

```

1 import pandas as pd
2 import numpy as np
3 from pprint import pprint
4 import sys
5
6 # Reads the data from CSV files, each attribute column can be obtained via its name,
  e.g., y = data['y']
7 def getDataframe(filePath):
8     data = pd.read_csv(filePath)
9     return data
10
11 # predicted_y and y are the predicted and actual y values respectively as numpy
  arrays
12 # function prints the accuracy
13 def compute_accuracy(predicted_y, y):
14     acc = 100.0
15     acc = np.sum(predicted_y == y)/predicted_y.shape[0]
16     return acc
17
18 #Compute entropy according to y distribution
19 def compute_entropy(y):
20     entropy = 0.0
21     elements,counts = np.unique(y, return_counts = True)
22     n = y.shape[0]
23
24     for i in range(len(elements)):
25         prob = counts[i]/n
26         if prob!= 0:
27             entropy -= prob * np.log2(prob)
28     return entropy
29
30 #att_name: attribute name; y_name: the target attribute name for classification
31 def compute_info_gain(data, att_name, y_name):
32     info_gain = 0.0
33
34     #Calculate the values and the corresponding counts for the select attribute
35     vals, counts = np.unique(data[att_name], return_counts=True)
36     total_counts = np.sum(counts)
37
38     #Calculate the conditional entropy
39     #=====#
40     # STRART YOUR CODE HERE #
41     #=====#
42     info = compute_entropy(data[y_name])
43     entropy = 0
44     for i in range(len(vals)):
45         probability = counts[i] / total_counts
46         val_data = data.loc[data[att_name] == vals[i]]
47         entropy += probability * compute_entropy(val_data[y_name])
48     info_gain = info - entropy
49     #=====#
50     # END YOUR CODE HERE #
51     #=====#
52
53     return info_gain
54
55
56 def comput_gain_ratio(data, att_name, y_name):
57     gain_ratio = 0.0
58     #Calculate the values and the corresponding counts for the select attribute

```

```

59     vals, counts = np.unique(data[att_name], return_counts=True)
60     total_counts = np.sum(counts)
61
62     #Calculate the information for the selected attribute
63     att_info = 0.0
64     #=====#
65     # STRART YOUR CODE HERE #
66     #=====#
67     for i in range(len(counts)):
68         probability = counts[i] / total_counts
69         att_info -= probability * np.log2(probability)
70     #=====#
71     # END YOUR CODE HERE #
72     #=====#
73     gain_ratio = 0.0 if np.abs(att_info) < 1e-9 else min(1, compute_info_gain(data,
att_name, y_name) / att_info)
74     return gain_ratio
75
76 # Class of the decision tree model based on the ID3 algorithm
77 class DecisionTree(object):
78     def __init__(self):
79         self.train_data = pd.DataFrame()
80         self.test_data = pd.DataFrame()
81
82     def load_data(self, train_file, test_file):
83         self.train_data = getDataframe(train_file)
84         self.test_data = getDataframe(test_file)
85
86     def train(self, y_name, measure, parent_node_class= None):
87         self.y_name = y_name
88         self.measure = measure
89         self.tree = self.make_tree(self.train_data, parent_node_class)
90
91     def make_tree(self, train_data, parent_node_class = None):
92         data = train_data
93         features = data.drop(self.y_name, axis = 1).columns.values
94         measure = self.measure
95         #Stopping condition 1: If all target_values have the same value, return this
value
96         if len(np.unique(data[self.y_name])) <= 1:
97             leaf_value = -1
98             #=====#
99             # STRART YOUR CODE HERE #
100            #=====#
101            leaf_value = data[self.y_name].values[0]
102            #=====#
103            # END YOUR CODE HERE #
104            #=====#
105            return leaf_value
106
107            #Stopping condition 2: If the dataset is empty, return the parent_node_class
108            elif len(data)== 0:
109                return parent_node_class
110
111            #Stopping condition 3: If the feature space is empty, return the majority
class
112            elif len(features) == 0:
113                return np.unique(data[self.y_name])
[ np.argmax(np.unique(data[y_name], return_counts=True)[1])]
114

```

```

115         # Not a leaf node, create an internal node
116         else:
117             #Set the default value for this node --> The mode target feature value of
the current node
118             parent_node_class = np.unique(data[self.y_name])
[ np.argmax(np.unique(data[self.y_name], return_counts=True)[1])]
119
120             #Select the feature which best splits the dataset
121             if measure == 'info_gain':
122                 item_values = [compute_info_gain(data, feature, self.y_name) for
feature in features] #Return the information gain values for the features in the
dataset
123             elif measure == 'gain_ratio':
124                 item_values = [comput_gain_ratio(data, feature, self.y_name) for
feature in features] #Return the gain_ratio for the features in the dataset
125             else:
126                 raise ValueError("kernel not recognized")
127
128             best_feature_index = np.argmax(item_values)
129             best_feature = features[best_feature_index]
130             print('best_feature is: ', best_feature)
131
132             #Create the tree structure. The root gets the name of the feature
(best_feature)
133             tree = {best_feature: {}}
134
135
136             #Grow a branch under the root node for each possible value of the root node
feature
137
138             for value in np.unique(data[best_feature]):
139                 #Split the dataset along the value of the feature with the largest
information gain and therwith create sub_datasets
140                 sub_data = data.where(data[best_feature] == value).dropna()
141
142                 #Remove the selected feature from the feature space
143                 sub_data = sub_data.drop(best_feature, axis = 1)
144
145                 #Call the ID3 algorithm for each of those sub_datasets with the new
parameters --> Here the recursion comes in!
146                 subtree = self.make_tree(sub_data, parent_node_class)
147
148                 #Add the sub tree, grown from the sub_dataset to the tree under the root
node
149                 tree[best_feature][value] = subtree
150
151             return tree
152
153
154         def test(self, y_name):
155             accuracy = self.classify(self.test_data, y_name)
156             return accuracy
157
158         def classify(self, test_data, y_name):
159             #Create new query instances by simply removing the target feature column from
the test dataset and
160             #convert it to a dictionary
161             test_x = test_data.drop(y_name, axis=1)
162             test_y = test_data[y_name]
163

```

```
164     n = test_data.shape[0]
165     predicted_y = np.zeros(n)
166
167     #Calculate the prediction accuracy
168     for i in range(n):
169         predicted_y[i] = DecisionTree.predict(self.tree, test_x.iloc[i])
170
171     output = np.zeros((n,2))
172     output[:,0] = test_y
173     output[:,1] = predicted_y
174     accuracy = compute_accuracy(predicted_y, test_y.values)
175     return accuracy
176
177 def predict(tree, query):
178     # find the root attribute
179     default = -1
180     for root_name in list(tree.keys()):
181         try:
182             subtree = tree[root_name][query[root_name]]
183         except:
184             return default ## root_name does not appear in query attribute list
185             (it is an error!)
186
187     ##if subtree is still a dictionary, recursively test next attribute
188     if isinstance(subtree,dict):
189         return DecisionTree.predict(subtree, query)
190     else:
191         leaf = subtree
192         return leaf
193
```