**Week 1**                                           **Varad Meru**

CS 273a - Introduction to Machine Learning (Winter '15)*          Student # 26648958

Prof. Alex Ihler                                      Due Date: 01/13/2015

# Homework 1[†][‡]

# Problem 1: Matlab & Data Exploration

In this problem, I looked at some constructs provided by Matlab to perform basic statistics and get visualizations of an example data set. I used the downloaded and loaded the *Fisher iris* dataset provided using code given in Listing 1.

Listing 1: Load iris data set and create necessary variables

```
1  % Fetching the dataset and separating it into X and Y.
2  iris=load('data/iris.txt');      % load the text file
3  y = iris(:,end);               % target value is last column
4  X = iris(:,1:end-1);           % features are other columns
5  features = char('Sepal length','Sepal width','Petal length','Petal ...
       width','Species');
6  features_short = char('SL','SW','PL','PW','SP');
7  disp(whos);
8
9  % Changed the shuffleData function to take a seed from the user. If not
10 % provided, then it generates it randomly.
11 [X, y] = shuffleData(X,y,s);
12
13 % Used in some sub-problems
14 X = X(:,1:2);
15 [Xtr,Xte,Ytr,Yte] = splitData(X,y, .75);
```

I have listed the approaches I used for solving the sub-parts of the problem. This question had six sub-parts.

(a) The size of the data matrix X is computed using the `size()` function provided by Matlab.

Listing 2: Displaying the size of input matrix X. output written in the next line of the command

```
1  % Part (a) - Use size(X,2) to get the number of features, and size(X,1) ...
       to get the number of data points.
2  disp(size(X, 1));
3      148
4  disp(size(X, 2));
5       2
```

(b) For plotting the histograms for each feature, I used the `histogram()` function provided by Matlab as shown in Listing 3. You can see the histograms in Figure 1. The saves function is used to store the produced histograms as image files. Figure 2 shows a unified view of all the features of the dataset overlapped.
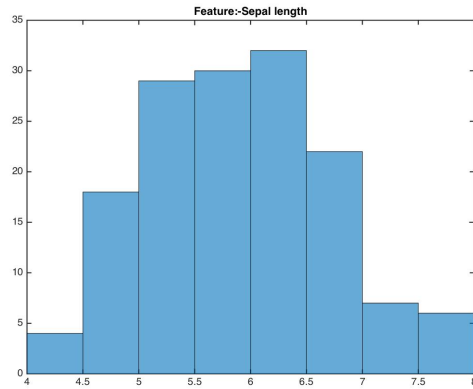
Listing 3: Matlab code for creating histograms.

```
1  %% Part (b)
2  % For each feature, plot a histogram ("hist") of the data values
3  for f=1:size(X, 2);
4      h=figure;
5      %subplot(1,4,f)
```
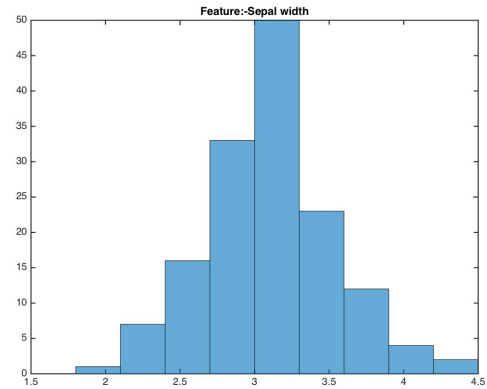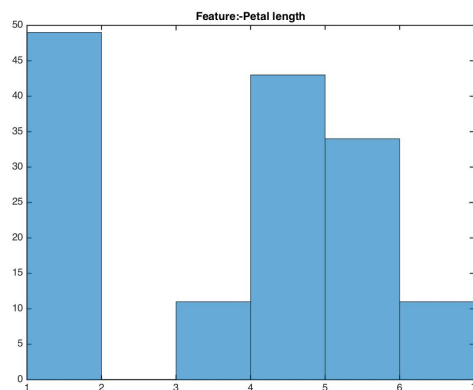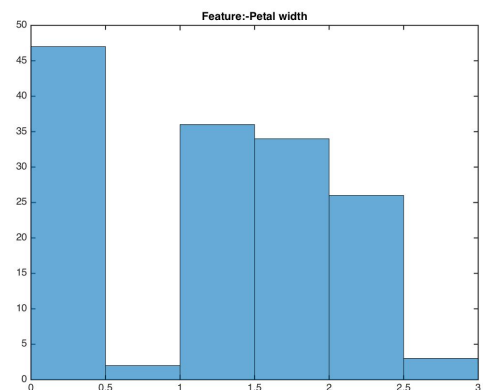
---

