

```

1  '''
2  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
8
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):
17     feature_labels = []
18     with open(file_path, 'r') as file:
19         for line in file:
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)
26 print("Data set: spambase.data")
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 probabilistic models for each class
58 #
59
60 # Determine each of the hypothesis ratios for each hypothesis (H), based on the test data
61 #  $P(H|D) = P(D|H) * P(H) / P(D)$ 
62 hypotheses, frequencies = np.unique(training_targets, return_counts=True)
63 hypothesis_ratios = frequencies / len(test_targets)
64
65 # Calculate the prior probabilities P(H)
66 #  $P(H) = P(D|H) * P(H) / P(D)$ 
67

```

```

68 for h, f, r in zip(hypotheses, frequencies, hypothesis_ratios):
69     print()
70     print(f"{'Hypothesis':<15} {'Frequency':<15} {'Ratio':<15}")
71     print(f"{h:<15} {f:<15} {r:<15.2f}")
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()
86     std_dev = group.std()
87     #print(f"\nGroup for target = {clas}")
88     #for col in range(len(mean)):
89     #    print(f"Column {col:<10} Mean: {mean[col]:<15.4f} \
90     # Standard Deviation: {std_dev[col]:<15.4f}")
91
92 # probabilistic Model
93 # build the variance table from the means for each column and the column variance,
94 # where the variance is the square of the standard deviation
95 print()
96 variance_table_data = []
97 for clas, group in grouped_data:
98     mean = group.mean()
99     std_dev = group.std()
100     std_dev[std_dev == 0.0000] = 0.0001
101     variance = std_dev ** 2
102     variance[variance == 0.0000] = 0.0001
103     for i, feature_label in enumerate(feature_labels):
104         print(f"class: {clas:<10} feature: {feature_label:<30} mean: {mean[i]:<10.4f} \
105             std_dev: {std_dev[i]:<10.4f} variance: {variance[i]:<10.4f}")
106         variance_table_data.append({'Class': clas, 'Feature': feature_label, \
107             'Mean': mean[i], 'Standard Deviation': std_dev[i], \
108             'Variance': variance[i]})
109 variance_table = pd.DataFrame(variance_table_data)
110 variance_table_class0 = variance_table[variance_table['Class'] == 0]
111 variance_table_class1 = variance_table[variance_table['Class'] == 1]
112
113 #print(np.shape(variance_table)    )
114
115 print("\nVariance Table")
116 print(f"{' ' class 0':<37} \
117     {' ' class 1':<25}")
118 print(f"{'Feature':<30} {'Spam_mean':<10} {'Spam_variance':<15} {'Not Spam_mean':<13} \
119     {'Not Spam_variance':<16}")
120 row_size = training_data.shape[1] - 1
121 for i in range(row_size):
122     print(f"{variance_table_class0.iloc[i, 1]:<30}", \
123         f"{variance_table_class0.iloc[i, 2]-1:<10.4f}", \
124         f"{variance_table_class0.iloc[i, 4]:<15.4f}", \
125         f"{variance_table_class1.iloc[i, 2]-1:<13.4f}", \
126         f"{variance_table_class1.iloc[i, 4]:<16.4f}")
127 print()
128
129 # step 3
130 # Gaussian Naive Bayes Algorithm `
131
132 #calculatethe gaussian probability density function
133
134 def gaussian_pdf(x, mean, variance):

```

```

135     if variance == 0:
136         variance = 0.0001
137     return (1 / np.sqrt(2 * np.pi * variance)) * \
138     np.exp(-((x - mean) ** 2) / (2 * variance))
139
140 def class_probabilites(data, variance_table_class0, variance_table_class1,
141 hypothesis_ratios):
142     # initialize the probabilities
143     probabilities = []
144     for i in range(len(data)):
145         # initialize the probabilities for each class
146         probabilities_class0 = 0
147         probabilities_class1 = 0
148         for j in range(len(data.iloc[i])):
149             # calculate the probabilities for each class
150             probabilities_class0 += np.log(gaussian_pdf(data.iloc[i, j], \
151                 variance_table_class0.iloc[j, 2], variance_table_class0.iloc[j, 4]
152             )))
153             probabilities_class1 += np.log(gaussian_pdf(data.iloc[i, j], \
154                 variance_table_class1.iloc[j, 2], variance_table_class1.iloc[j, 4]
155             )))
156         # calculate the probabilities for each class
157         probabilities_class0 += np.log(hypothesis_ratios[0])
158         probabilities_class1 += np.log(hypothesis_ratios[1])
159         # append the probabilities for each class
160         probabilities.append([probabilities_class0, probabilities_class1])
161     return probabilities
162
163 # calculate the class probabilities for the training data
164 probabilities = class_probabilites(training_data, variance_table_class0,
165 variance_table_class1, \
166                                     hypothesis_ratios)
167
168 print()
169 print("Training_data class probabilities")
170 print("class 0", "      class 1")
171 for i in range(10):
172     print(f"{probabilities[i][0]:10.4f} {probabilities[i][1]:10.4f}")
173 print()
174
175 # calculate the class probabilities for the test data
176 probabilities = class_probabilites(test_data, variance_table_class0,
177 variance_table_class1, \
178                                     hypothesis_ratios)
179
180 print("Test_data class probabilities")
181 [print("class 0", "      class 1")]
182 for i in range(10):
183     print(f"{probabilities[i][0]:10.4f} {' ' * 6} {probabilities[i][1]:10.4f}")
184 print()
185
186 # calculate the accuracy, precision, and recall of the model on the test_data \
187 # and the test_targets
188 def accuracy_precision_recall(test_targets, probabilities):
189     # initialize the accuracy, precision, and recall
190     accuracy = 0
191     precision = 0
192     recall = 0
193     for i in range(len(test_targets)):
194         # determine the predicted class
195         predicted_class = 0 if probabilities[i][0] > probabilities[i][1] else 1
196         # determine the actual class
197         actual_class = test_targets.iloc[i]
198         # increment the accuracy, precision, and recall
199         accuracy += 1 if predicted_class == actual_class else 0
200         precision += 1 if predicted_class == 1 and actual_class == 1 else 0
201         recall += 1 if actual_class == 1 else 0
202     # calculate the accuracy, precision, and recall

```

```

197     accuracy /= len(test_targets)
198     precision /= recall
199     recall /= len(test_targets)
200     return accuracy, precision, recall
201
202 print()
203 print("Accuracy    Precision    Recall")
204 accuracy, precision, recall = accuracy_precision_recall(test_targets, probabilities)
205 print(f"{{accuracy:.4f}}    {{precision:.4f}}    {{recall:.4f}}")
206 print()
207
208 # generate a confusion matrix for the test data given the test_targets
209 def confusion_matrix(test_targets, probabilities):
210     # initialize the confusion matrix
211     confusion_matrix = np.zeros((2, 2))
212     for i in range(len(test_targets)):
213         # determine the predicted class
214         predicted_class = 0 if probabilities[i][0] > probabilities[i][1] else 1
215         # determine the actual class
216         actual_class = test_targets.iloc[i]
217         # increment the confusion matrix
218         confusion_matrix[actual_class][predicted_class] += 1
219     return confusion_matrix
220
221 print()
222 print("Confusion Matrix")
223 print(confusion_matrix(test_targets, probabilities))
224

```