```
1.1.1
    Machine Language Programming assignment 2
 3
     see the assignment description in the README.md file
 4
 5
 6
     import numpy as np
 7
     import pandas as pd
9
    print("Gaussian Naive Bayes Algorithm")
10
    print ("Machine Language Programming assignment 2")
11
   print("Author: Larry Bilodeau")
12
   print()
13
    print()
14
15
    # define a function to process the feature labels.
16 def read and truncate file (file path):
         feature labels = []
17
18
         with open (file path, 'r') as file:
             for line in file:
19
20
                 feature label = line.split(':')[0]
21
                 feature labels.append(feature label)
22
         return feature labels
23
24
   # Load the raw data
25
    rawdata = pd.read csv('...\datasets\spambase.data', header=None)
    print("Data set: spambase.data")
26
27
    print("rawdata shape", rawdata.shape)
28
    # load the feature labels for the rawdata
29
    feature_labels = read_and_truncate_file('..\\datasets\\feature_labels.txt')
30
31
32
    # step 1: split the data into training and testing sets
33
34
    # Define the labels
35
    labels = rawdata.iloc[:, -1]
36
37
    # Split the data into training and testing sets
38
    # Select 2300 random rows with 40% labels 1 and 60% labels 0
39
     test data = rawdata.sample(n=2300, weights=labels.map(\{0: 0.6, 1: 0.4\}), random state=42)
40
    print("test data shape", test data.shape)
41
42
    # Define the labels for test data
43 test targets = test data.iloc[:, -1]
44
    test data = test data.iloc[:, :-1]
45
46
    # Remove the selected test data rows from rawdata to get training data
47
    training data = rawdata.drop(test data.index)
48
    # Define the labels for training_data
49
    training targets = training data.iloc[:, -1]
50
     # trim off the targets from the training data
51
    training data = training data.iloc[:, :-1]
52
    print("training data shape", training data.shape)
53
54
55
56
57
     # step 2: create probalistic models for each class
58
59
60
     # Determine each of the hypothesis ratios for each hypotesis (H), based on the test data
61
     \# P(H|D) = P(D|H) * P(H) / P(D)
    hypotheses, frequencies = np.unique(training targets, return counts=True)
63
    hypothesis ratios = frequencies / len(test targets)
64
65
     # Calculate the prior probabilities P(H)
     \# P(H) = P(D|H) * P(H) / P(D)
66
67
```

```
68
      for h, f, r in zip(hypotheses, frequencies, hypothesis ratios):
 69
          print()
 70
          print(f"{'Hypothesis':<15} {'Frequency':<15} {'Ratio':<15}")</pre>
 71
          print(f"{h:<15} {f:<15} {r:<15.2f}")</pre>
 72
          print()
 73
          print("Prior Probability of each hypothesis")
 74
          print("hyposis: spam, P(h1):", hypothesis ratios[1])
 75
          print("hyposis: not spam, P(h0):", hypothesis_ratios[0])
 76
          print(" P(+|h1):", 1 - hypothesis ratios[1], "P(-|h1):", hypothesis ratios[1])
 77
          print(" P(+|h0):", 1 - hypothesis ratios[0], "P(-|h0):", hypothesis ratios[0])
 78
 79
      # Filter rows for each unique value in training targets
 80
      grouped data = training data.groupby(training targets)
 81
 82
      # Calculate the mean and standard deviation for each group
 83
 84
      for clas, group in grouped data:
 85
          mean = group.mean()
          std dev = group.std()
 86
 87
          #print(f"\nGroup for target = {clas}")
          #for col in range(len(mean)):
 88
 89
               print(f"Column {col:<10} Mean: {mean[col]:<15.4f} Standard Deviation:</pre>
          {std dev[col]:<15.4f}")
 90
 91
      # probalistic Model
 92
      # build the variance table from the means for each column and the column variance,
 93
      # where the variance is the square of the standard deviation
 94
     print()
 95
     variance table data = []
 96
     for clas, group in grouped_data:
 97
          mean = group.mean()
 98
          std dev = group.std()
 99
          std dev[std dev == 0.0000] = 0.0001
          variance = std dev ** 2
100
101
          variance[variance == 0.0000] = 0.0001
102
          for i, feature label in enumerate(feature labels):
              print(f"class: {clas:<10} feature: {feature label:<30} mean: {mean[i]:<10.4f}</pre>
103
              std dev: {std dev[i]:<10.4f} \
104
                     variance: {variance[i]:<10.4f}")</pre>
              variance table data.append({'Class': clas, 'Feature': feature label, 'Mean': mean
105
              [i], 'Standard Deviation': std dev[i],\
106
                                             'Variance': variance[i]})
107
      variance table = pd.DataFrame(variance table data)
108
      variance table class0 = variance table[variance table['Class'] == 0]
109
      variance table class1 = variance table[variance table['Class'] == 1]
110
111
      #print(np.shape(variance table)
112
     print("\nVariance Table")
113
114
    print(f"{'
                                                      class 0':37}{'
      class 1':<25}")
115
      print(f"{'Feature':<30} {'Spam mean':<10} {'Spam variance':<15} {'Not Spam mean':<13} {</pre>
