

June 18, 2024

1 Zadatak-Opcija 2

U timu predmeta na MSTeams platformi dostupna je baza slikanih šaka koje pokazuju papir“, kamen“ i makaze“. Projektovati inovativni sistem za prepoznavanje pokazanih znakova zasnovan na testiranju hipoteza.

a) Detaljno opisati algoritam za obradu slike i odabir obeležja koji prethodi samoj klasifikaciji. Algoritam treba da bude što robusniji (na različite osvetljaje, položaje šaka, načine pokazivanja znakova itd).

b) Izvršiti podelu na trening i test skup. Rezultate klasifikacije test skupa prikazati u obliku matrice konfuzije.

c) Odabrati dva znaka i dva obeležja takva da su odabrani znakovi što separabilniji u tom prostoru. Prikazati histogram obeležja za oba slova i prokomentarisati njihov oblik.

d) Za slova i obeležja pod c) projektovati parametarski klasifikator po izboru i iscertati klasifikacionu liniju.

Po jedan primer učitanih slika iz svake klase izgledaju na sledeći način:

Kamen



Figure 1: Kamen

Papir



Figure 2: Papir



Figure 3: Makaze

Za početak je potrebno sve slike binarizovati. To je odrađeno prebacivanjem slike iz RGB color sistema u HSV color sistem. Primeri H, S i V komponenti za sve tri klase kao i histogrami intenziteta dobijenih komponenti su sledeći:



Figure 4: Kamen u HSV color sistemu



Figure 5: Papir u HSV color sistemu

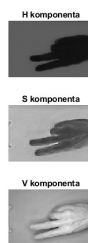


Figure 6: Makaze u HSV color sistemu

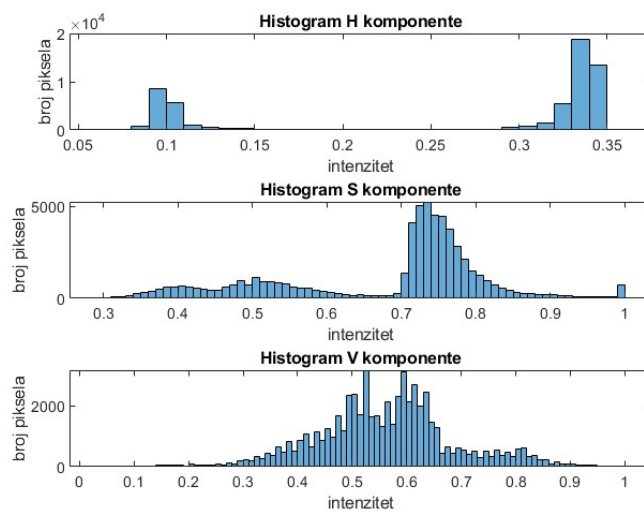


Figure 7: Histogram HSV komponenti kamena

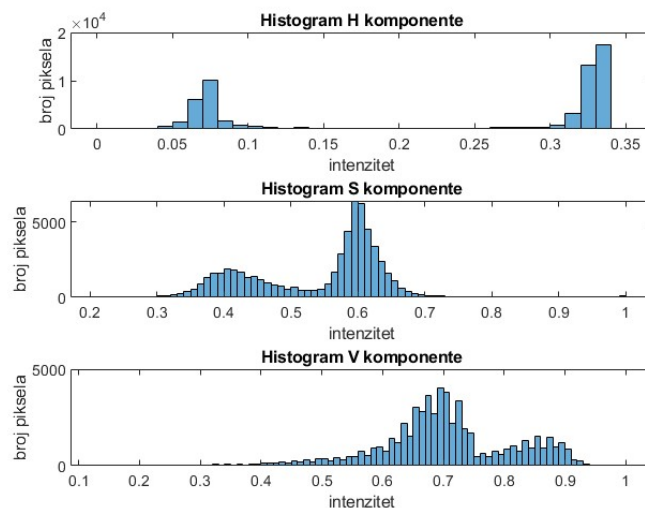


Figure 8: Histogram HSV komponenti papira

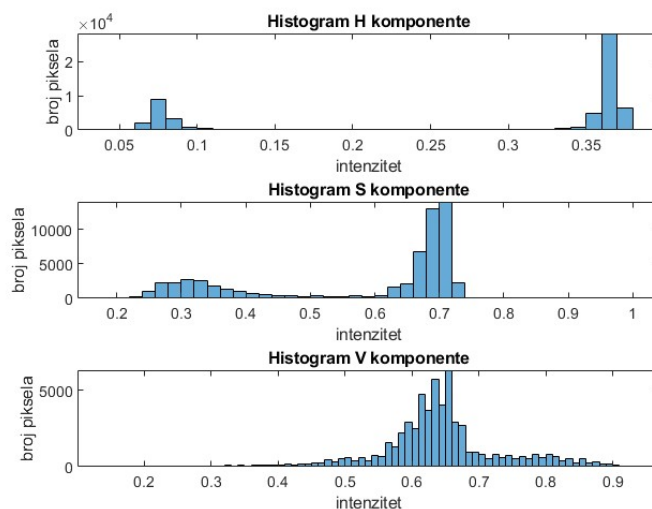


Figure 9: Histogram HSV komponenti makaza

Na histogramima možemo videti da je za sve tri hlase intenzitet H komponente najseparabilniji pa će na osnovu nje biti odrađena binarizacija. Uzet je prag 0.2 i na osnovu njega dobijene binarizovane slike na koje su primenjene morfološke operacije, diletacija i erozija, kao i inverzija kako bi pozadina bila crna i konačne binarizovane slike su sledeće:

Binarizovana slika nakon diletacije, erozije i inverzije



Figure 10: Binarizovani kamen

Binarizovana slika nakon diletacije, erozije i inverzije



Figure 11: Binarizovani papir

Binarizovana slika nakon diletacije, erozije i inverzije



Figure 12: Binarizovane makaze

Sada su binarizovane slike pogodne za izdvajanje obeležja koja bi trebalo da što više separatišu klase. Uzeta obeležja su dužina od prvog levog belog piksela do maksimalne širine ruke, kao i dužina ivica tog istog dela slike:

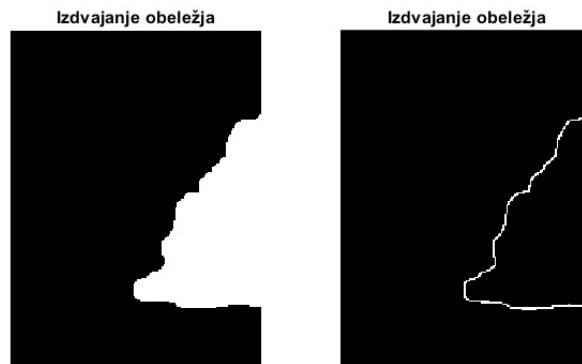


Figure 13: Obeležja kamena

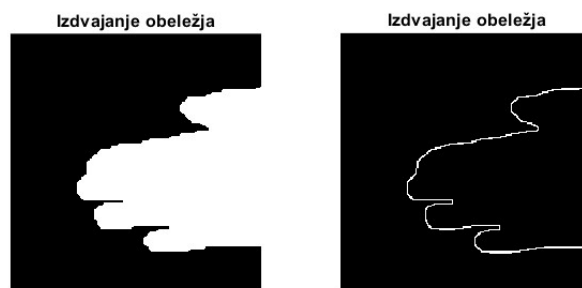


Figure 14: Obeležja papira



Figure 15: Obeležja makaza

Odbirci su podeljeni na trening i test skup, gde je trening skup prvih 600 odbiraka svake klase, a test skup je ostatak. Konačno dobijena obeležja sve tri klase trening skupa su sledeća:

