

EE 516  
Pattern Recognition  
Prof. Jane Zhang

# Project 5

May 31, 2023

Nicholas BRUNET  
Nathan JAGGERS  
Jordan PERLAS

# Introduction

For this project, we were tasked with working with two separate data sets. The first is the Iris Plant classification data set. This set consists of 50 samples, containing four features, from 3 species of Iris. The other data set is almost 600 samples of benign and malignant breast masses with 30 logged features. Our objective is to parse through the data and extract relevant and robust features to develop and train a Bayes and Naive Bayes classifier for both data sets. We are also tasked with evaluating the classifier and reporting results.

## Iris Plant

### Data

The Iris Data Set used is a popular data set and built into Matlab so it can be called using a single load instruction. The data set is hosted by UC Irvine and can be found at <http://archive.ics.uci.edu/ml/datasets/Iris>.

As stated before it is a dataset of three types of Iris plants (Setosa, Versicolour and Virginica) with 50 samples each, totaling 150 samples. Each sample is described by four features: sepal length, sepal width, petal length and petal width.

### Procedure

To start, we loaded the data into Matlab using the load command and that gave us access to the 4 features mentioned previously and the associated class for each sample. Using this data set we created histograms of the different plants' features (Figures 1 – 4 in the Result section) which allowed us to analyze the separability of the different plant classes. The features that showed high separability were chosen and would later be used when creating the Bayes and Naive Bayes classifier.

LISTING 1: Separating Iris Data into Classes

---

```
load("fisheriris.mat")

setosa_X = meas(1:50,:); % w1
versicolor_X = meas(51:100,:); % w2
virginica_X = meas(101:150,:); % w3

setosa_Y = ones(50, 1);
versicolor_Y = ones(50, 1) * 2;
virginica_Y = ones(50, 1) * 3;
Y = [setosa_Y; versicolor_Y; virginica_Y];
```

---

The separable features were used in training the Bayes Classifier. Training consisted of finding the mean and covariance matrices of the features and using them to develop discriminant functions. The challenge with this data set is that there are not a lot of samples

which makes it difficult to separate the data into a training and test set. To remedy this we employ a Leave-One-Out (LOO) technique which allows us to use the entire data set for training and testing by rotating out samples used for either set. This is accomplished by doing 150 iterations of training (the same amount of samples we have) where each iteration, we leave a single sample out and then test if our classifier correctly distinguishes the sample. To evaluate the effectiveness of the classifier we determine its accuracy when classifying samples in the LOO testing and also construct a confusion matrix to display the performance per class.

LISTING 2: Selection of Test Sample

---

```

for i = 1:150
    % create train/test split
    setosa_X_train = setosa_X;
    versicolor_X_train = versicolor_X;
    virginica_X_train = virginica_X;
    if i <= 50
        x = setosa_X(i,:);
        setosa_X_train(i,:) = [];
    elseif i <= 100
        x = versicolor_X(i-50,:);
        versicolor_X_train(i-50,:) = [];
    else
        x = virginica_X(i-100,:);
        virginica_X_train(i-100,:) = [];
    end
    y = Y(i);
    ...
end

```

---

Lastly we explore the effectiveness of using a Naive Bayes Classifier for this data set. The procedure stays relatively constant when comparing to the Bayes Classifier. The only significant difference is during training. The assumption of the Naive Bayes Classifier is that the features are independent of each other. This means we are able to independently find and multiply the different features probability distributions when determining what class a sample belongs to. So instead of finding a mean and covariance for the whole class, we end up finding them for the features of the different classes. Using this training method, we gather results the same as described for the Bayes Classifier.

## Results

Figures 1 through 4 show the histograms for the dataset. Here we can see that the best features for classification are petal length and petal width because of their high degree of separability.

Table 1 shows the Bayes and Naive Bayes error for the classifier. Figures 5 and 6 show the confusion matrix for the Leave-One-Out process using Bayes and Naive Bayes Classifiers.

