# Homework 5 Report

Brian King

## Question 1

Using 5 folds with the Naïve Bayes algorithm obtains the following accuracies for the validation sets:

```
Scores: [93.33333333333, 96.66666666666667, 100.0, 93.33333333333, 93.33333333333]
Mean Accuracy: 95.333%
```

## Question 2

The following functions were modified in order to implement a Multivariate Bayes classifier in place of the Naive Bayes classifier:

**Calculate_class_probabilities function**

```python
145     # Calculate the probabilities of predicting each class for a given row
146     def calculate_class_probabilities(summaries, row):
147         del(row[-1])
148         total_rows = sum([summaries[label][2] for label in summaries])
149         probabilities = dict()
150         for class_value, class_summaries in summaries.items():
151             probabilities[class_value] = summaries[class_value][2]/float(total_rows)
152             mean, cov, _ = class_summaries
153             probabilities[class_value] *= calculate_probability(row, mean, cov)  # discriminant function
154         return probabilities
```

**Calculate_probability function**

```python
136     # Calculate the Gaussian probability distribution function for x
137     def calculate_probability(x, mean, cov):
138         inv_cov = np.linalg.inv(cov)
139         multiply1 = np.dot(np.transpose(x-mean), inv_cov)
140         multiply2 = np.dot(multiply1, (x-mean))
141         exponent = exp(-.5*multiply2)
142         return (1 / ((2 * pi)**(4/2) * (np.linalg.det(cov))**.5)) * exponent
```

**Summarize_dataset function**

```python
114     # Calculate the mean, cov and count for each dataset
115     def summarize_dataset(dataset):
116         columns = [column for column in zip(*dataset)]
117         del(columns[-1])     # remove classification from array
118         mean = means(columns)
119         length = len(columns[1])
120         cov = covariance_creator(columns, mean)
121         summaries = [mean, cov, length]
122         return summaries
```

**Covariance function**

```python
103     # Calculate covariance arrays
104     def covariance_creator(datalist, mean):
105         covariances = np.zeros((4, 4))
106         for i in range(len(datalist)):
107             for j in range(len(datalist)):
108                 for k in range(len(datalist[1])):
109                     r = ((datalist[i][k] - mean[i]) * (datalist[j][k] - mean[j]))
110                     covariances[i][j] = covariances[i][j] + r/len(datalist[1])
111         return covariances
```

**Mean function**

```python
95      # Calculate mean arrays
96      def means(datalist):
97          mean = np.array([])
98          for i in range(len(datalist)):
99              this_mean = sum(datalist[i])/len(datalist[i])
100             mean = np.append(mean, this_mean)
101         return mean
```

# Question 3

Utilizing the aforementioned Multivariate Bayes classifier, the following accuracies were obtained.

```
Scores: [96.66666666666667, 96.66666666666667, 100.0, 100.0, 93.33333333333333]
Mean Accuracy: 97.333%
```

This is an improvement of 2% over the Naïve Bayes classifier. However, with large amounts of data (greater than the previous 150 data points used in the report), it is more efficient to use the Naïve Bayes classifier as it computationally less intensive.