(a) Sepal Length

(b) Sepal Width

(c) Petal Length

(d) Petal Width

Figure 1: Histogram of features in the 'Fisher iris' data set

```
6        h1=histogram(X(:,f)); % histogram is preferred over hist - matlab ...
             documentation
7        title(strcat('Feature:- ',features(f,:)))
8        %hold on;
9        saveas(h,strcat('histogram',num2str(f)),'jpg');
10   end;
11   %hold off;
```

Listing 4: Overlapping histograms to show a unified view of 'Fisher iris' dataset.

```
1   % All in one
2   h=figure;
3   for f=1:size(X, 2);
4       histogram(X(:,f));
5       hold on;
6   end;
7   hold off;
8   saveas(h,'histogram5.jpg','jpg');
```

(c) The mean is computed using the `mean()` method provided by matlab (Listing 6)

Listing 5: Mean of features of 'Fisher iris' data set.

```
1   %% Part (c)
2   % Compute the mean of the data points for each feature (mean)
3   disp('Mean');
```
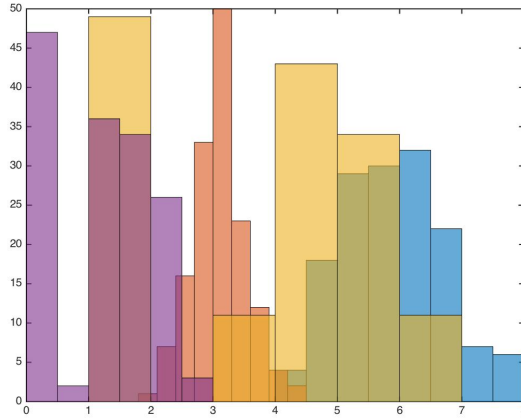
Figure 2: Histogram of features in the 'Fisher iris' data set.

```
4  meanX = mean(X);
5  for f=1:size(X, 2);
6      disp(strcat({'Feature:- '}, features(f,:),{': '} , num2str(meanX(:,f))));
7  end;
8
9  %Output
10 Mean
11     'Feature:- Sepal length: 5.9001'
12     'Feature:- Sepal width: 3.0989'
13     'Feature:- Petal length: 3.8196'
14     'Feature:- Petal width: 1.2526'
```

(d) The mean is computed using the `mean()` method provided by matlab (Listing 6)

Listing 6: Mean of features of 'Fisher iris' data set.

```
1  %% Part (d)
2  % Compute the variance and standard deviation of the data points for each ...
       feature
3  disp('Standard Deviations')
4  stdX = std(X);
5  for f=1:size(X, 2);
6      disp(strcat({'Feature:- '}, features(f,:),{': '} , num2str(stdX(:,f))));
7  end;
8
9  disp('Variances')
10 varX = var(X);
11 for f=1:size(X, 2);
12     disp(strcat({'Feature:- '}, features(f,:),{': '} , num2str(varX(:,f))));
13 end;
14
15 %Output
16 Standard Deviations
17     'Feature:- Sepal length: 0.83623'
18     'Feature:- Sepal width: 0.43777'
19     'Feature:- Petal length: 1.76'
20     'Feature:- Petal width: 0.76135'
21
22 Variances
23     'Feature:- Sepal length: 0.69928'
24     'Feature:- Sepal width: 0.19165'
25     'Feature:- Petal length: 3.0976'
26     'Feature:- Petal width: 0.57965'
```

3

(e) I normalized the dataset using the `bsxfun()`, which applies element-by-element binary operations on the specified function. The functions I used were `@minus` and `@rdivide`. The code described in Listing 9 shows the use. One more approach to achieve the same is seen in Listing 8.