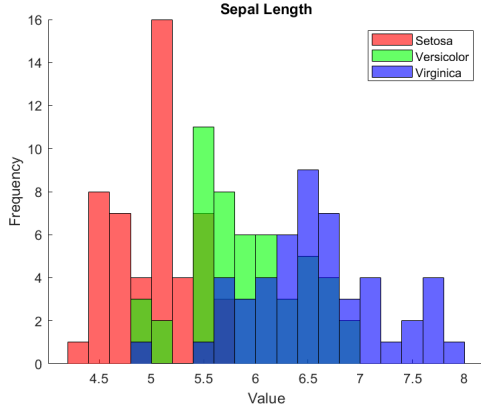


FIGURE 1: Sepal Length

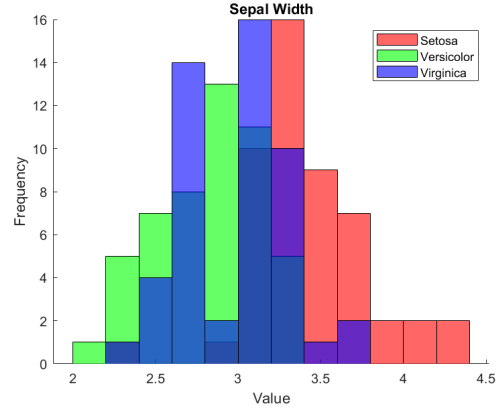


FIGURE 2: Sepal Width

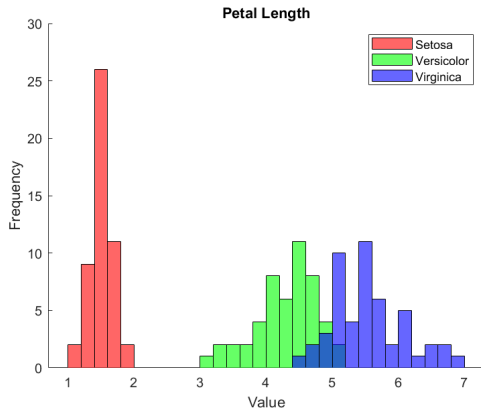


FIGURE 3: Petal Length

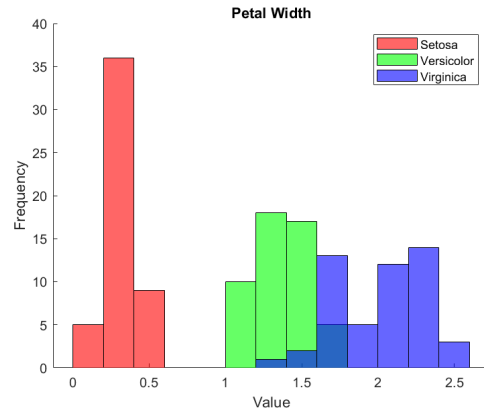


FIGURE 4: Petal Width

TABLE 1: Error Result

Classifier	Empirical Error (%)
Bayes	3.3
Naive Bayes	4

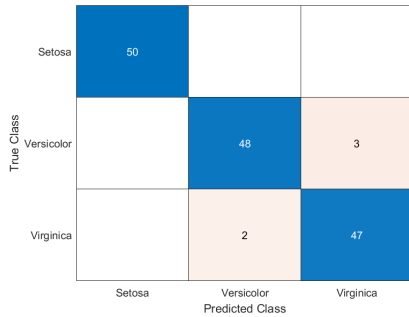


FIGURE 5: Bayes LOO Confusion Matrix

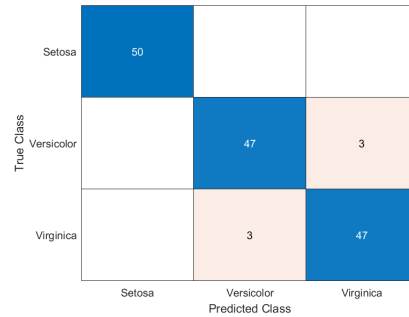


FIGURE 6: Naive Bayes LOO Confusion Matrix

In the confusion matrices, we see perfect classification for Setosa. If we reference the histograms of the features we are using, we can clearly see why. In figures 3 and 4 we see

that the features for Setosa are completely separated from the other two irises. We see some misclassification for Versicolour and Virginica because there is some overlap in the features distribution. When comparing the results between the Bayes and Naive classifier, we see that they perform similarly. The Naive classifier only misclassifies one more sample than the Bayes classifier.

# Breast Cancer

## Data

The data for this section was taken from the Breast Cancer Wisconsin (Diagnostic) data set. It includes 357 benign samples and 212 malignant samples totaling 569 samples, each with 30 features. This data set can be found at <http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/>.

## Procedure

Once the data set was downloaded into Matlab, our first task was to split the data into training and test data sets. From the entire data set, 80% was set aside for training data and the other 20% was stored to be used for testing after the Bayes Classifier was trained. The steps we took to do this are outlined in Listing 3.

LISTING 3: Breast Cancer Initial Processing

---

```
% load data set
wdbc = importdata("wdbc.data");

%separate benign and malignant
pos_data = wdbc.data((wdbc.textdata(:,2))=="M",:);
neg_data = wdbc.data((wdbc.textdata(:,2))=="B",:);
pos80 = round(length(pos_data)*0.8);

neg80 = round(length(neg_data)*0.8);

pos_Train = (pos_data(1:pos80,:));
neg_Train = (neg_data(1:neg80,:));

pos_Test = (pos_data(pos80:length(pos_data),:));
neg_Test = (neg_data(neg80:length(neg_data),:));

train_data = [pos_Train; neg_Train];
test_data = [pos_Test; neg_Test];
```

---

The next step in creating the Bayes classifier is finding the features that will be used. To do this we iterated through all 30 features and created histograms comparing the benign distributions with the malignant distributions. We visually checked the histograms for which features were the most invariant and robust. We also retroactively calculated the errors between the histograms to ensure we were using the most accurate and separate data. This process can be seen in Listing 4. From the 30 features, we chose Feature 23 and Feature 28.

Next, a classifier was created based on the Bayes Decision Rule. To make the process simpler we created multiple functions that can be seen implemented in Listing 5. The first one was tasked with determining the different characteristic values of our data. We input our feature data into the function and were returned means, covariances, priors, as well as

LISTING 4: Breast Cancer Histograms

---

```

%create histograms to compare positive and negative data to find good
%discriminant features
feature_overlap = size(neg_data,2);
for feature=1:30
    figure;
    hold on;
    histogram(pos_data(:,feature));
    histogram(neg_data(:,feature));
    xlabel("Value");
    ylabel("Frequency");
    %figtitle = sprintf("Feature %d");
    title("Feature ",feature);
    hold off;
    feature_overlap(feature) = hist_overlap(pos_data(:,feature), neg_data(:,feature));
end

```

---

a whole data set and classification table. The next function created was to develop our dichotomizer. We input the data set and the characteristic values of our data. The function calculates the terms needed for a Bayes classifier and creates a prediction list depending on the classifier.

LISTING 5: Full Bayes

---

```

%make test set feature vectors
pos_feat_test = [pos_Test(:,feature_1) pos_Test(:,feature_2)];
neg_feat_test = [neg_Test(:,feature_1) neg_Test(:,feature_2)];

%get dataset and augmented matrix from test features
[X_test, Y_test, ~] = feat_details(pos_feat_test,neg_feat_test);

%classify test set
[predict,~] = dichotomizer(X_test, Y_test, mu_pos, cov_pos, prior_pos, mu_neg, cov_neg,
    ↪ prior_neg);

```

---

Next, we tested the accuracy of our dichotomizer, using both confusion maps and ROC curves. As seen in Listing 6 we were able to get the confusion map by using the built-in `confusionmat()` function from Matlab. To get the ROC curve, we swept the priors of the classes and therefore changing our dichotomizer. And for every prior pair we got a new confusion map and calculated the values that would be plotted on the ROC.

