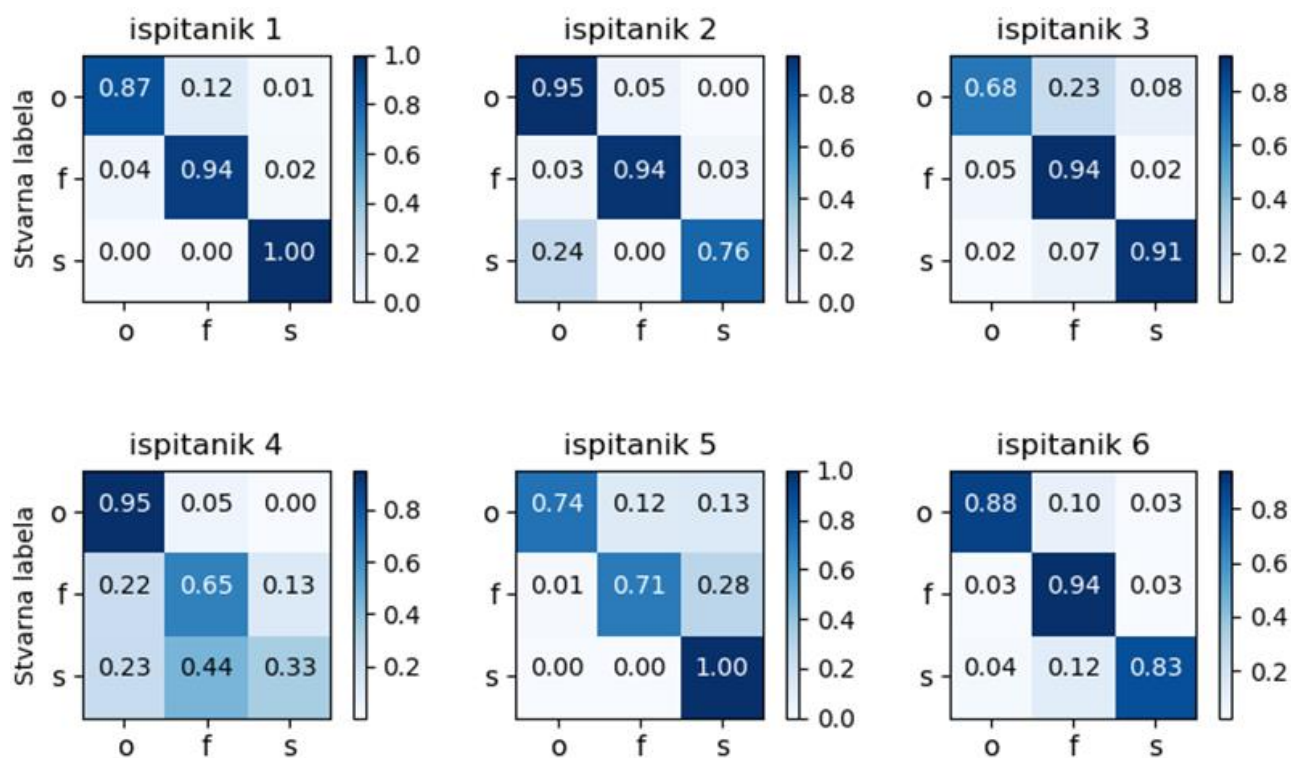
#### 4.4.3 Rezultati primene SVM metode

Projektuje se model SVM klasifikatora sa rbf kernelskom funkcijom, opisan u pododeljku 3.3.2. Na slici 4.4.6 prikazan je rezultat klasifikacije celokupnog, združenog testirajućeg skupa.