Listing 7: Normalizing the data set

```
1  % bsxfun: applies the element-by-element binary operation specified by ...
      the function handle fun to arrays A and B, with singleton expansion ...
      enabled.
2  normX = bsxfun(rdivide, bsxfun(minus, X, mean(X)), std(X));
```

Listing 8: Normalizing the data set using `repmat()`

```
1  %One More way to do this (more understandable way) -
2  meanXMat = repmat(meanX, size(X,1),1); stdXMat = repmat(stdX, size(X,1),1);
3  norX1 = X - meanXMat; normX = norX1 ./ stdXMat;
```

(f) The plot for pairs of features (1,2), (1,3), and (1,4) is created using a for-loop and the `scatter()`. The colors are set by passing the class values to the scatter function. The code that performs the operation can be seen in Listing 9 and the resulting set of plots generated can be seen in Figure 3

Listing 9: Generating scatter plots for feature pairs - (1,2), (1,3), and (1,4).

```
1  h=figure;
2  i = 1;
3
4  for f=2:size(X, 2);
5      subplot(1,3,i);
6      h1 = scatter(X(:,1), X(:,f), 50, y, 'filled');
7      i = i+1;
8      title(strcat({'X: '},features_short(1,:), {',Y:'}, features_short(f,:)))
9      hold on;
10 end;
11
12 hold off;
13 saveas(h,'plots.jpg','jpg');
```

## Problem 2: kNN predictions

In this problem, I looked at knnClassify class and the features provided by it to classify data in the Iris dataset. The data is shuffle and split into training and test subsets. The question had 2 sub-questions:

(a) I modified the code snippet provided in the homework description and created the classifier. The code after modification is given in Listing 10. The classes for various values of $k$ can be seen in Figure 4. The `Randstream()` function helps in giving a seed to the randomizer and helps in reproducing the results. The `shuffleData()` was modified for accommodating the same with an option of being empty, in which the results would be generated randomly.

Listing 10: Computes the classes using kNNClassify class and generates the scatter plot.

```
1  i = 1;
2  for k=[1, 5, 10, 50];
3      h=figure;
4      %subplot(2,2,i);
5      i = i+1;
6      knn = knnClassify( Xtr, Ytr, k );
7      YteHat = predict( knn, Xte );
8      plotClassify2D( knn, Xtr, Ytr );
9      title(strcat({'K= '},num2str(k)))
```
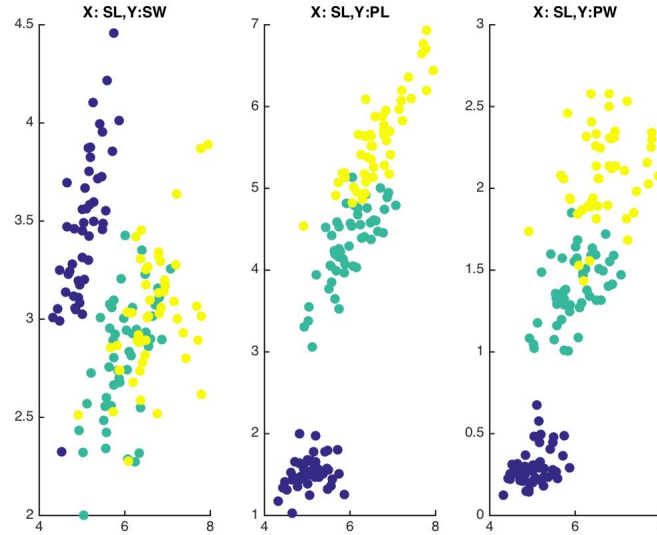
Figure 3: Plots generated with the feature pairs - (1,2), (1,3), and (1,4).

```
10      saveas(h,strcat('2a-',num2str(k),'.jpg'),'jpg');
11      hold on;
12  end;
13  hold off;
```

(b) The error is computed using a different function `errorTrain()`, which compares the expected and actual assigned classes and checks for the matches. `errorTrain()` is described in Listing 11. Listing 12 describes the whole approach to solve the problem given. I used the `semilogx` based graph to present the test-vs-train error plotting for various values of $k$ and got $k = 50$ to be a good value, as also shown by Figure 5. Also to understand the graph in a better way, I changed the value of $k$ to a range of 1 to 100. The results of this are portrayed in Figure 6. It is very similar to the one showed in the class and portrays the overfitting and underfitting zones well.