Another task we had to complete was the Naive Bayes classifier. The Naive Bayes assumes independence between features making it easier to compute. The Naive Bayes code can be found in Listing 7.

LISTING 6: Confusion Map and ROC Code

---

```

%showing results through confusion matrix (TP, TN, FP, FN)
C = confusionmat(Y_test, predict);
confusionchart(C,["Benign","Malignant"]);

%ROC curve
step = 0.05;
prior_list = 0:step:1;
tp_list = length(prior_list);
fp_list = length(prior_list);

for i = 1:length(prior_list)
%gather predictions for different priors
[predict,~] = dichotomizer(X_test, Y_test, mu_pos, cov_pos, prior_list(i), mu_neg,
↪ cov_neg, 1-prior_list(i));

%collect TP and FP info for each set of priors
C = confusionmat(Y_test, predict);
% tp_list(i) = C(2,2)/sum(C(:));
% fp_list(i) = C(1,2)/sum(C(:));
tp_list(i) = C(2,2)/(C(2, 2)+ C(2, 1));
fp_list(i) = C(1,2)/(C(1, 2)+ C(1, 1));
end

plot(fp_list,tp_list);
xlabel("False Alarms (FP)");
ylabel("Correct Detections (TP)");

```

---

LISTING 7: Naive Bayes Code

---

```

%Naïve Bayes classifier
cov_pos_naive = [cov(pos_feat(:,1));cov(pos_feat(:,2))];
cov_neg_naive = [cov(neg_feat(:,1));cov(neg_feat(:,2))];

[predict,~] = naive_bayes(X_test, Y_test, mu_pos, cov_pos_naive, prior_pos, mu_neg,
↪ cov_neg_naive, prior_neg);

```

---

## Results

From the 30 different features we could use from the breast cancer data set, we chose to work with the features 23, 28, 21, and 8 because they had the most separability between distributions. In other words, they had the least amount of area overlap between the two classes. The histograms for these features can be seen in figures 7 to 10.

With the features selected, we can move on to developing and training a classifier. Once both classifiers were trained, we used the test data that we had set aside from earlier. To test the accuracy of the classifiers, we created ROC curves. To do this, we swept the priors of the classes and created the confusion matrix at every instance. This can be seen in Figures



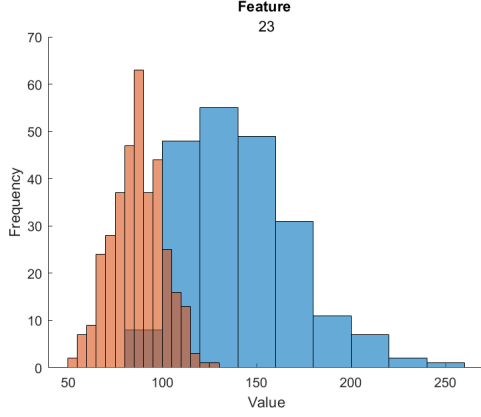


FIGURE 7: Breast Cancer Feature 23

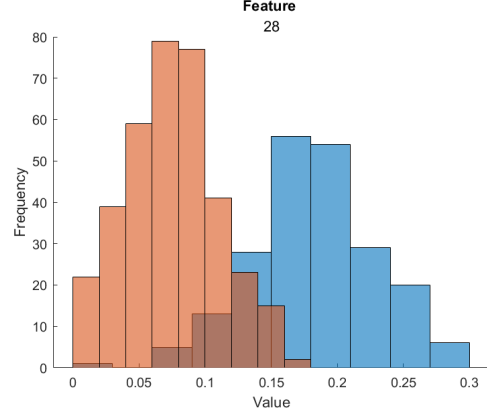


FIGURE 8: Breast Cancer Feature 28

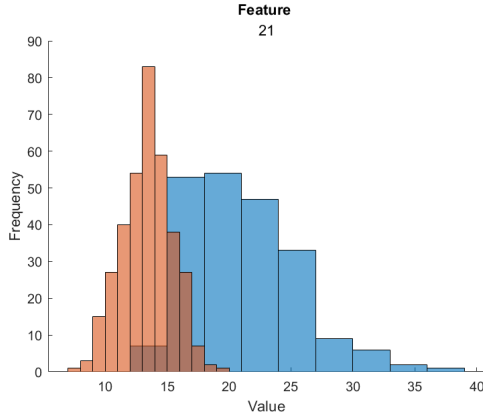


FIGURE 9: Breast Cancer Feature 21

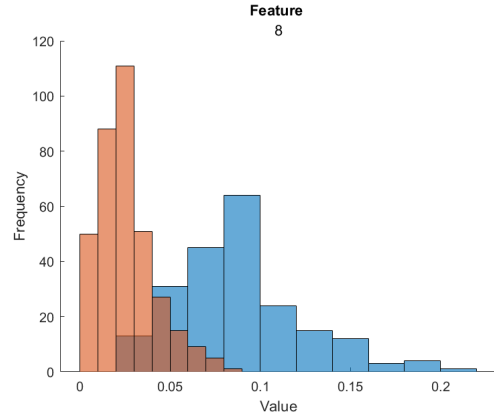


FIGURE 10: Breast Cancer Feature 8

11 to 14.

The results can be found in Table 2.

TABLE 2: Breast Cancer Accuracy Results

Features	Training (%)	Test (%)	Naive Test (%)
28v8	91	90	95
28v23	94	95	88
28v23v8	93	95	91

We can see from the table and the figures that using features 28 and 23 give us the best results. This makes sense because these are the two features with the least overlap area. When comparing the Bayes and the Naive classifier, we see a drop in accuracy. When observing the difference in performance of Bayes and Naive classifiers for features 28 and 8 this is not the case. The Naive classifier in this case sees an increase in accuracy. It is also important to point out that in the table, when increasing the features used to classify data, we do not see an increase upon our best performance. Our Best classifier turns out to be using just soley the two features 28 and 23.

