

# CT4101 Assignment 2

Name: Aideen McLoughlin  
Name: Louise Kilheaney

Student ID: 17346123  
Student ID: 16100463

Class: 4th year CS  
Class: 4th year CS

## Description of Algorithm

For this assignment, we chose to develop an implementation of the C4.5 Algorithm.

Our algorithm starts by gathering the user input. The user must type the filepath of the file from which to fetch the data, and the train to test split fraction, selected from one of 3 options -  $\frac{2}{3}$ ,  $\frac{1}{3}$  or  $\frac{1}{5}$ . i.e. if a user picks  $\frac{2}{3}$  then  $\frac{2}{3}$  is training and  $\frac{1}{3}$  is testing. If the file path to the data is wrong, an error window is loaded which will prompt the user to try again. Once a valid file path has been passed, the gather\_data function returns the file data, and the split fraction. At this point, createTreeWhileShowingLoadingScreen is called, this function uses python multiprocessing to render an animation of a loading screen while the tree is being built, so that the user can know that there have been no errors, and the tree is building. At the same time as this loading animation function, renderLoadingWindow, is called, the createTree function is called. createTree takes the file data and split value. createTree calls the function split\_data\_training\_testing, which splits the data into training and testing data in the ratio of the split value. The function then saves these datasets so that the Weka implementation can run with the same data, then the datasets of train\_data and test\_data are returned to createTree.

Once the training and testing data has been returned, the buildTree function is called. buildTree takes the training data and the list of attributes as input.

Build\_tree allows for 3 bases cases

1. Check that all data is not in the same class and return a failure leaf node fail if it is
2. Check that the training data set is not empty and return a failure leaf node fail if it is
3. Otherwise, Find the attribute with the highest information gain and return that attribute
  - a. find\_best\_attribute is called, which takes the data and calls spilt\_into\_subsets which splits the data in half around a threshold value
  - b. Then information gain is called, which passes the subset to the entropy function and uses the returned value to calculate the information gain of the attribute. The information gain value is then returned.
4. If the best\_attribute function returns an empty string, then there are no more attributes to analyse, and a node can be created representing the majority class from the data.
5. Otherwise, a new node can be created for the best attribute, and the dataset is then split in 2 at the threshold value for that attribute. Those values are then appended as child nodes to the new best attribute node by recursively calling the createTree function, by calling build tree with in the subset and the unused attribute
6. Return the root node, which will return have connections to all the child nodes through the recursive calling of createTree

Once the tree has been created, the root\_node is passed to Print\_tree, which uses a Graphviz Digraph to visualise the tree and save it as a png so it can be displayed in the GUI.

The Accuracy of the tree is calculated by passing the root\_node and the testing data from the tree to the test\_tree function. The test\_tree function removes the target column from the dataset (style) and stores those values so they can be used in evaluation. The test\_tree function then calls the test\_data function, which loops through all the test data values, and calls test\_lr for each of them. Test\_lr takes a row of data and the root node of the tree. It starts at the root node of the tree, and works its way down to the relevant predicted leaf node for the row data passed to it. The leaf node reached in the test\_lr function is the predicted target for the row. Once all test rows have been evaluated, the test\_data function returns a list of all the predictions, and then loops through them alongside the actual values to calculate the accuracy of the tree. Test\_data then saves the test output to a file and returns the accuracy of the test.

The tree is then built using the Python Weka Wrapper module. The Weka C4.5 model is trained with the training dataset which was saved as a file in the training\_testing\_split function, and tested using the test data which was saved in that same function. The tree is then rendered as a png and saved, and the accuracy of the tree of the testing data is calculated. The results of the tests are stored in a csv file as was done in the python implementation.

The time it took to build the tree using the python code (built from scratch) and the weka code are formatted nicely, and then the final pop-up window will render, showing both the python tree and the weka tree, along with the accuracies of both trees built from the same data set on the same testing data set, and the time it took for each implementation to build the tree.