Listing 11: `errorTrain()` to compute the error in the classes assigned.

```
1  function [error] = errorTrain(Yte, YteHat)
2  % Y = from1ofK(Y1k [,values]) : convert 1-of-K valued Y into discrete ...
       representation
3  %  optional "values" specifies the possible values of Y (default 1..K)
4
5  error = 0;
6  for i=1:size(Yte,1);
7      if(Yte(i) ≠ YteHat(i));
8          error = error + 1;
9      end;
10 end;
11
12 error = error/size(Yte,1);
```
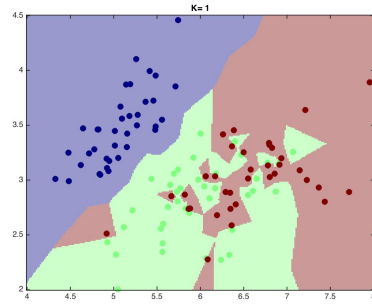
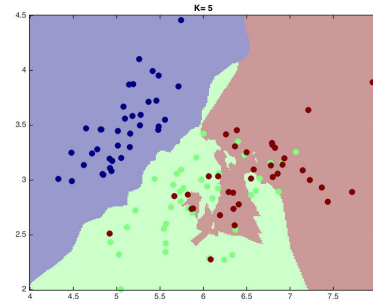Listing 12: Code to generate the `semilogx` classifier error graphs

```
1  %% Part (b)
2  % Training with Xtr and Ytr
3  K=[1,2,5,10,50,100,200, 300];
4  errorsTr = zeros(size(K));
5  errorsTe = zeros(size(K));
6  for i=1:size(K,2);
```
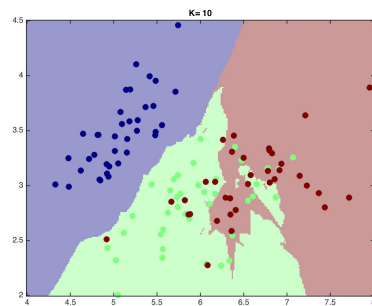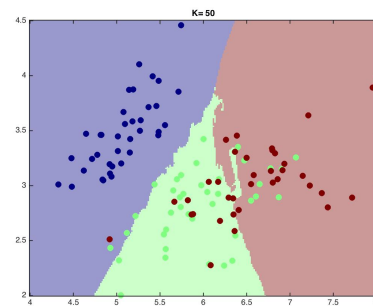
(a) kNN boundaries for K=1

(b) kNN boundaries for K=5

(c) kNN boundaries for K=10

(d) kNN boundaries for K=50

Figure 4: kNN Classifier output for various values of k.

```
7       knn = knnClassify( Xtr, Ytr, K(:,i));
8       YtrHat = predict( knn, Xtr );
9       errorsTr(1,i) = errorTrain(Ytr, YtrHat);
10      YteHat = predict( knn, Xte );
11      errorsTe(1,i) = errorTrain(Yte, YteHat);
12  end;
13
14  h=figure;
15  semilogx(log(K), errorsTr);
16  hold on;
17  semilogx(log(K), errorsTe);
18  hold off;
19
20  disp(K(find(errorsTe == min(errorsTe))));
21      50
```