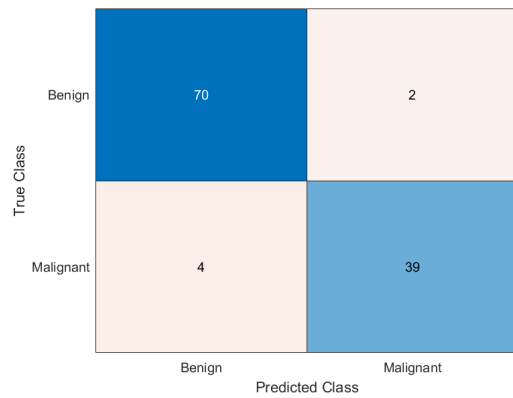


FIGURE 11: Bayes Confusion Matrix: Feature 28 vs 23

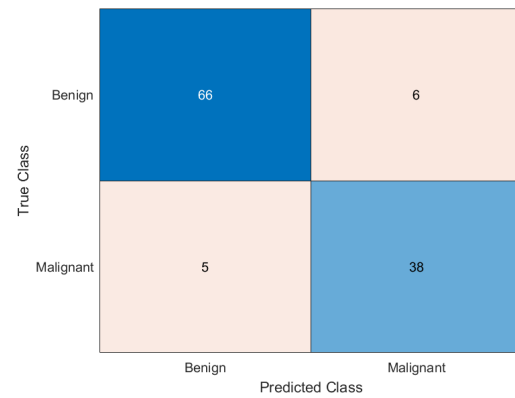


FIGURE 12: Bayes Confusion Matrix: Feature 28 vs 8

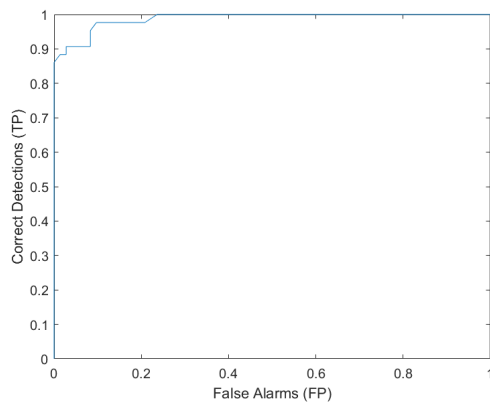


FIGURE 13: Bayes ROC: Feature 28 vs 23

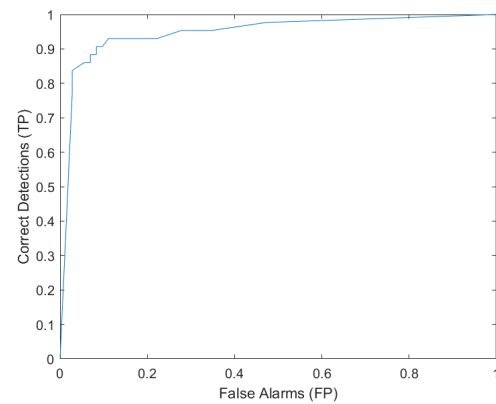


FIGURE 14: Bayes ROC: Feature 28 vs 8

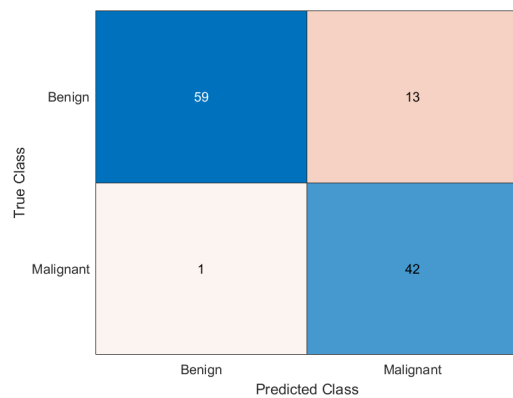


FIGURE 15: Naive Confusion Matrix: Feature 28 vs 23

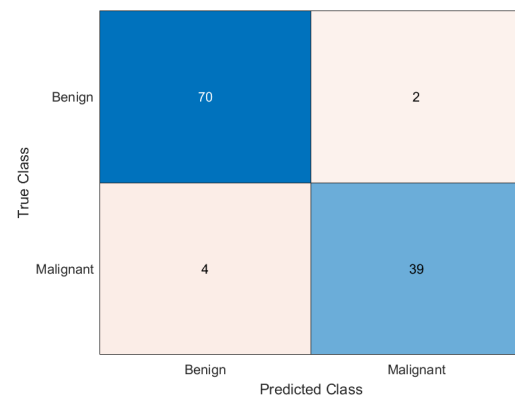


FIGURE 16: Naive Confusion Matrix: Feature 28 vs 8

## Conclusion

**Nathan Jagers**

In this project we explored the different ways of training and testing Bayes and Naive Bayes classifiers for different sized data sets. For the smaller data set we employed the Leave-One-

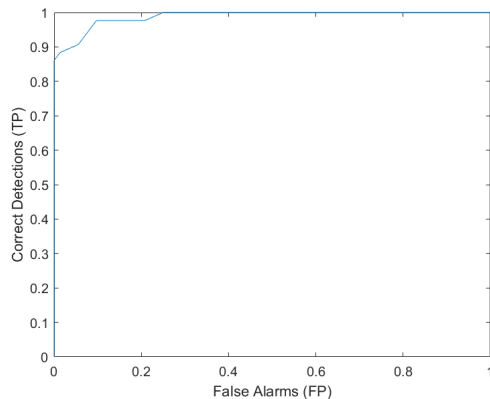


FIGURE 17: Naive ROC: Feature 28 vs 23

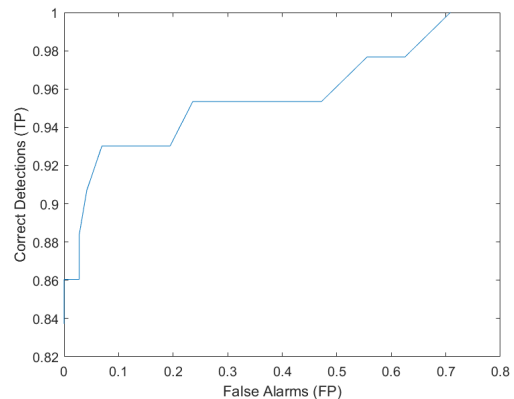


FIGURE 18: Naive ROC: Feature 28 vs 8

Out procedure so we could still have a significant testing and training pool of data. For the larger data set we used a more traditional 80% and 20% split of the data for training and testing respectively. This project was a great way to explore these concepts and the effectiveness of Naive Bayes in different situations.

