

ЕЛЕКТРОТЕХНИЧКИ ФАКУЛТЕТ

ИЗВЕШТАЈ ДОМАЋЕГ

ИЗ ПРЕПОЗНАВАЊА ОБЛИКА

Сви домаћи

Аутор

Ђорђе Живановић, 3013/2017

Ментор

Професор Др Жељко Ђуровић

20. фебруар 2018.

Садржај

1	Први домаћи задатак	2
1.1	Опис система	2
1.2	Избор обележја	6
1.3	Пример класификација	7
2	Други домаћи задатак	8
2.1	Генерисање одбирака	8
2.2	Функција густине вероватноће	10
2.3	Бајесов класификатор минималне грешке	11
2.4	Бајесово правило одлучивања минималне цене	13
2.5	Волдов секвенцијални тест	14
3	Трећи домаћи задатак	18
3.1	Линеарни класификатор	18
3.1.1	Генерисање класа	18
3.1.2	Пројектовање линеарног класификатора методом ресупституције	18
3.1.3	Метод жељеног излаза	21
3.2	Квадратни класификатор	24
4	Четврти домаћи задатак	25
4.1	C-Mean алгоритам	25
4.1.1	Генерисање одбирака	25
4.1.2	Алгоритам	26
4.1.3	Анализа резултата алгоритма	30
4.2	Метод квадратне декомпозиције	31
4.2.1	Генерисање одбирака	31
4.2.2	Алгоритам	32
4.2.3	Анализа резултата	35

1 Први домаћи задатак

За руком писаних цифара 6, 8 и 9 која је доступна на сајту предмета испројектовати иновативни систем за препознавање цифара заснован на тестирању хипотеза. Резултате приказати у облику матрице конфузија. Извештај треба да садржи кратки опис испројектованог система, образложени избор обележја као и карактеристичне примере правилно и неправилно класификованих цифара

1.1 Опис система

На располагању налазило се 360 руком писаних цифара које су достављене у .bmp формату. Било је по 120 цифара 6, 8 и 9. Задатак је одрађен у програмском пакету *Matlab vR2017b*. Сlike су обрађене у неколико фаза. У првој фази одрађено је бинаризовање слике, тј. вредност свим пикселима је преведена да буде 255 или 0, с обзиром да је слика представљена као низ црних и белих пиксела различитих јачина. Бинаризација је одрађена функцијом **1**.

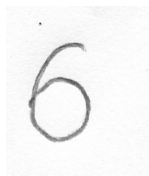
Одсечак кода 1: Бинаризација слике

```
1 function [binImage] = binarizeImage(image, ratio)
2     val = 255 - double(image);
3     val(val < ratio * max(max(val))) = 0;
4     val(val >= ratio * max(max(val))) = 255;
5     binImage = uint8(255 - val);
6 end
```

Функција прима *image* и *ratio* који представљају дату слику, односно праг бељења-колика црнина се прихвата као белина.

Друга фаза омогућава одсецање неискоришћеног простора, односно белина. Поред белина занемарене су све мрље унутар белина. Кад се каже мрља мисли се тачка величине бар 5 пиксела (в. сл. **1a**, **1b**, **1c** као пример).

Слика 1: Тачке



(a) Тачка у 6



(b) Тачка у 8



(c) Тачка у 9

Код за уклањање белина и тачки је одрађен функцијом **2**.

Одсечак кода 2: Одсецање слике

```
1 function [X] = cropImage(image)
2     ratio = 0.20;
3     X = binarizeImage(image, ratio);
4
5     [numRow, numCol] = size(X);
6     begin = 1;
7     while (begin < numRow) && (sum(X(begin, 1 : numCol)) / numCol > 250
8         || ...
9         sum(X(begin + 1, 1 : numCol)) / numCol > 250 || ...
10        sum(X(begin + 2, 1 : numCol)) / numCol > 250 || ...
11        sum(X(begin + 3, 1 : numCol)) / numCol > 250 || ...
12        sum(X(begin + 4, 1 : numCol)) / numCol > 250 || ...
13        sum(X(begin + 5, 1 : numCol)) / numCol > 250)
14         begin = begin + 1;
15
16     endV = numRow;
```

```

17 while(endV > begin)&&...
18     (sum(X(endV, 1 : numCol)) / numCol > 250 ||...
19     sum(X(endV - 1, 1 : numCol)) / numCol > 250 ||...
20     sum(X(endV - 2, 1 : numCol)) / numCol > 250 ||...
21     sum(X(endV - 3, 1 : numCol)) / numCol > 250 ||...
22     sum(X(endV - 4, 1 : numCol)) / numCol > 250 ||...
23     sum(X(endV - 5, 1 : numCol)) / numCol > 250)
24     endV = endV - 1;
25 end
26
27 left = 1;
28 while (left < numCol)&&...
29     (sum(X(1 : numRows, left)) / numRows > 250 ||...
30     sum(X(1 : numRows, left + 1)) / numRows > 250 ||...
31     sum(X(1 : numRows, left + 2)) / numRows > 250 ||...
32     sum(X(1 : numRows, left + 3)) / numRows > 250 ||...
33     sum(X(1 : numRows, left + 4)) / numRows > 250 ||...
34     sum(X(1 : numRows, left + 5)) / numRows > 250)
35     left = left + 1;
36 end
37
38 right = numCol;
39 while (right > 1)&&...
40     (sum(X(1 : numRows, right)) / numRows > 250 ||...
41     sum(X(1 : numRows, right - 1)) / numRows > 250 ||...
42     sum(X(1 : numRows, right - 2)) / numRows > 250 ||...
43     sum(X(1 : numRows, right - 3)) / numRows > 250 ||...
44     sum(X(1 : numRows, right - 4)) / numRows > 250 ||...
45     sum(X(1 : numRows, right - 5)) / numRows > 250)
46     right = right - 1;
47 end
48
49 X = X(begin : endV, left : right);
50 end

```

Овде је коришћен $ratio = 0,2$ јер даје најбољу тачност за одговарајући скуп података. Међутим и поред свега дође до лошег одсецања (в. сл. 2a, 2b).

Слика 2: Лоше одсецање



(a) Лоше одсецање 6



(b) Лоше одсецање 8

У осталим исправним случајевима добију се бројеви као на сл. 3a, 3b, 3c.

Крајња, трећа фаза рачуна "центар масе" цифара добијених из друге фазе и формира функцију густине расподеле вероватноће за сваку од центара масе цифара 6, 8, 9. Центар маса се рачуна функцијом 3. Функција помера координатни почетак слике тако да буде у центру слике и онда сабира све вредности пиксела померене са одговарајућим вектором положаја. На крају се добијена вредност подели са бројем пиксела на слици и то је центар масе дате слике.

Слика 3: Добро одсечене



(a) Добро одсечена цифра 6 (b) Добро одсечена цифра 8 (c) Добро одсечена цифра 9

Одсечак кода 3: Центар масе

```
1 function [center] = calculateCenter(X)
2     [numRow, numCol] = size(X);
3     centerX = 0;
4     centerY = 0;
5     for row = 1 : numRow
6         for col = 1 : numCol
7             val = floor(double(X(row, col)));
8             centerX = centerX + val * (row - floor(numRow / 2));
9             centerY = centerY + val * (col - floor(numCol / 2));
10        end
11    end
12    center = [centerX; centerY] / (numRow * numCol);
13 end
```

Обрада цифара за све 3 фазе одрађена је кодом 4. Овде се може приметити да је узето 300 од 360 слика за тренирање, а 60 за проверу, што је 83% – 17% однос у подели скупа на обучавајући и тренирајући скуп.

Одсечак кода 4: Обрада цифара

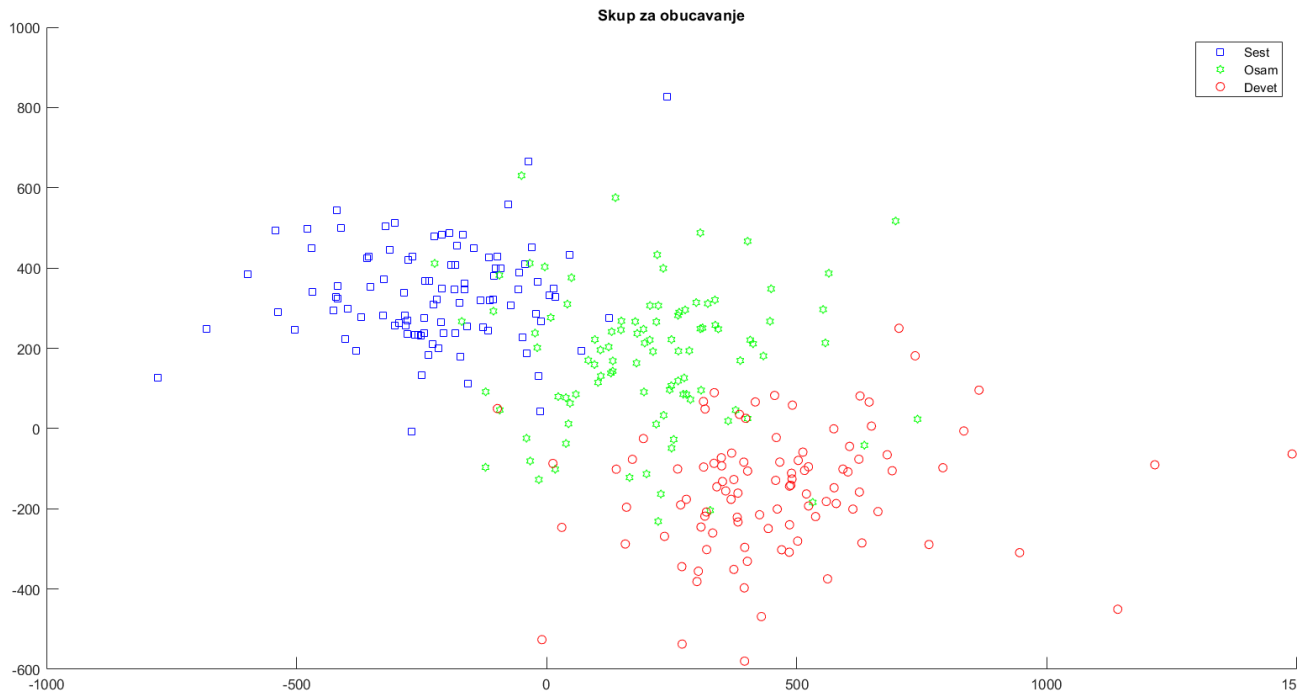
```
1 N = 100;
2 for i = 1 : N
3     name = ['BazaCifara6/baza6' num2str(i, '%03d')];
4     X = imread(name, 'bmp');
5     nameCropped = ['BazaCifaraCropped6/baza6' num2str(i, '%03d') '.bmp'];
6     croppedImage = cropImage(X);
7     imwrite(croppedImage, nameCropped, 'bmp');
8     features6(:, i) = calculateCenter(croppedImage);
9
10    name = ['BazaCifara8/baza8' num2str(i, '%03d')];
11    X = imread(name, 'bmp');
12    nameCropped = ['BazaCifaraCropped8/baza8' num2str(i, '%03d') '.bmp'];
13    croppedImage = cropImage(X);
14    imwrite(croppedImage, nameCropped, 'bmp');
15    features8(:, i) = calculateCenter(croppedImage);
16
17    name = ['BazaCifara9/baza9' num2str(i, '%03d')];
18    X = imread(name, 'bmp');
19    nameCropped = ['BazaCifaraCropped9/baza9' num2str(i, '%03d') '.bmp'];
20    croppedImage = cropImage(X);
21    imwrite(croppedImage, nameCropped, 'bmp');
22    features9(:, i) = calculateCenter(croppedImage);
23 end
24
25 figure;
```

```

26 hold on
27 plot(features6(1, :), features6(2, :), 'sb');
28 plot(features8(1, :), features8(2, :), 'hg');
29 plot(features9(1, :), features9(2, :), 'ro');
30 legend('Sest', 'Osam', 'Devet');
31 title('Skup za obucavanje');

```

Приказ расподеле скупа тренирајућих цифара је приказано на слици 4.



Слика 4: Обучавајући скуп

Функција густине вероватноће за сваку цифру се рачуна тако што се израчунају математичка очекивања и варијансе за сваки од 3 скупа цифара (в. одсечак 5).

Одсечак кода 5: Варијансе и математичка очекивања

```

1 M1 = mean(features6')';
2 S1 = cov(features6');
3 M2 = mean(features8')';
4 S2 = cov(features8');
5 M3 = mean(features9')';
6 S3 = cov(features9');

```

И потом се израчунају 3 функције густине вероватноће за сваку цифру при чему се густина рачуна користећи функцију 6, при чему се проследи одговарајућа математичка очекивања и коваријациона матрица.

Одсечак кода 6: Гаусова расподела

```

1 function [Y] = gaussianMultivariate(X, Mx, Sx)
2 n = size(Sx, 1);
3 Y = (1 / ((2 * pi) ^ (n / 2) * sqrt(det(Sx)))) * exp((X - Mx)' * inv(Sx) *
4 (X - Mx) / -2);
end

```

У зависности која функција густине вероватноће је већа, тој класи припада цифра која се препознаје. Код за одлучивање којој класи припада је дат одсечком 7.

Одсечак кода 7: Одлучивање припадности класи

```

1  for k = 1:3
2      if k == 1
3          T = featuresTest6;
4      elseif k == 2
5          T = featuresTest8;
6      else
7          T = featuresTest9;
8      end
9
10     for i = 1 : size(T, 2)
11         f1 = gaussianMultivariate(T(:, i), M1, S1);
12         f2 = gaussianMultivariate(T(:, i), M2, S2);
13         f3 = gaussianMultivariate(T(:, i), M3, S3);
14         if (f1 > f2) && (f1 > f3)
15             CM(1, k) = CM(1, k) + 1;
16             if (k ~= 1)
17                 fprintf("Predicted: 1 Actual: %d %d\n", k, i + N);
18             end
19         elseif f2 > f1 && f2 > f3
20             CM(2, k) = CM(2, k) + 1;
21             if (k ~= 2)
22                 fprintf("Predicted: 2 Actual: %d %d\n", k, i + N);
23             end
24         elseif f3 > f2 && f3 > f1
25             CM(3, k) = CM(3, k) + 1;
26             if (k ~= 3)
27                 fprintf("Predicted: 3 Actual: %d %d\n", k, i + N);
28             end
29         end
30     end
31 end

```

Добијена матрица конфузије је у табели 1. Можемо да приметимо да је 88,33 прецизност оваквог препознавања. Такође можемо да приметимо да се 8 најчешће погрешно класификује (4 пута, а 6 и 9 збирно 3 пута су погрешно класификоване), што је за очекивати.

Табела 1: Матрица конфузије

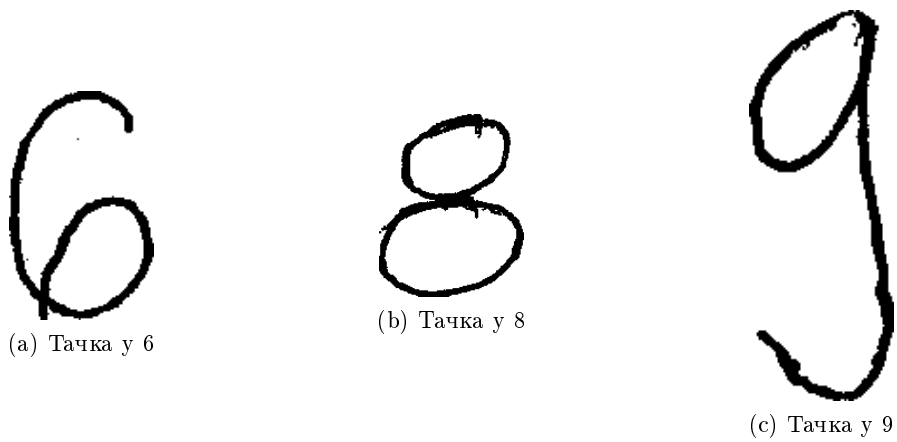
		Стварно		
		6	8	9
Класификовано	6	18	3	0
	8	2	16	1
	9	0	1	19

1.2 Избор обележја

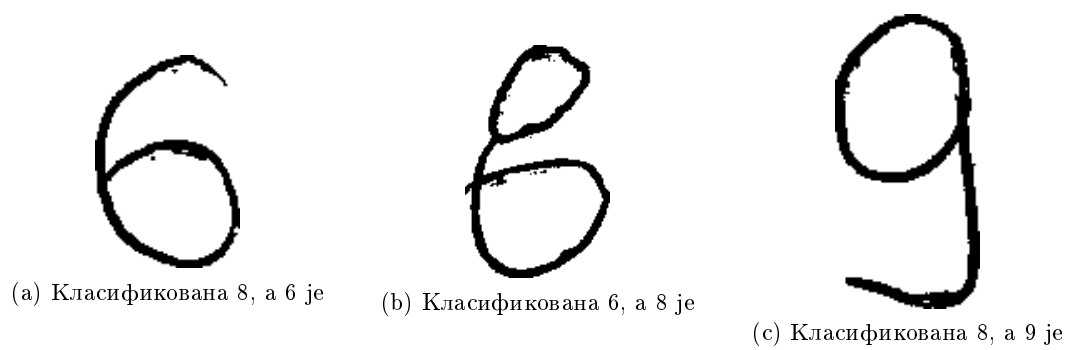
Избор "центра масе" као обележја за које ће се посматрати функција расподеле је разуман јер ако би се 6, 8, 9 и посматрали као објекти, могло би се приметити да је маса 6 пребачена у доњи део, док је за 8 у средњем, а за 9 у горњем, што се може приметити и на слици 4 јер осмице су груписане око координатног центра, деветке су груписане у горњем делу (негативно за y због тога што за слике је координатни систем обрнут код y) и шестице су груписане у доњем делу. Занимљиво је да ће груписање по x оси за шестице и деветке бити симетрично у односу на координатни почетак, јер маса девет је углавном у "десном" крају, а шест у "левом" крају слике.