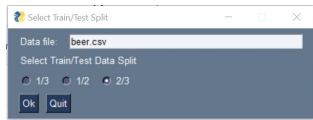
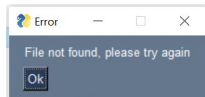
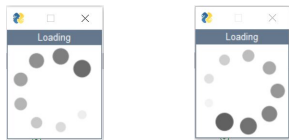
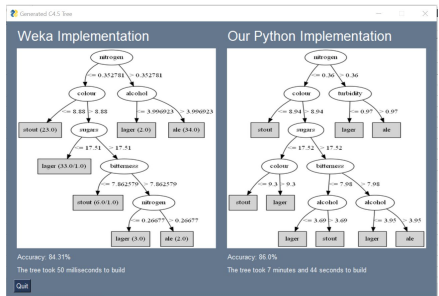
## Design Decisions

### Algorithm Design

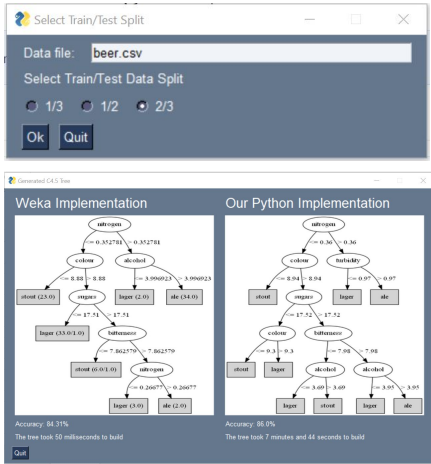
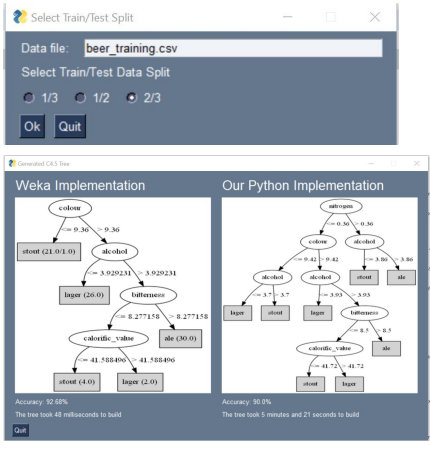
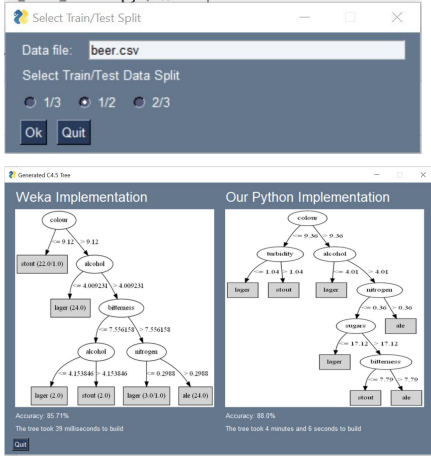
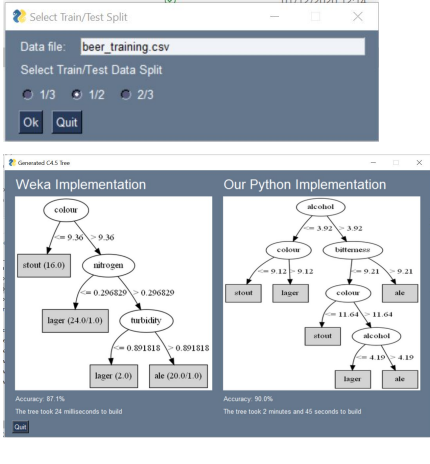
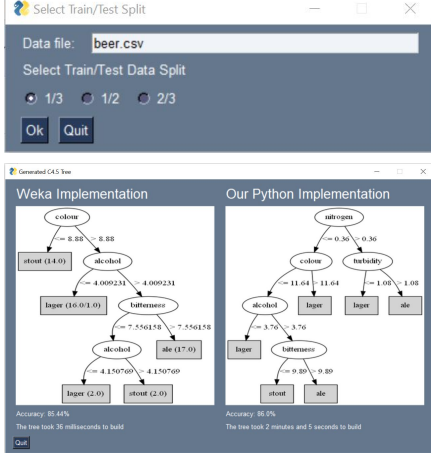
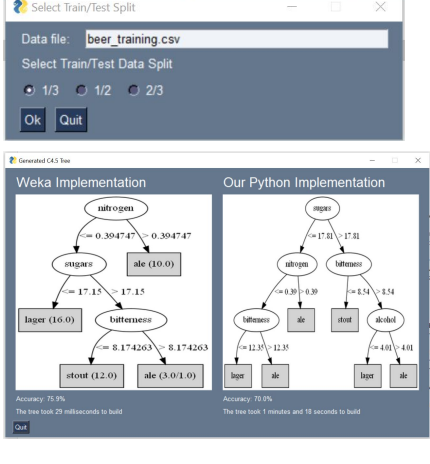
- We chose to use Python for this assignment as python facilitates the complex data operations required by the C4.5 algorithm, while still allowing for relatively easy implementation of elements such as the GUI.
- We chose to develop an GUI (A User Interface) for our algorithm implementation to make our tree building and testing clearer and easier to use
- We chose to input and process all of our data in CSV format as it made it easier for us to use and store.
- We chose to style the graphviz graph for our tree similarly to the weka tree for the sake of consistency, so that any graphical differences between the trees would be minimal
- We chose the split value options of  $\frac{1}{3}$ ,  $\frac{1}{2}$  and  $\frac{2}{3}$  to use in our GUI as they are all reasonable enough to allow for a nicely representative test and train data set, while being different enough to show how the split function can change the overall tree outcome.