## Nicholas Brunet

In this lab, we learned several extensions to the simple Bayes classifier explored in the previous lab. For the Iris dataset, we worked with a 3-class classifier with the leave-one-out cross-validation method. We explored the effects of varying the number of parameters on the breast cancer dataset. We also implemented the Naive Bayes classifier on both datasets. My primary work was implementing the Bayes classifier with leave-one-out for the Iris dataset.

## Jordan Perlas

During this lab, I learned a lot more about the different variations of the Bayes Classifier and how they work. The two data sets offer a great insight into all the different ways pattern recognition can be applied in the practical world. The breast cancer data especially showed me the impact of misclassification in a real-world setting. I enjoyed using the confusion maps and seeing how they can be used to create the ROC curve.

# Complete Code Listings

## Part A

---

```
1 %% Project 5
2 %
3 % EE 516 - Pattern Recognition
4 % Spring 2023
5 %
6 % Group 4: Nathan Jaggers, Nicholas Brunet, Jordan Rubio Perlas
7 %
8 % Description: See corresponding document <Can add description later>
9 %% Part 1
10 close all;
11 clear;
12 clc;
13
14 %%
15 % load data set
16 load("fisheriris.mat")
17
18 %%
19 % Features:
20 % 1. sepal length
21 % 2. sepal width
22 % 3. petal length
23 % 4. petal width
24
25 setosa_X = meas(1:50,:); % w1
26 versicolor_X = meas(51:100,:); % w2
27 virginica_X = meas(101:150,:); % w3
28
29 setosa_Y = ones(50, 1);
30 versicolor_Y = ones(50, 1) * 2;
31 virginica_Y = ones(50, 1) * 3;
32 Y = [setosa_Y; versicolor_Y; virginica_Y];
33
34 % BAD
35 hold on
36 histogram(setosa_X(:,1), 'FaceColor', 'r', 'BinWidth', .2);
37 histogram(versicolor_X(:,1), 'FaceColor', 'g', 'BinWidth', .2);
38 histogram(virginica_X(:,1), 'FaceColor', 'b', 'BinWidth', .2);
39 xlabel("Value");
40 ylabel("Frequency");
41 title("Sepal Length");
42 legend("Setosa", "Versicolor", "Virginica");
43 hold off
44
45 % BAD
46 figure;
47 hold on
48 histogram(setosa_X(:,2), 'FaceColor', 'r', 'BinWidth', .2);
```

```

49 histogram(versicolor_X(:,2), 'FaceColor', 'g', 'BinWidth', .2);
50 histogram(virginica_X(:,2), 'FaceColor', 'b', 'BinWidth', .2);
51 xlabel("Value");
52 ylabel("Frequency");
53 title("Sepal Width");
54 legend("Setosa", "Versicolor", "Virginica");
55 hold off
56
57 % GOOD
58 figure;
59 hold on
60 histogram(setosa_X(:,3), 'FaceColor', 'r', 'BinWidth', .2);
61 histogram(versicolor_X(:,3), 'FaceColor', 'g', 'BinWidth', .2);
62 histogram(virginica_X(:,3), 'FaceColor', 'b', 'BinWidth', .2);
63 xlabel("Value");
64 ylabel("Frequency");
65 title("Petal Length");
66 legend("Setosa", "Versicolor", "Virginica");
67 hold off
68
69 % GOOD
70 figure;
71 hold on
72 histogram(setosa_X(:,4), 'FaceColor', 'r', 'BinWidth', .2);
73 histogram(versicolor_X(:,4), 'FaceColor', 'g', 'BinWidth', .2);
74 histogram(virginica_X(:,4), 'FaceColor', 'b', 'BinWidth', .2);
75 xlabel("Value");
76 ylabel("Frequency");
77 title("Petal Width");
78 legend("Setosa", "Versicolor", "Virginica");
79 hold off
80
81 % Prior probabilities
82 w1_P = size(setosa_X,1)/size(meas,1); % 1/3 % each have 50 samples
83 w2_P = size(versicolor_X,1)/size(meas,1); % 1/3 % total samples is 150
84 w3_P = size(virginica_X,1)/size(meas,1); % 1/3 % 1/3 = 50/150
85
86 %% Data cleaning
87
88 % Remove sepal features
89 setosa_X = setosa_X(:,3:4);
90 versicolor_X = versicolor_X(:,3:4);
91 virginica_X = virginica_X(:,3:4);
92
93
94 %% Analysis
95
96 % PROCESS
97 % for s in samples
98 %   leave s out
99 %   train with dataset
100 %   test on s => record score
101 % average scores
102