1.3 Пример класификација

Слика 5: Правилна класификација



Слика 6: Неправилна класификација



2 Други домаћи задатак

По угледу на пример 4.3 из документа Предавање 2, генерисати по $N = 500$ одбирака из двеју дводимензионалних класа.

- На дијаграму приказати одбирке
- Генерисати геометријско место тачака са константном вредношћу функција густина вероватноће па их приказати на дијаграму у простору облика.
- Испројектовати Бајесов класификатор минималне грешке и на дијаграму, заједно са одбирцима, скицирати класификациону линију, па проценити вероватноћу грешке.
- Поновити претходну тачку за неки други класификатор по избору.
- а класе облика генерисаних у претходним тачкама, испројектовати Валдов секвенцијални тест, па скицирати зависност броја потребних одбирака од усвојене вероватноће грешке првог, односно другог типа.

2.1 Генерисање одбирака

За почетак генерисаћемо одбирке двеју дводимензионалних класа чије су функције густине вероватноће у облику бимодалних Гаусовских расподела:

$$f_1(X) = P_{11}N(M_{11}, \Sigma_{11}) + P_{12}N(M_{12}, \Sigma_{12})$$
$$f_2(X) = P_{21}N(M_{21}, \Sigma_{21}) + P_{22}N(M_{22}, \Sigma_{22}),$$

при чему вероватноће, средње вредности и коваријационе матрице имају следеће вредности:

$$P_{11} = 0,6, M_{11} = \begin{bmatrix} -1,5 \\ 4,5 \end{bmatrix}, S_{11} = \begin{bmatrix} 3,5 & -1 \\ -1 & 2,2 \end{bmatrix}, P_{12} = 0,4, M_{12} = \begin{bmatrix} -0,5 \\ 0,5 \end{bmatrix}, S_{12} = \begin{bmatrix} 1,3 & 0,9 \\ 0,9 & 2 \end{bmatrix}$$
$$P_{21} = 0,45, M_{21} = \begin{bmatrix} 7,5 \\ -3,5 \end{bmatrix}, S_{21} = \begin{bmatrix} 1,5 & 1,1 \\ 1,1 & 1,5 \end{bmatrix}, P_{22} = 0,55, M_{22} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, S_{22} = \begin{bmatrix} 3 & -0,8 \\ -0,8 & 3 \end{bmatrix}$$

За генерисање одбирака коришћене су две функције. Одсечак кода прве функције која је коришћена је дат на одсечку 8. Овде се користи својство да се "обична" Гаусова расподела са варијансама 1 и очекивањем 0 кад се помножи са $\Phi\Lambda^{\frac{1}{2}}$, при чему су Φ матрица вектора сопствених вредности и Λ матрица сопствених вредности од S_x , и на крају дода M_x , при чему је M_x математичко очекивање вектора за који хоће да се генерише дата расподела, добије се Гаусова расподела са математичким очекивањем M_x и S_x .

Одсечак кода 8: Вишеваријабална Гаусова расподела-генерисање

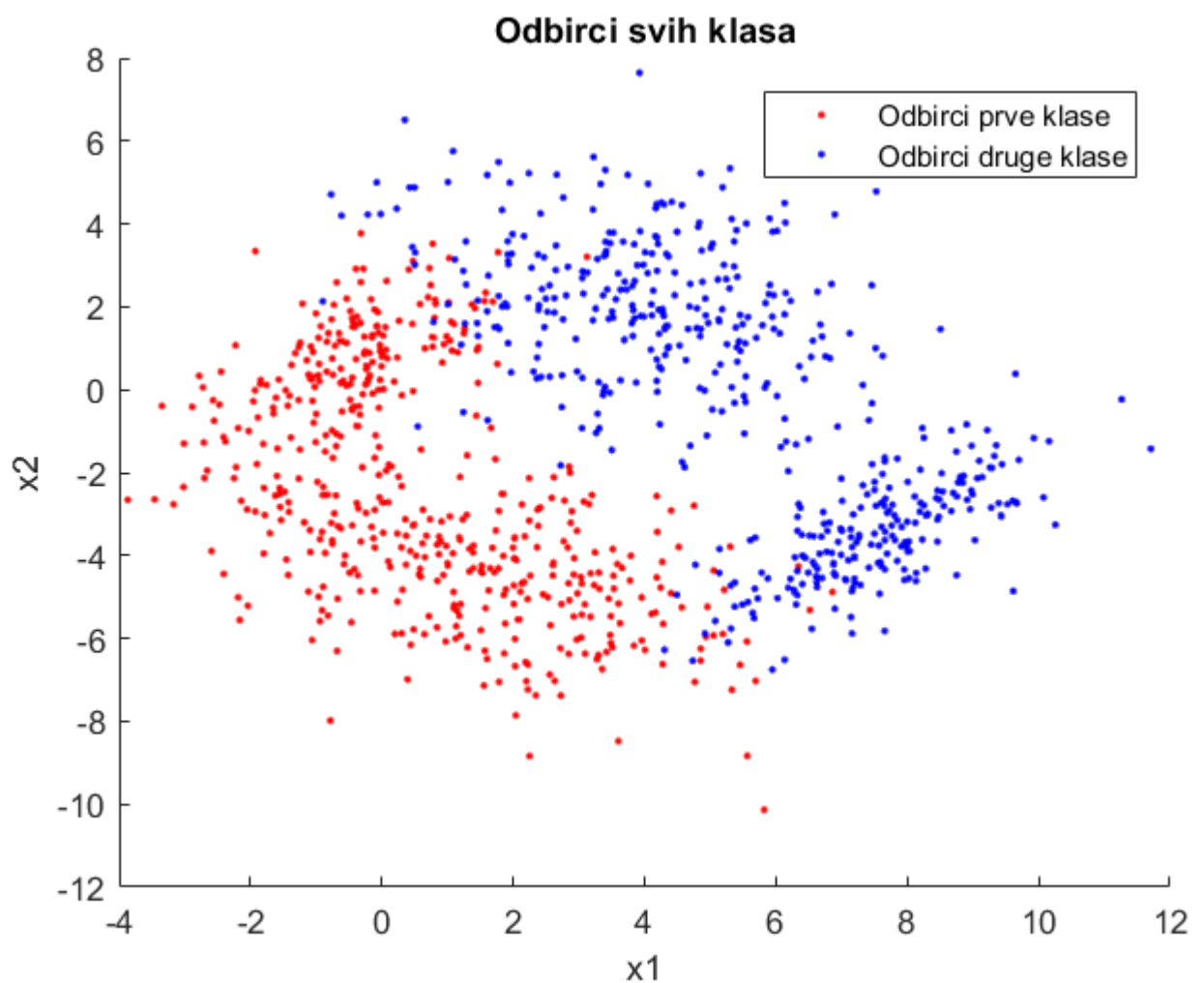
```
1 function [X] = gaussianMultivariateGenerate(Mx, Sx, num_points)
2 % Generates a matrix whose size is num_points, also it has
3 % gaussian multivariate distribution with Mx mean and
4 % Sx covariance matrix.
5
6 [F, L] = eig(Sx);
7
8 n = size(Sx, 1);
9 X = randn(n, num_points);
10
11 for i = 1 : num_points
12     X(1:n, i) = F * sqrt(L) * X(1:n, i) + Mx;
13 end
14
15 end
```

Друга коришћена функција је дата у одсечку 9. Она користи својство да је бимодална Гаусова расподела збир две Гаусове са одређеним вероватноћама и тако генерише потребну расподелу.

Одсечак кода 9: Вишеваријабална мултимодална Гаусова расподела-генерисање

```
1 function [X] = gaussianMultimodalGenerate(P1,P2, M1, M2, S1, S2, num_points)
2 X = zeros(num_points, 2);
3 for i = 1 : num_points
4     probability = rand(1, 1);
5     if (probability < P1)
6         X(i, :) = gaussianMultivariateGenerate(M1, S1, 1);
7     else
8         X(i, :) = gaussianMultivariateGenerate(M2, S2, 1);
9     end
10 end
11 end
```

Позивањем *gaussianMultimodalGenerate* са одговарајућим параметрима за обе класе добија се расподела као на слици 20.



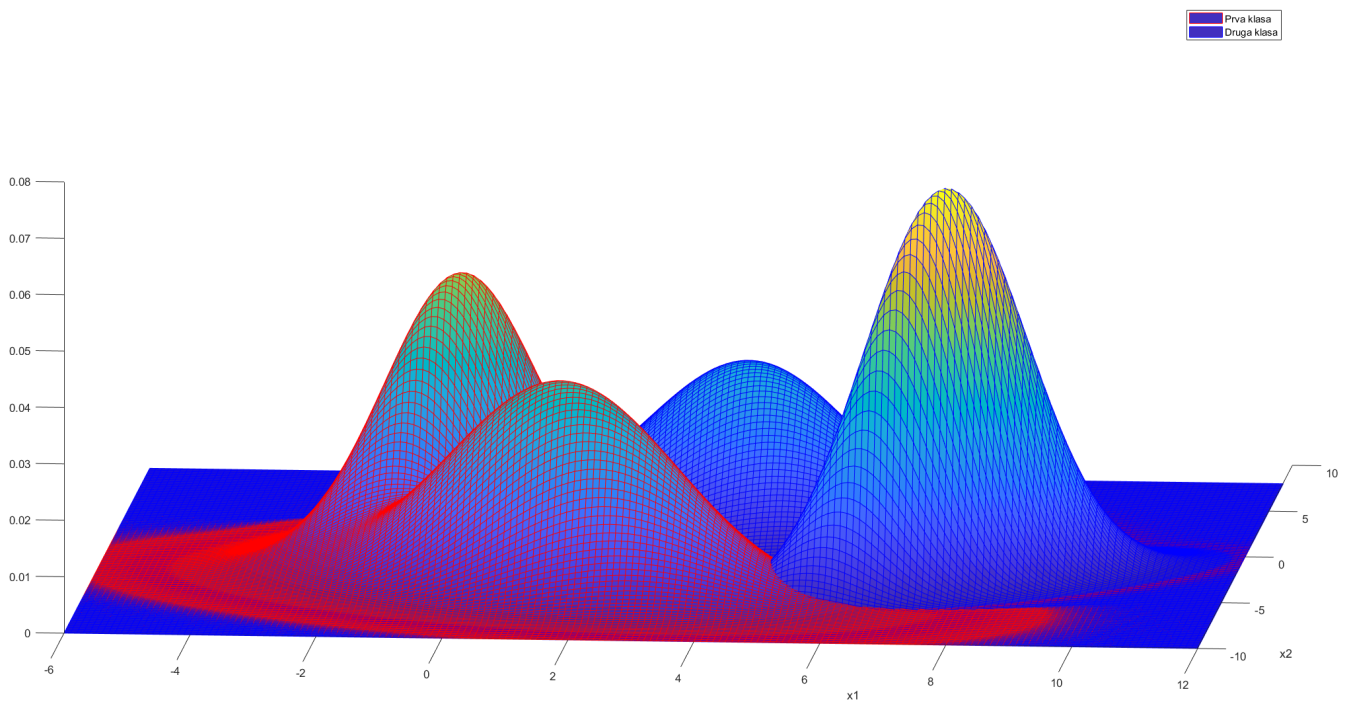
Слика 7: Одбирци

2.2 Функција густине вероватноће

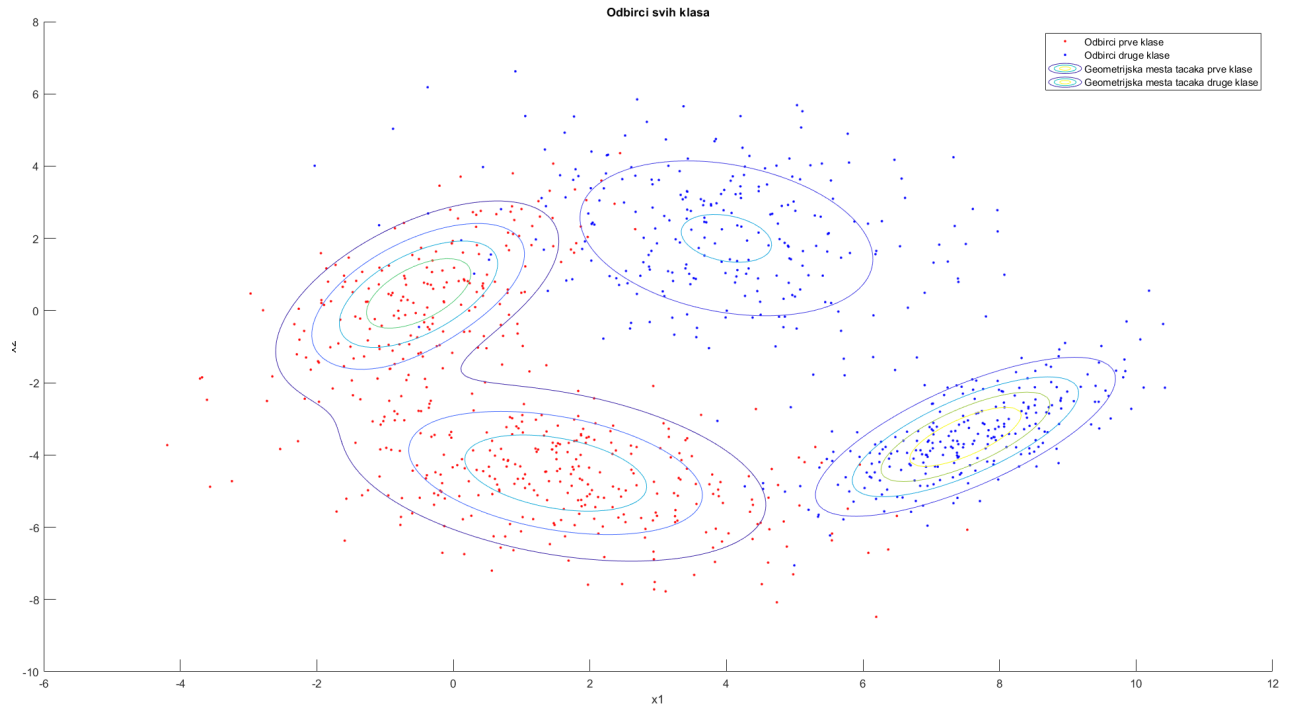
Функција густине вероватноће је израчуната за сваки "део" тако што је за сваку тачку у размаку од 0,1 у интервалу $[-6, -12]$ на x оси и на интервалу $[-10, 8]$ на y оси позивана функција која је дата одсечком кода 10. Дата функција густине вероватноће може да се види у 3D на слици 8. На слици 9 налазе се геометријска места тачака функције густине вероватноће од сваке класа, и то $0,8 * \max(f(X))$, $0,6 * \max(f(X))$, $0,4 * \max(f(X))$, $0,2 * \max(f(X))$. Може да се примети да они делови "класа" који су "оштрији" имају неке вредности које "тупљи" делови немају.