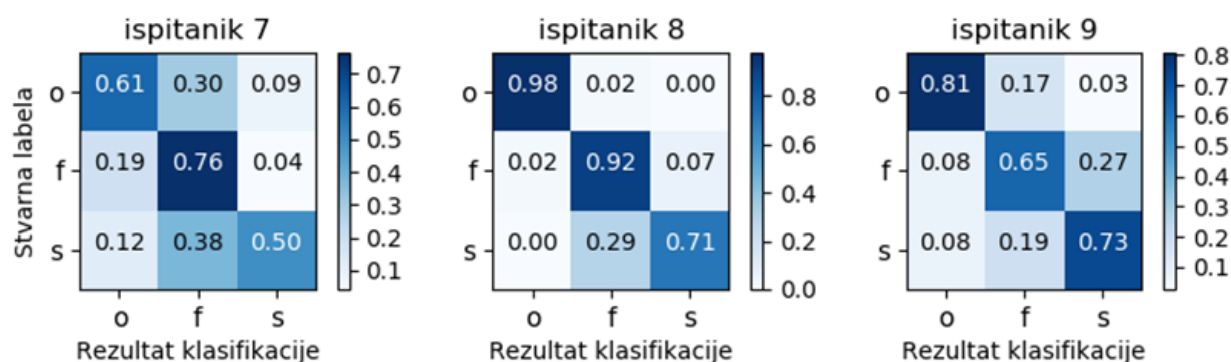


**Slika 4.4.6** Rezultat interpersonalne klasifikacije nelinearnog SVM klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah

Vidi se značajan napredak u odnosu na linearni model obučavan u prethodnom odeljku. Na slikama 4.4.7 i 4.4.8 prikazan je rezultat klasifikacije na individualnim testirajućim skupovima ispitanika.



**Slika 4.4.7** Rezultat *rbf* SVM klasifikatora na pojedinačnim test skupovima ispitanika 1-6

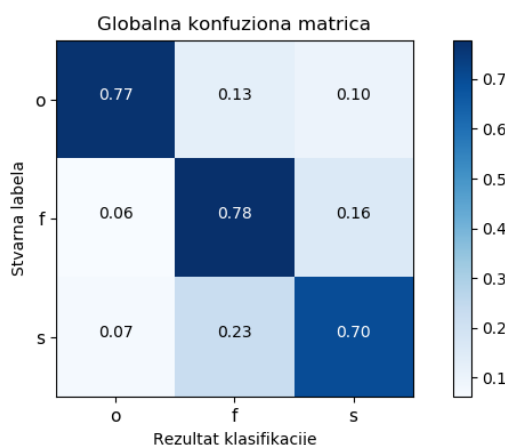


Slika 4.4.8 Rezultat *rbf* SVM klasifikatora na pojedinačnim test skupovima ispitanika 7-9

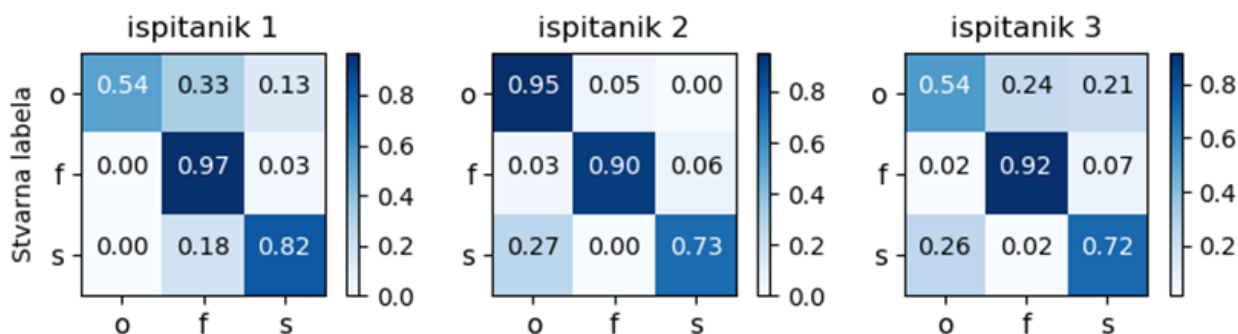
Konstanta regularizacije je uzeta kao  $c=5$ , parametar kernelske funkcije  $a=0.008$  a parametri balansiranja su 0.6, 0.3 i 10 za opuštenost, fokusiranost i strah respektivno. Pored boljih globalnih performansi uočavaju se i bolje individualne performanse. Ipak, postoje ispitanici kod kojih ovako projektovan model nije konzistentan. Ispitanici 4 i 7 imaju nešto lošije rezultate od ostalih ispitanika.

