```python
import csv
import random
import sys
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
from tkinter import filedialog
import tkinter as tk
import ntpath
import ctypes


# Used to load in the file data into the dataset
# Also returns the number of entries in the file
def load_data(file):
    d = []
    att = []
    with open(file, 'r') as f:
        r = csv.reader(f, delimiter='\t')
        i = 0
        for row in r:
            if i == 0:
                att = row
                i += 1
            else:
                d.append(row)
                i += 1
    return d, i, att


# Used to shuffle and split the data into 2/3s training data and 1/3 testing data

def split_shuffle(ds, parts):
    random.shuffle(ds)
    p = (len(ds)) // parts
    test = ds[:p]
    train = ds[p:]
    return test, train


# Used to returns all values for a set column

def get_column(rows, col):
    cols = []

    for row in rows:
        cols.append(row[col])
    return cols


# Used to count the number of each label/feature (beer_style) in the rows passed to the
function

def y_count(r):
    y_num = {}
    for row in r:
        y = row[label_col]
        if y not in y_num:
            y_num[y] = 0
        y_num[y] += 1
    return y_num


# Used to compare and test if the current row is greater than or equal to the test value
# in order to split up the data

def compare(r, test_c, test_val):
```

```python
67          if r[test_c].isdigit():
68              return r[test_c] == test_val
69
70          elif float(r[test_c]) >= float(test_val):
71              return True
72
73          else:
74              return False
75
76
77      # Splits the data into two lists for the true/false results of the compare test
78      def fork(r, c, test_val):
79          true = []
80          false = []
81
82          for row in r:
83
84              if compare(row, c, test_val):
85                  true.append(row)
86              else:
87                  false.append(row)
88
89          return true, false
90
91
92      # Used to calculate the Gini Index/Impurity of the rows inputted (of beer style)
93
94      def gini_index(r):
95          stylesNum = y_count(r)
96          impurity = 1
97
98          for style in stylesNum:
99              style_prob = stylesNum[style] / float(len(r))
100             impurity -= style_prob ** 2
101         return impurity
102
103
104     # Used to calculate the Information gain, incorporates the gini index (impurity)
105
106     def gain(left, right, impurity):
107         p = float(len(left)) / (len(left) + len(right))
108         ig = impurity - p * gini_index(left) - (1 - p) * gini_index(right)
109         return ig
110
111
112     # Used to find the best split for data among all attributes
113
114     def split(r):
115         max_ig = 0
116         max_att = 0
117         max_att_val = 0
118
119         # calculates gini for the rows provided
120         curr_gini = gini_index(r)
121         no_att = len(r[0])
122
123         # Goes through the different attributes
124
125         for c in range(no_att):
126
127             # Skip the label column (beer style)
128
129             if c == label_col:
130                 continue
131             column_vals = get_column(r, c)
132
133             i = 0
```

```python
134            while i < len(column_vals):
135                # value to compare
136                att_val = r[i][c]
137
138                # Use the attribute value to fork the data to true and false streams
139                true, false = fork(r, c, att_val)
140
141                # Calculate the information gain
142                ig = gain(true, false, curr_gini)
143
144                # If this gain is the highest found then mark this as the best choice
145                if ig > max_ig:
146                    max_ig = ig
147                    max_att = c
148                    max_att_val = r[i][c]
149                i += 1
150
151        return max_ig, max_att, max_att_val


# Used to recursively go through the tree in order to find the optimal attribute to
split the tree with

def rec_tree(r):
    ig, att, curr_att_val = split(r)

    if ig == 0:
        return Leaf(r)

    true_rows, false_rows = fork(r, att, curr_att_val)

    true_branch = rec_tree(true_rows)
    false_branch = rec_tree(false_rows)

    return Node(att, curr_att_val, true_branch, false_branch)


# Defines the classifications of the leaf

class Leaf:
    def __init__(self, rows):
        self.predictions = y_count(rows)


# Defines a split node - contains the primary attribute its value and the two child
branches

class Node:
    def __init__(self, att, att_value, true_branch, false_branch):
        self.att = att
        self.att_value = att_value
        self.true_branch = true_branch
        self.false_branch = false_branch


# Confidence is used in order to determine what is each value

def confidence(r, node):
    if isinstance(node, Leaf):
        return node.predictions

    c = node.att
    att_value = node.att_value

    if compare(r, c, att_value):
        return confidence(r, node.true_branch)
    else:
```

```python
199                 return confidence(r, node.false_branch)
200
201
202    # Prints and formats the tree based on the branches and questions
203
204    def build_tree(node, spacing=""):
205        # If you've reached the terminal state then predict
206        if isinstance(node, Leaf):
207            print(spacing + "Predict", node.predictions)
208            return
209
210        print(spacing + "Is " + attributes[node.att] + " > " + str(node.att_value) + " ?")