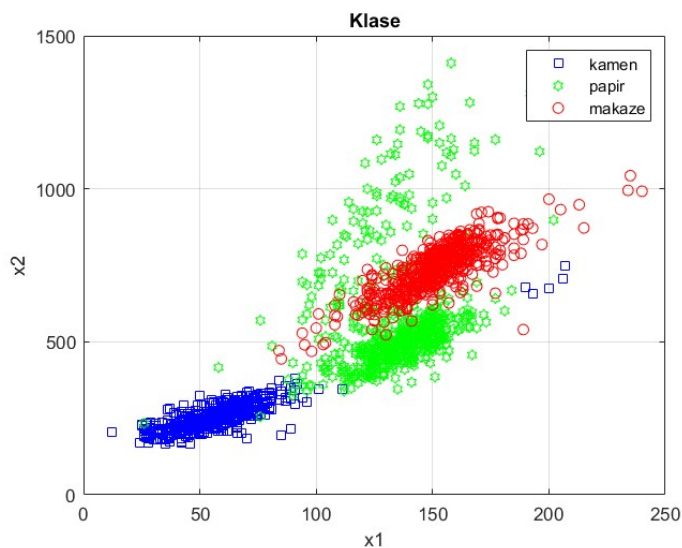


Figure 16: Obeležja klasa

Sada je potrebno projektovati klasifikator. Za klasifikaciju sve tri klase je projektovan Bayes-ov klasifikator, to jest proveravano je koja bi funkcija gustine verovatnoće bila najveća, za svaki odbirak, a zatim je odbirak dodeljivan toj klasi. Izgled klasifikacije, kao i konfuziona matrica nad test skupom izgledaju na sledeći način:

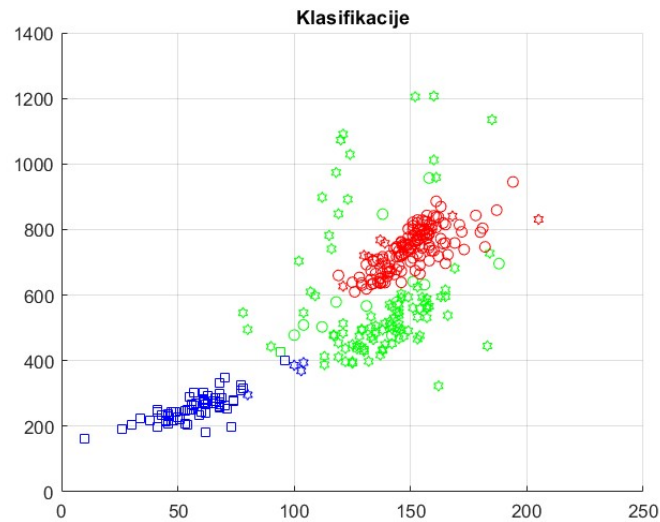


Figure 17: Klasifikacija test skupa

True Class	1	68	1	
	2	4	98	10
	3		12	138
		1	2	3
		Predicted Class		

Figure 18: Konfuziona matrica sve tri klase

Tačnost ovakve procene je 91.84%, što je prihvatljiva složenost. Međutim da se primetiti da je papir najmanje separabilan od ostale dve klase i da bi klasifikacija samo makaza i kamena bila dosta jednostavnija. Stoga se na dalje vrši upravo ta klasifikacija.

Histogram obeležja makaza i kamena je sledeći:

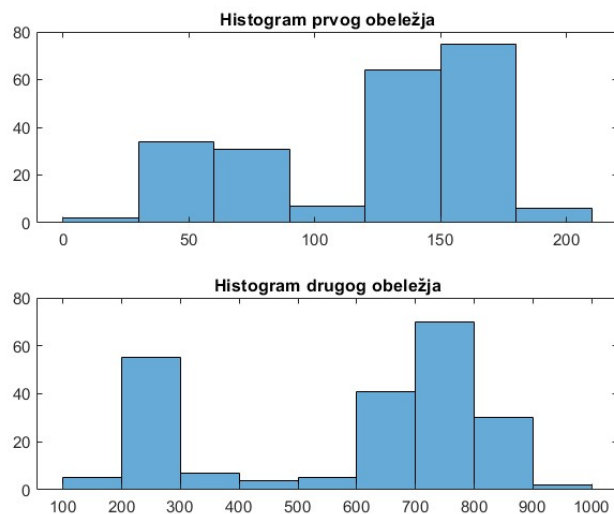


Figure 19: Histogram obeležja makaza i kamena

Vidimo da su obeležja poprilično separabilna i za njihovu klasifikaciju je korišćen linearni klasifikator na bazi željenog izlaza, kod koga je data malo veća težina makazama, radi bolje klasifikacije. Diskriminaciona kriva, kao i konfuzionna matrica su dati na sledećim graficima:

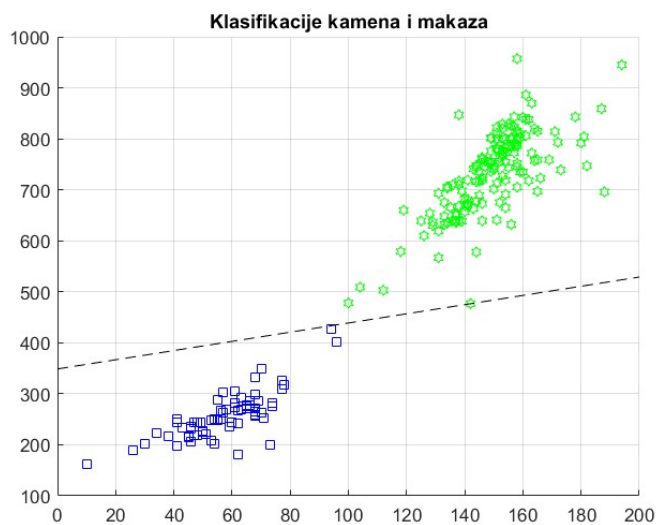


Figure 20: Diskriminaciona kriva

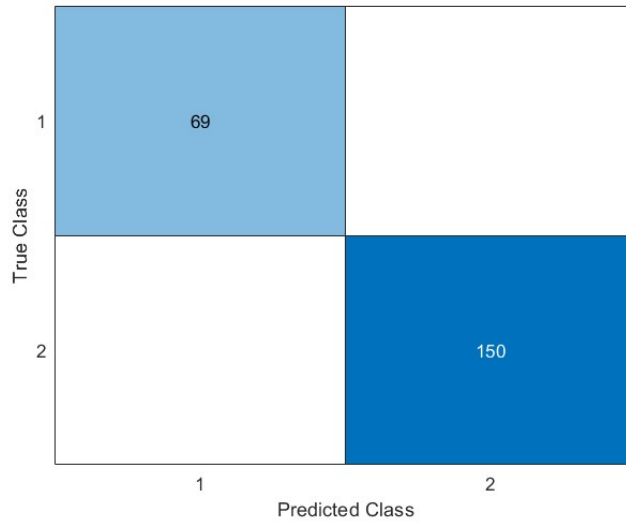


Figure 21: Konfuzionna matrica

Vidimo da je tačnost klasifikacije 100%.

2 Zadatak

Generisati po $N = 500$ odbiraka iz dveju dvodimenzionih bimodalnih klasa:

$$\Omega_1 \sim P_{11}N(M_{11}, \Sigma_{11}) + P_{12}N(M_{12}, \Sigma_{12}) \quad (1)$$

$$\Omega_2 \sim P_{21}N(M_{21}, \Sigma_{21}) + P_{22}N(M_{22}, \Sigma_{22}) \quad (2)$$

Parametre klasa samostalno izabrati.

a) Na dijagramu prikazati odbirke.

b) Is crtati kako teorijski izgledaju funkcije gustine verovatnoće za raspodele klasa i uporediti ih sa histogramom generisanih odbiraka.

c) Projektovati Bajesov klasifikator minimalne greške i na dijagramu, zajedno sa odbircima, skicirati klasi fikacionu liniju. Uporediti grešku klasifikacije konkretnih odbiraka sa teorijskom greškom klasifikacije prve i druge vrste za datu postavku.

d) Projektovati klasifikator minimalne cene tako da se više penalizuje pogrešna klasifikacija odbiraka iz prve klase.

e) Ponoviti prethodnu tačku za Neuman-Pearson-ov klasifikator. Obrazložiti izbor $\epsilon_1 = \epsilon_0$.

f) Za klase oblika generisanih u prethodnim tačkama, projektovati Wald-ov sekvencijalni test pa skicirati zavisnost broja potrebnih odbiraka od usvojene verovatnoće grešaka prvog, odnosno drugog tipa.