#### 4.4.4 Rezultati primene Random Forest metode

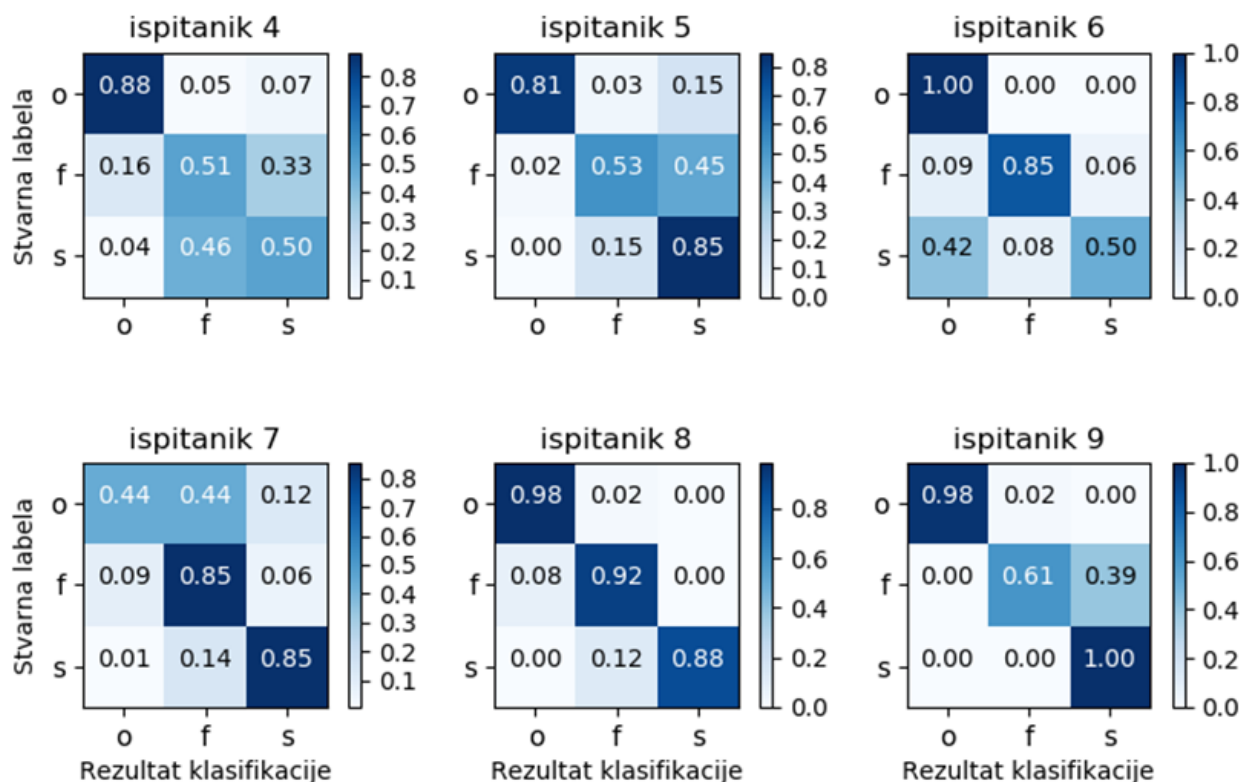
Projektovan je klasifikator “slučajnih šuma”, opisan u odeljku 3.3.3. Na slici 4.4.9 prikazan je rezultat klasifikacije celokupnog, združenog testirajućeg skupa, dok su na slikama 4.4.10 i 4.4.11 prikazani individualni rezultati.



Slika 4.4.9 Rezultat interpersonalne klasifikacije *Random Forest* klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah



Slika 4.4.10 Rezultat *Random Forest* klasifikatora na pojedinačnim testirajućim skupovima ispitanika 1-3



Slika 4.4.11 Rezultat *Random Forest* klasifikatora na pojedinačnim testirajućim skupovima ispitanika 4-9

Hiper-parametri su izabrani kao:

*hiper – parametri*: {*max\_depth* = 8, *n\_estimators* = 120, *max\_features* = 0.2,  
*min\_samples\_split* = 24, *min\_samples\_leaf* = 12, *random\_state* = 42}.