```

```

103
104
105 C = loo_bayes(Y, setosa_X, w1_P, versicolor_X, w2_P, virginica_X, w3_P);
106 %confusionchart(C);
107 confusionchart(C,["Setosa", "Versicolor", "Virginica",]);
108
109 % xvalues = {'Setosa', 'Versicolor', 'Virginica'};
110 % yvalues = {'Setosa', 'Versicolor', 'Virginica'};
111 % h = heatmap(xvalues, yvalues, conf_mat);
112 % h.XLabel = 'Actual';
113 % h.YLabel = 'Predicted';
114
115 %%
116 C = naive_loo_bayes(Y, setosa_X, w1_P, versicolor_X, w2_P, virginica_X, w3_P);
117 %confusionchart(C);
118 confusionchart(C,["Setosa", "Versicolor", "Virginica",]);
119
120
121 %%
122 function conf_mat = loo_bayes(Y, setosa_X, w1_P, versicolor_X, w2_P, virginica_X, w3_P)
123
124     correct = 0;
125     conf_mat = zeros(3, 3); % confusion matrix
126
127     for i = 1:150
128
129         % create train/test split
130         setosa_X_train = setosa_X;
131         versicolor_X_train = versicolor_X;
132         virginica_X_train = virginica_X;
133         if i <= 50
134             x = setosa_X(i,:);
135             setosa_X_train(i,:) = [];
136         elseif i <= 100
137             x = versicolor_X(i-50,:);
138             versicolor_X_train(i-50,:) = [];
139         else
140             x = virginica_X(i-100,:);
141             virginica_X_train(i-100,:) = [];
142         end
143         y = Y(i);
144
145         x = x'; % make x a column vector
146
147         w1_mean = mean(setosa_X_train)';
148         w2_mean = mean(versicolor_X_train)';
149         w3_mean = mean(virginica_X_train)';
150
151         w1_cov = cov(setosa_X_train);
152         w2_cov = cov(versicolor_X_train);
153         w3_cov = cov(virginica_X_train);
154
155         % Test
156

```

```

157     g1 = g(x, w1_mean, w1_cov, w1_P);
158     g2 = g(x, w2_mean, w2_cov, w2_P);
159     g3 = g(x, w3_mean, w3_cov, w3_P);
160
161     if g1 >= g2 && g1 >= g3
162         result = 1;
163     elseif g2 >= g1 && g2 >= g3
164         result = 2;
165     else
166         result = 3;
167     end
168
169     conf_mat(result, y) = conf_mat(result, y) + 1;
170     if(result == y)
171         correct = correct + 1;
172     end
173
174
175
176 end
177
178 fprintf("Accuracy: %.3f\n", correct / 150);
179
180 end
181
182 function conf_mat = naive_loo_bayes(Y, setosa_X, w1_P, versicolor_X, w2_P, virginica_X,
↵ w3_P)
183
184     correct = 0;
185     conf_mat = zeros(3, 3); % confusion matrix
186
187     for i = 1:150
188
189         % create train/test split
190         setosa_X_train = setosa_X;
191         versicolor_X_train = versicolor_X;
192         virginica_X_train = virginica_X;
193         if i <= 50
194             x = setosa_X(i,:);
195             setosa_X_train(i,:) = [];
196         elseif i <= 100
197             x = versicolor_X(i-50,:);
198             versicolor_X_train(i-50,:) = [];
199         else
200             x = virginica_X(i-100,:);
201             virginica_X_train(i-100,:) = [];
202         end
203         y = Y(i);
204
205         x = x'; % make x a column vector
206
207         %Training
208         w1_mean = mean(setosa_X_train)';
209         w2_mean = mean(versicolor_X_train)';

```

```

210     w3_mean = mean(virginica_X_train)';
211
212     w1_cov = zeros(size(setosa_X_train,2),1);
213     w2_cov = zeros(size(versicolor_X_train,2),1);
214     w3_cov = zeros(size(virginica_X_train,2),1);
215
216     for cols = 1:size(setosa_X_train,2)
217         w1_cov(cols,1) = cov(setosa_X_train(:,cols));
218         w2_cov(cols,1) = cov(versicolor_X_train(:,cols));
219         w3_cov(cols,1) = cov(virginica_X_train(:,cols));
220     end
221
222     % Test
223
224     %multiplication pdfs of of first class
225     pdf_per_class =
↪ (1./(sqrt(2.*pi()).*w1_cov))).*exp(-0.5.*((x-w1_mean).^2)./w1_cov);
226     pdf_pi = prod(pdf_per_class);
227     %discriminant of first class
228     g1 = log(pdf_pi)*w1_P;
229
230     %multiplication pdfs of of second class
231     pdf_per_class =
↪ (1./(sqrt(2.*pi()).*w2_cov))).*exp(-0.5.*((x-w2_mean).^2)./w2_cov);
232     pdf_pi = prod(pdf_per_class);
233     %discriminant of second class
234     g2 = log(pdf_pi)*w2_P;
235
236     %multiplication pdfs of of third class
237     pdf_per_class =
↪ (1./(sqrt(2.*pi()).*w3_cov))).*exp(-0.5.*((x-w3_mean).^2)./w3_cov);
238     pdf_pi = prod(pdf_per_class);
239     %discriminant of second class
240     g3 = log(pdf_pi)*w3_P;
241
242     if g1 >= g2 && g1 >= g3
243         result = 1;
244     elseif g2 >= g1 && g2 >= g3
245         result = 2;
246     else
247         result = 3;
248     end
249
250     conf_mat(result, y) = conf_mat(result, y) + 1;
251     if(result == y)
252         correct = correct + 1;
253     end
254
255 end
256
257 fprintf("Accuracy: %.3f\n", correct / 150);
258
259 end
260

```



```

261 % Calculate parameters once in separate function if this takes too long
262 function result = g(x, mean, cov, P)
263     cov_i = inv(cov);
264     W = -0.5 * cov_i;
265     w = cov_i * mean;
266     w0 = -0.5 * (mean' * cov_i * mean) - 0.5 * log(det(cov)) + log(P);
267     result = x'*W*x + w'*x + w0;
268 end
269