      'Not Spam variance':<16}")
116
      row size = training data.shape[1] -1
117
      for i in range(row size):
118
          print(f"{variance table class0.iloc[i, 1]:<30}", f"{variance table class0.iloc[i, 2]-</pre>
          1:<10.4f}", \
                f"{variance table class0.iloc[i, 4]:<15.4f}", f"{variance table class1.iloc[i,
119
                2]-1:<13.4f}", \
120
                f"{variance table class1.iloc[i, 4]:<16.4f}")</pre>
121
      print()
122
123
      # step 3
124
      # Gaussian Naive Bayes Algorithm `
125
126
      #calculatethe gaussian probability density function
127
```

```
128
      def gaussian pdf(x, mean, variance):
129
          if variance == 0:
130
              variance = 0.0001
          return (1 / np.sqrt(2 * np.pi * variance)) * np.exp(-((x - mean) ** 2) / (2 *
131
          variance))
132
133
      def class probabilites (data, variance table class0, variance table class1,
      hypothesis ratios):
134
          # initialize the probabilities
135
          probabilities = []
136
          for i in range(len(data)):
137
              # initialize the probabilities for each class
138
              probabilities class0 = 0
139
             probabilities class1 = 0
140
              for j in range(len(data.iloc[i])):
141
                  # calculate the probabilities for each class
142
                  probabilities class0 += np.log(gaussian pdf(data.iloc[i, j],
                  variance table class0.iloc[j, 2], variance table class0.iloc[j, 4]))
143
                  probabilities class1 += np.log(gaussian pdf(data.iloc[i, j],
                  variance table class1.iloc[j, 2], variance table class1.iloc[j, 4]))
              # calculate the probabilities for each class
144
145
              probabilities class0 += np.log(hypothesis ratios[0])
              probabilities class1 += np.log(hypothesis ratios[1])
146
147
              # append the probabilities for each class
148
              probabilities.append([probabilities class0, probabilities class1])
149
          return probabilities
150
151
      # calculate the class probabilities for the training data
152
     probabilities = class probabilites (training data, variance table class0,
     variance_table_class1, hypothesis_ratios)
153
     print()
154
     print("Training data class probabilities")
     print("class 0", "
155
                             class 1")
156
     for i in range(10):
157
          print(f"{probabilities[i][0]:10.4f} {probabilities[i][1]:10.4f}")
158
     print()
159
160
161
      # calculate the class probabilities for the test data
162
      probabilities = class probabilites (test data, variance table class0,
      variance table class1, hypothesis ratios)
163
     print("Test data class probabilities")
164
     [print("class 0", "
                                     class 1")]
165
     for i in range (10):
166
          print(f"{probabilities[i][0]:10.4f} { ' ' * 6} {probabilities[i][1]:10.4f}")
167
     print()
168
169
      # calculate the accuracy, precision, and recall of the model on the test data and the
      test targets
170
      def accuracy precision recall(test targets, probabilities):
          # initialize the accuracy, precision, and recall
171
172
          accuracy = 0
173
          precision = 0
174
          recall = 0
175
          for i in range(len(test targets)):
176
              # determine the predicted class
177
              predicted class = 0 if probabilities[i][0] > probabilities[i][1] else 1
178
              # determine the actual class
179
              actual class = test targets.iloc[i]
180
              # increment the accuracy, precision, and recall
181
              accuracy += 1 if predicted class == actual class else 0
182
              precision += 1 if predicted class == 1 and actual class == 1 else 0
183
              recall += 1 if actual class == 1 else 0
184
          # calculate the accuracy, precision, and recall
185
          accuracy /= len(test targets)
          precision /= recall
186
          recall /= len(test targets)
187
```

```
188
         return accuracy, precision, recall
189
190 print()
191 print("Accuracy Precision Recall")
accuracy, precision, recall = accuracy_precision_recall(test_targets, probabilities)
193 print(f"{accuracy:.4f} {precision:.4f} {recall:.4f}")
194
    print()
195
196
    # generate a confusion matrix for the test data given the test targets
197
     def confusion matrix(test targets, probabilities):
198
         # initialize the confusion matrix
199
         confusion matrix = np.zeros((2, 2))
200
         for i in range(len(test targets)):
201
             # determine the predicted class
            predicted class = 0 if probabilities[i][0] > probabilities[i][1] else 1
202
             # determine the actual class
203
             actual class = test targets.iloc[i]
204
205
             # increment the confusion matrix
206
             confusion matrix[actual class][predicted class] += 1
207
         return confusion matrix
208
209 print()
210 print("Confusion Matrix")
211
     print(confusion matrix(test targets, probabilities))
212
```