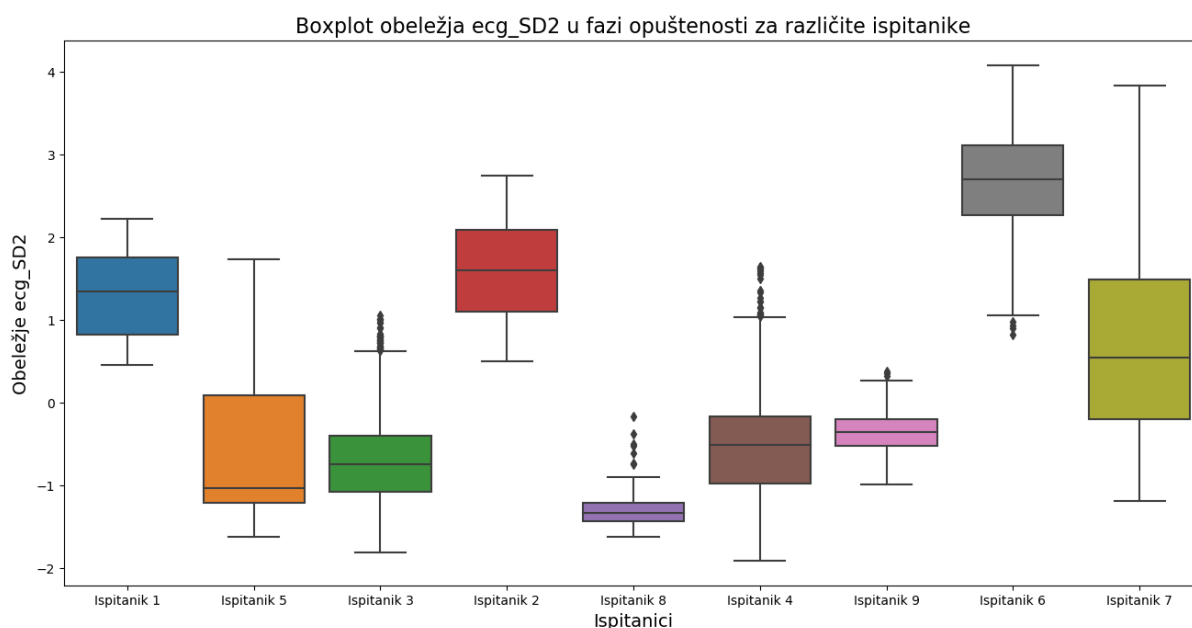
Dok su parametri balansiranja dobili vrednosti 1, 0.42 i 6 za klase opuštenost, fokusiranost, strah respektivno.

U pogledu globalnih performansi, sa slike 4.4.9, *Random Forest* klasifikator daje bolje balansirane performanse od nelinearnog SVM modela. Međutim, kada uporedimo individualne rezultate RFC klasifikatora sa SVM modelom uočavamo da su performanse RFC značajno slabije konzistentne. Samo ispitanici broj 8, 9 i 2 imaju zadovoljavajuće rezultate, dok drugi svedoče o jako lošim performansama *Random Forest* modela.

## 5. Zaključak

U poglavlju o klasifikaciji opisan je postupak projektovanja kako personalnog tako i interpersonalnog modela predikcije emocija. U potpoglavlju 4.3, koje se bavi personalizovanim modelima, je pokazano da ukoliko se klasifikator projektuje i obučava za konkretnog ispitanika postižu se relativno zadovoljavajući rezultati. Dok je za neke ispitanike finalni model predikcije usvojen kao linearni *softmax* ili SVM klasifikator, za druge ispitanike je bilo potrebno projektovati kompleksnije modele. Pored toga, dok usvojeni modeli jednih ispitanika postižu ujednačene performanse na različitim klasama, finalni modeli za druge ispitanike, uprkos našim težnjama, imaju disbalansirane performanse. Odavde se naziru naznake individualnih razlika između ispitanika, tj. razlika u načinu ispoljavanja emocija kroz elektrofiziološke signale. Uprkos ovome, individualno podešavani modeli su u većoj ili manjoj meri sposobni da vrše predikciju emocija ispitanika na kojima su obučavani.