211
212        print(spacing + '--> True:')
213        build_tree(node.true_branch, spacing + "  ")
214
215        print(spacing + '--> False:')
216        build_tree(node.false_branch, spacing + "  ")
217
218
219    # Prints out the leaf (the beer style)
220
221    def print_leaf(counts):
222        total = sum(counts.values())
223        probs = {}
224        for lbl in counts.keys():
225            probs[lbl] = str(int(counts[lbl] / total * 100)) + "%"
226        return probs
227
228    # Ntpath is used in order to retrieve the name of the file from the file path
229
230    def path_name(path):
231        head, tail = ntpath.split(path)
232        return tail or ntpath.basename(head)
233
234
235    if __name__ == "__main__":
236        #TKinter is used in order to open file dialog to get the training and testing data
237        root = tk.Tk()
238        root.withdraw()
239        #ctypes is used in order to print out a message box to tell the user which files
240        are being asked of them
        ctypes.windll.user32.MessageBoxW(0, "Select your training + testing data", "File
        Selection", 0)
241
242        file_path = filedialog.askopenfilename()
243        print(file_path)
244        filename = path_name(file_path)
245        filename="beer.txt"
246
247        #Label col in this case is beer style and is adjustable to whichever attritbute you
        choose
248        label_col = 3
249        avg_acc = 0
250        avg_ref_acc = 0
251        i = 0
252        data, classes, attributes = load_data(filename)
253
254        # ----------------------------------------------------------------#
255        # This is for the reference implementation of the decision tree classifier
256
257
258        featuredCols = ['calorific_value', 'nitrogen', 'turbidity', 'alcohol', 'sugars',
        'bitterness', 'beer_id',
259                        'colour', 'degree_of_fermentation']
260        ref_attributes = ['calorific_value', 'nitrogen', 'turbidity', 'beer_style',
        'alcohol', 'sugars', 'bitterness',
```

```python
261                             'beer_id', 'colour', 'degree_of_fermentation']
262
263         ctypes.windll.user32.MessageBoxW(0, "Select your reference algorithm training
            data", "File Selection", 0)
264         ref_train_path = filedialog.askopenfilename()
265         train_path_filename = path_name(ref_train_path)
266
267         ctypes.windll.user32.MessageBoxW(0, "Select your reference algorithm testing data",
            "File Selection", 0)
268         ref_test_path = filedialog.askopenfilename()
269         test_path_filename = path_name(ref_test_path)
270
271         trainingData = pd.read_csv("training.txt", sep='\t', names=ref_attributes)
272         testData = pd.read_csv("test.txt", sep='\t', names=ref_attributes)
273         sys.stdout = open('output.txt', 'wt')
274
275         # ----------------------------------------------------------------#
276         # Main random divisions of the algorithm. Each time the testing and training data is
277         # shuffled and split randomly
278
279         while i < 10:
280             testing, training = split_shuffle(data, 3)
281             tree = rec_tree(training)
282             build_tree(tree)
283
284             correct = 0
285             incorrect = 0
286             for r in testing:
287                 print("Actual: %s. Predicted: %s" % (r[label_col], print_leaf(confidence(r,
                    tree))))
288                 for key, value in confidence(r, tree).items():
289                     if r[label_col] == key:
290                         correct += 1
291                     else:
292                         incorrect += 1
293         print('Percentage Correctly Classified')
294         print(correct / (correct + incorrect) * 100)
295         print('Percentage Incorrectly Classified')
296         print(incorrect / (correct + incorrect) * 100)
297
298             i += 1
299             avg_acc += correct / (correct + incorrect)
300
301             # ----------------------------------------------------------------#
302             # REFERENCE IMPLEMENTATION
303
304             x_train = trainingData[featuredCols]
305             y_train = trainingData.beer_style
306             x_test = testData[featuredCols]
307             y_test = testData.beer_style
308
309             dtc = DecisionTreeClassifier()
310             dtc = dtc.fit(x_train, y_train)
311             y_predict = dtc.predict(x_test)
312
313             avg_ref_acc += metrics.accuracy_score(y_test, y_predict)
314             print('Reference Algorithms Percentage Accuracy')
315             print(metrics.accuracy_score(y_test, y_predict)*100)
316             # ----------------------------------------------------------------#
317
318         print("\nThe Average Accuracy across 10 iterations: ")
319         acc10 = avg_acc / 10 * 100
320         print(acc10)
321
322         refAcc = ((avg_ref_acc / 10) * 100)
323         print("\nThe Average Accuracy for the reference decision tree classifier across 10
            iterations: ")
```

```
324        print(refAcc)
325
```