### Interface and Graphics Design

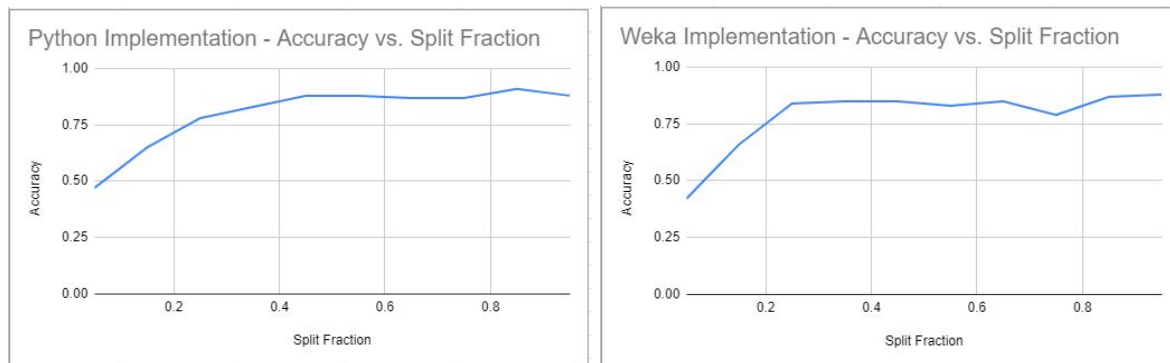
Our implementation of the C4.5 Algorithm runs with a simple Graphical User Interface built using the PySimpleGUI python library.

Action	Response	User Interface
Run the program	A pop up window appears, prompting the user to enter the data file location, and the split between training and testing data to use for that data	
Enter invalid data file location	Popup appears to notify you that an invalid file path was entered. Once User clicks OK, the initial window will pop up again so that the user can enter valid data	
Valid inputs in first popup, User clicks 'OK'	The program will begin building the C4.5 tree in the background. While this happens, a loading gif will be animated on the screen to show that the program is still working and not frozen.	
Once the tree has built	The loading popup will close, and a new window will pop up to show the Weka tree and the python tree built by our code.	

# Tests and Results

Data	Output	Data	
<p>Dataset: <b>beer.csv</b></p> <p>Train/test split: <b>2/3</b></p> <p>Accuracy: <b>86%</b> Weka Accuracy: <b>88.24%</b></p>		<p>Dataset: <b>beer_training.csv</b></p> <p>Train/test split: <b>2/3</b></p> <p>Accuracy: <b>90%</b> Weka Accuracy: <b>85.37%</b></p>	
<p>Dataset: <b>beer.csv</b></p> <p>Train/test split: <b>1/2</b></p> <p>Accuracy: <b>88%</b> Weka Accuracy: <b>87.01%</b></p>		<p>Dataset: <b>beer_training.csv</b></p> <p>Train/test split: <b>1/2</b></p> <p>Accuracy: <b>90%</b> Weka Accuracy: <b>87.1%</b></p>	
<p>Dataset: <b>beer.csv</b></p> <p>Train/test split: <b>1/3</b></p> <p>Accuracy: <b>86%</b> Weka Accuracy: <b>85.44%</b></p>		<p>Dataset: <b>beer_training.csv</b></p> <p>Train/test split: <b>1/3</b></p> <p>Accuracy: <b>70%</b> Weka Accuracy: <b>75.9%</b></p>	

In a Separate file - 10tests.py - We ran the build tree functions for our implementation and the weka implementation 10 times with different split fractions, and stored the results in files. We then plotted these results generated by these runs on graphs, Mapping the split fraction to the accuracy.



## Conclusions and Observations

From the Graphs shown above, we can see a definite positive correlation between Accuracy and split fraction. This does seem to level out at a Split Fraction of about 0.5 (or  $\frac{1}{2}$ ) for both the python implementation and the Weka implementation. From this we can gather that, for the beer dataset of 154 attributes, roughly 77 random attributes are enough to generate a reliable accuracy of 80% to 85%.