# Problem 3: Bayes Classifiers

(a) I computed the

Listing 13: Calculating the necessary probabilities for Naive Bayes

```
1  email_data = [0 0 1 1 0 -1;
2      1 1 0 1 0 -1;
3      0 1 1 1 1 -1;
4      1 1 1 1 0 -1;
5      0 1 0 0 0 -1;
6      1 0 1 1 1 1;
7      0 0 1 0 0 1;
8      1 0 0 0 0 1;
9      1 0 1 1 0 1;
```
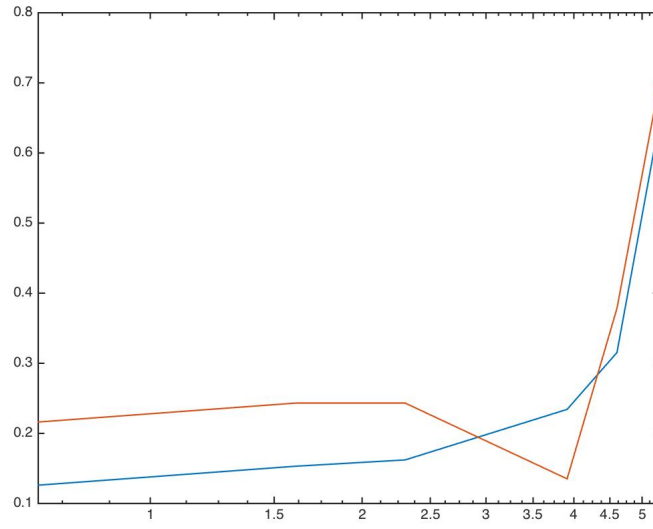
Figure 5: Resulting error rate functions using a semi-log plot with training errors in green and testing error in red for $k = [1,2,5,10,50,100,200]$.
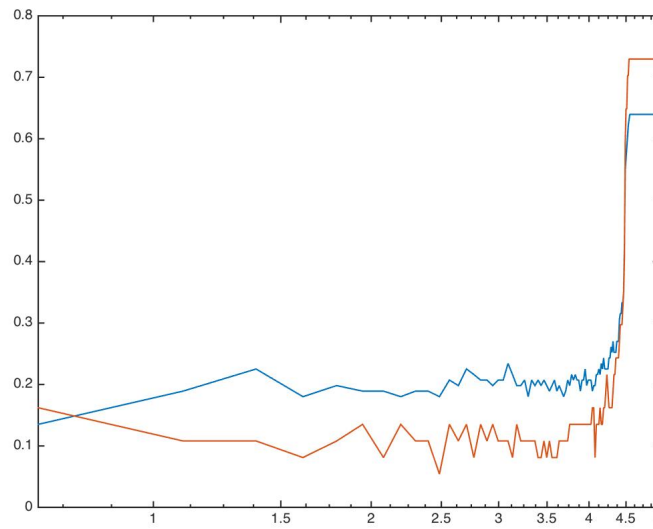


Figure 6: Resulting error rate functions using a semi-log plot with training errors in green and testing error in red for $k = 1{:}100$.

```
10          1 1 1 1 1 -1];
11
12   y = email_data(:,end);
13   X = email_data(:,1:end-1);
14   classes = unique(y);
15
16   [uniqClasses, numUniqClasses] = count_unique(y);
17   class_prob = zeros(size(uniqClasses)); % We'll calculate it later.
18   % class_prob(i, :) = numUniqClasses (i,:) / sum(numUniqClasses);
19
20   class_count = zeros(size(uniqClasses,1), size(X,2));
21
22   % Counting the occurances of Xi for different values of the classes
23   % For each feature -
24   for xi = 1:size(X,2);
25       [uniqX, numUniqX] = count_unique(X(:,xi));
26       for xn = 1:size(X,1);
27           if X(xn,xi) == 1;
28               class_count(find(uniqClasses == y(xn)) , xi) = ...
                     class_count(find(uniqClasses == y(xn)) , xi) + 1;
29           end;
30       end;
31   end;
32
33   % Gets the probabilities of the ones cases, where the data is 1.
34   probabilities_of_one_cases = zeros(size(class_count));
35   for col = 1:size(probabilities_of_one_cases, 2);
36       for row = 1:size(probabilities_of_one_cases, 1);
37           probabilities_of_one_cases(row,:) = class_count(row,:) ./ ...
                 numUniqClasses (row);
38       end;
39   end;
40
41   % Gets the probabilities of the zero cases, where the data is 0.
42   probabilities_of_zero_cases = 1 - probabilities_of_one_cases;
43
44   disp('probabilities_of_one_classes');
45   [probabilities_of_one_cases, uniqClasses]
46   disp('probabilities_of_zero_classes');
47   [probabilities_of_zero_cases, uniqClasses]
48
49   ---------------------------------------
50   probabilities_of_one_classes
51   ans =
52       0.5000    0.8333    0.6667    0.8333    0.3333   -1.0000
53       0.7500         0    0.7500    0.5000    0.2500    1.0000
54
55   probabilities_of_zero_classes
56   ans =
57       0.5000    0.1667    0.3333    0.1667    0.6667   -1.0000
58       0.2500    1.0000    0.2500    0.5000    0.7500    1.0000
```