Za samostalno izabrane parametre, dobijene klase su sledeće:

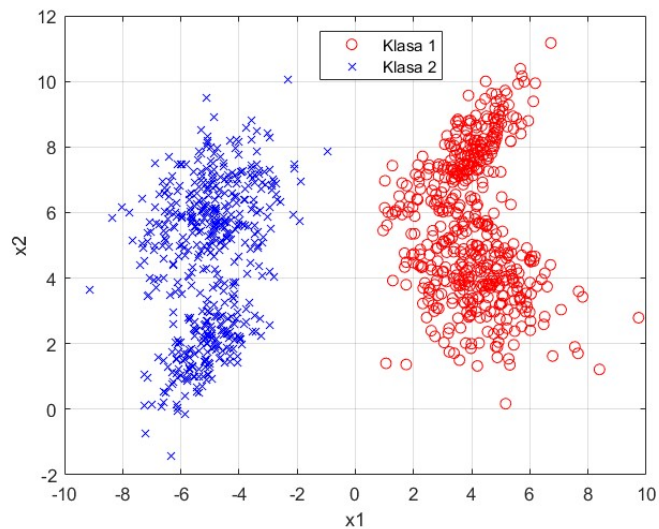


Figure 22: Odbirci klasa

Konture teorijskih funkcija gustina verovatnoća izgledaju na sledeći način: Nakon dobijenih odbiraka

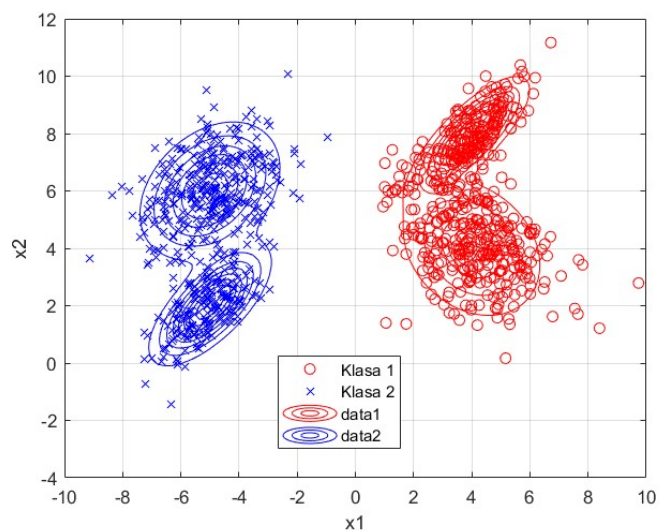


Figure 23: Odbirci klasa sa konturama teorijske funkcije gustine verovatnoće

nad njima je za početak potrebno primeniti Bayesov klasifikator, kao i odrediti teorijske i eksperimentalne verovatnoće greške prvog i drugog tipa. Dobijeni Bayesov klasifikator izgleda na sledeći način:

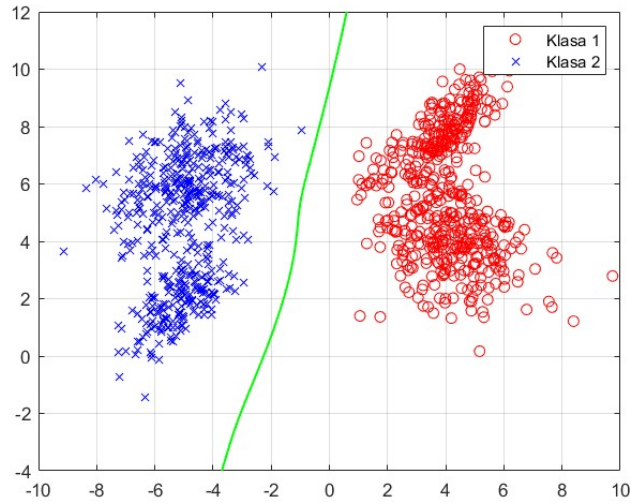


Figure 24: Bayesov klasifikator

Dobijene eksperimentalne greške prvog i drugog tipa su obe 0, dok su teorijske $e_1 = 1.6656e - 04$ i $e_2 = 1.1668e - 04$.

Nakon toga je potrebno primeniti test minimalne cene takav da se više penalizuje kada je odbirak iz prve klase klasifikovan kao odbirak iz druge klase pa c_{21} treba da bude veće od c_{12} , i očekivano bi bilo da se diskriminaciona kriva pomeri ka klasi 2:

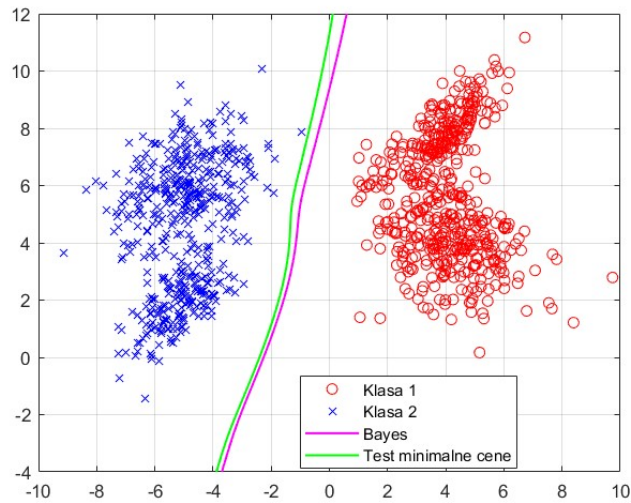


Figure 25: Klasifikator minimalne cene

Nakon toga je primenjen Neuman-Pearson-ov test koji se zasniva na tome da jednu od verovatnoća greške fiksiram, dok druga pokušava da se minimizuje, pa se uglavnom fiksira manje rizična greška, i to na takav način da i ona bude što manja. Kako se klasifikacija vrši na osnovu formule:

$$-ln\left(\frac{f_1(x)}{f_2(x)}\right) \begin{cases} > -ln(\mu) & x \in \omega_2 \\ < -ln(\mu) & x \in \omega_1 \end{cases} \quad (3)$$

a i za $\epsilon_2 = \epsilon_0$ važi:

$$\epsilon_0 = \int_{-\infty}^{-ln(\mu)} f(h|\omega_2)dh \quad (4)$$

Odatle za željeno ϵ_0 možemo naći μ na osnovu sledećeg grafika:

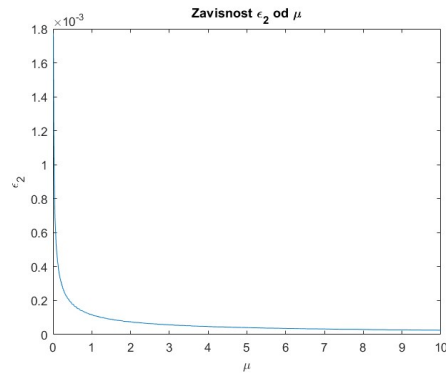


Figure 26: Grafik zavisnosti ϵ_0 od μ

Na osnovu grafika je uzeto $\epsilon_0 = 0.00003$ i na osnovu toga dobijeno odgovarajuće μ . Konačna diskriminaciona kriva dobijena na ovaj način je sledeća:

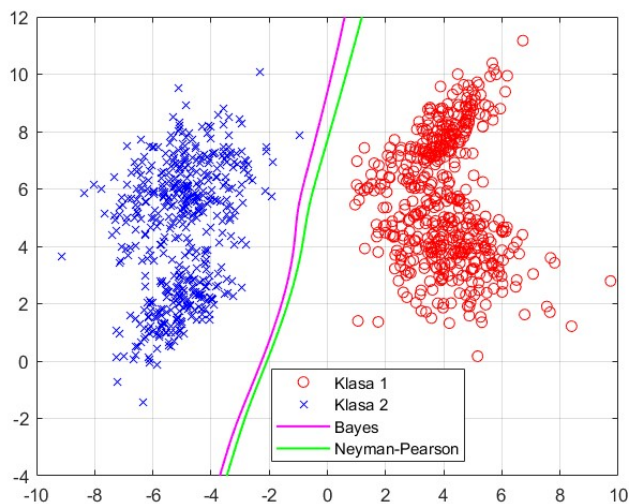


Figure 27: Neuman-Pearson-ov test

Konačno je ideja da se primeni Wald-ov sekvencijalni test. Ideja sekvencijalnog testa je da se odluka ne donese odmah već nakon nekoliko odbiraka, to jest nakon minimalno odbiraka za željene verovatnoće greške prvog i drugog tipa. Grafik na kome je prikazan proces odabiranja je sledeći:

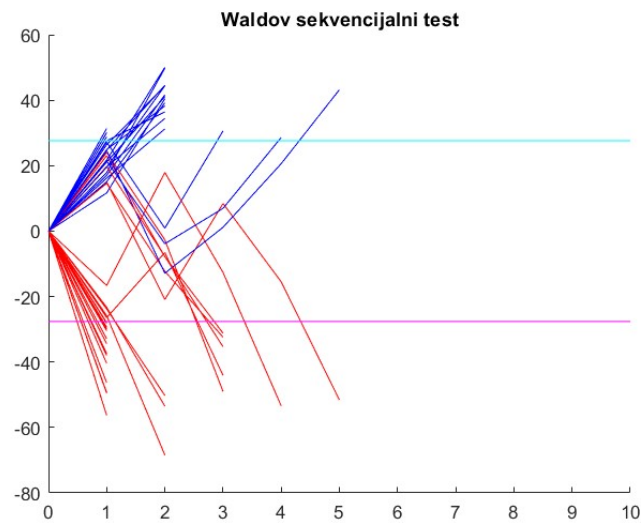


Figure 28: Wald-ov test

Zavisnosti broja iteracija od verovatnoća grešaka prvog i drugog tipa su sledeće:

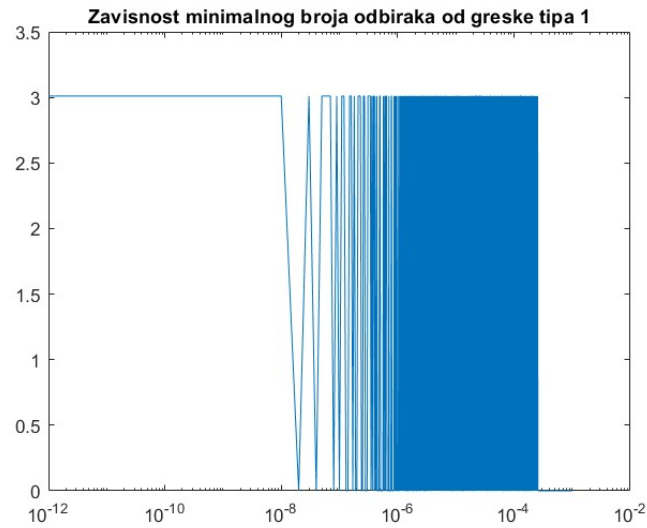


Figure 29: Zavisnost broja iteracija od verovatnoće greške prvog tipa

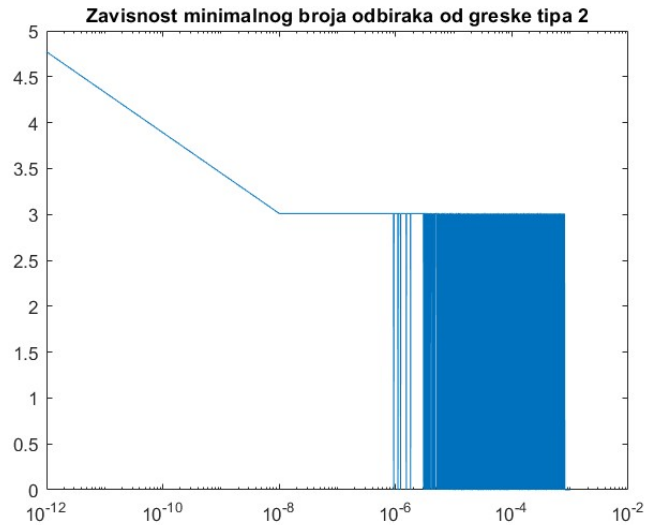


Figure 30: Zavisnost broja iteracija od verovatnoće greške drugog tipa

Očekivano je da u logaritamskoj razmeri jedan od grafika ima linearnu dok drugi ima konstantnu zavisnost, međutim za veće vrednosti grešaka gotovo da se stalno već posle prvog ili drugog odbirka odredi klasa, i zbog toga na kraju oba grafika postoji oscilovanje.

3 Zadatak-Opcija 1

1. Generisati tri klase dvodimenzionalnih oblika. Izabrati funkciju gustine verovatnoće oblika tako da klase budu linearno separabilne.
 - a) Za tako generisane oblike izvršiti projektovanje linearnog klasifikatora jednom od tri iterativne procedure. Rezultate prikazati u obliku matrice konfuzije. Detaljno opisati postupak klasifikacije.
 - b) Ponoviti prethodni postupak korišćenjem metode željenog izlaza. Analizirati uticaj elemenata u matrici željenih izlaza na konačnu formu linearnog klasifikatora.
2. Generisati dve klase dvodimenzionalnih oblika koje jesu separabilne, ali ne linearno, pa isprojektovati kvadratni klasifikator metodom po želji.

Odbirci tri generisane linearno separabilne klase su sledeći:

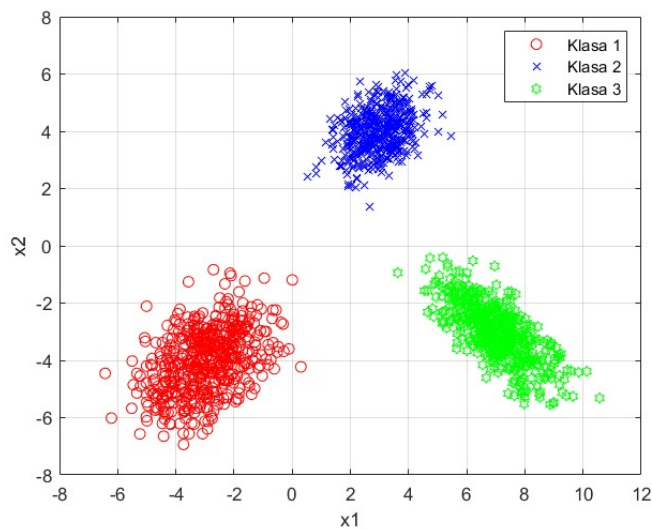


Figure 31: Odbirci klasa

Za početak je nad njima potrebno primeniti jednu od numeričkih metoda za pronalaženje optimalnog linearnog klasifikatora. Izabrana je metoda resupstitucije. Kako bi se klasifikovale tri klase linearnim klasifikatorima, za početak se dve klase spoje u jednu pa se klasifikuju u odnosu na treću, a tek onda se preostale dve međusobno klasifikuju. Klasifikator nakon klasifikacije samo jedne klase izgleda na sledeći način:

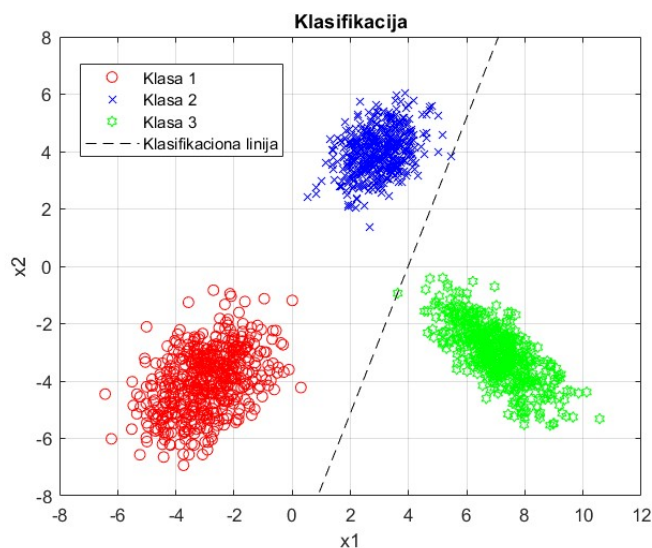


Figure 32: Linearni klasifikator

Nakon konačne klasifikacije izgled diskriminacionih krivih i konfuziona matrica su sledeći:

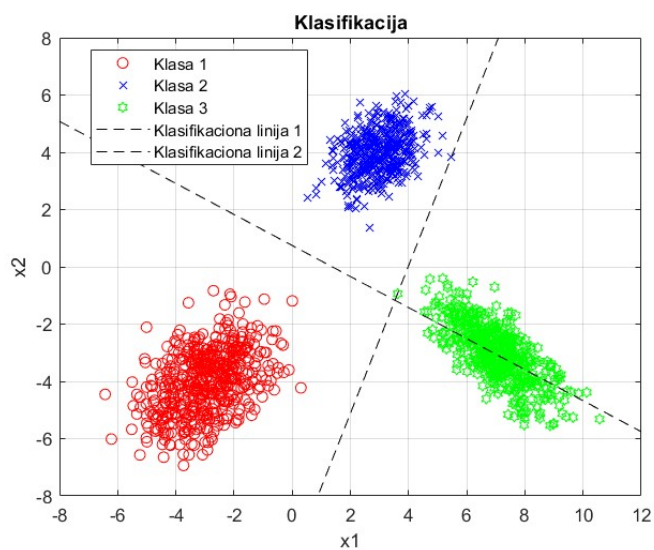


Figure 33: Linearni klasifikator

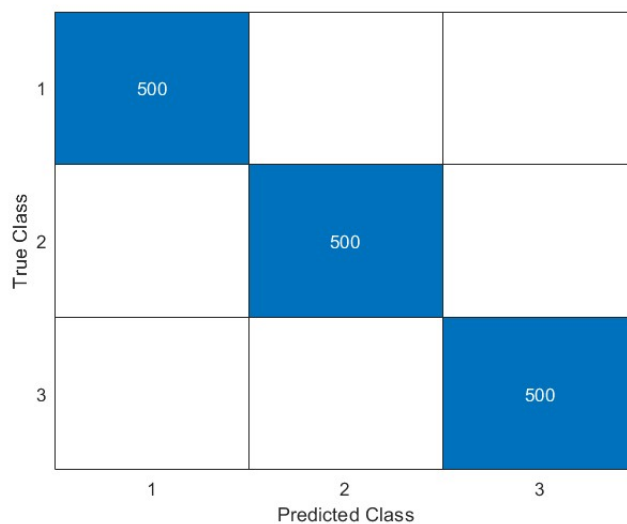


Figure 34: Linearni klasifikator-konfuziona matrica

Vidimo da je tačnost ovakve klasifikacije 100%.

Sada je potrebno isti ovaj postupak ponoviti linearnim klasifikatorom na bazi željenog izlaza. Klasifikacija je ponovo rađena u dva koraka na isti način, i konačne diskriminacione krive kao i konfuziona matrica su sledeći:

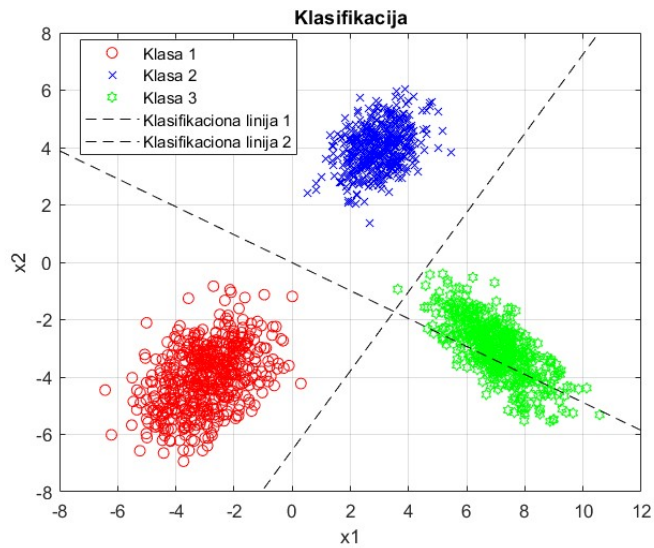


Figure 35: Linearni klasifikator na bazi željenog izlaza

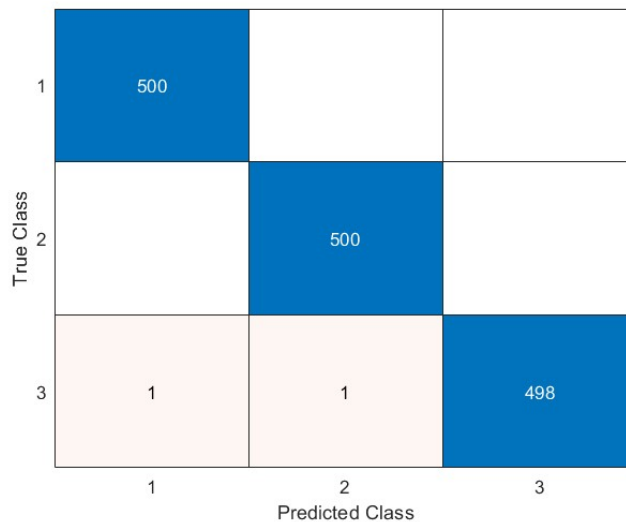


Figure 36: Linearni klasifikator na bazi željenog izlaza-konfuziona matrica

Vidimo da su u ovom slučaju dva odbirka klase 3 pogrešno klasifikovana, što ne narušava previše tačnost, ali bi moglo da se popravi ako bi se klasi 3 dala nešto veća težina, to jest povećalo γ uz klasu 3. To je u sledećem delu i urađeno i dobijeni rezultati su sledeći:

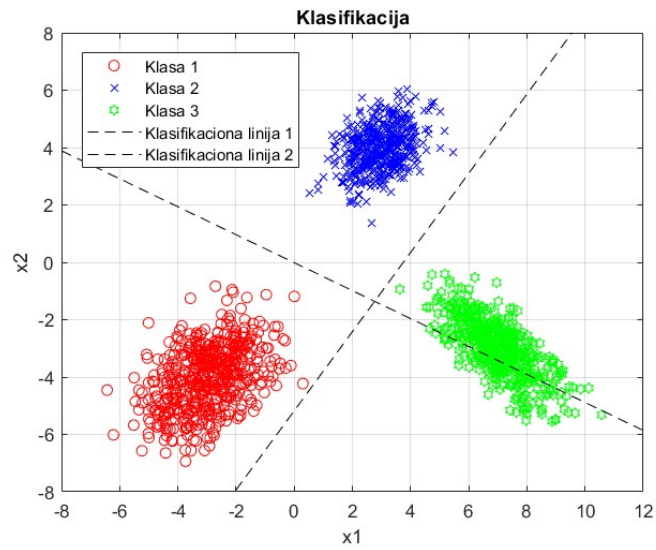


Figure 37: Linearni klasifikator na bazi željenog izlaza

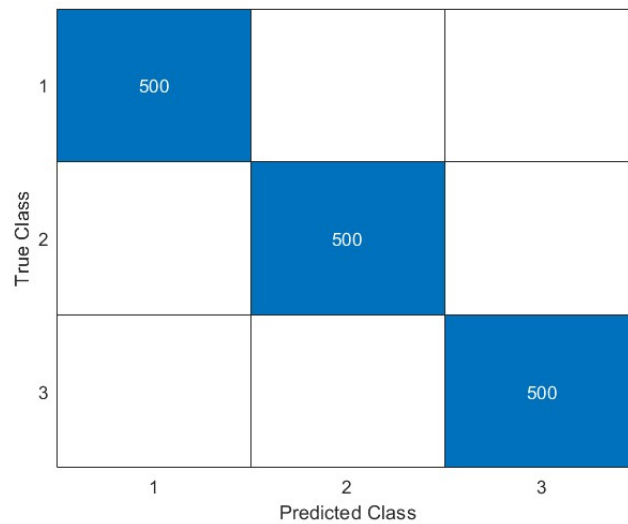


Figure 38: Linearni klasifikator na bazi željenog izlaza-konfuziona matrica

Uz ovakvu popravku vidimo da je tačnost 100%.