U potpoglavlju 4.4 vršeno je projektovanje interpersonalnog klasifikatora. U odeljku 4.4.1 vršena je normalizacija HRV i BRV signala u cilju dobijanja uporedivih obeležja. Normalizacija koja je primenjena nije opštepoznata u literaturi već je rezultat našeg istraživanja i rezonovanja. U odeljku 4.4.2 projektovani linearni *softmax* klasifikator je pokazao znatno lošije performanse nego za slučaj personalizovane klasifikacije iz odeljka 4.3.2. Ovo nagoveštava da je problem interpersonalne klasifikacije izrazito nelinearne prirode. U odeljcima 4.4.3 i 4.4.4 projektovani su nelinearni klasifikatori SVM model sa *rbf* kernelskom funkcijom i model slučajnih šuma (*Random Forest*). Među navedenim modelima, za interpersonalnu klasifikaciju, najbolje rezultate u globalnom smislu kao i na nivou ispitanika je dao SVM model. Međutim, postignute performanse na konkretnim ispitanicima su znatno lošije od individualno projektovanih modela. Navedena pojava je prvenstveno rezultat individualnih razlika između ispitanika u pogledu elektrofizioloških reakcija. Posledicu ove pojave vidimo sa slike 5.1 koja predstavlja *boxplot* obeležja *ecg\_SD1* tokom faze opuštanja za različite ispitanike.



Slika 5.1 Boxplot obeležja *ecg\_SD2* u fazi opuštenosti za različite ispitanike

Na slici 5.1 se vidi da za različite ispitanike a za istu emociju isto obeležje ima drastično različite raspodele. Slična pojava se vidi i na drugim obeležjima. Iako je primenjena normalizacija dovela HRV i BRV signale na uporedive nivoe, ovakva normalizacija nije dovela do uporedivosti dinamike ovih signala što se odrazilo na ekstrahovana obeležja.

Za uspešnu interpersonalnu klasifikaciju potrebno je pronaći adekvatniji vid normalizacije. U cilju ovoga bi bilo svrsishodno da postoji matematički model HRV i BRV signala kako bi se mogle modelovati i normalizovati ne samo amplituda ovih signala već i njihove frekvencijske karakteristike i nelinearna dinamika.

Još jedno potencijalno poboljšanje pristupa bi bilo da se eksperiment izvrši na većem broju ispitanika pri čemu bi se mogao više puta ponoviti na jednom ispitaniku. Ponavljanje na jednom ispitaniku bi nam obezbedilo veći obim skupa podataka nad kojima bi se mogao projektovati kvalitetniji i robusniji personalizovani model predikcije. Ovakav pristup bi doprineo i rešavanju problema nebalansiranih klasa. Sa druge strane, veći broj ispitanika bi nam omogućio bolju uvid u raznolikost elektrofizioloških promena u populaciji. Tu bi se otvorio prostor za kategorizaciju ljudi u grupe sa sličnim načinom ispoljavanja emocija. Nad ovakvim grupama bi se mogla izvršiti kvalitetnija interpersonalna klasifikacija ■

## 6. Spisak referenci

- [1] H. Salih and L. Kulkarni, "Study of video based facial expression and emotions recognition methods," 2017 International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), Palladam, 2017, pp. 692-696
- [2] P. Tarnowski, M. Kołodziej, A. Majkowski, and R. J. Rak, "Emotion recognition using facial expressions," International Conference on Computational Science, pp. 1175–1184, June 2017
- [3] A. Bandrabur, L. Florea, C. Florea and M. Mancas, "Emotion identification by facial landmarks dynamics analysis," *2015 IEEE International Conference on Intelligent Computer Communication and Processing (ICCP)*, Cluj-Napoca, 2015, pp. 379-382.
- [4] J. Kim and E. Andre, "Emotion recognition based on physiological changes in music listening," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 12, pp. 2067–2083, 2008.
- [5] Kim, K.H., Bang, S.W. & Kim, S.R. *Med. Biol. Eng. Comput.* (2004) 42: 419.
- [6] <http://www.smartex.it/en/our-products/232-wearable-wellness-system-www>
- [7] <https://mbraintrain.com/smaring/>
- [8] Li, R., Principe, J.C., 2006. Blinking artifact removal in cognitive EEG data using ICA. *Conf. Proc. IEEE Eng. Med. Biol. Soc.* 1, 5273–5276.
- [9] Task Force of the European Society of Cardiology and the North American Society of Pacing and Electrophysiology. Heart rate variability: standards of measurement, physiological interpretation, and clinical use. *Circulation* 1996; 93: 1043-65.
- [10] [https://en.wikipedia.org/wiki/Fast\\_Fourier\\_transform](https://en.wikipedia.org/wiki/Fast_Fourier_transform)
- [11] Khandoker et al, *Poincaré Plot Methods for Heart Rate Variability Analysis*, - Springer US – 2013
- [12] P. Guzik, J. Piskorski, T. Krauze, A. Wykretowicz, H. Wysocki, Heart rate asymmetry by Poincaré plots of RR intervals. *Biomed. Tech.* 51, 530–537 (2006)
- [13] A. Porta, K.R. Casali, A.G. Casali, T. Gneccchi-Ruscone, E. Tovaldini, N. Montano, S. Lange, D. Geue, D. Cysarz, P. Van Leeuwen, Temporal asymmetries of short-term heart period variability are linked to autonomic regulation. *Am. J. Physiol. Regul. Integr. Comp. Physiol.* 295, R550–R557 (2008)
- [14] Kuusela, Tom. (2012). *Methodological Aspects of Heart Rate Variability Analysis*. 9-42. 10.1201/b12756-4.