(b) For computing the class using the Naive Bayes classifier, we used the basic code used in the previous question (Listing 13) and computing it by calculating the chance of the class being one of the given classes (+1 and -1 in our example). The code to compute the associated class is given in Listing 14.

- For input = [0, 0, 0, 0, 0]: The computed class of the email is 1 with a probability of 0.8351. The system would recommended this email to be read.

- For input = [1, 1, 0, 1, 0]: The computed class of the email is -1 with a probability of 1. The system would recommended this email to not be read.

Listing 14: For some input

```matlab
1  %
2  % P(y|X) = p(X|y) p(y) / P(X)
3  % P(y = -1 |X) = p(X|y = -1) p(y = -1) / P(X)
4
5  % P(y = -1 |X) =
6  % p(X|y = -1) p(y = -1)
7  % ------------------------------------
8  %(p(X|y=-1) p(y=-1) + p(X|y=1) p(y=1))
9
10 % P(y = 1 |X) =
11 % p(X|y = 1) p(y = 1)
12 % ------------------------------------
13 %(p(X|y=-1) p(y=-1) + p(X|y=1) p(y=1))
14
15 input = [0, 0, 0, 0, 0];
16 % input = [1, 1, 0, 1, 0];
17
18 % finding the classes for -1
19 productClass1 = 1;
20 for i = 1:size(input,2);
21     if input(:,i) == 1;
22         productClass1 = productClass1 *  probabilities_of_one_cases(1,i);
23     else
24         productClass1 = productClass1 *  probabilities_of_zero_cases(1,i);
25     end;
26 end;
27
28 % finding the classes for 1
29 productClass2 = 1;
30 for i = 1:size(input,2);
31     if input(:,i) == 1;
32         productClass2 = productClass2 *  probabilities_of_one_cases(2,i);
33     else
34         productClass2 = productClass2 *  probabilities_of_zero_cases(2,i);
35     end;
36 end;
37
38 probClass1 = productClass1 * class_prob(1,:) / (productClass1 * ...
        class_prob(1,:) + productClass2 * class_prob(2,:));
39 probClass2 = productClass2 * class_prob(2,:) / (productClass1 * ...
        class_prob(1,:) + productClass2 * class_prob(2,:));
```

(c) The posterior probability of y = +1 given X = [1, 1, 0, 1, 0] is 0. This corroborates with the fact that X was also a data point which the data was trained on, and its class with -1.

(d) Using a Bayes classifier for such kind of data where there is high cost associated with false positives and false negatives can have high impact. Also, for very high number of features, not assuming independent features can be computationally very expensive as it would be a combinatorial explosion.

## Problem 4: Gaussian Bayes Classifiers

I used some of the features provided in the `gaussBayesClassify` class and was able to create

(a) As before, I split data vertically to get the first 2 features of Iris dataset. After getting this data, I computed the number of unique classes and each ones count. This was done using the `count_unique()` function. This is previously used in Problem 3. Once I got the unique classes, I fetched all the indices from the training set which matched the class value

and found the mean and covariance of the temporary matrices using the `mean()` and `cov()` functions of Matlab.