When the size of the training data falls below 20% of the total file (~46 lines), the created tree will be much less accurate, as the training dataset does not provide a clear enough picture of the entire data set.

We can see that when the algorithm is run on the beer\_training.csv dataset from Assignment 1, the trees generated for the  $\frac{2}{3}$  and  $\frac{1}{2}$  data split values are similarly accurate to the trees generated for those same split values from the beer.csv dataset. The tree generated for the  $\frac{1}{3}$  data split is significantly less accurate with the beer\_training data set than it was with the beer data set, and this is likely due to the fact that the beer\_training dataset is smaller than the beer dataset (124 lines). Thus, the training dataset would be smaller and less accurate when generated with this smaller dataset.

The Weka implementation performed similarly to the python-from-scratch implementation in all of the tested scenarios, with less than a 6% difference in accuracy for all tests run through the GUI and through the 10tests.py file. For all scenarios however, the Weka implementation was much faster than our python-from-scratch implementation as the Weka implementation was built with a greater focus on efficiency. The python implementation took between 2 and 5 seconds per training data row when calculating this test data. This time is understandable given the complexity of the algorithm, and the number of recursive calls required to make it work, along with the fact that the tree was being built as a multiprocessing.

In every scenario run, the python implementation builds a different tree to the tree built by our implementation from scratch, despite building from identical datasets. This can be attributed to the fact that the Weka implements a more optimised c4.5 algorithm called j48 [9], while our code builds a more basic implementation of the C4.5 algorithm. Despite the trees having different structures, they perform similarly and successfully on the testing data set.

## Dependencies

Pip Modules - numpy, pandas, graphviz, javabridge, python-weka-wrapper3, PySimpleGUI, weka

To initialize graphviz for use in generating graphs, open command prompt as admin, type 'dot -c', and then restart your code env

## Work Division between Team Members

We helped each other out, and paired to solve a lot of issues we were facing in the development of this C4.5 implementation. We did each lead the development for different aspects of the code, as shown in the breakdown below:

function	Written by	function	Written by
main()	louise	split_data_styles(data)	Aideen
processTimes(time)	Aideen	split_data_training_testing(data, ratio)	Aideen
createTreeWhileShowingLoadingWindow(data, split)	Aideen	gather_data()	louise
createTree(data, data_split, quit, queue)	louise	entropy(dataset)	Aideen
renderLoadingWindow(quit)	Aideen	information_gain(train_target, subsets)	louise
getInputData()	Aideen	test_tree(root_node, testing_data, split)	Aideen
errorWindow(text)	Aideen	test_data(data, node, test_results)	Aideen
resize_images()	Aideen	test_lr(node, row)	louise
DisplayTreesPopup(python_accuracy, weka_accuracy, p_time, w_time)	Aideen	class Node	louise
build_tree(data, attributes)	louise	if __name__ == '__main__':	louise
print_tree(root_node)	louise	build_weka_tree(split)	louise
addEl(node, g, rootname)	louise	build_weka(split)	louise
find_best_attribute(train_data, attributes)	louise	<a href="#">10tests.py</a>	Aideen
getMajorityClass(data)	louise		
split_into_subsets(column_header, training_data)	louise		

