```python
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     #Calculate the conditional entropy
38     #=======================#
39     # STRART YOUR CODE HERE  #
40     #=======================#
41     total_info = compute_entropy(data[y_name])
42     info_A = 0.0
43     for i in range(len(vals)):
44         info_A += (counts[i]/total_counts) *
   compute_entropy(data.loc[data[att_name] == vals[i]][y_name])
45     #=======================#
46     #   END YOUR CODE HERE   #
47     #=======================#
48     info_gain = total_info - info_A
49     return info_gain
50
51
52 def comput_gain_ratio(data, att_name, y_name):
53     gain_ratio = 0.0
54     #Calculate the values and the corresponding counts for the select
```

```python
55      vals, counts = np.unique(data[att_name], return_counts=True)
56      total_counts = np.sum(counts)
57
58      #Calculate the information for the selected attribute
59      att_info = 0.0
60      #========================#
61      # STRART YOUR CODE HERE  #
62      #========================#
63      for i in range(len(vals)):
64          p = counts[i] / total_counts
65          att_info -= p * np.log2(p)
66      #========================#
67      #   END YOUR CODE HERE   #
68      #========================#
69      gain_ratio = 0.0 if np.abs(att_info) < 1e-9 else min(1,
    compute_info_gain(data, att_name, y_name) / att_info)
70      return gain_ratio
71
72 # Class of the decision tree model based on the ID3 algorithm
73 class DecisionTree(object):
74      def __init__(self):
75          self.train_data = pd.DataFrame()
76          self.test_data = pd.DataFrame()
77
78      def load_data(self, train_file, test_file):
79          self.train_data = getDataframe(train_file)
80          self.test_data = getDataframe(test_file)
81
82      def train(self, y_name, measure, parent_node_class = None):
83          self.y_name = y_name
84          self.measure = measure
85          self.tree = self.make_tree(self.train_data, parent_node_class)
86
87      def make_tree(self, train_data, parent_node_class = None):
88          data = train_data
89          features = data.drop(self.y_name, axis = 1).columns.values
90          measure =  self.measure
91          #Stopping condition 1: If all target_values have the same value,
    return this value
92          if len(np.unique(data[self.y_name])) <= 1:
93              leaf_value = -1
94              #========================#
95              # STRART YOUR CODE HERE  #
96              #========================#
97              leaf_node = np.unique(data[self.y_name])[0]
98              #========================#
99              #   END YOUR CODE HERE   #
100             #========================#
101             return leaf_node
102
103         #Stopping condition 2: If the dataset is empty, return the
    parent_node_class
104         elif len(data)== 0:
105             return parent_node_class
106
107         #Stopping condition 3: If the feature space is empty, return the
    majority class
108         elif len(features) == 0:
109             return np.unique(data[self.y_name])
```

```python
110
111          # Not a leaf node, create an internal node
112          else:
113              #Set the default value for this node --> The mode target feature
     value of the current node
114              parent_node_class = np.unique(data[self.y_name])
     [np.argmax(np.unique(data[self.y_name],return_counts=True)[1])]
115
116              #Select the feature which best splits the dataset
117              if measure == 'info_gain':
118                  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
119              elif measure == 'gain_ratio':
120                  item_values = [comput_gain_ratio(data, feature, self.y_name)
     for feature in features] #Return the gain_ratio for the features in the
     dataset
121              else:
122                  raise ValueError("kernel not recognized")
123
124              best_feature_index = np.argmax(item_values)
125              best_feature = features[best_feature_index]
126              print('best_feature is: ', best_feature)
127
128              #Create the tree structure. The root gets the name of the feature
     (best_feature)
129              tree = {best_feature:{}}
130
131
132          #Grow a branch under the root node for each possible value of the
     root node feature
133
134          for value in np.unique(data[best_feature]):
135              #Split the dataset along the value of the feature with the
     largest information gain and therwith create sub_datasets
136              sub_data = data.where(data[best_feature] == value).dropna()
137
138              #Remove the selected feature from the feature space
139              sub_data = sub_data.drop(best_feature, axis = 1)
140
141              #Call the ID3 algorithm for each of those sub_datasets with the
     new parameters --> Here the recursion comes in!
142              subtree = self.make_tree(sub_data, parent_node_class)
143
144              #Add the sub tree, grown from the sub_dataset to the tree under
     the root node
145              tree[best_feature][value] = subtree
146
147          return tree
148
149
150      def test(self, y_name):
151          accuracy = self.classify(self.test_data, y_name)
152          return accuracy
153
154      def classify(self, test_data, y_name):
155          #Create new query instances by simply removing the target feature
     column from the test dataset and
156          #convert it to a dictionary
```

```python
            test_y = test_data[y_name]

            n = test_data.shape[0]
            predicted_y = np.zeros(n)

            #Calculate the prediction accuracy
            for i in range(n):
                predicted_y[i] = DecisionTree.predict(self.tree, test_x.iloc[i])

            output = np.zeros((n,2))
            output[:,0] = test_y
            output[:,1] = predicted_y
            accuracy = compute_accuracy(predicted_y, test_y.values)
            return accuracy

    def predict(tree, query):
        # find the root attribute
        default = -1
        for root_name in list(tree.keys()):
            try:
                subtree = tree[root_name][query[root_name]]
            except:
                return default ## root_name does not appear in query
    attribute list (it is an error!)

                ##if subtree is still a dictionary, recursively test next
    attribute
            if isinstance(subtree,dict):
                return DecisionTree.predict(subtree, query)
            else:
                leaf = subtree
                return leaf
```