Konačno je potrebno napraviti nove dve klase koje nisu linearno separabilne i projektovati kvadratni klasifikator, koji se projektuje tako što se transformiše u linearni, a zatim ponovo radi na bazi željenog izlaza. Generisane klase i diskriminaciona kriva su sledeći:

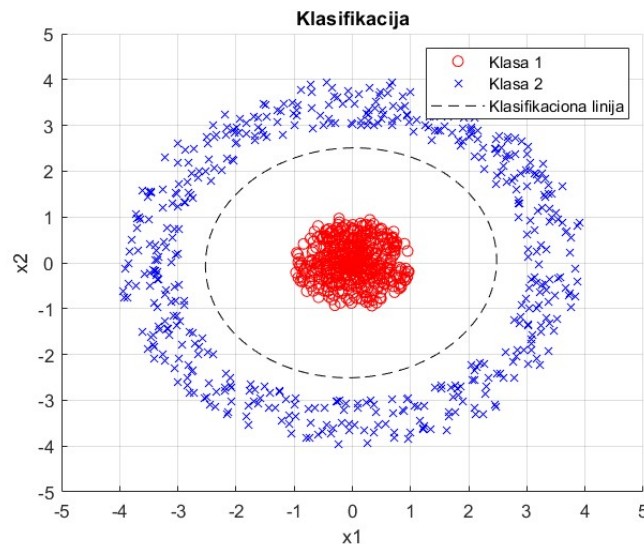


Figure 39: Kvadratni klasifikator

4 Zadatak

1. Generisati po $N = 500$ dvodimenzionih odbiraka iz četiri klase koje će biti linearno separabilne. Preporuka je da to budu Gausovski raspodeljeni dvodimenzioni oblici. Izabrati jednu od metoda za klasterizaciju (cmean metod, metod kvadratne dekompozicije) i primeniti je na formirane uzorke klase. Izvršiti analizu osetljivosti izabranog algoritma na početnu klasterizaciju kao i srednji broj potrebnih iteracija. Takođe izvršiti analize slučaja kada se apriorno ne poznaje broj klase.
2. Na odbircima iz prethodne tačke izabrati jednu od metoda klasterizacije (metod maksimalne verodostojnosti ili metod grana i granica) i primeniti je na formirane uzorke klase. Izvršiti analizu osetljivosti izabranog algoritma na početnu klasterizaciju kao i srednji broj potrebnih iteracija. Takođe izvršiti analize slučaja kada se apriorno nepoznaje broj klase.
3. Generisati po $N = 500$ dvodimenzionih odbiraka iz dve klase koje su nelinearno separabilne. Izabrati jednu od metoda za klasterizaciju koje su primenjive za nelinearno separabilne klase (metod kvadratne dekompozicije ili metod maksimalne verodostojnosti) i ponoviti analizu iz prethodnih tačaka.

Početne klase izgledaju na sledeći način:

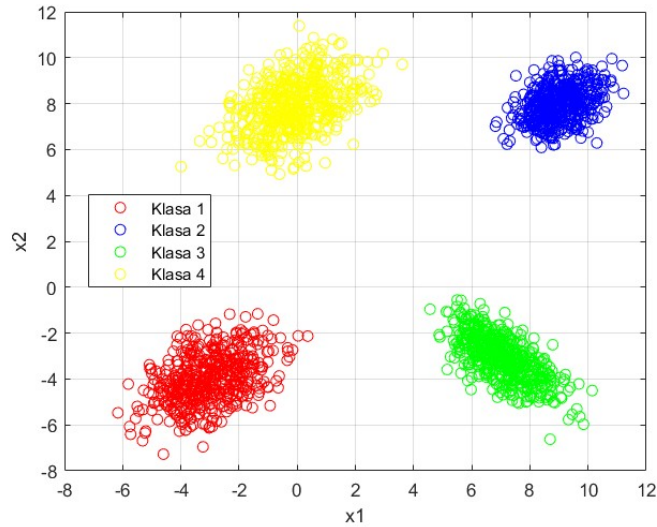


Figure 40: Početne klase

Nad njima je za početak primenjen C-mean metod klasterizacije koji se sastoji od nekoliko koraka:

1. Napravi se početna klasterizacija slučajno i odrede se početne verovatnoće klastera, matematička očekivanja i kovarijacione matrice.
2. Za svaki odbirak se odredi Euklidsko rastojanje do svakog klastera i prebaci se u najbliži klaster.
3. Ponavlja se proces sve dok i dalje ima promena ili do maksimalnog broja iteracija.

Za početak je pokušano rešenje zadatka sa slučajnom početnom klasterizacijom, i na taj način je često u prvom koraku izbacivan jedan ili dva klastera i zato se početna klasterizacija ne radi na potpuno slučajan način i prvih 20 odbiraka je smešteno u svoje predodređene klase. Početna klasterizacija izgleda na sledeći način:

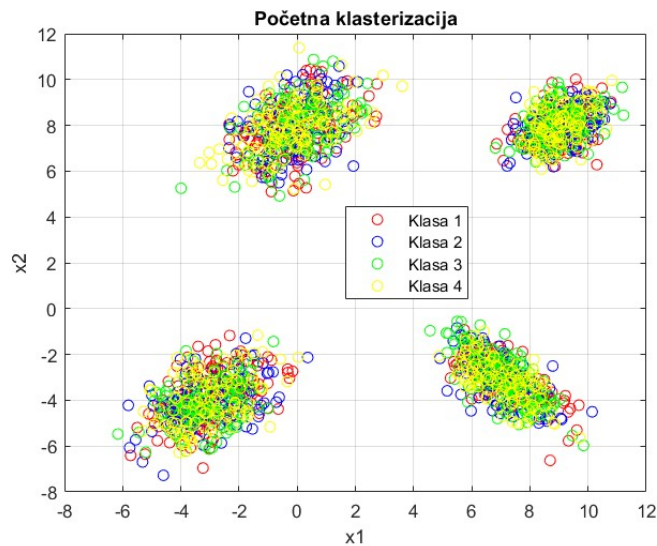


Figure 41: Početna klasterizacija

Prvi korak jedne klasterizacije bi izgledao na sledeći način:

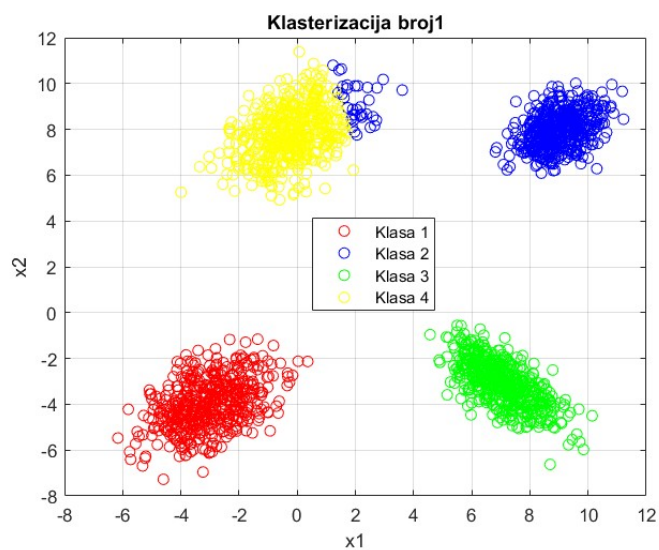


Figure 42: Prva klasterizacija

Vidimo da je većina odbiraka smeštena u odgovarajući klaster i već nakon sledeće iteracije će svi biti smešteni:

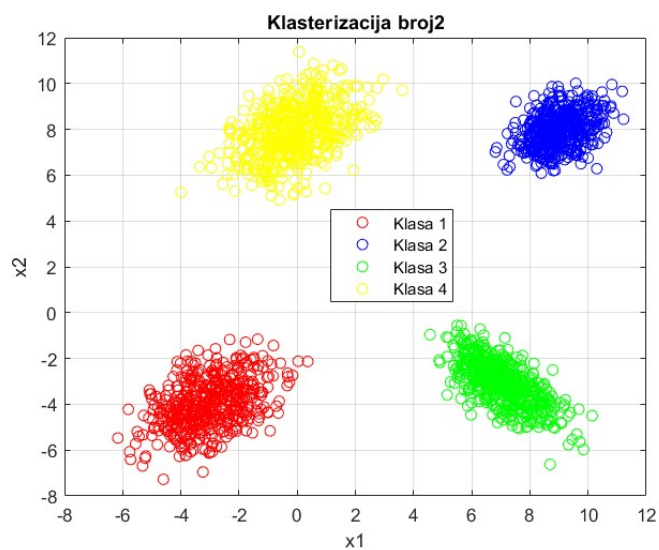


Figure 43: Druga klasterizacija

Nakon ovoga bi u sledećoj iteraciji bio završen čitav proces jer ne bi bilo premeštanja iz klastera u druge klastere. Prosečan broj iteracija u ovom slučaju je 3.9.