Одсечак кода 10: Вишеваријабална мултимодална Гаусова расподела-функција густине

```
1 function [f] = gaussianMultimodal(X, P1, P2, M1, M2, S1, S2)
2 %GAUSIANMULTIMODAL Calculates multimodal probability density function.
3 f = P1 * gaussianMultivariate(X, M1, S1) + P2 * gaussianMultivariate(X, M2,
   S2);
4 end
```



Слика 8: Функције густине вероватноће



Слика 9: Функције густине вероватноће-геометријска места тачака

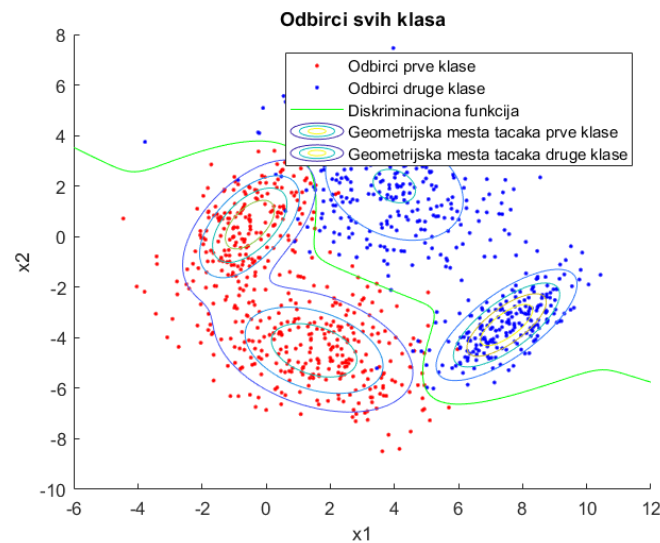
2.3 Бајесов класификатор минималне грешке

Дати класификатор има следеће правило одлучивања:

$$h(X) = -\ln(f_1(X) + \ln(f_2(X))) < \ln\left(\frac{P_1}{P_2}\right) \implies X \in \omega_1$$

$$h(X) = -\ln(f_1(X) + \ln(f_2(X))) > \ln\left(\frac{P_1}{P_2}\right) \implies X \in \omega_2$$

где су P_1, P_2 априорне вероватноће појављивања класа. Како је број генерисаних тачака 500, важи $P_1 = P_2 = 0,5$, стога $h(X) = 0$. Геометријска места тачака, као и дискриминациона функција је дата на слици 10.



Слика 10: Функције густине вероватноће, геометријска места тачака и дискриминациона линија

На одсечку кода 11 је приказано како је генерисана дискриминациона функција.

Одсечак кода 11: Генерисање дискриминационе функције

```

1 X1 = gaussianMultimodalGenerate(P11, P12, M11, M12, S11, S12, num_points);
2 X2 = gaussianMultimodalGenerate(P21, P22, M21, M22, S21, S22, num_points);
3
4 step = 0.1;
5 xIt = -6 : step : 12;
6 yIt = -10 : step : 8;
7
8 m = length(xIt);
9 n = length(yIt);
10
11 f1 = zeros(m, n);
12 f2 = zeros(m, n);
13 h = zeros(m, n);
14
15 [XItGrid, YItGrid] = meshgrid(xIt, yIt);
16
17 for i = 1 : m
18     for j = 1 : n
19         tempInput = [XItGrid(i, j) YItGrid(i, j)]';
20         f1(i, j) = gaussianMultimodal(tempInput, P11, P12, M11, M12, S11,
21                                     S12);
22         f2(i, j) = gaussianMultimodal(tempInput, P21, P22, M21, M22, S21,
23                                     S22);
24         h(i, j) = log(f2(i, j)) - log(f1(i, j));
25     end
26 end
27
28 f1Max = max(max(f1));
29 f2Max = max(max(f2));
30
31 figure(1);
32 hold on;
33 plot(X1(:, 1), X1(:, 2), 'r. ');
34 plot(X2(:, 1), X2(:, 2), 'b. ');
35
36 xlabel('x1');
37 ylabel('x2');
38
39 title('Odbirci svih klasa');
40
41 contour(XItGrid, YItGrid, h, [0 0], 'g');
42 contour(XItGrid, YItGrid, f1, [0.8 * f1Max f1Max * 0.6 f1Max * 0.4 f1Max *
43     0.2]);
44 contour(XItGrid, YItGrid, f2, [0.8 * f2Max f2Max * 0.6 f2Max * 0.4 f2Max *
45     0.2]);
46
47 legend('Odbirci prve klase', 'Odbirci druge klase', 'Diskriminaciona
48     funkcija', ...
49     'Geometrijska mesta tacaka prve klase', 'Geometrijska mesta tacaka
50     druge klase');

```

Вероватноћа грешке првог и другог типа се рачунају по формули:

$$\varepsilon_1 = \int_{L_2} f_1(X) dX, \varepsilon_2 = \int_{L_1} f_2(X) dX,$$

При чему су $L1, L2$ области на које дели дискриминациона функција, при чему је $L1$ она област за коју густина вероватноће припадност класи 1 има већу вредност, за $L2$ обратно. Пошто интеграл није решив, рачунао се методом правоугаоника (квадара). Функција из одсечка кода 12 је имплементација методе правоугаоника(квадара). Вредност грешке првог типа је 3,7% а другог типа је 3,8%.

Одсечак кода 12: Процена грешке

```

1 function [epsilon1, epsilon2] = errorEstimation(f1, f2, disFunction,
    num_rows, num_cols, step)
2 epsilon1 = 0;
3 epsilon2 = 0;
4
5 for i = 1 : num_rows
6     for j = 1 : num_cols
7         if (disFunction(i, j) < 0)
8             epsilon2 = epsilon2 + f2(i, j) * step^2;
9         else
10            epsilon1 = epsilon1 + f1(i, j) * step^2;
11        end
12    end
13 end
14
15 end

```

2.4 Бајесово правило одлучивања минималне цене

Класификацију ћемо поново извршити, али овога пута помоћу Бајесовог правила одлучивања минималне цене. Тест је сличан Бајесовом, с тим да одговарајућих коефицијенти кажњавају грешке једног или другог типа, односно "дозвољавају" их. c_{ij} цена одлуке $X \in \omega_i$ кад је заправо $X \in \omega_j$. Условна цена одлуке износи

$$r_i(X) = c_{i1}q_1(X) + c_{i2}q_2(X),$$

тј.

$$r_1(X) < r_2(X) \implies X \in \omega_1$$

$$r_1(X) > r_2(X) \implies X \in \omega_2$$

Тада априорни ризик добија облик: $r(X) = \min(r_1(X), r_2(X))$. Минимизацијом израза добија се:

$$\frac{f_1(X)}{f_2(X)} > \frac{(c_{12} - c_{22})P_2}{(c_{21} - c_{11})P_1} \implies X \in \omega_1$$

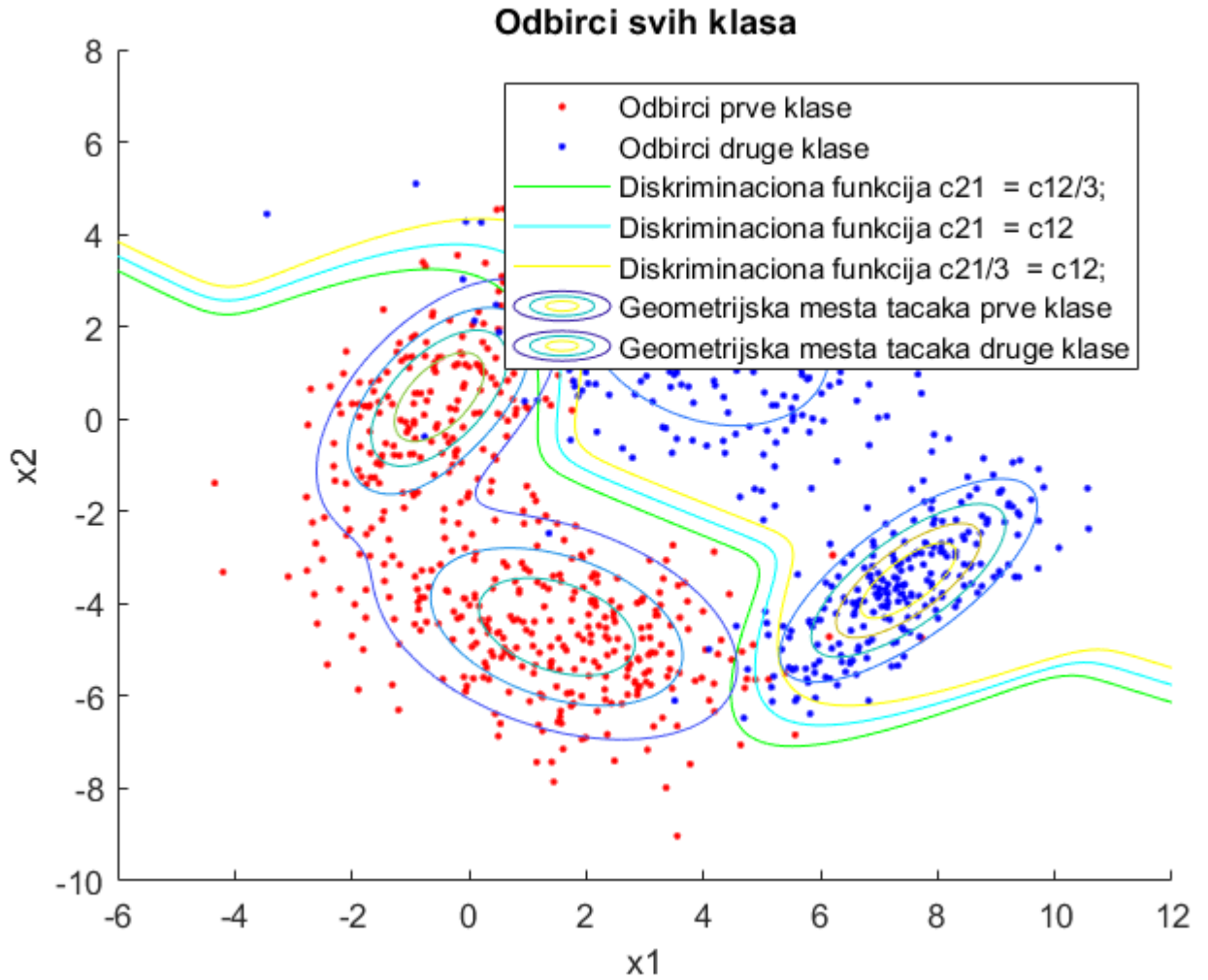
$$\frac{f_1(X)}{f_2(X)} < \frac{(c_{12} - c_{22})P_2}{(c_{21} - c_{11})P_1} \implies X \in \omega_2$$

Како је $P_2 = P_1$, и $h(X) = -\ln(\frac{f_1(X)}{f_2(X)})$ добија се:

$$h(X) = \ln(f_2(X)(c_{12} - c_{22})) - \ln(f_1(X)(c_{21} - c_{11})) < 0 \implies X \in \omega_1$$

$$h(X) = \ln(f_2(X)(c_{12} - c_{22})) - \ln(f_1(X)(c_{21} - c_{11})) > 0 \implies X \in \omega_2$$

Ако ставимо $c_{11} = c_{22} = 0$ постављањем c_{12}, c_{21} на одговарајуће вредности добијамо да се кажњава више одговарајући тип грешке (в. сл. 11). У случају $c_{12} = 3, c_{21} = 1$ добијамо да је дискриминациона функција примакнута првој класи, што има смисла, јер се кажњава промашај ако је се претпоставило да је прва класа, а треба да је друга. Обратно, важи само за другу класу кад је $c_{21} = 3, c_{12} = 1$. Кад је $c_{21} = c_{12}$ тест постаје Бајесов класификатор минималне грешке (што се и види на слици). Грешке су рачунају истом функцијом као и Бајесов класификатор, при чему у првом случају је грешка прве врсте 0,072, а друге 0,018, а у другом 0,017 односно 0,072. Претходно је такође логично јер је у првом случају дискриминациона функција близу прве класе и више ће бити погрешно класификованих њених одбирака у другу класу.



Слика 11: Тест минималне цене

2.5 Волдов секвенцијални тест

У претходно обрађеним класификаторима одлука је доношена само на основу информација које тренутно се поседују. Али у многим практичним применама информације стижу секвенцијално и њихов број се увећава. Фамилија секвенцијалних класификатора у озбир узима и претходна мерења. Један од њих је Волдов тест. Он доноси одлуку након коначног броја мерења када вероватноћа грешке спадне на жељене нивое ε_1 и ε_2 .

Дефинишемо s који представља негативни логаритам здружене функције густине вероватноће свих мерења и узимамо претпоставку да су оне независне:

$$s_m = -\ln \frac{f_1(X_1, \dots, X_m)}{f_2(X_1, \dots, X_m)} = \sum_{i=1}^m \left(-\ln \frac{f_1(X_i)}{f_2(X_i)} \right) = \sum_{i=1}^m h(X_i)$$

Волдов тест функционише:

- 1) $s_m \leq a \implies X \in \omega_1$
- 2) $a < s_m < b \implies$ узети следеће мерење
- 3) $s_m \geq b \implies X \in \omega_2$

Границе a и b су директно повезане са ε_1 и ε_2 и износе: $a = -\ln \frac{1-\varepsilon_1}{\varepsilon_2}$, $b = -\ln \frac{\varepsilon_1}{1-\varepsilon_2}$

Код који имплементира је одсечак кода [13](#). Можемо да видимо да код понавља *testCases* пута Волдов тест да би усредњио број итерација неопходан за детекцију класе којој припада јер се насумично узимају тачке из одређене класе X .

Одсечак кода 13: Одређивање припадностне

```

1  function [classIdx , i] = waldosTest(X, epsilon1 , epsilon2 , P11, P12, P21,
    P22, ...
2                                     M11, M12, M21, M22, S11, S12, S21,
    S22)
3  a = -log((1 - epsilon1) / epsilon2);
4  b = -log(epsilon1 / (1 - epsilon2));
5
6  num_points = size(X, 1);
7  testCases = 25;
8  sumI = 0;
9
10 for testIdx = 1 : testCases
11     idx = randperm(num_points);
12
13     i = 1;
14
15     s = -log(gaussianMultimodal(X(idx(i), :)', P11, P12, M11, M12, S11, S12)
16             ) + ...
17           log(gaussianMultimodal(X(idx(i), :)', P21, P22, M21, M22, S21,
18             S22)) ;
19
20     while ((i + 1 <= num_points) && (s < b) && (s > a))
21         i = i + 1;
22         s = s + -log(gaussianMultimodal(X(idx(i), :)', P11, P12, M11, M12,
23             S11, S12)) + ...
24             log(gaussianMultimodal(X(idx(i), :)', P21, P22, M21, M22, S21,
25             S22));
26     end
27
28     sumI = sumI + i;
29
30     if (s <= a)
31         classIdx = 1;
32     elseif (s >= b)
33         classIdx = 2;
34     else
35         classIdx = 0;
36     end
37 end
38
39 i = sumI / testCases;
40 end

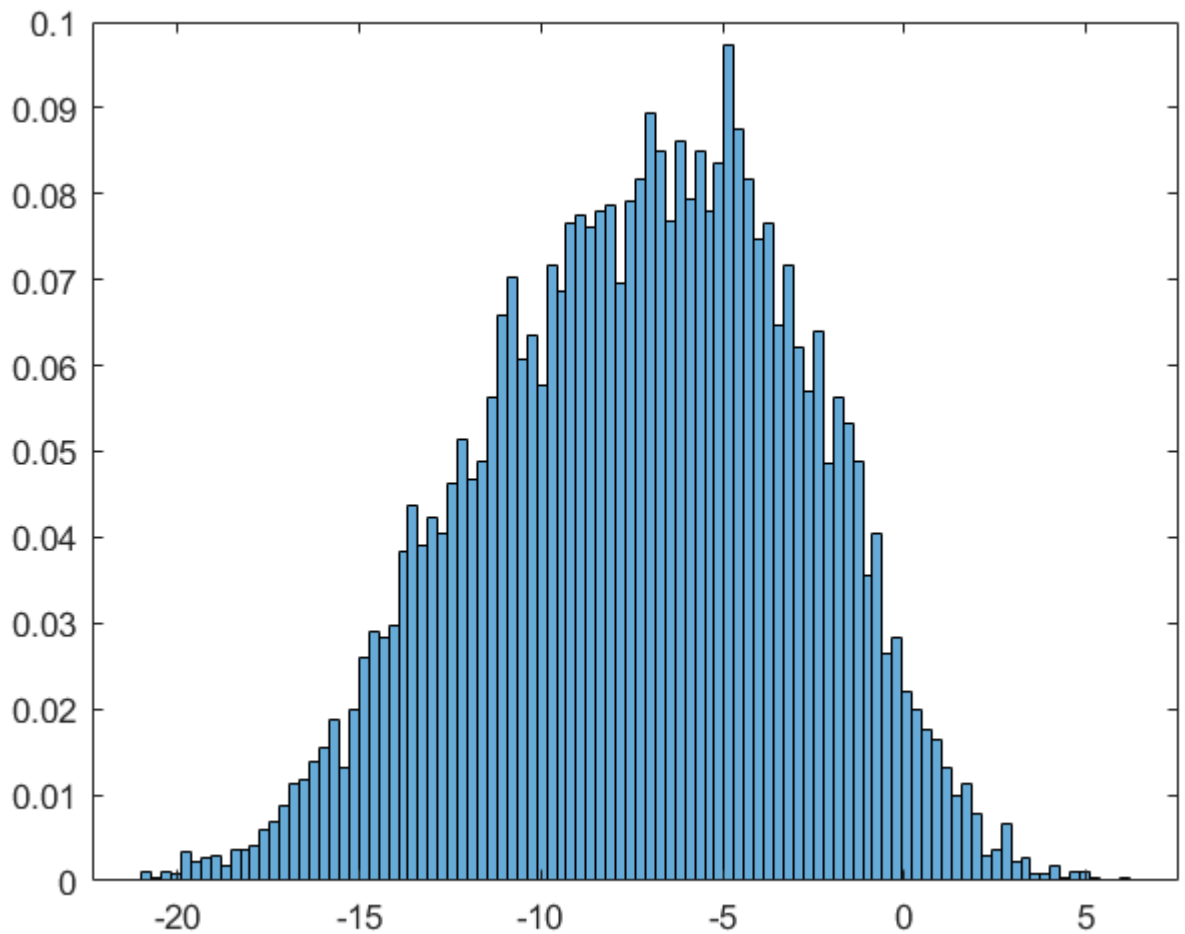
```