## References

- [1] Explanation of C4.5 Algorithm which we studied before beginning to construct our code [Accessed 13 November 2020].  
<https://towardsdatascience.com/what-is-the-c4-5-algorithm-and-how-does-it-work-2b971a9e7db0>
- [2] More In-depth explanation of Information gain studied while writing the code [Accessed 13 November 2020].  
<https://machinelearningmastery.com/information-gain-and-mutual-information/>
- [3] Python Weka Wrapper documentation used in building the reference algorithm [Accessed 15 November 2020].  
<https://fracpete.github.io/python-weka-wrapper3/index.html>
- [4] PySimpleGUI documentation used in building the Python User Interface [Accessed 16 November 2020].  
<https://pysimplegui.readthedocs.io/en/latest/>
- [5] Python Multiprocessing documentation used when adding the multiprocessing component to animate the loading gif while the tree is built [Accessed 16 November 2020].  
<https://docs.python.org/2/library/multiprocessing.html>
- [6] Machine learning with python: Decision Trees in python [Accessed 14 November 2020].  
[https://www.python-course.eu/Decision\\_Trees.php](https://www.python-course.eu/Decision_Trees.php)
- [7] Examples - graphviz 0.15 documentation [Accessed 17 November 2020].  
<https://graphviz.readthedocs.io/en/stable/examples.html#btree-py>
- [8] Decision Trees - C4.5 | octavians blog [Accessed 19 November 2020].  
<https://octaviansima.wordpress.com/2011/03/25/decision-trees-c4-5/>
- [9] ResearchGate. 2020. *Figure 1. The Typical Decision Tree Model. J48 Are The Improved....* [online] Available at: <https://www.researchgate.net/figure/The-typical-decision-tree-model-J48-are-the-improved-versions-of-C45-algorithms-or-can-fig1-288807267#:~:text=decision%20tree%20model.-,J48%20are%20the%20improved%20versions%20of%20C4.,J48%20is%20the%20Decision%20tree.>> [Accessed 2 December 2020].

```

1  # implementation_from_scratch.py
2  # Main file for the implementation of the C4.5 algorithm
3  import pandas as pd
4  import numpy as np
5  import graphviz
6  from graphviz import Digraph
7  from math import log2
8  from copy import deepcopy
9  from weka_implementation import build_weka_tree
10 import PySimpleGUI as sg
11 from PIL import Image, ImageTk, ImageSequence
12 from multiprocessing import Process, Queue
13 import multiprocessing as mp
14 import time
15 import os.path
16 from os import path
17 import PIL.Image
18 import io
19 import base64
20 import math
21
22
23 # louise Kilheeney -16100463
24 def main():
25
26     # Get the data, the train/test split percentage and the filepath of the data location
27     data, split = gather_data()
28
29     # While animate a 'loading' gif to show that the process is running,
30     # Build a C4.5 tree from the data.
31     # Return the root node of the tree, the dataset to use in testing and the time it
32     # took to build the tree
33     root_node, testing_data, python_time_to_build =
34     createTreeWhileShowingLoadingWindow(data, split)
35
36     # Draw the tree and store it in png format
37     print_tree(root_node)
38
39     # Calculate the accuracy of the built tree using the testing data
40     python_accuracy = test_tree(root_node, testing_data, split)
41
42     # Build a tree using the same data in weka.
43     # Draw the weka tree and store it in png format
44     # Get the accuracy of the weka tree, and the time it took to construct
45     weka_accuracy, weka_time_to_build = build_weka_tree(split)
46
47     # Scale and format the times it took to run each implementation
48     python_time = processTimes(python_time_to_build)
49     weka_time = processTimes(weka_time_to_build)
50
51     # Display both trees side to side, with their accuracies, and the time it took to
52     # build them
53     DisplayTreesPopup(python_accuracy, weka_accuracy, python_time, weka_time)
54
55 # Aideen McLoughlin - 17346123
56 # Take in the time it took to run a process in seconds
57 # and return a nicely scaled representation of that value
58 def processTimes(time):
59     if time < 1:
60         scaledtime = round(time * 1000)
61         timemsg = str(scaledtime) + " milliseconds"
62     elif time <= 60:
63         timemsg = str(round(time)) + " seconds"
64     else:
65         minutes = math.floor(time/60)
66         seconds = round(time - minutes*60)

```

```

65         timemsg = str(minutes) + " minutes and "+str(seconds)+" seconds"
66     return "The tree took " + timemsg + " to build"
67
68
69 # Aideen McLoughlin - 17346123
70 # Using python multiprocessing, build a tree while showing a 'loading' animation
71 def createTreeWhileShowingLoadingWindow(data, split):
72
73     # Create an 'Event' to indicate when the tree is built
74     quit = mp.Event()
75     # Create a 'Queue' to store generated values in
76     Q = Queue()
77
78     # Define both processes in the multiprocessing - the tree creation and the loading
    animation
79     p1 = Process(target = createTree, args=(data, split, quit, Q,))
80     p2 = Process(target = renderLoadingWindow, args=(quit, ))
81
82     # Store the time before starting to build the tree
83     starttime = time.time()
84
85     # Start both processes
86     p1.start()
87     p2.start()
88
89     # Wait for the quit event to be set, which will happen once the tree is built
90     quit.wait()
91
92     # Store the time once the tree has been built
93     endtime = time.time()
94
95     # Get the data stored in the Queue
96     queue_data = Q.get()
97
98     # Return the Queue data and the time to build
99     return queue_data[0], queue_data[1], endtime-starttime
100
101
102 # louise Kilheeney -16100463
103 def createTree(data, data_split, quit, queue):
104
105     #call function to split the data into training and test data.
106     train_data, test_data = split_data_training_testing(data, (data_split))
107
108     #list of attributes in the data
109     attributes =
110     ['calorific_value','nitrogen','turbidity','style','alcohol','sugars','bitterness','co
    lour','degree_of_fermentation']
111
112     #calling function to build tree with the training data and list of attributes
113     root_node = build_tree(train_data, attributes)
114
115     queue.put([root_node, test_data])
116     queue.cancel_join_thread()
117     quit.set()
118     return True
119
120 # Aideen McLoughlin - 17346123
121 def renderLoadingWindow(quit):
122     # Declare the PySimpleGUI layout for the popup window
123     layout_loading = [[sg.Text("Loading")],[sg.Image(r'loading.gif', key='-IMAGE-')]]
124
125     # Create the popup window
126     window = sg.Window("Building C4.5 Tree", layout_loading, element_justification='c',
    margins=(0,0), element_padding=(0,0), finalize=True)
127