Listing 15: Visualizeing the classifier and the boundaries

```matlab
%% Part(a)
% Splitting by class
[tempUniq, numTemp] = count_unique(Ytr);
mean_matrix = zeros(size(tempUniq,1), size(Xtr,2));
for i = 1:size(tempUniq,1);
    temp = Xtr(find(Ytr==tempUniq(i,:)),:); % All the indexes of
    disp(strcat('Class:',num2str(i)))
    disp('-------------------------------')
    mean_matrix (i,:) = mean(temp); % Mean Stored for future use
    disp('Mean of the class')
    disp(mean_matrix (i,:));
    disp('Covariance matrix')
    disp(cov(temp)); % Creating and displaying the Covariance matrix of ...
        the class data.
end;
```

Listing 16: Output of the code from Listing 18

```
Class:1
-------------------------------
Mean of the class
    5.0721    3.4765
Covariance matrix
    0.1194    0.0880
    0.0880    0.1159

Class:2
-------------------------------
Mean of the class
    5.9581    2.8093
Covariance matrix
    0.2967    0.0939
    0.0939    0.1138

Class:3
-------------------------------
Mean of the class
    6.5659    3.0264
Covariance matrix
    0.3730    0.1039
    0.1039    0.1039
```

(b) The scatter plot generated by Listing 18 is given in Figure 7.

Listing 17: Visualizeing the classifier and the boundaries

```matlab
h = figure;
scatter(Xtr(:,1), Xtr(:,2), 50, Ytr, 'filled');
saveas(h,'scatter-bayes.jpg','jpg');
```

(c) The scatter plot along with the contours based on the covariance matrix and the mean of data is generated by Listing 16 is given in Figure 8.

Listing 18: Visualizeing the classifier and the boundaries

```matlab
% Evaluate each point of feature space and predict the class
[tempUniq, numTemp] = count_unique(Ytr);
mean_matrix = zeros(size(tempUniq,1), size(Xtr,2));
color = {'blue','green','yellow'};
figure;
```
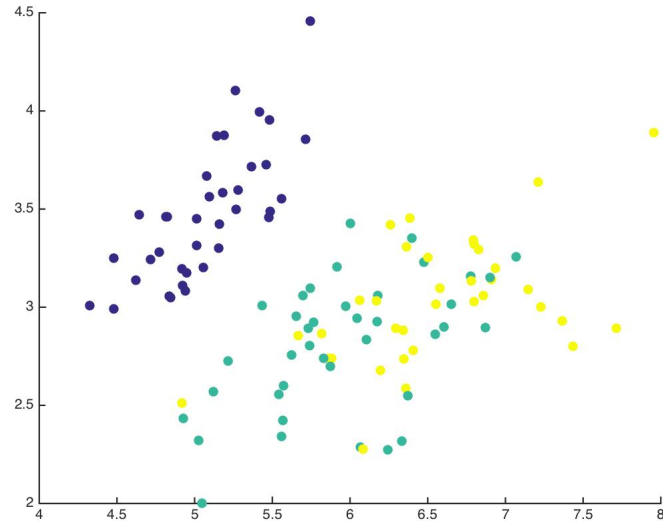
Figure 7: Scatter plot of Iris data set used Gaussian Bayes Learner.

```
6   scatter(Xtr(:,1),Xtr(:,2), 50, Ytr, 'filled');
7   hold on;
8   for i = 1:size(tempUniq,1);
9       temp = Xtr(find(Ytr==tempUniq(i,:)),:); % All the indexes of
10      mean_matrix (i,:) = mean(temp); % Mean Stored for future use
11      plotGauss2D(mean(temp),cov(temp), color{i});
12      hold on;
13  end;
14  hold off;
```

(d) The code given in Listing 19 generates the plot given in Figure 9.

Listing 19: Visualizeing the classifier and the boundaries

```
1   %% Part(d)
2   h = figure;
3   bc = gaussBayesClassify( Xtr, Ytr );
4   plotClassify2D(bc, Xtr, Ytr);
5   saveas(h,'bayes-2.jpg','jpg');
```

(e) The code given in Listing 20 computes the training error rate and the testing error rate. The code uses errorTrainBayes() to compute the number of misclassifications and the rate. It is a variant of Listing 11.