Зависност броја потребних одбарака из прве, односно друге класе је:

$$E\{m|\omega_1\} = \frac{a(1 - \varepsilon_1) + b\varepsilon_1}{\eta_1} = \frac{\ln \frac{\varepsilon_2}{1-\varepsilon_1}(1 - \varepsilon_1) + \ln \frac{1-\varepsilon_2}{\varepsilon_1}\varepsilon_1}{\eta_1}$$

$$E\{m|\omega_2\} = \frac{b(1 - \varepsilon_2) + a\varepsilon_2}{\eta_2} = \frac{\ln \frac{1-\varepsilon_2}{\varepsilon_1}(1 - \varepsilon_2) + \ln \frac{\varepsilon_2}{1-\varepsilon_1}\varepsilon_2}{\eta_2},$$

при чему је $\eta_i = E\{h(X)|\omega_i\} = \int_{-\infty}^{+\infty} h f_h(h|\omega_i)dh$. Проценићемо израз само један израз $E\{m|w_i\}$ јер је други аналоган. Прво морамо да проценимо η_1 , то је обављено коришћењем исечка кода 14, који се своди на метод хистограма тако што množимо одговарајућу функцију густине вероватноће са одговарајућом вредношћу. Добија се да је за ω_1 вредност $\eta_1 = -7,263$. Хистограм је дат на слици 12. Добијена зависност усредњеног броја одбарака од усвојене вероватноће грешке је на сликама 13a и 13b.



Слика 12: Хистограм процењене густине вероватноће

Одсечак кода 14: Процена условног h

```

1 function [expValCond, hApprox] = condProbDiscApprox(P11, P12, P21, P22, M11
    , M12, M21, M22, S11,...
2                                     S12, S21, S22)
3 num_points = 10000;
4 X1Approx = gaussianMultimodalGenerate(P11, P12, M11, M12, S11, S12,
    num_points);
5 m = size(X1Approx, 1);
6 hApprox = zeros(m, 1);
7
8 for i = 1 : m
9     tempInput = X1Approx(i, :)';
10    f1Approx = gaussianMultimodal(tempInput, P11, P12, M11, M12, S11, S12);
11    f2Approx = gaussianMultimodal(tempInput, P21, P22, M21, M22, S21, S22);
12    hApprox(i) = log(f2Approx) - log(f1Approx);
13 end
14
15 [condProb, bins] = histcounts(hApprox, 100, 'Normalization', 'pdf');
16
17 expValCond = 0;
18 step = bins(2) - bins(1);

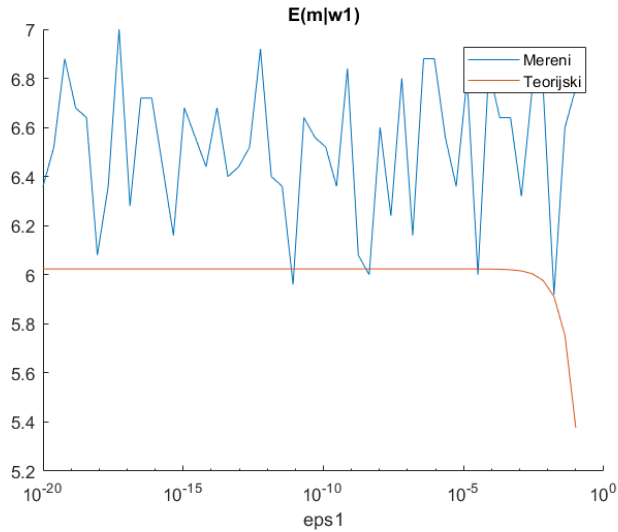
```

```

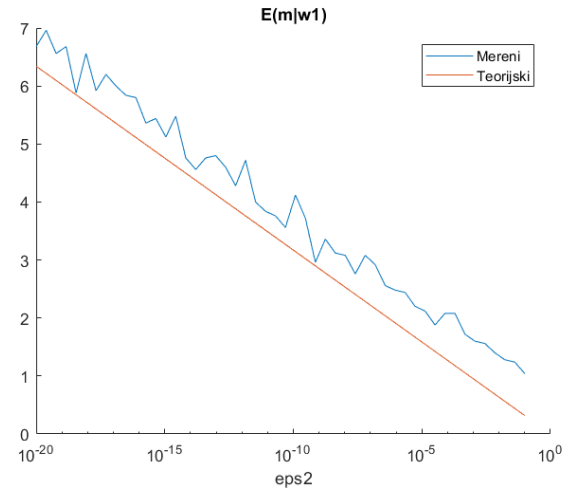
19 for i = 1 : length(condProb)
20     x = (bins(i) + bins(i + 1)) / 2;
21     expValCond = expValCond + x * condProb(i) * step;
22 end
23
24 end

```

Слика 13: Графици зависности



(a) $\varepsilon_2 = 10^{-20}$, зависност у односу на ε_1



(b) $\varepsilon_1 = 10^{-20}$, зависност у односу на ε_2

Са последње две слике се уочава да очекивани број одбирака из прве класе за доношење одлуке зависи од вероватноће грешке другог типа, док вероватноћа грешке првог типа скоро па нема утицај. Аналогно за другу класу, потребан број одбирака из друге класе зависи од ε_1 , док скоро уопште не зависи од ε_2

3 Трећи домаћи задатак

1. Генерисати две класе дводимензионалних облика. Изабрати функцију густине вероватноће облика тако да класе буду линеарно сепарабилне.
 - а) За тако генерисане облике извршити пројектовање линеарног класификатора једном од три итеративне процедуре.
 - б) Поновити претходни поступак методом жељеног излаза. Анализирати утицај елемената у матрици жељених излаза на коначну форму линеарног класификатора.
2. Генерисати две класе дводимензионалних облика које су сепарабилне, али не линеарно, па испројектовати квадратни класификатор методом по жељи.

3.1 Линеарни класификатор

3.1.1 Генерисање класа

За почетак генерисаћемо одбирке двеју дводимензионалних класа чије су функције густине вероватноће у облику бимодалних Гаусових расподела:

$$f_1(X) = P_{11}N(M_{11}, \Sigma_{11}) + P_{12}N(M_{12}, \Sigma_{12})$$

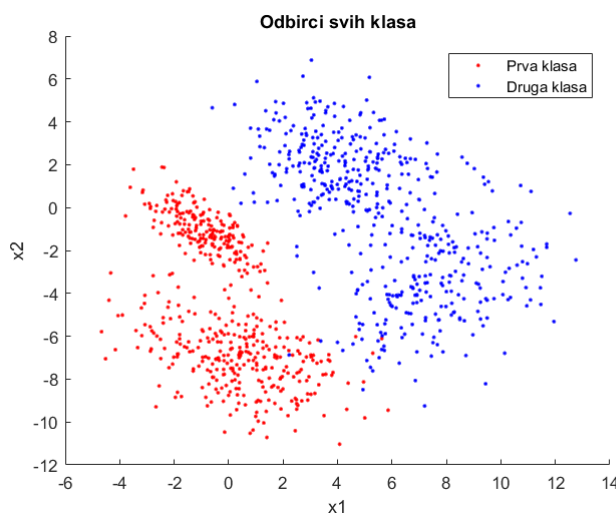
$$f_2(X) = P_{21}N(M_{21}, \Sigma_{21}) + P_{22}N(M_{22}, \Sigma_{22}),$$

при чему вероватноће, средње вредности и коваријационе матрице имају следеће вредности:

$$P_{11} = 0,6, M_{11} = \begin{bmatrix} 0,5 \\ -7 \end{bmatrix}, S_{11} = \begin{bmatrix} 3,5 & -1 \\ -1 & 2,2 \end{bmatrix}, P_{12} = 0,4, M_{12} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, S_{12} = \begin{bmatrix} 1,3 & -0,9 \\ -0,9 & 2 \end{bmatrix}$$

$$P_{21} = 0,45, M_{21} = \begin{bmatrix} 7,5 \\ -3,5 \end{bmatrix}, S_{21} = \begin{bmatrix} 4 & 1,1 \\ 1,1 & 4 \end{bmatrix}, P_{22} = 0,55, M_{22} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, S_{22} = \begin{bmatrix} 3 & -0,8 \\ -0,8 & 3 \end{bmatrix}$$

Добијене класе су на слици 14.



Слика 14: Одбирци

3.1.2 Пројектовање линеарног класификатора методом ресупституције

Линеарни класификатор је један од најједноставнијих класификатора, Иако у већини случајева није оптималан јер класе не морају да буду линеарно сепарабилне, често се користи због једноставности. Дискриминациона функција је линеарна и правило одлучивања гласи:

$$h(X) = V^T + v_0 < 0 \implies X \in \omega_1$$

$$h(X) = V^T + v_0 > 0 \implies X \in \omega_2$$

Метод који је коришћен је метод ресупституције. Првобитно су процењени вектори средњих вредности $M1Approx$, $M2Approx$ и коваријационе матрице $S1Approx$, $S2Approx$, то је коришћено одсечком кода 15.

Одсечак кода 15: Процена математичког очекивања и коваријационе матрице

```

1 M1Approx = mean(X1)';
2 M2Approx = mean(X2)';
3
4 S1Approx = cov(X1);
5 S2Approx = cov(X2);

```

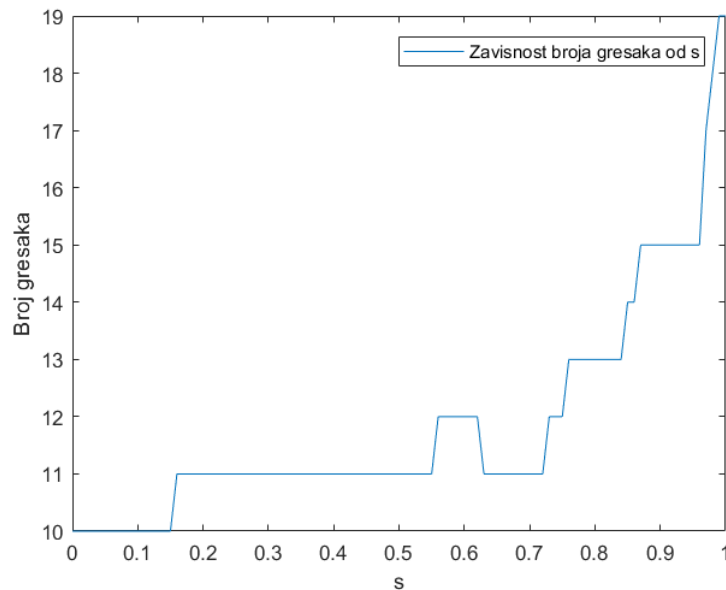
Потом се у сваком кораку мења вредност параметра $s \in [0, 1]$, и за свако s одреди се вектор :

$$V = [s * S1Approx + (1 - s) * S2Approx]^{-1} (M2Approx - M1Approx)$$

Потом се вектор X пројектује на вектор V :

$$y_j^{(i)} = V^T X_j^{(i)}, i = 1, 2, j = 1, \dots, N_i$$

где је N_i број облика из i -те класе, а $X_j^{(i)}$ j -ти обучавајући вектор из i -те класе. Затим се v_0 мења да се добије најмања грешка. Мења се у опсегу $[-\max(\max(y_j^{(i)}), \max(y_j^{(2)})), -\min(\min(y_j^{(1)}), \min(y_j^{(2)}))]$. На крају процедуре бира се оно s за које је број погрешно класификованих објеката (односно грешка) најмањи. На слици 15 имамо како се s мења у интервалу $[0, 1]$. Можемо да приметимо да је минимална вредност било која између 0 и 0,15, стога бирамо 0.



Слика 15: Зависност броја грешака од s

За овај случај се добија $V = [1, 52541.1170]^T$, $v_0 = -0,6212$ што се види на слици 16. Такође треба приметити да је број грешака 10 што значи да је грешка 1%.

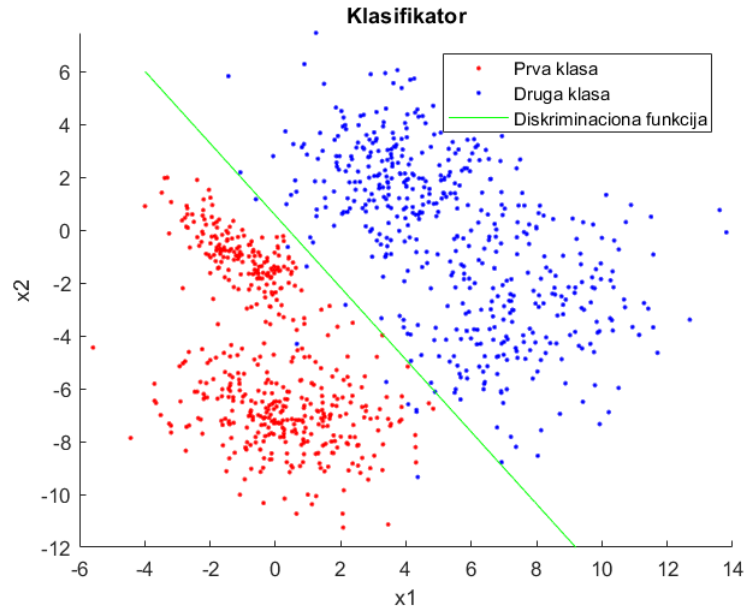
Код којим је одрађено претходно је одсечак кода 16.