```

```

128 # Animate the loading gif for the duration of time that 'quit' is not set
129 interframe_duration = Image.open(r'loading.gif').info['duration']
130 while not quit.is_set():
131     event, values = window.read(timeout=interframe_duration)
132     if event == sg.WIN_CLOSED:
133         exit()
134         break
135
136     window.FindElement("-IMAGE-").UpdateAnimation("loading.gif",time_between_frames=i
137         nterframe_duration)
138
139 # Close the popup window
140 window.close()
141 return True
142
143 # Aideen McLoughlin - 17346123
144 def getInputData():
145     # Declare the PySimpleGUI layout for the popup window
146     layout = [
147         [sg.Text('Data file: '), sg.InputText("beer.csv", key="file")],
148         [sg.Text("Select Train/Test Data Split")],
149         [sg.Radio('1/3',"1", key="1/3"),
150          sg.Radio('1/2',"1", key="1/2"),
151          sg.Radio('2/3',"1", key="2/3", default=True)],
152         [sg.Button('Ok'), sg.Button('Quit')]]
153
154 # Create the popup window
155 window = sg.Window("Select Train/Test Split", layout)
156
157 # Loop until the window is closed with 'x' 'OK' or 'Quit'
158 # Set split to be the selected train/test split value
159 # And set the filepath to be the contents of the InputText box (default beer.csv)
160 while True:
161     event, values = window.read()
162     if event == sg.WIN_CLOSED or event == 'Quit':
163         break
164     if event == 'Ok':
165         if values["1/3"] == True:
166             split = (1/3)
167         elif values["1/2"] == True:
168             split = (1/2)
169         elif values["2/3"] == True:
170             split = (2/3)
171         break
172     filepath = values["file"]
173
174 # Close the popup window
175 window.close()
176 return split, filepath
177
178 # Aideen McLoughlin - 17346123
179 def errorWindow(text):
180     # Declare the PySimpleGUI layout for the popup window
181     layout = [[sg.Text(text)],[sg.Button('Ok')]]
182
183 # Create the popup window
184 window = sg.Window("Error", layout)
185
186 # Display a popup window with the text passed as a function param, until the user
187 # closes the window
188 while True:
189     event, values = window.read()
190     if event == sg.WIN_CLOSED or event == 'Ok':
191         break

```



```

192
193     # Close the popup window
194     window.close()
195
196
197 # Aideen McLoughlin - 17346123
198 # Resize the tree png images to be equal and fit nicely into the popup window
199 def resize_images():
200     for image in (r'weka-test.gv.png', r'test.gv.png'):
201         img = PIL.Image.open(image)
202         img = img.resize((400, 400), PIL.Image.ANTIALIAS)
203         img.save(image, format="PNG")
204
205
206 # Aideen McLoughlin - 17346123
207 def DisplayTreesPopup(python_accuracy, weka_accuracy, p_time, w_time):
208
209     # Resize the tree images to be the right size for the popup
210     resize_images()
211
212     # Define 2 columns for the layout
213     # The right one for the custome python implementation from scratch
214     # The left one for the implementation with weka
215     weka_column = [
216         [sg.Text("Weka Implementation",font=('Helvetica 20'))],
217         [sg.Image(r'weka-test.gv.png',key='-IMAGE-')],
218         [sg.Text("Accuracy: "+str(weka_accuracy)+"%")],
219         [sg.Text(w_time)]
220     ]
221     python_column = [
222         [sg.Text("Our Python Implementation",font=('Helvetica 20'))],
223         [sg.Image(r'test.gv.png',key='-IMAGE-')],
224         [sg.Text("Accuracy: "+str(python_accuracy*100)+"%")],
225         [sg.Text(p_time)]
226     ]
227
228     # Declare the PySimpleGUI layout for the popup window with the two columns, and a
229     # QUIT button
230     layout = [
231         [
232             sg.Column(weka_column),
233             sg.Column(python_column)
234         ],
235         [sg.Button('Quit')],
236     ]
237
238     # Create the popup window
239     window = sg.Window("Generated C4.5 Tree", layout)
240
241     # Display the popup window until the user closes it
242     while True:
243         event, values = window.read()
244         if event == sg.WIN_CLOSED or event == 'Quit': # if user closes window or clicks
245             # cancel
246             break
247
248     # Close the popup window
249     window.close()
250
251 # louise Kilheeneey - 16100463
252 def build_tree(data, attributes):
253     #1. check the base cases
254     #• All the examples from the training set belong to the same class ( a tree
255     #   leaf labeled with that class is returned ).
256     if data['style'].nunique() == 1:
257         return Node(True, data['style'][0], None)