```

---

## Part B

```

1     %% Project 5
2 %
3 % EE 516 - Pattern Recognition
4 % Spring 2023
5 %
6 % Group 4: Nathan Jagers, Nicholas Brunet, Jordan Rubio Perlas
7 %
8 % Description: See corresponding document <Can add description later>
9 %% Part 2
10 close all;
11 clear;
12 clc;
13
14 %%
15 % load data set
16 wdbc = importdata("wdbc.data");
17
18 %separate benign and malignant
19 pos_data = wdbc.data((wdbc.textdata(:,2))=="M",:);
20 neg_data = wdbc.data((wdbc.textdata(:,2))=="B",:);
21
22 %%
23 %create histograms to compare positive and negative data to find good
24 %discriminant features
25 feature_overlap = size(neg_data,2);
26 for feature=1:30
27     figure;
28     hold on;
29     histogram(pos_data(:,feature));
30     histogram(neg_data(:,feature));
31     xlabel("Value");
32     ylabel("Frequency");
33     %figtitle = sprintf("Feature %d");
34     title("Feature ",feature);
35     hold off;
36     feature_overlap(feature) = hist_overlap(pos_data(:,feature), neg_data(:,feature));
37
38 end
39

```

```

40 %%
41 %use overlap to find best features
42 best_features = zeros(5,2);
43 min_Fill = max(feature_overlap);
44 tempMat = feature_overlap;
45 for i=1:size(best_features,1);
46 [best_features(i,1), best_features(i,2)] = min(tempMat);
47 tempMat(best_features(i,2)) = min_Fill;
48 fprintf("Minimums at %d: %f\n",best_features(i,2),best_features(i,1));
49 end
50
51 %%
52 % create data set for training and test
53
54 pos80 = round(length(pos_data)*0.8);
55 neg80 = round(length(neg_data)*0.8);
56
57 pos_Train = (pos_data(1:pos80,:));
58 neg_Train = (neg_data(1:neg80,:));
59
60 pos_Test = (pos_data(pos80:length(pos_data),:));
61 neg_Test = (neg_data(neg80:length(neg_data),:));
62
63 train_data = [pos_Train; neg_Train];
64 test_data = [pos_Test; neg_Test];
65
66 %getting prior probabilities from ratios in training data
67 % pos_prior = pos80/(pos80+neg80);
68 % neg_prior = neg80/(pos80+neg80);
69
70 %%
71 %ones that we like: 28, 23, 21ish, 8, 3
72
73 %trying to use features 8 and 28
74 %make bayesian characteristics mu and sigma
75 feature_1 = 23;
76 feature_2 = 28;
77 feature_3 = 8;
78 % pos_feat = [pos_Train(:,feature_1) pos_Train(:,feature_2)];
79 % neg_feat = [neg_Train(:,feature_1) neg_Train(:,feature_2)];
80
81 pos_feat = [pos_Train(:,feature_1) pos_Train(:,feature_2) pos_Train(:,feature_3)];
82 neg_feat = [neg_Train(:,feature_1) neg_Train(:,feature_2) neg_Train(:,feature_3)];
83
84 %%
85 %training and set up for dichotomizer
86 [X, Y, mu_pos, cov_pos, prior_pos, mu_neg, cov_neg, prior_neg] =
    ↪ feat_details(pos_feat,neg_feat);
87
88 %using dichotomizer
89 [predict,~] = dichotomizer(X, Y, mu_pos, cov_pos, prior_pos, mu_neg, cov_neg, prior_neg);
90
91 %showing results through confusion matrix
92 C = confusionmat(Y, predict);

```

```

93  confusionchart(C,["Benign","Malignant"]);
94
95  %%
96  %using previous training on test set
97
98  %make test set feature vectors
99  % pos_feat_test = [pos_Test(:,feature_1) pos_Test(:,feature_2)];
100 % neg_feat_test = [neg_Test(:,feature_1) neg_Test(:,feature_2)];
101
102 pos_feat_test = [pos_Test(:,feature_1) pos_Test(:,feature_2) pos_Test(:,feature_3)];
103 neg_feat_test = [neg_Test(:,feature_1) neg_Test(:,feature_2) neg_Test(:,feature_3)];
104
105 %get dataset and augmented matrix from test features
106 [X_test, Y_test, ~] = feat_details(pos_feat_test,neg_feat_test);
107
108 %classify test set
109 [predict,~] = dichotomizer(X_test, Y_test, mu_pos, cov_pos, prior_pos, mu_neg, cov_neg,
    ↪ prior_neg);
110
111 %showing results through confusion matrix (TP, TN, FP, FN)
112 C = confusionmat(Y_test, predict);
113 confusionchart(C,["Benign","Malignant"]);
114
115 %%
116 %ROC curve
117 step = 0.05;
118 prior_list = 0:step:1;
119 tp_list = length(prior_list);
120 fp_list = length(prior_list);
121
122 for i = 1:length(prior_list)
123 %gather predictions for different priors
124 [predict,~] = dichotomizer(X_test, Y_test, mu_pos, cov_pos, prior_list(i), mu_neg,
    ↪ cov_neg, 1-prior_list(i));
125
126 %collect TP and FP info for each set of priors
127 C = confusionmat(Y_test, predict);
128 tp_list(i) = C(2,2)/(C(2, 2)+ C(2, 1));
129 fp_list(i) = C(1,2)/(C(1, 2)+ C(1, 1));
130 end
131
132 plot(fp_list,tp_list);
133 xlabel("False Alarms (FP)");
134 ylabel("Correct Detections (TP)");
135
136 %%
137 %Naïve Bayes classifier
138 % cov_pos_naive = [cov(pos_feat(:,1));cov(pos_feat(:,2))];
139 % cov_neg_naive = [cov(neg_feat(:,1));cov(neg_feat(:,2))];
140
141 cov_pos_naive = [cov(pos_feat(:,1));cov(pos_feat(:,2));cov(pos_feat(:,3))];
142 cov_neg_naive = [cov(neg_feat(:,1));cov(neg_feat(:,2));cov(neg_feat(:,3))];
143