Одсечак кода 16: Пројектовање линеарног класификатора

```

1 step = 0.01;
2 s = 0 : step : 1;
3 num_s = length(s);
4 for i = 1 : num_s
5     V = (s(i) * S1Approx + (1 - s(i)) * S2Approx) \ (M2Approx - M1Approx);

```



Слика 16: Одбирци

```

6     Y1 = X1 * V;
7     Y2 = X2 * V;
8
9     downLimit = -max(max(Y1), max(Y2));
10    upLimit = -min(min(Y1), min(Y2));
11
12    V0 = downLimit : step : upLimit;
13
14    minError = size(Y1, 1) + size(Y2, 1);
15    minIdx = 0;
16
17    V0Num = length(V0);
18
19    for j = 1 : V0Num
20        v0 = V0(j);
21        wrongClass1 = nnz(Y1 > -v0);
22        wrongClass2 = nnz(Y2 < -v0);
23        wrongClass = wrongClass1 + wrongClass2;
24        if (wrongClass < minError)
25            minError = wrongClass;
26            minIdx = j;
27        end
28    end
29    epsilon(i) = minError;
30    V0Min(i) = V0(minIdx);
31
32 end
33
34 [epsilon, idx] = min(epsilon);
35
36 sVal = s(idx);
37
38 V = (sVal * S1Approx + (1 - sVal) * S2Approx) \ (M2Approx - M1Approx)
39 v0 = V0Min(idx)

```

3.1.3 Метод жељеног излаза

Поновићемо поступак пројектовања линеарног класификатора над датим класама, али овога пута методом жељеног излаза. Код ове методе линија класификациона линија се формира тако да за исправно класификован објекат увек буде позитивна:

$$h(X) = -V^T - v_0 > 0 \implies X \in \omega_1$$

$$h(X) = V^T + v_0 > 0 \implies X \in \omega_2$$

Да би се то постигло уводи се нови улазни вектор:

$$Z = [-1 - X_1 \dots - X_n]^T, X \in \omega_1$$

$$Z = [1 X_1 \dots X_n]^T, X \in \omega_2$$

Проблем се своди на одређивање вектора $W_{(n+1) \times 1}$ за који ће што више одбирака Z да испуне релацију $h(Z) = W^T Z > 0$. Критеријумска функција која је најзахвалнија за минимизацију је:

$$\overline{\varepsilon^2} = \frac{1}{N} \sum_{j=1}^N (W^T Z_j - \gamma(Z_j))^2$$

Применом парцијалног извода на њу, и изједначавањем са 0, добија се $U^T W = \Gamma$. Дата једначина се решава псеудоинверзијом:

$$W = (U U^T)^{-1} U \Gamma,$$

где је $U = [Z_1 Z_2 \dots Z_n]$ матрица узорака, а $\Gamma = [\gamma(Z_1) \dots \gamma(Z_n)]^T$ матрица жељених излаза. Често се за матрицу жељених излаза узима јединични вектор. Одсечак кода 17 имплементира методу жељеног излаза.

Одсечак кода 17: Метода очекиваног излаза

```

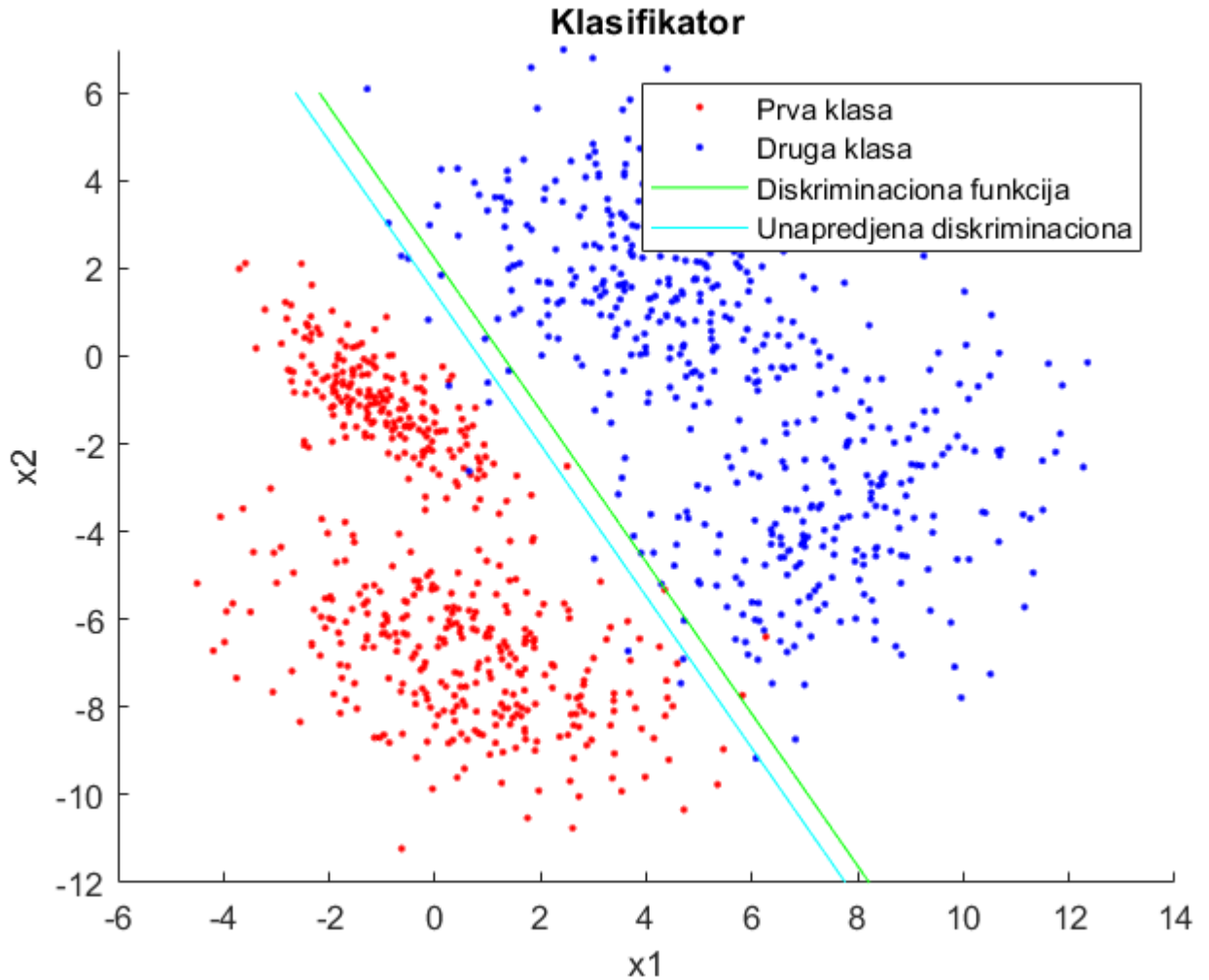
1 U = zeros(3, num_points * 2);
2 G = ones(num_points * 2, 1);
3 for i = 1 : num_points
4     U(:, 2 * i - 1) = [-1 -X1(i, :)];
5     U(:, 2 * i) = [1 X2(i, :)];
6 end
7
8 W = U * U' \ U * G;
9
10 xIt = -6 : 0.01 : 14;
11 yIt = -12 : 0.01 : 6;
12 hExpectedOutput = zeros(length(xIt), length(yIt));
13 for i = 1 : length(xIt)
14     for j = 1 : length(yIt)
15         hExpectedOutput(i, j) = [1 xIt(i) yIt(j)] * W;
16     end
17 end
```

Добијена функција је $h(X) = -0,2581 + 0,20417_1 + 0,1166x_2$ је приказан на слици 17. Грешка се рачуна функцијом чији је одсечак кода 18.

Одсечак кода 18: Метода очекиваног излаза - грешка

```

1 function [epsilon1, epsilon2, epsilon] = errorEstimationEO(U, W, num_points)
2     Y = U' * W;
3     epsilon = numel(Y(Y < 0)) / (2 * num_points);
4     epsilon1 = 0;
5     epsilon2 = 0;
6     for i = 1 : num_points
7         epsilon1 = (Y(2 * i - 1) < 0) + epsilon1;
```



Слика 17: Дискриминациона функција код методе очекиваног излаза

```

8         epsilon2 = (Y(2 * i) < 0) + epsilon2;
9     end
10    epsilon1 = epsilon1 / num_points;
11    epsilon2 = epsilon2 / num_points;
12 end

```

Добијене грешке су $\varepsilon_2 = 0,034$, $\varepsilon_1 = 0,004$, $\varepsilon = 0,019$. Можемо да приметимо да је ε_2 много већа него ε_1 . Ако бисмо желели да је смањимо на уштрб грешке ε_1 онда бисмо очекиване излазе за елементе који припадају другој класи морали да повећамо. Када то урадимо добија се $\varepsilon_1 = 0,008$, $\varepsilon = 0,024$, $\varepsilon = 0,016$. На слици 17 се налази и новодобијена функција. Уопштено ако имамо неку грешку која је већа и желимо да је умањимо методом жељеног излаза, само треба да повећамо жељене излазе за тип грешке који желимо да смањимо.

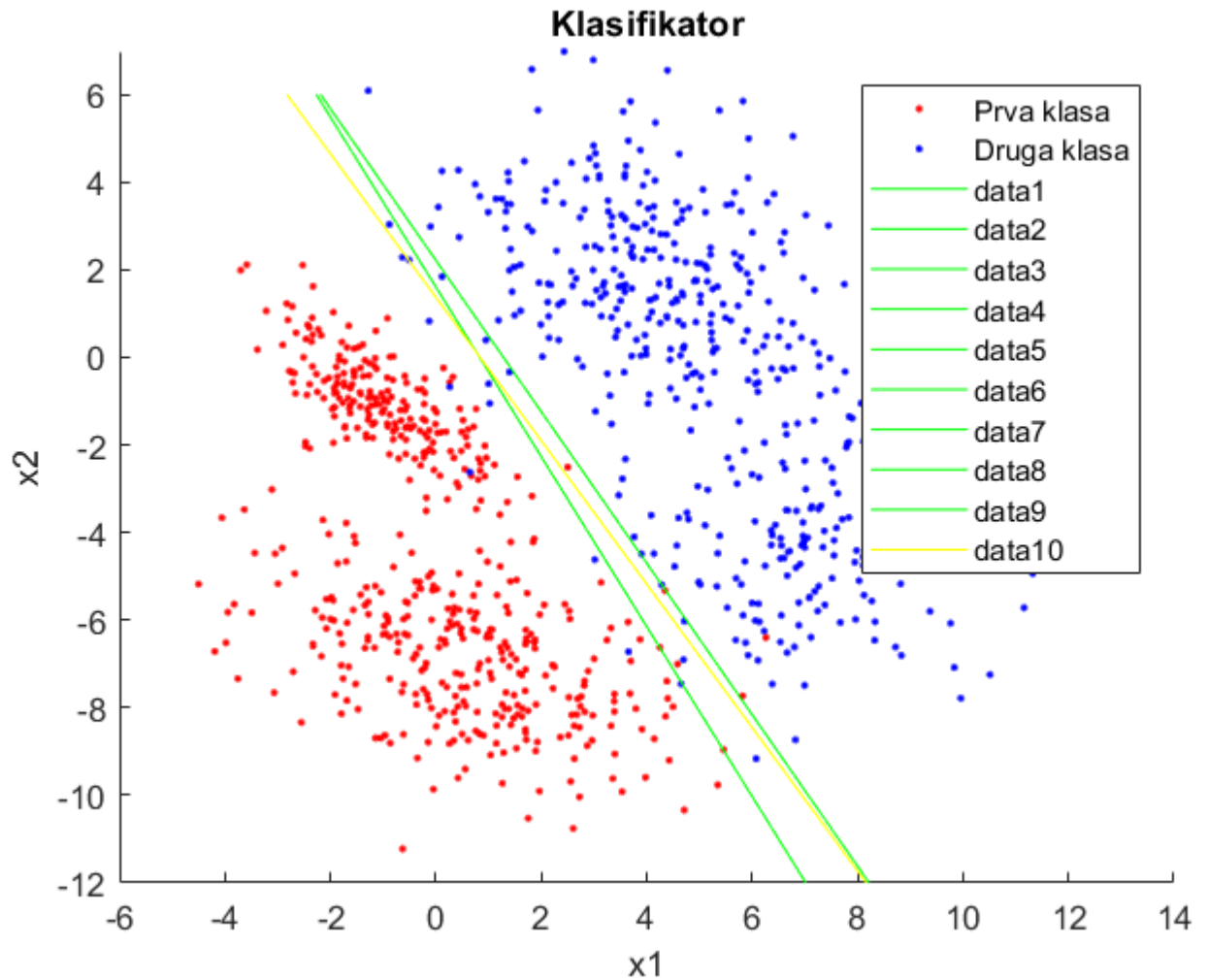
Ако бисмо желели да минимизујемо свеукупну грешку тад би свим тачкама у близини дискриминационе функције требали да повећамо вредност. Ако итеративно понављамо дату операцију повећавања добијемо низ дискриминационих функција које конвергирају као на слици 18 жутој линији. Грешка се смањује и $\varepsilon = 0,015$ након ове операције, што је мање него првобитна свеукупна грешка. Одсечак кода 19 приказује имплементацију претходно наведеног.

Одсечак кода 19: Метода очекиваног излаза - конвергирање

```

1  U = zeros(3, num_points * 2);
2  G = ones(num_points * 2, 1);
3  for i = 1 : num_points
4      U(:, 2 * i - 1) = [-1 -X1(i, :)] ;

```



Слика 18: Низ дискриминационих функција код методе очекиваног излаза

```

5     U(:, 2 * i) = [1 X2(i, :)] ;
6 end
7 for t = 1 : 10
8
9     W = U * U' \ U * G;
10    xIt = -6 : 0.01 : 14;
11    yIt = -12 : 0.01 : 6;
12    hExpectedOutput = zeros(length(xIt), length(yIt));
13    for i = 1 : length(xIt)
14        for j = 1 : length(yIt)
15            hExpectedOutput(i, j) = [1 xIt(i) yIt(j)] * W;
16        end
17    end
18
19    [e1, e2, e] = errorEstimationEO(U, W, num_points);
20
21    fprintf("ExpClose : Greska %.3f Greska1 : %.3f Greska 2 %.3f\n", e, e1,
22            e2);
23    if (t == 10)
24        contour(xIt, yIt, hExpectedOutput', [0 0], 'y');
25    else
26        contour(xIt, yIt, hExpectedOutput', [0 0], 'g');

```



```

26     end
27     Y = U' * W;
28     idx = find(Y(1:num_points) < 0.1);
29     for i = 1:length(idx)
30         G(idx(i)) = 10;
31     end
32 end

```

3.2 Квадратни класификатор

Представља сложељнију форму класификатора. Користи се кад две класе нису линеарно сепарабилне. Општа форма квадратног класификатора је:

$$h(X) = X^T Q X + V^T X + v_0 < 0 \implies X \in \omega_1$$

$$h(X) = X^T Q X + V^T X + v_0 > 0 \implies X \in \omega_2$$

Да би се могао пројектовати неком од постојећих метода, треба да се изврши првидна линеаризација квадратног класификатора. То ћемо учинит на следећи начин:

$$h(X) = q_{11}x_1^2 + q_{22}x_2^2 + q_{12}x_1x_2 + v_1x + v_2x_2 + v_0 = W^T Z > 0$$

где је улазни вектор

$$Z = [-1 - x_1 - x_2 - x_1x_2 - x_1^2 - x_2^2]^T, X \in \omega_1$$

$$Z = [1 \ x_1 \ x_2 \ x_1x_2 \ x_1^2 \ x_2^2]^T, X \in \omega_2$$

вектор параметара $W = [v_0 v_1 v_2 q_{12} q_{11} q_{22}]$

За почетак генерисаћемо одбирке двеју дводимензионалних класа чије су функције густине вероватноће у облику бимодалних Гаусових расподела:

$$f_1(X) = P_{11}N(M_{11}, \Sigma_{11}) + P_{12}N(M_{12}, \Sigma_{12})$$

$$f_2(X) = P_{21}N(M_{21}, \Sigma_{21}) + P_{22}N(M_{22}, \Sigma_{22}),$$

при чему вероватноће, средње вредности и коваријационе матрице имају следеће вредности:

$$P_{11} = 0,6, M_{11} = \begin{bmatrix} 0,5 \\ -7 \end{bmatrix}, S_{11} = \begin{bmatrix} 3,5 & -1 \\ -1 & 2,2 \end{bmatrix}, P_{12} = 0,4, M_{12} = \begin{bmatrix} 12 \\ 6 \end{bmatrix}, S_{12} = \begin{bmatrix} 1,3 & -0,9 \\ -0,9 & 1,2 \end{bmatrix}$$