```

```

256     #• The training set is empty ( returns a tree leaf called failure ).
257     if len(data) == 0:
258         return Node(True, "Fail", None)
259
260     #2. find attribute with highest info gain, retrun best_attribute
261     # calling function find-best-attribute which retruns the best attribute, the
        attribute subsets and the threshold divisor
262     best_attribute, attribute_subsets, threshold_divisor = find_best_attribute(data,
        attributes)
263
264     #if best attribute is empty
265     if best_attribute == "":
266         #calling function get majorityclass to return the majority class
267         majClass = getMajorityClass(data)
268         return Node(True, majClass, None)
269     else:
270         #3. split the set (data) in subsets arrcording to value of best_attribute
271         #attribute_subsets = split_into_subsets(best_attribute, data)
272         remainColumns = deepcopy(data)
273         remainColumns = data.drop(columns=[best_attribute])
274
275         #4. repeat steps for each subset
276         node = Node(False, best_attribute, threshold_divisor)
277         for attr_subset in attribute_subsets:
278             node.children.append(build_tree(attr_subset, remainColumns))
279
280         return node
281
282
283 # Louise Kilheeney - 16100463
284 # Generate a png of the tree from the root node
285 def print_tree(root_node):
286
287     # Create a digraph in which to store the tree data
288     g = Digraph('python_tree_implementation')
289
290     # Add the root node, and all its children recursively
291     addEl(root_node, g, 'a')
292
293     # Format the graph as a png, and save it
294     g.format = "png"
295     g.render('test.gv', view=False)
296
297
298 # Louise Kilheeney - 16100463
299 # Add a node to the tree
300 def addEl(node, g, rootname):
301
302     # If the node is not a leaf
303     if not node.isLeaf:
304         #Create the node
305         g.node(name=str(rootname), label=node.label)
306
307         # Create an edge from the node to its left child
308         nodename1 = rootname+'b'
309         g.edge(rootname, nodename1, label="<= "+str(round(node.divisor,2)))
310         # Recursively add the nodes left child
311         addEl(node.children[0],g,nodename1)
312
313         # Create an edge from the node to its right child
314         nodenamec = rootname + 'c'
315         g.edge(rootname, nodenamec, label="> "+str(round(node.divisor,2)))
316         # Recursively add the nodes right child
317         addEl(node.children[1],g,nodenamec)
318     else:
319         # If the node is a leaf, add it while styling it as a leaf
320         g.node(name=rootname, label=node.label, shape='box', style='filled')