- [15] A.Bhavani Sankar, D.Kumar, K.Seethalakshmi, “Enhanced Method for Extracting Features of Respiratory Signals and Detection of Obstructive Sleep Apnea Using Threshold Based Automatic Classification Algorithm”, International Journal of Computer Science & Emerging Technologies (E-ISSN: 2044-6004) 38, Volume 1, Issue 4, December 2010
- [16] Dengyong Zhou, Olivier Bousquet, Thomas Navin Lal, Jason Weston, Bernhard Schoelkopf. Learning with local and global consistency (2004)
- [17] <http://www.holehouse.org/mlclass/>
- [18] [https://en.wikipedia.org/wiki/Pearson\\_correlation\\_coefficient](https://en.wikipedia.org/wiki/Pearson_correlation_coefficient)
- [19] [https://en.wikipedia.org/wiki/Mutual\\_information](https://en.wikipedia.org/wiki/Mutual_information)
- [20] Smola, A.J. & Schölkopf, B. Statistics and Computing (2004) 14: 199.
- [21] L.Breiman, “Random Forests”, Machine Learning, 45(1), 5-32, 2001
- [22] Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011

## 7. Spisak skraćenica

- 1) VR – virtuelna realnost
- 2) 3D – trodimenzionalni
- 3) EKG - elektrokardiogram
- 4) EEG - elektroencefalogram
- 5) EMG - elektromiogram
- 6) HRV - *Heart rate variability*
- 7) BB - *Breath to Breath*
- 8) BRV - *Breathing Rate Variability*
- 9) ICA - *Independent Component Analysis*
- 10) FFT - *Fast Fourier Transform*
- 11) ULF – *Ultra low frequency*
- 12) VLF – *Very low frequency*
- 13) LF – *Low frequency*
- 14) HF – *High frequency*
- 15) CCM - *complex correlation measure*
- 16) ApEn – *Approximate entropy*
- 17) SampEn – *Sample entropy*
- 18) UHF – *Ultra high frequency*
- 19) LDA - Linearna diskriminantna analiza
- 20) PCA - *Principal component analysis*
- 21) PCC – Pearsonov koeficijent korelacije
- 22) MI - *Mutual information*
- 23) RFC – *Random Forest Classifier*
- 24) HR - *Heart rate signal*
- 25) SVM – *Support Vector Machine*
- 26) oob – *out-of-bag* skup odbiraka