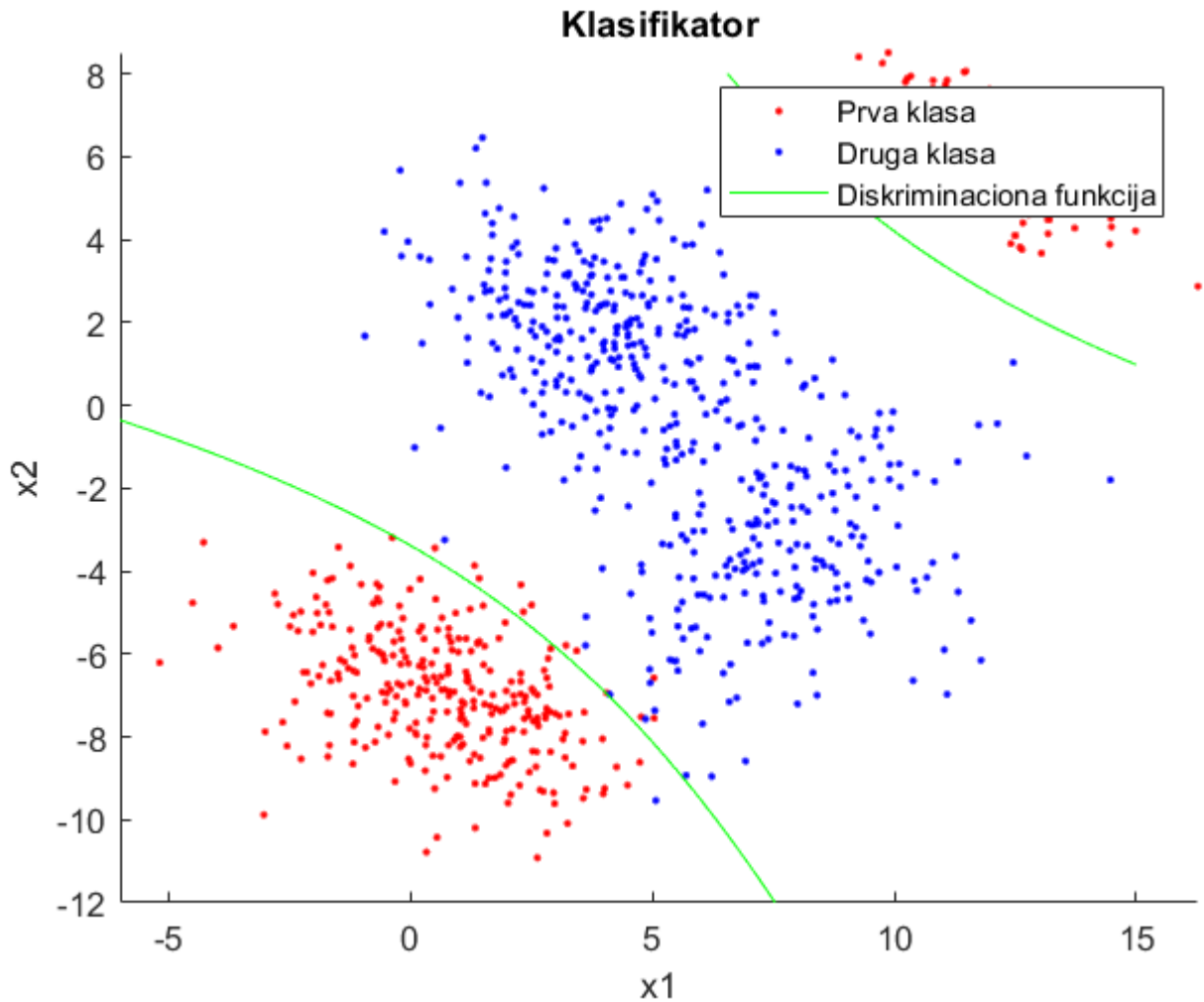
$$P_{21} = 0,45, M_{21} = \begin{bmatrix} 7,5 \\ -3,5 \end{bmatrix}, S_{21} = \begin{bmatrix} 4 & 1,1 \\ 1,1 & 4 \end{bmatrix}, P_{22} = 0,55, M_{22} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, S_{22} = \begin{bmatrix} 3 & -0,8 \\ -0,8 & 3 \end{bmatrix}$$

Применом методе псудоинверзије добијамо да ус вектор параметара и дискриминациона функција:

$$W = [0,7240, 0,0622 \ 0,1785 \ -0,0310 \ -0,0061 \ 0,0105]$$

$$h(X) = -0,7240 + 0,0622x_1 + 0,1785x_2 - 0,0310x_1x_2 - 0,0061x_1^2 + 0,0105x_2^2$$

Добијају се грешке $\varepsilon = 0,013, \varepsilon_1 = 0,024, \varepsilon_2 = 0,002$. Дискриминациона функција и одбирци су приказани на слици 19.



Слика 19: Квадратни класификатор

4 Четврти домаћи задатак

1. Генерисати по $N = 500$ дводимензионалних одбирака из четири класе које ће бити линеарно сепарабилне. Препоручује се да то буду гаусовски расподељени дводимензионални облици. Изабрати једну од метода за кластеризацију (С-mean метод, метод квадратне декомпозиције, метод максималне веродостојности или метод грана и граница) и применити је на формиране узорке класа. Извршити анализу осетљивости изабраног алгорита на почетну кластеризацију, као и средњи број потребних итерација. Такође извршити анализе случаја када се априорно не познаје број класа.
2. Генерисати по $N = 500$ дводимензионалних одбирака из две класе које су нелинеарно сепарабилне. Изабрати једну од метода за кластеризацију које су применљиве за нелинеарно сепарабилне класе (метод квадратне декомпозиције или метод максималне веродостојности) и поновити анализу из претходне тачке.

4.1 С-Mean алгоритам

4.1.1 Генерисање одбирака

За почетак генерисаћемо одбирке четири дводимензионалних класа чије су функције густине вероватноће у облику гаусовских расподела:

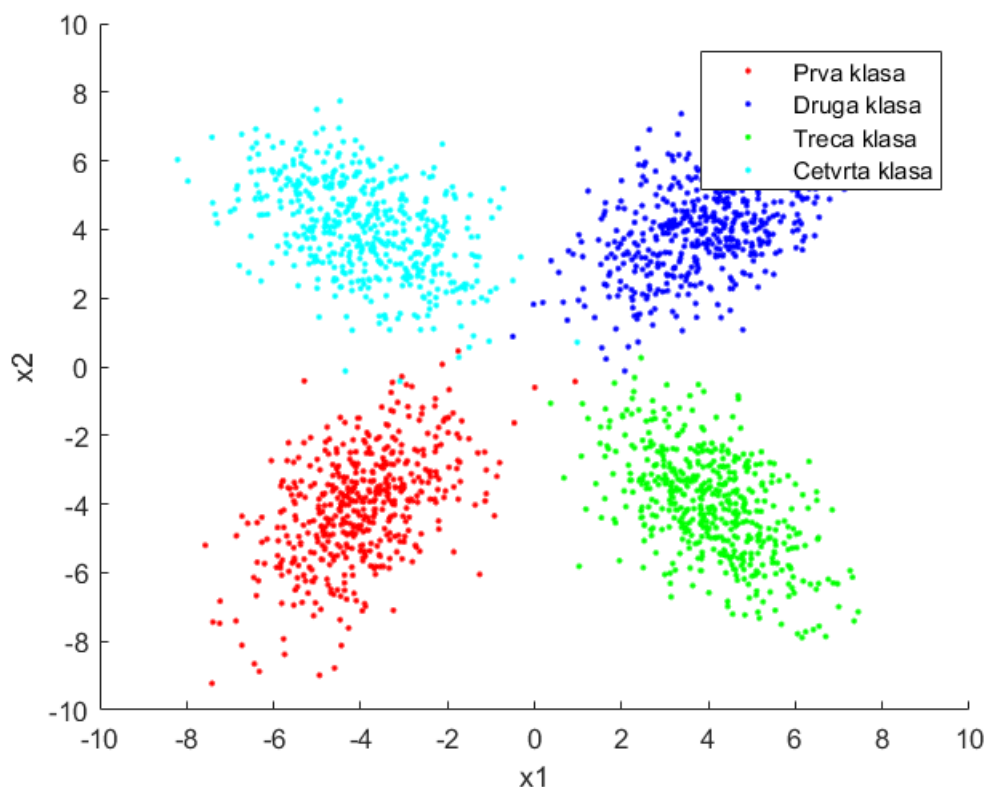
$$f_1(X) = N(M_1, \Sigma_1), f_2(X) = N(M_2, \Sigma_2), f_3(X) = N(M_3, \Sigma_3), f_4(X) = N(M_4, \Sigma_4)$$

при чему вероватноће, средње вредности и коваријационе матрице имају следеће вредности:

$$M_1 = \begin{bmatrix} -4 \\ -4 \end{bmatrix}, M_2 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}, M_3 = \begin{bmatrix} 4 \\ -4 \end{bmatrix}, M_4 = \begin{bmatrix} -4 \\ 4 \end{bmatrix}, M_5 = \begin{bmatrix} -4 \\ 4 \end{bmatrix},$$

$$S_1 = \begin{bmatrix} 1,5 & 1 \\ 1 & 2,5 \end{bmatrix}, S_2 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}, S_3 = \begin{bmatrix} 1,5 & -1 \\ -1 & 1,5 \end{bmatrix}, S_4 = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

Добијају се 4 класе одбирака као на слици 20



Слика 20: Одбирци

4.1.2 Алгоритам

Кластеризацију над датим скупом података ћемо извршити уз помоћ C-mean алгоритма. Једна од предности ове методе су њена једноставност, могућност стохастичке почетне кластеризације. Међутим, има и одређених ограничења. Кластери су поелени део по део линеарним линијама, што ограничава примену ове методе на линеарно сепарабилне класе. Такође, захтева априорно знање о броју кластера и не гарантује глобалну конвергенцију.

Одсечак кода 20 имплементира c-Mean алгоритам.

Одсечак кода 20: C-Mean функција

```

1 function [num_it, critFun] = cMean(X1, X2, X3, X4, num_points, num_classes,
   apPer, drawFigure)
2
3 chunkRand = floor((1 - apPer) * num_points);
4 chunkApp = num_points - chunkRand;
5
6 Z = [X1(1 : chunkRand, :); X2(1 : chunkRand, :); ...
7      X3(1 : chunkRand, :); X4(1 : chunkRand, :)];
8
```

```

9  num_points_all = 4 * chunkRand;
10
11  idx = randperm(num_points_all);
12  chunk_size = floor(num_points_all / num_classes);
13
14  Y1 = Z(idx(1 : chunk_size), :);
15  Y1App = X1(chunkRand + 1 : num_points, :);
16  Y1 = [Y1App; Y1];
17  Y2 = Z(idx(chunk_size + 1:2 * chunk_size), :);
18  Y2App = X2(chunkRand + 1 : num_points, :);
19  Y2 = [Y2App; Y2];
20  Y3 = [];
21  Y3App = [];
22  Y4 = [];
23  Y4App = [];
24  Y5 = [];
25  Y5App = [];
26  if (num_classes > 2)
27      Y3 = Z(idx(2 * chunk_size + 1:3 * chunk_size), :);
28      Y3App = X3(chunkRand + 1 : num_points, :);
29      Y3 = [Y3App; Y3];
30  end
31  if (num_classes > 3)
32      Y4 = Z(idx(3 * chunk_size + 1:4 * chunk_size), :);
33      Y4App = X4(chunkRand + 1 : num_points, :);
34      Y4 = [Y4App; Y4];
35  end
36  if (num_classes > 4)
37      Y5 = Z(idx(4 * chunk_size + 1:5 * chunk_size), :);
38  end
39  run = 1;
40
41  it = 0;
42
43  while (run && it < 50)
44      it = it + 1;
45      M1 = mean(Y1);
46      M2 = mean(Y2);
47      M3 = mean(Y3);
48      M4 = mean(Y4);
49      M5 = mean(Y5);
50
51      if (size(M1, 2) == 1)
52          M1 = [realmax realmax];
53      end
54      if (size(M2, 2) == 1)
55          M2 = [realmax realmax];
56      end
57      if (size(M3, 2) == 1)
58          M3 = [realmax realmax];
59      end
60      if (size(M4, 2) == 1)
61          M4 = [realmax realmax];
62      end
63      if (size(M5, 2) == 1)
64          M5 = [realmax realmax];
65      end
66      clear X1 X2 X3 X4;

```

```

67
68 M = [M1; M2; M3; M4; M5];
69
70 [X11, X21, X31, X41, X51, changed1] = reCluster(Y1(chunkApp + 1: end,
71      :), M, 1);
72 [X12, X22, X32, X42, X52, changed2] = reCluster(Y2(chunkApp + 1: end,
73      :), M, 2);
74 [X13, X23, X33, X43, X53, changed3] = reCluster(Y3(chunkApp + 1: end,
75      :), M, 3);
76 [X14, X24, X34, X44, X54, changed4] = reCluster(Y4(chunkApp + 1: end,
77      :), M, 4);
78 [X15, X25, X35, X45, X55, changed5] = reCluster(Y5(chunkApp + 1: end,
79      :), M, 5);
80 run = changed1 | changed2 | changed3 | changed4 | changed5;
81 Y1 = [Y1App; X11; X12; X13; X14; X15];
82 Y2 = [Y2App; X21; X22; X23; X24; X25];
83 Y3 = [Y3App; X31; X32; X33; X34; X35];
84 Y4 = [Y4App; X41; X42; X43; X44; X45];
85 Y5 = [Y5App; X51; X52; X53; X54; X55];
86 end
87
88 num_it = it;

```

На почетку алгоритма имамо насумично кластеризовање података. На сликама [21a](#), [21b](#), [21c](#), [21d](#) видимо почетне кластеризације за различит број почетних класа.

Критеријум минимизује критеријумску функцију:

$$J = \frac{1}{N} \sum_{r=1}^L \sum_{j=1}^{N_r} \|X_j^{(r)} - M_r\|^2$$

где је N -укупан број одбирака, L -укупан број кластера, M_r -вектор средње вредности r -тог кластера, N_r -број елемената у r -том кластеру. С обзиром да критеријум минимизује средње квадратно одступање од центара кластера, закључује се да ће се сваки одбирак придружити најближем центру кластера. Међутим, овде се може закључити да С-mean узима у обзир само математичко очекивање, али не и коваријационе матрице одбирака. Сређивање израза добија се да је правило одлучивања:

$$\|X_i - M_t(l)\| = \min_{1 \leq j \leq L} \|X_i - M_j(l)\| \implies X_i \in \omega_t$$

Дакле прво се обави рекласификација, потом поново се израчуна $M_j(l)$ и то се извршава док има промене изгледа кластера. Кад нема завршава се. Одсечци кода [22](#) и [21](#) су функције које имплементирају логику налажења минималне вредности критеријумске функције.