```

```

144 [predict,~] = naive_bayes(X_test, Y_test, mu_pos, cov_pos_naive, prior_pos, mu_neg,
    ↪ cov_neg_naive, prior_neg);
145
146 %showing results through confusion matrix (TP, TN, FP, FN)
147 C = confusionmat(Y_test, predict);
148 confusionchart(C,["Benign","Malignant"]);
149
150 %%
151 %ROC curve
152 step = 0.05;
153 prior_list = 0:step:1;
154 tp_list = length(prior_list);
155 fp_list = length(prior_list);
156
157 for i = 1:length(prior_list)
158 %gather predictions for different priors
159 [predict,~] = naive_bayes(X_test, Y_test, mu_pos, cov_pos_naive, prior_list(i), mu_neg,
    ↪ cov_neg_naive, 1-prior_list(i));
160
161 %collect TP and FP info for each set of priors
162 C = confusionmat(Y_test, predict);
163 tp_list(i) = C(2,2)/(C(2, 2)+ C(2, 1));
164 fp_list(i) = C(1,2)/(C(1, 2)+ C(1, 1));
165 end
166
167 plot(fp_list,tp_list);
168 xlabel("False Alarms (FP)");
169 ylabel("Correct Detections (TP)");
170
171 %%
172 %set up for dichotomizer
173 function [dataset, classification, w1_mean, w1_cov, w1_Prior, w2_mean, w2_cov, w2_Prior]
    ↪ = feat_details(w1_features, w2_features)
174 %calculate mean and cov to train dichotomizer
175 w1_mean = mean(w1_features)';
176 w2_mean = mean(w2_features)';
177
178 w1_cov = cov(w1_features);
179 w2_cov = cov(w2_features);
180
181 %priors section
182 w1_samples = size(w1_features,1);
183 w2_samples = size(w2_features,1);
184 w1_Prior = w1_samples/(w1_samples+w2_samples);
185 w2_Prior = w2_samples/(w1_samples+w2_samples);
186
187 %create input matrix and augmented matrix
188 dataset = [w1_features; w2_features]; %input
189 classification = [ones(size(w1_features, 1), 1); zeros(size(w2_features, 1), 1)];
    ↪ %augmented (Y)
190 end
191
192 %general dichotomizer

```

```

193 function [prediction, accuracy] = dichotomizer(dataset, classification, w1_mean, w1_cov,
↪ w1_Prior, w2_mean, w2_cov, w2_Prior)
194     %dataset - input data to be classified
195     %classification - true classification of samples
196     %w1_mean - mean for class 1
197     %w1_cov - covariance for class 2
198     %w1_Prior - prior for class 1
199     %w2_mean - mean for class 2
200     %w2_cov - covariance for class 2
201     %w2_Prior - prior for class 2
202
203     %create class prediction matrix
204     prediction = [zeros(size(classification, 1), 1)];
205
206     %initialize correct counter
207     correct = 0;
208
209     for i = 1:size(dataset, 1)
210         x = dataset(i,:)' ;
211         y = classification(i);
212         g1_result = g(x, w1_mean, w1_cov, w1_Prior);
213         g2_result = g(x, w2_mean, w2_cov, w2_Prior);
214         prediction(i) = g1_result - g2_result > 0;
215         correct = correct + (prediction(i) == y);
216     end
217     accuracy = correct / size(dataset, 1);
218     fprintf("Accuracy: %.2f\n", accuracy);
219 end
220
221 % Calculate parameters once in separate function if this takes too long
222 function result = g(x, mean, cov, P)
223     cov_i = inv(cov);
224     W = -0.5 * cov_i;
225     w = cov_i * mean;
226     w0 = -0.5 * (mean' * cov_i * mean) - 0.5 * log(det(cov)) + log(P);
227     result = x'*W*x + w'*x + w0;
228 end
229
230 %Naïve Bayes classifier
231 function [prediction, accuracy] = naive_bayes(dataset, classification, w1_mean, w1_cov,
↪ w1_Prior, w2_mean, w2_cov, w2_Prior)
232     %dataset - input data to be classified
233     %classification - true classification of samples
234     %w1_mean - mean for class 1
235     %w1_cov - naive covariance for class 2 (one variance for each feature; column
↪ vector)
236     %w1_Prior - prior for class 1
237     %w2_mean - mean for class 2
238     %w2_cov - covariance for class 2 (one variance for each feature; column vector)
239     %w2_Prior - prior for class 2
240
241     %create class prediction matrix
242     prediction = [zeros(size(classification, 1), 1)];
243

```

```

244     %initialize correct counter
245     correct = 0;
246
247     for i = 1:size(dataset, 1)
248         x = dataset(i,:)';
249         y = classification(i);
250         %multiplication pdfs of of first class
251         pdf_per_class =
↪ (1./(sqrt(2.*pi()).*w1_cov))).*exp(-0.5.*((x-w1_mean).^2)./w1_cov);
252         pdf_pi = prod(pdf_per_class);
253         %discriminant of first class
254         g1_result = log(pdf_pi)*w1_Prior;
255
256         %multiplication pdfs of of second class
257         pdf_per_class =
↪ (1./(sqrt(2.*pi()).*w2_cov))).*exp(-0.5.*((x-w2_mean).^2)./w2_cov);
258         pdf_pi = prod(pdf_per_class);
259         %discriminant of second class
260         g2_result = log(pdf_pi)*w2_Prior;
261
262         prediction(i) = g1_result - g2_result > 0;
263         correct = correct + (prediction(i) == y);
264     end
265     accuracy = correct / size(dataset, 1);
266     fprintf("Accuracy: %.2f\n", accuracy);
267 end
268
269 %% Compare Histograms
270 % edited from Dylan Baxter's and Nicholas Brunet's code
271 function [overlap] = hist_overlap(feature1, feature2)
272     numBins = 50;
273     % Extract histogram Information
274     max_val = max([feature1;feature2], [], "all");
275     min_val = min([feature1;feature2], [], "all");
276     bins = linspace(min_val, max_val, numBins+1);
277     hist1 = histcounts(feature1, bins);
278     hist2 = histcounts(feature2, bins);
279
280     % Iterate over histograms and calculate overlap
281     overlap = 0;
282     for m = 1:length(hist1)
283         overlap = overlap + min(hist1(m), hist2(m));
284     end
285     overlap = overlap / sum(hist1, "all");
286 end
287

```

---