Listing 20: Calculating the error rate.

```
1   %% Part(e)
2   % Training error
3   bc = gaussBayesClassify( Xtr, Ytr );
4   YtrHat = predict( bc, Xtr );
5   [errorsCount, errorRate] = errorTrainBayes(Ytr, YtrHat);
6   >> [errorsCount, errorRate]
7   ans =
8      23.0000    0.2072
9
10  % Testing error
11  bcTe = gaussBayesClassify( Xte, Yte );
12  YteHat = predict( bcTe, Xte );
13  [errorsCount, errorRate] = errorTrainBayes(Yte, YteHat);
```
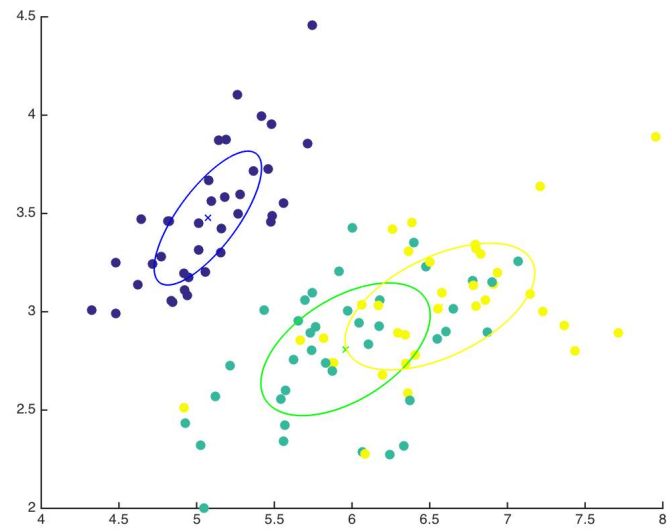
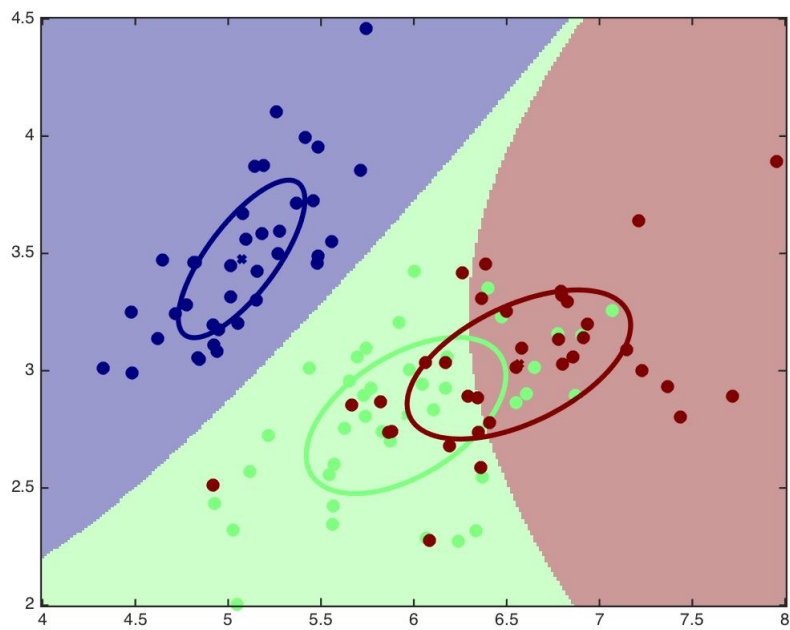Figure 8: Scatter plot with counters of Iris data set.



Figure 9: Bayes classifier in action: Generating the decision boundaries using Bayes classifier.

```
14  >> [errorsCount, errorRate]
15  ans =
16      9.0000    0.2432
```

(f) When training the data with all the features of the data set, I saw tremendous improvement in the working of the classifier. It predicted all the test data results with 100% accuracy. The code of the same is given in the Listing 21. I've put the output within the code itself.

Listing 21: Bayesian classifier for all the 4 features of the Iris data set.

```
1   %% Part(f)
2   bc = gaussBayesClassify( Xtr, Ytr );
3   YtrHat = predict( bc, Xtr );
4   [errorsCountTr, errorRateTr] = errorTrainBayes(Ytr, YtrHat);
5   [errorsCountTr, errorRateTr]
6   % Output
7   ans =
8       2.0000    0.0180
9
10  bcTe = gaussBayesClassify( Xte, Yte );
11  YteHat = predict( bcTe, Xte );
12  [errorsCountTe, errorRateTe] = errorTrainBayes(Yte, YteHat);
13  [errorsCountTe, errorRateTe]
14  % Output
15  ans =
16       0      0
```