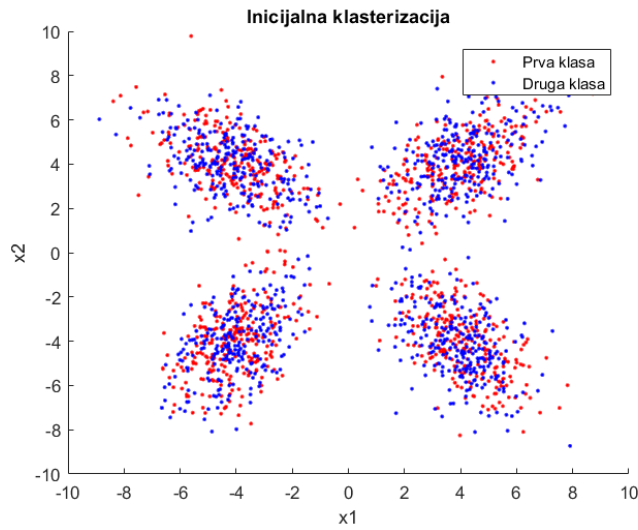
Одсечак кода 21: Функција налажења кластера за дате одбирке

```

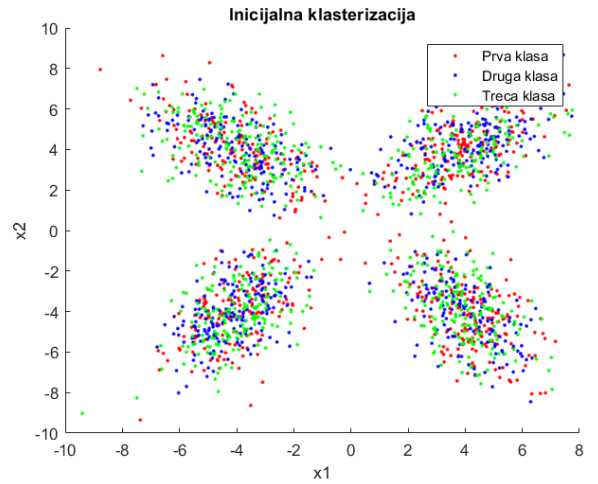
1 function [X1, X2, X3, X4, X5, changed] = reCluster(Y,M, curId)
2 changed = 0;
3 X1 = [];
4 X2 = [];
5 X3 = [];
6 X4 = [];
7 X5 = [];
8 for i = 1 : size(Y, 1)
9     clusterId = findClosestDist(Y(i, :), M);
10    if clusterId ~= curId
11        changed = 1;
12    end
13    switch clusterId
14        case 1
15            X1 = [X1; Y(i, :)];
16        case 2

```

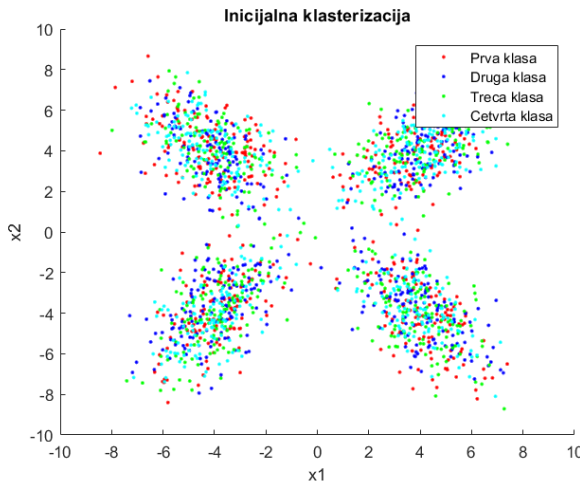
Слика 21: Иницијална кластеризација



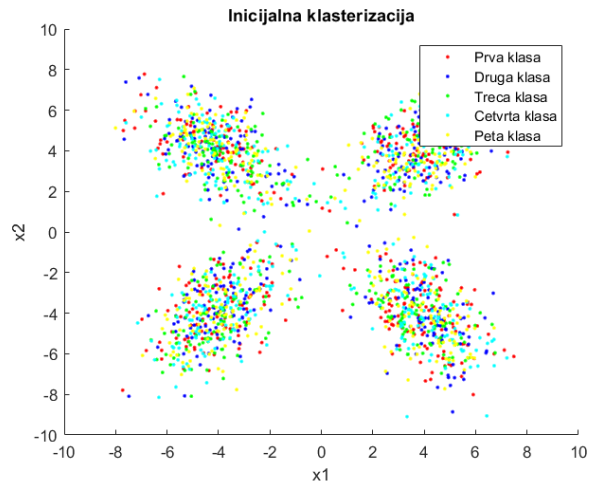
(a) Случај кад је $L = 2$



(b) Случај кад је $L = 3$



(c) Случај кад је $L = 4$



(d) Случај кад је $L = 5$

```

17         X2 = [X2; Y(i, :) ] ;
18         case 3
19             X3 = [X3; Y(i, :) ] ;
20         case 4
21             X4 = [X4; Y(i, :) ] ;
22         case 5
23             X5 = [X5; Y(i, :) ] ;
24         end
25     end
26 end

```

Одсечак кода 22: Функција налажења кластера за дати одбирак

```

1 function [clusterId] = findClosestDist(elem, center)
2 for i = 1 : distLen
3     curCenter = center(i, :) ;
4     curDist = sqrt(sum((elem - curCenter).^2));
5     if (curDist < distMin)

```

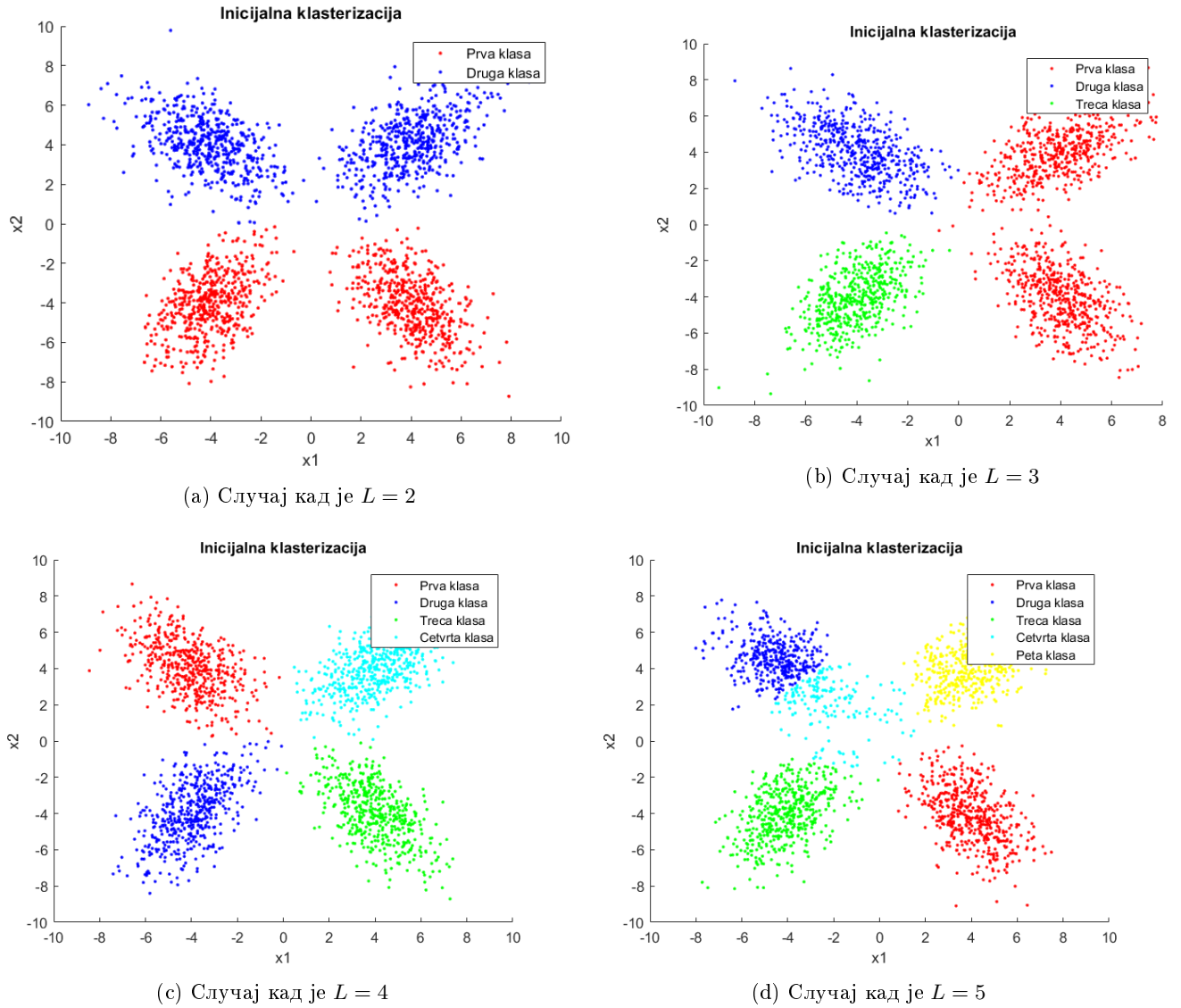
```

6         distMin = curDist;
7         distMinIdx = i;
8     end
9 end
10 clusterId = distMinIdx;

```

Крајњи резултат извршавања се може видети на сликама [22a](#), [22b](#), [22c](#), [22d](#).

Слика 22: Крајња кластеризација



4.1.3 Анализа резултата алгорита

Алгоритам није превише осетљив на почетну кластеризацију, али не гарантује да ће се достићи оптимална кластеризација (у зависности од почетне расподеле у случају $L = 4$ добије се изглед као [22b](#)). Просечан број итерација за иницијално претпостављене величине кластера $L = 2, L = 3, L = 4, L = 5$ варира, али усредњене вредности су редом 3, 7; 5, 93; 6, 42; 9, 13. Ово је и у неку руку јасно јер су ове 4 класе на неки начин подељена на 2 скупа одбирака, односно 4, стога просечан број итерација расте драстично са случаја $L = 2$ на случај $L = 3$ и са $L = 4$ на $L = 5$, док са $L = 3$ на $L = 4$ не расте. Ако је међутим познат број неких тачака које припадају одређеним кластерима, ово може да убрза рад алгоритма значајно. У случају $L = 4$ добија се да за априорно вероватноћу

познавања броја одбирака одређеног кластера 0,1 просечан број одбирака је 3,7 а за случај 0,25 спада на 3, а за 0,5 чак је спао на 2,3, што је и логично јер су кластери већ у неку руку одређени.

Уколико априорно није познат број кластера неопходно је одрадит минимизацију коришћењем неког критеријума. Мало измењена верзија Акаикевог критеријума $AIC(L) = J(L) + \ln(L)$ где је $J(L)$ минимизирана критеријумска функција за конкретан број кластера. Редом, за различит број класа 2, 3, 4, 5 су добијене 20, 10; 12.8; 5, 194; 5, 68. Што и каже оно што видимо на слици, да постоје 4 класе.

4.2 Метод квадратне декомпозиције

4.2.1 Генерисање одбирака

За почетак генерисаћемо одбирке двеју дводимензионалних класа чије су функције густине вероватноће у облику бимодалних гаусовских расподела:

$$f_1(X) = P_{11}N(M_{11}, \Sigma_{11}) + P_{12}N(M_{12}, \Sigma_{12})$$

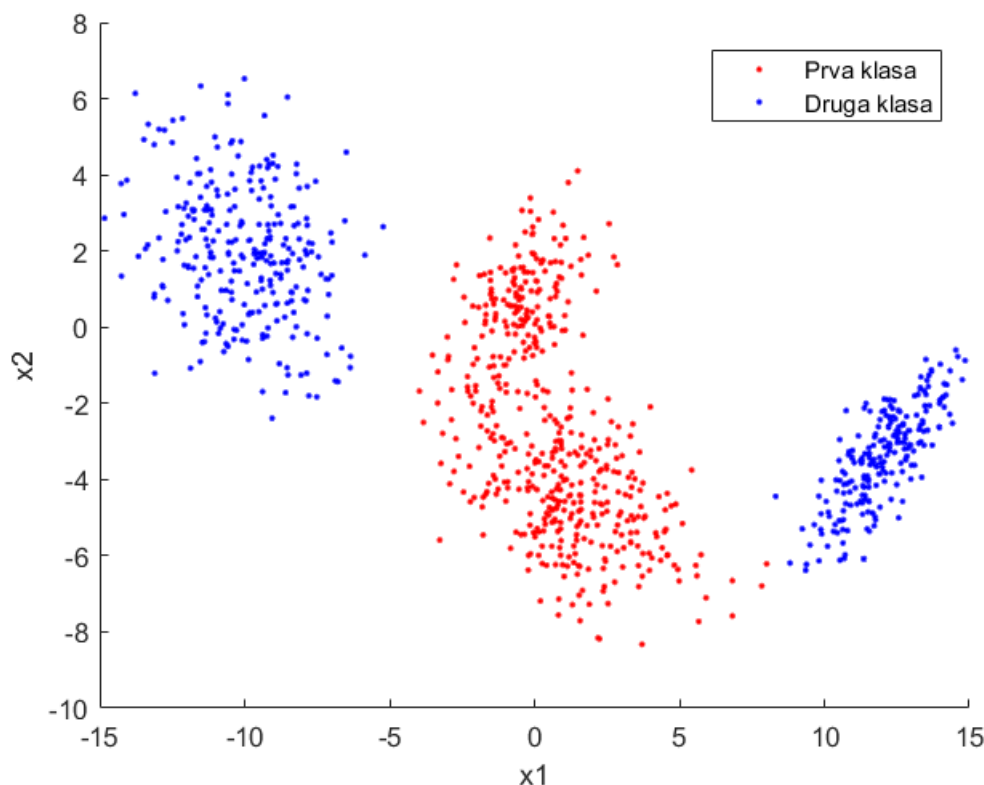
$$f_2(X) = P_{21}N(M_{21}, \Sigma_{21}) + P_{22}N(M_{22}, \Sigma_{22}),$$

при чему вероватноће, средње вредности и коваријационе матрице имају следеће вредности:

$$P_{11} = 0,6, M_{11} = \begin{bmatrix} 1,5 \\ -4,5 \end{bmatrix}, S_{11} = \begin{bmatrix} 3,5 & -1 \\ -1 & 2,2 \end{bmatrix}, P_{12} = 0,4, M_{12} = \begin{bmatrix} -0,5 \\ 0,5 \end{bmatrix}, S_{12} = \begin{bmatrix} 1,3 & 0,9 \\ 0,9 & 2 \end{bmatrix}$$

$$P_{21} = 0,45, M_{21} = \begin{bmatrix} 12 \\ -3,5 \end{bmatrix}, S_{21} = \begin{bmatrix} 1,5 & 1,1 \\ 1,1 & 1,5 \end{bmatrix}, P_{22} = 0,55, M_{22} = \begin{bmatrix} -10 \\ 2 \end{bmatrix}, S_{22} = \begin{bmatrix} 3 & -0,8 \\ -0,8 & 3 \end{bmatrix}$$

Добијају се 2 класе одбирака као на слици [23](#).



Слика 23: Одбирци

4.2.2 Алгоритам

Овај метод је сложенији и захтевнији за израчунавање. Али његова предност је што границе изђему кластера не морају да буду само линеарне, него су квадратне криве. Стога се он употребљава кад је потребна кластеризација нелинеарлно сепарабилних класа. Одсечак кода представља **23** имплементацију алгоритма.

На почетку се изабере иницијална кластеризација $\Omega(0)$ И на основу ње процене априорне вероватноће појава $P_i(0)$ вектори средњих вредности класа $M_i(0)$, Као и коваријационе матриц $\Sigma_i(0)$. У свакој итерацији одбирак X_j се придружује класи t за коју је дати израз најмањи:

$$(X_j - M_t(l))^T \Sigma_t^{-1}(l) (X_j - M_t(l)) + \ln|\Sigma_t(l)| - \ln(P_t(l))$$

Логика испитивања минималне вредност критеријумске функције имплементирана је функцијом чији је одсечак кода је **24** и функцијом чији одсечак кода је **25**.