Sa druge strane, u prethodnom slučaju je bilo pretpostavljeno da je poznat broj klastera, međutim ni to ne mora uvek biti tačno. Ukoliko bismo pretpostavili da postoji samo dva klastera početna klasterizacija bi izgledala na sledeći način:

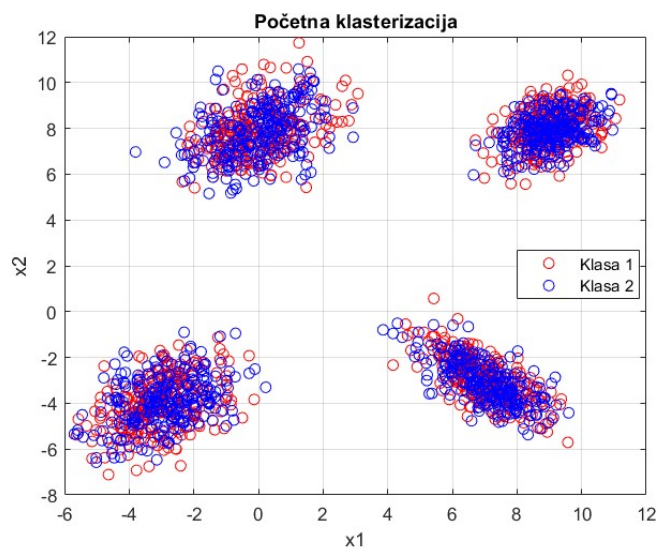


Figure 44: Početna klasterizacija

Već u prvoj iteraciji bi klasteri bili potpuno odvojeni:

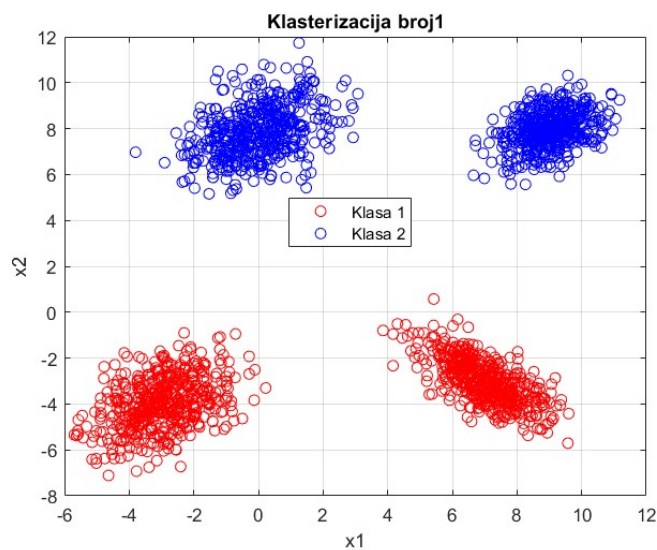


Figure 45: Prva klasterizacija

U ovom slučaju je prosečan broj iteracija 2.

U sledećem koraku je potrebno generisati nove dve klase koje nisu linearno separabilne i nad njima primeniti kvadratnu dekompoziciju koja se razlikuje od C-mean klasterizacije isključivo po tome što ne traži Euklidsko rastojanje od klastera već uzima u obzir i kovarijacionu matricu klastera. Početni klasteri bi izgledali na sledeći način:

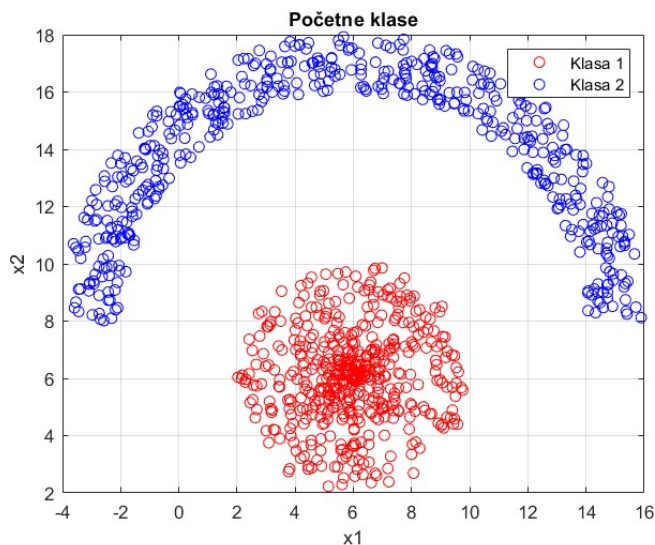


Figure 46: Početni klasteri

I u ovom slučaju početna klasterizacija nije potpuno slučajno urađena, već je određeni broj odbiraka na početku smešten u svoje klaster.

Početna klasterizacija izgleda na sledeći način:

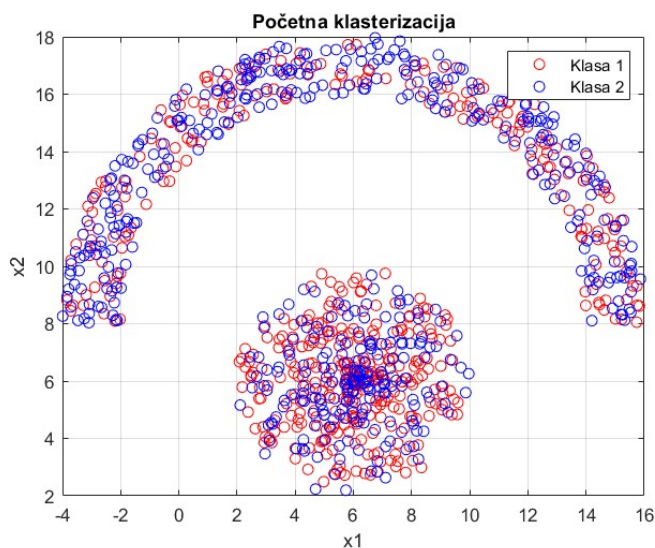


Figure 47: Početna klasterizacija

Prva tri koraka klasterizacije izgledaju na sledeći način:

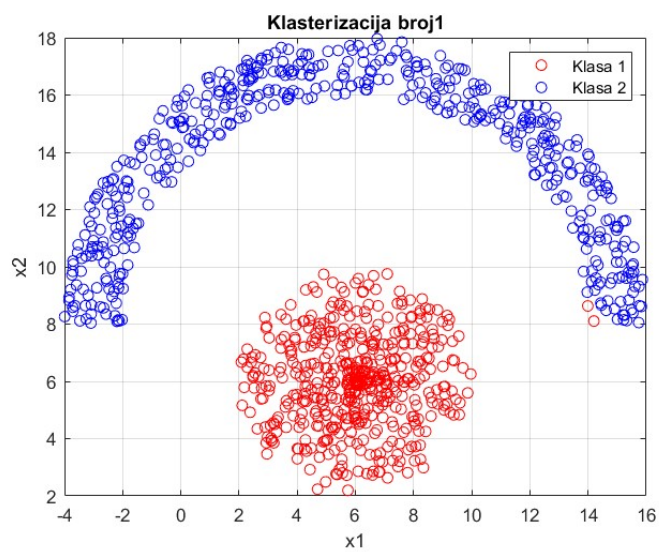


Figure 48: Prva klasterizacija

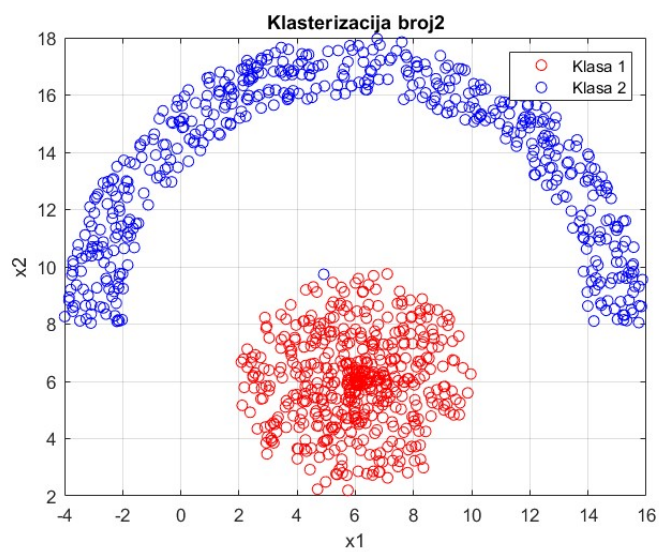


Figure 49: Druga klasterizacija

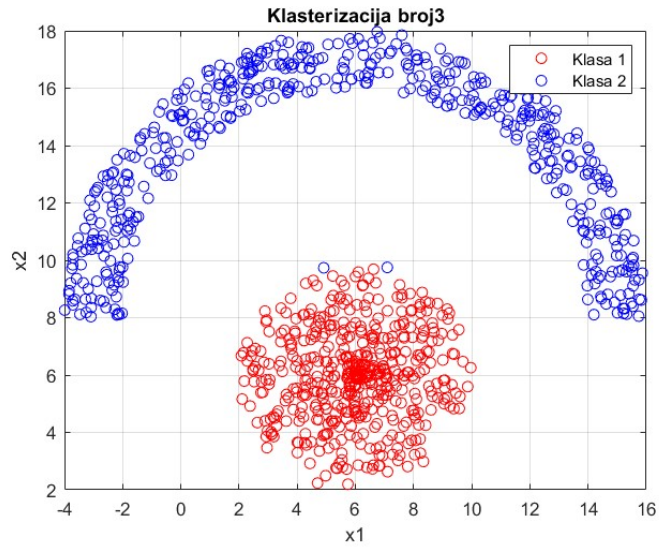


Figure 50: Treća klasterizacija

U ovom slučaju vidimo da zapravo i ne bude ispunjena potpuna klasterizacija i da postoje dva odbirka koja su pogrešno klasifikovana. Prosečan broj iteracija u ovom slučaju je 4.2.

Konačno i u ovom slučaju je moguće da nam je početna pretpostavka o broju klastera bila pogrešna. Ukoliko bismo pretpostavili da postoji četiri klastera, početna klasterizacija bi izgledala na sledeći način:

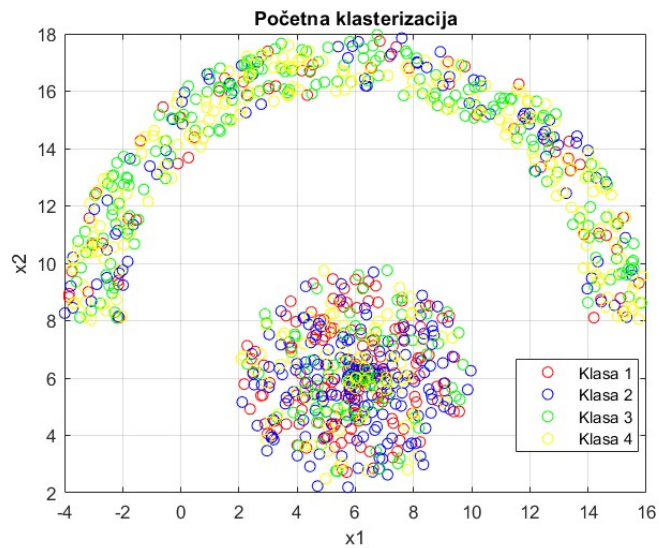


Figure 51: Početna klasterizacija

U ovom slučaju se svakako ne može naći dovoljno dobra klasterizacija i prve 3 iteracije će izgledati na sledeći način:

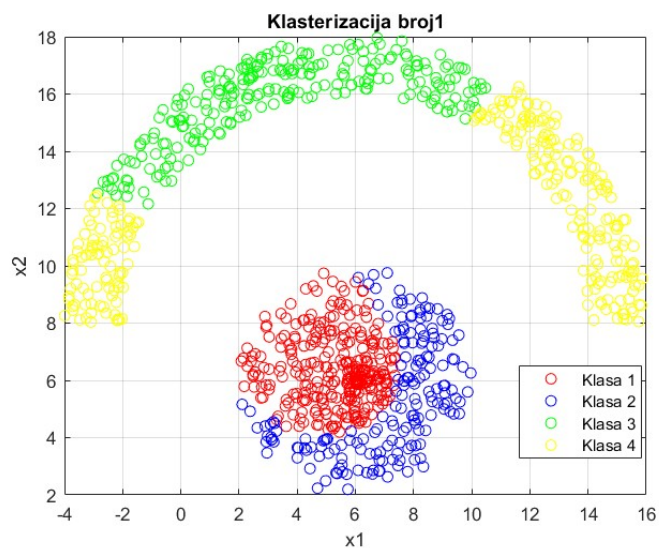


Figure 52: Prva klasterizacija

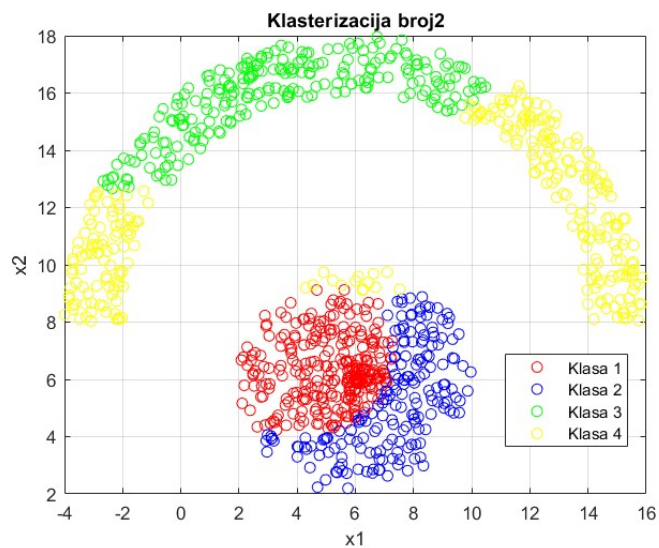


Figure 53: Druga klasterizacija

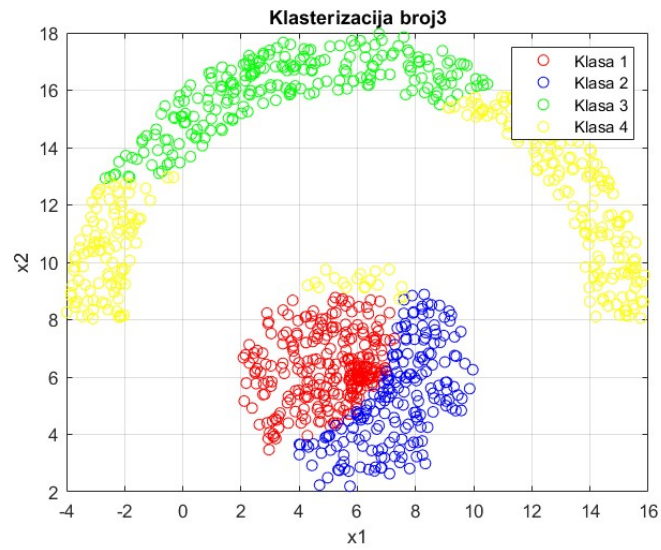


Figure 54: Treća klasterizacija

Ovaj proces bi se nastavio ili do maksimalnog broja iteracija ili dok se ne bi sasvim slučajno desilo da ne dođe ni do jedne reklasterizacije. Prosečan broj iteracija u ovom slučaju je 12.9.