## 8.Spisak slika

Slika 2.1.1 Postavka eksperimenta sa ispitanikom .....	6
Slika 2.2.1 Filtrirani EKG signal .....	7
Slika 2.2.2 Signal HRV .....	7
Slika 2.2.3 Filtriran signal respiracije .....	8
Slika 2.2.4 BRV ( <i>Breathing Rate Variability</i> ) signal .....	8
Slika 2.2.5 EEG segment posle filtracije i odstranjivanja očnih atrefakata .....	9
Slika 4.1 Dijagram toka opšte procedure.....	20
Slika 4.1.1 Ručno formirana referenca .....	22
Slika 4.1.2 LDA prikaz podataka labeliranih ručnom metodom bez izuzimanja odbiraka niske pouzdanosti .....	22
Slika 4.1.3 LDA prikaz podataka labeliranih ručnom metodom sa izuzimanjem odbiraka niske pouzdanosti .....	23
Slika 4.1.4 LDA prikaz podataka labeliranih polu-supervizijskom klasterizacijom .....	23
Slika 4.1.5 Rezultat klasterizacije u vidu verovatnoća klasne pripadnosti .....	24
Slika 4.2.1 Udeo sadržane varijanse u odnosu na broj PCA komponenti.....	25
Slika 4.2.2 Informativnost EKG obeležja .....	27
Slika 4.2.3 Informativnost obeležja respiracije.....	27
Slika 4.2.4 Važnost EKG obeležja prema RFC modelu .....	28
Slika 4.2.5 Važnost obeležja respiracije prema RFC modelu.....	29
Slika 4.3.1 Podela eksperimenta na obučajući, validirajući i testirajući skup: crvena šrafura predstavlja validirajući skup, plava i ljubičasta šrafura predstavlja testirajući skup, nešrafirani delovi eksperimenta predstavljaju obučavajući skup .....	30
Slika 4.3.2 Rezultati klasifikacije za <i>softmax</i> regresiju na testirajućem skupu, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah .....	32
Slika 4.3.3 Rezultati klasifikacije za linearni <i>SVM</i> klasifikator na testirajućem skupu za ispitanike 1-3, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah .....	33
Slika 4.3.4 Rezultati klasifikacije za linearni <i>SVM</i> klasifikator na testirajućem skupu za ispitanike 4-9, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah .....	34
Slika 4.3.5 Rezultati klasifikacije za nelinearni <i>SVM</i> klasifikator na testirajućem skupu, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah .....	35
Slika 4.3.6 Rezultati klasifikacije za <i>Random Forest</i> klasifikator na testirajućem skupu, značenje oznaka na slici je: o-opuštenost, f-fokusiranost, s-strah .....	36
Slika 4.4.1 Razlika vrednosti <i>meanRR</i> obeležja tokom eksperimenta kod dva ispitanika.....	37
Slika 4.4.2 Razlika vrednosti <i>stdRR</i> obeležja tokom eksperimenta kod dva ispitanika.....	37

Slika 4.4.3 Razlika vrednosti <i>stdRR</i> obeležja tokom eksperimenta kod dva ispitanika.....	38
Slika 4.4.4 Rezultat interpersonalne klasifikacije softmax linearnog klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah.....	38
Slika 4.4.5 Rezultat softmax linearnog klasifikatora na pojedinačnim testirajućim skupovima .....	39
Slika 4.4.6 Rezultat interpersonalne klasifikacije nelinearnog SVM klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah.....	40
Slika 4.4.7 Rezultat <i>rbf</i> SVM klasifikatora na pojedinačnim test skupovima ispitanika 1-6 .....	40
Slika 4.4.8 Rezultat <i>rbf</i> SVM klasifikatora na pojedinačnim test skupovima ispitanika 7-9 .....	41
Slika 4.4.9 Rezultat interpersonalne klasifikacije <i>Random Forest</i> klasifikatora na celokupnom testirajućem skupu, na slici o-opuštenost, f-fokusiranost, s-strah.....	41
Slika 4.4.10 Rezultat <i>Random Forest</i> klasifikatora na pojedinačnim testirajućim skupovima ispitanika 1-3.....	41
Slika 4.4.11 Rezultat <i>Random Forest</i> klasifikatora na pojedinačnim testirajućim skupovima ispitanika 4-9.....	42
Slika 5.1 Boxplot obeležja <i>ecg_SD2</i> u fazi opuštenosti za različite ispitanike .....	43

## 9.Spisak tabela

Tabela 2.3.1 Pregled obeležja dobijenih iz HRV signala .....	9
Tabela 2.3.2 Pregled obeležja dobijenih iz BRV signala.....	10
Tabela 4.2.1 Koeficijent korelacije prve grupe obeležja .....	25
Tabela 4.2.2 Koeficijent korelacije druge grupe obeležja .....	26
Tabela 4.3.1 Hiper-parametri <i>softmax</i> klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika .....	31
Tabela 4.3.2 Parametri linearnog <i>SVM</i> klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika.....	33
Tabela 4.3.3 Parametri nelinearnog <i>SVM</i> klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika.....	34
Tabela 4.3.4 Parametri <i>random forest</i> klasifikatora za personalizovanu klasifikaciju kod različitih ispitanika .....	35