Одсечак кода **23**: Метод квадратне декомпозиције

```

1 function [num_it, critFun] = quadratic(X1, X2, num_points, num_classes,
    apPer, drawFigure)
2
3 chunkRand = floor((1 - apPer) * num_points);
4 chunkApp = num_points - chunkRand;
5
6 Z = [X1(1 : chunkRand, :); X2(1 : chunkRand, :)];
7
8 num_points_all = 2 * chunkRand;
9
10 idx = randperm(num_points_all);
11 chunk_size = floor(num_points_all / num_classes);
12
13 Y1 = Z(idx(1 : chunk_size), :);
14 Y1App = X1(chunkRand + 1 : num_points, :);
15 Y1 = [Y1App; Y1];
16 Y2 = Z(idx(chunk_size + 1:2 * chunk_size), :);
17 Y2App = X2(chunkRand + 1 : num_points, :);
18 Y2 = [Y2App; Y2];
19 Y3 = [];
20 Y3App = [];
21 Y4 = [];
22 Y4App = [];
23 Y5 = [];
24 Y5App = [];
25 if (num_classes > 2)
26     Y3 = Z(idx(2 * chunk_size + 1:3 * chunk_size), :);
27 end
28 if (num_classes > 3)
29     Y4 = Z(idx(3 * chunk_size + 1:4 * chunk_size), :);
30 end
31 if (num_classes > 4)
32     Y5 = Z(idx(4 * chunk_size + 1:5 * chunk_size), :);
33 end
34
35 run = 1;
36
37 it = 0;
38
39 while (run && it < 50)
40     it = it + 1;
41     M1 = mean(Y1);
42     M2 = mean(Y2);

```

```

43 M3 = mean(Y3);
44 M4 = mean(Y4);
45 M5 = mean(Y5);
46
47 S1 = cov(Y1);
48 S2 = cov(Y2);
49 S3 = cov(Y3);
50 S4 = cov(Y4);
51 S5 = cov(Y5);
52 P1 = size(Y1, 1) / (2 * num_points);
53 P2 = size(Y2, 1) / (2 * num_points);
54 P3 = size(Y3, 1) / (2 * num_points);
55 P4 = size(Y4, 1) / (2 * num_points);
56 P5 = size(Y5, 1) / (2 * num_points);
57
58 if (size(M1, 2) == 1)
59     M1 = [realmax realmax];
60     S1 = [realmax realmax; realmax realmax];
61 end
62 if (size(M2, 2) == 1)
63     M2 = [realmax realmax];
64     S2 = [realmax realmax; realmax realmax];
65 end
66 if (size(M3, 2) == 1)
67     M3 = [realmax realmax];
68     S3 = [realmax realmax; realmax realmax];
69 end
70 if (size(M4, 2) == 1)
71     M4 = [realmax realmax];
72     S4 = [realmax realmax; realmax realmax];
73 end
74 if (size(M5, 2) == 1)
75     M5 = [realmax realmax];
76     S5 = [realmax realmax; realmax realmax];
77 end
78 clear X1 X2 X3 X4;
79
80 M = [M1; M2; M3; M4; M5];
81 S = [S1; S2; S3; S4; S5];
82 P = [P1; P2; P3; P4; P5];
83
84 [X11, X21, X31, X41, X51, changed1] = reClusterQuadr(Y1(chunkApp + 1:
85     end, :), M, S, P, 1);
86 [X12, X22, X32, X42, X52, changed2] = reClusterQuadr(Y2(chunkApp + 1:
87     end, :), M, S, P, 2);
88 [X13, X23, X33, X43, X53, changed3] = reClusterQuadr(Y3(chunkApp + 1:
89     end, :), M, S, P, 3);
90 [X14, X24, X34, X44, X54, changed4] = reClusterQuadr(Y4(chunkApp + 1:
91     end, :), M, S, P, 4);
92 [X15, X25, X35, X45, X55, changed5] = reClusterQuadr(Y5(chunkApp + 1:
93     end, :), M, S, P, 5);
94 run = changed1 | changed2 | changed3 | changed4 | changed5;
95 Y1 = [Y1App; X11; X12; X13; X14; X15];
96 Y2 = [Y2App; X21; X22; X23; X24; X25];
97 Y3 = [Y3App; X31; X32; X33; X34; X35];
98 Y4 = [Y4App; X41; X42; X43; X44; X45];
99 Y5 = [Y5App; X51; X52; X53; X54; X55];
100 end

```

96

97 num_it = it;

Одсечак кода 24: Функција проналажења одговарајућег кластера

```

1 function [X1, X2, X3, X4, X5, changed] = reClusterQuadr(Y, M, S, P, curId)
2 changed = 0;
3 X1 = [];
4 X2 = [];
5 X3 = [];
6 X4 = [];
7 X5 = [];
8 for i = 1 : size(Y, 1)
9     clusterId = findClosestDistQuadr(Y(i, :), M, S, P);
10    if clusterId ~= curId
11        changed = 1;
12    end
13    switch clusterId
14        case 1
15            X1 = [X1; Y(i, :)];
16        case 2
17            X2 = [X2; Y(i, :)];
18        case 3
19            X3 = [X3; Y(i, :)];
20        case 4
21            X4 = [X4; Y(i, :)];
22        case 5
23            X5 = [X5; Y(i, :)];
24    end
25 end
26 end

```

Одсечак кода 25: Функција проналажења одговарајућег кластера

```

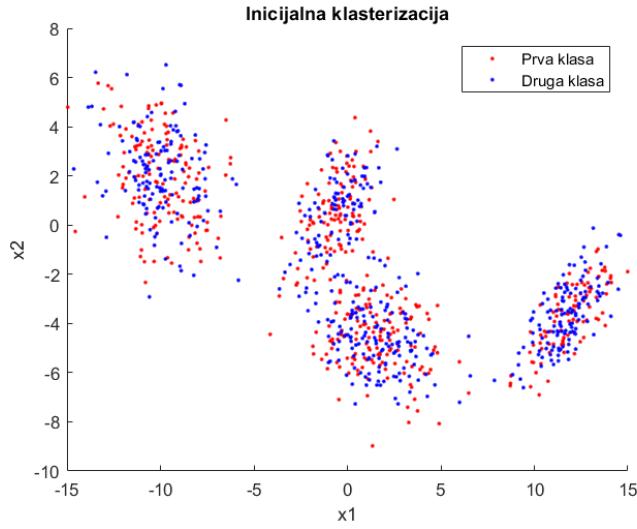
1 function [clusterId] = findClosestDistQuadr(elem, M, S, P)
2 distLen = size(M, 1);
3 distMinIdx = 0;
4 distMin = realmax;
5 for i = 1 : distLen
6     curM = M(i, :);
7     curS = S(2 * i - 1 : 2 * i, :);
8     curP = P(i, :);
9     if (S(2*i - 1, 1) > 10000)
10        j = realmax;
11    else
12        j = (elem - curM) * inv(curS) * (elem - curM)' + log(det(curS)) -
            log(curP);
13    end
14    if (j < distMin)
15        distMin = j;
16        distMinIdx = i;
17    end
18 end
19 clusterId = distMinIdx;

```

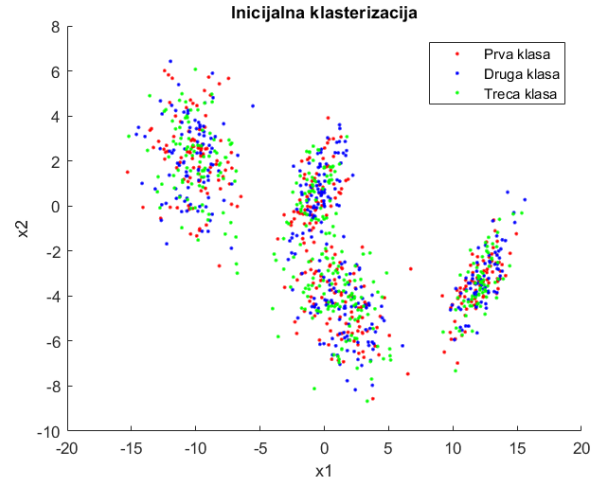
На почетку алгоритма имамо насумично кластеризовање података. На сликама 24a, 24b, 24c, 24d видимо почетне кластеризације за различит број почетних класа.

Крајњи резултат извршавања се може видети на сликама 25a, 25b, 25c, 25d. Дакле алгоритам уопште не кластеризује добро.

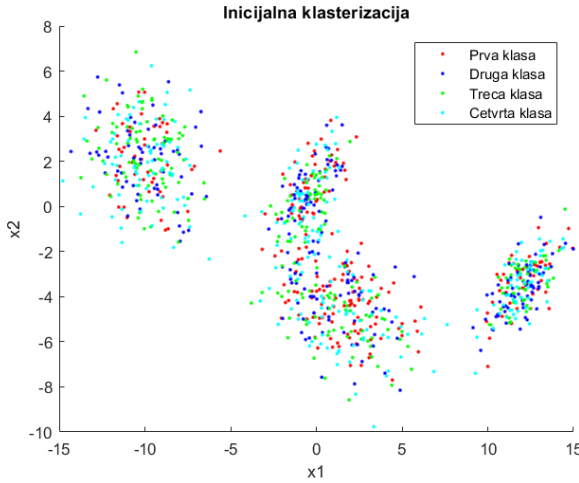
Слика 24: Иницијална кластеризација



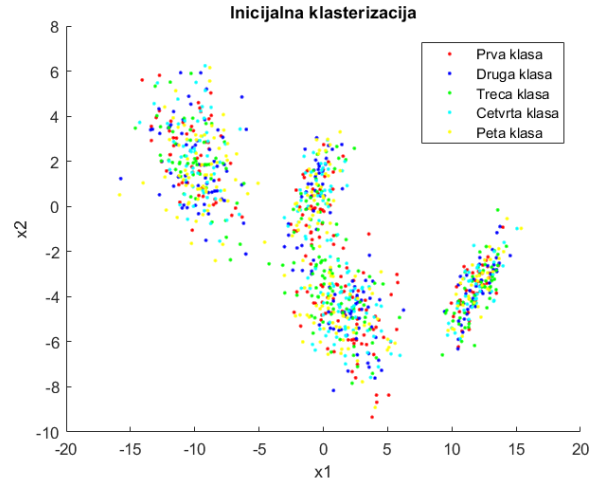
(a) Случај кад је $L = 2$



(b) Случај кад је $L = 3$



(c) Случај кад је $L = 4$



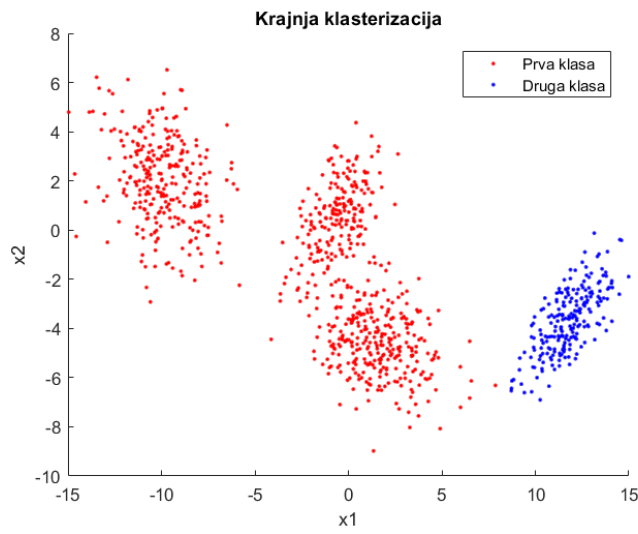
(d) Случај кад је $L = 5$

4.2.3 Анализа резултата

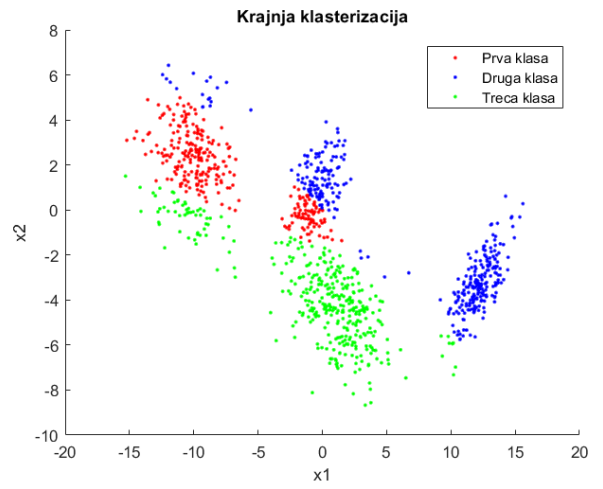
Алгоритам је осетљив на почетну кластеризацију, и не гарантује да ће се достићи оптимална кластеризација. Просечан број итерација за иницијално претпостављене величине кластера $L = 2, L = 3, L = 4, L = 5$ варира, али усредњене вредности су редом 10.2; 17.9; 20.4; 20.1. Ово је и у неку руку јасно јер су ове 2 класе на неки начин подељена на 2 скупа одбирака, стога просечан број итерација расте драстично са случаја $L = 2$ на случај $L = 3$ и релативно, драстично са $L = 3$ на $L = 4$ док са $L = 4$ на $L = 5$ чак и опада. Ако је међутим познат број неких тачака које припадају одређеним кластерима, ово може да убрза рад алгоритма значајно и да се постигне тачност. У случају $L = 4$ добија се да за априорно вероватноћу познавања броја одбирака одређеног кластера 0, 1 просечан број одбирака је 7, 1 а за случај 0, 25 спада на 6 а за 0, 5 чак је спао на 4, 6. Такође, са порастом априорне вероватноће познавања броја одбирака одређеног кластера, добија се пуно већа тачност, што за вероватноће 0, 15 и 0, 25 можемо да видимо на сликама 26а и 26б да је релативно добро одредио кластере.

Уколико априорно није познат број кластера неопходно је одрадит минимизацију коришћењем неког критеријума. Мало измењена верзија Акаикевог критеријума $AIC(L) = J(L) + 3\ln(L)$ где је $J(L)$ минимизирана критеријумска функција за конкретан број кластера. Редом, за различит број класа 2, 3, 4, 5 су добијене 8, 778; 8, 76; 8, 62; 9, 16. Што и каже оно што видимо на слици, да постоје

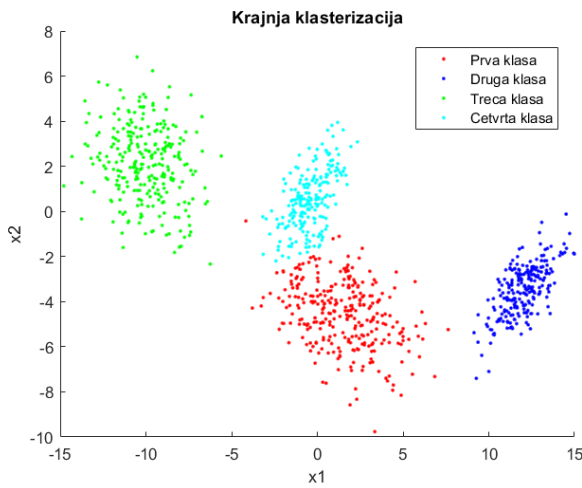
Слика 25: Крајња кластеризација



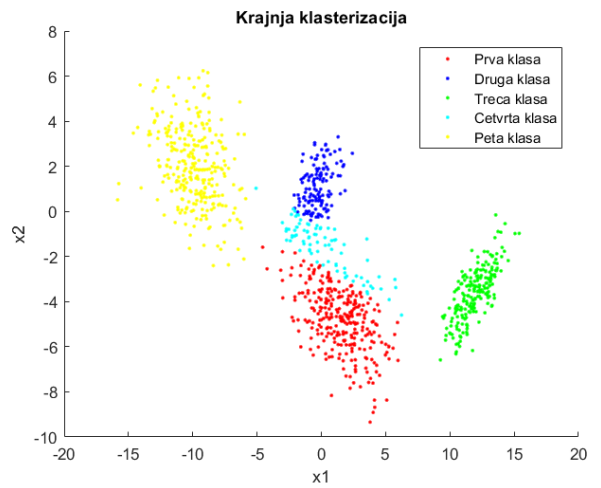
(a) Случај кад је $L = 2$



(b) Случај кад је $L = 3$



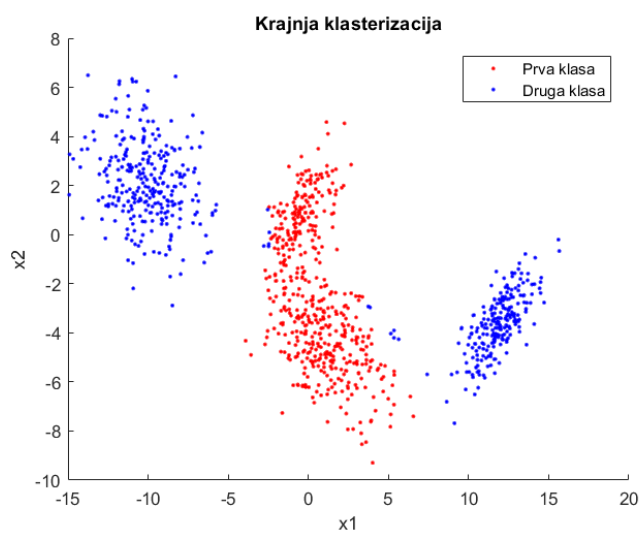
(c) Случај кад је $L = 4$



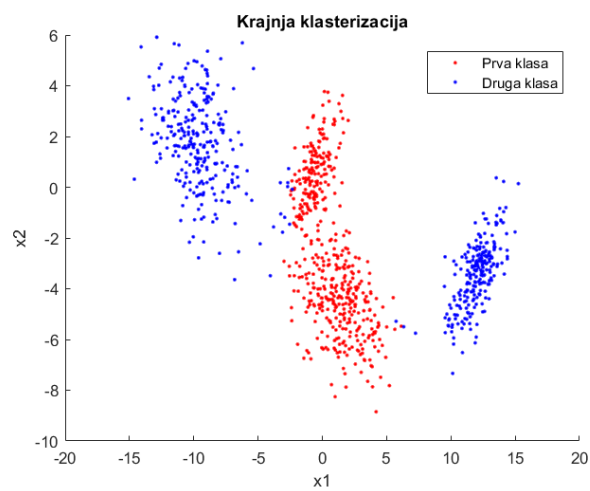
(d) Случај кад је $L = 5$

2 класе, али подједнаке шансе постоје и да постоје 3 и 4.

Слика 26: Крајња кластеризација



(a) Случај кад је априорна вероватноћа 15%



(b) Случај кад је априорна вероватноћа 25%