```

```

321
322
323 # Louise Kilheeneey - 16100463
324 def find_best_attribute(train_data, attributes):
325     # Returns the best attribute from all
326     best_information_gain = 0
327     best_attribute = ""
328     threshold_divisor = ""
329     subsets = []
330     for attribute in attributes:
331         #making sure not to include style
332         if attribute != 'style':
333
334             #calling function split_into_subsets
335             temp_subsets, temp_divisor = split_into_subsets(attribute, train_data)
336
337             # temp gain is equal to the information gain function for the train_data
338             # and the subsets.
339             temp_gain = information_gain(train_data, temp_subsets)
340
341             #check for the best attribute
342             if temp_gain > best_information_gain:
343                 best_attribute = attribute
344                 best_information_gain = temp_gain
345                 subsets = temp_subsets
346                 threshold_divisor = temp_divisor
347             # return the best attribute , subsets and the threshold divisor
348             return best_attribute, subsets, threshold_divisor
349
350 # Louise Kilheeneey - 16100463
351 # Function to get the majority class of the data been passed in.
352 def getMajorityClass(data):
353     #find the majority class in the data with the data - style
354     grouped = data.groupby(data['style'])
355     #return majority class
356     return max(grouped.groups)
357
358
359 # Louise Kilheeneey - 16100463
360 def split_into_subsets(column_header, training_data):
361     split_values = []
362     maxEnt = -1*float("inf")
363     best_threshold = ""
364     sorted_data = training_data.sort_values(by=[column_header])
365     for item in range(0, len(training_data) - 1):
366         if type(sorted_data.iloc[item][column_header]) != 'style':
367             if sorted_data.iloc[item][column_header] !=
368                 sorted_data.iloc[item+1][column_header]:
369                 threshold = (sorted_data.iloc[item][column_header] +
370                     sorted_data.iloc[item+1][column_header]) / 2
371                 smaller_than_threshold = pd.DataFrame()
372                 bigger_than_threshold = pd.DataFrame()
373                 for index, row in sorted_data.iterrows():
374                     if(row[column_header] > threshold):
375                         bigger_than_threshold = bigger_than_threshold.append(row,
376                             ignore_index = True)
377                     else:
378                         smaller_than_threshold = smaller_than_threshold.append(row,
379                             ignore_index = True)
380
381                 igain = information_gain(training_data, [smaller_than_threshold,
382                     bigger_than_threshold])
383
384                 if igain >= maxEnt:
385                     split_values = [smaller_than_threshold, bigger_than_threshold]
386                     best_threshold = threshold

```

```

382         maxEnt = igain
383     return split_values, best_threshold
384
385
386 # Aideen McLoughlin - 17346123
387 # Split the python DataFrame object into 3 DataFrame objects
388 # One storing all the values with the style 'ale'
389 # One storing all the values with the style 'lager'
390 # and One storing all the values with the style 'stout'
391 def split_data_styles(data):
392
393     # Declare empty subsets array
394     subsets = {}
395
396     # Group the data passed to the function by its style
397     grouped = data.groupby(data['style'])
398
399     # For each style name, add the values in that group to the subsets array as a new
400     # DataFrame object
401     for index, beer_style in enumerate(['ale', 'lager', 'stout']):
402         if beer_style in grouped.groups.keys():
403             subsets[index] = grouped.get_group(beer_style)
404         else:
405             subsets[index] = {}
406
407     # return the subsets array of DataFrame objects
408     return subsets
409
410 # Aideen McLoughlin - 17346123
411 # Split the data into training and testing datasets
412 def split_data_training_testing(data, ratio):
413
414     # Drop the beer_id column as it is not relevant to the beer style
415     data = data.drop(columns=['beer_id'])
416
417     # Get a random sample from the data file as the training data
418     # This data will be ratio% of the initial dataset
419     train = data.sample(frac=ratio, random_state=5)
420
421     # Get the rest of the dataset values as the testing data
422     test = data.merge(train, how='left', indicator=True)
423     test = test[(test['_merge'] == 'left_only')].copy()
424     test = test.drop(columns='_merge').copy()
425
426     # Save the divided dataset so that it can be used in the weka implementation
427     train.to_csv('train_data_generated.csv', index=False, header=True)
428     test.to_csv('test_data_generated.csv', index=False, header=True)
429
430     # Create a folder to store the results in if it does not already exist
431     if not path.exists('results'):
432         print("Create results folder")
433         os.mkdir('results')
434
435     # Return the training and testing data
436     return train, test
437
438 # Louise Kilheeney - 16100463
439 # Get the filepath of the data file, and the train/test data split fro user input in a
440 # PySimpleGUI popup
441 # If a filepath provided is not valid, prompt the user to input a new filepath.
442 # Repeat until a valid filepath is provided
443 def gather_data():
444
445     # Get the train/test split and the filepath of the data file
446     split, filepath = getInputData()

```

```

447     # Create an empty pandas dataframe element
448     data = pd.DataFrame()
449
450     # While the dataframe element remains empty
451     while data.empty:
452
453         # Check If the filepath is valid
454         if path.isfile(filepath):
455             # If it is, set the data to be the csv data at that filepath
456             data = pd.read_csv(filepath)
457         else:
458             # If it is not valid, Display an error pop-up and prompt the user to input
459             # the split and filepath again
460             errorWindow("File not found, please try again")
461             split, filepath = getInputData()
462
463     # Once the dataframe element is filled, return the data, the train/test split
464     # percentage and the filepath (For use in the weka implementation)
465     return data, split
466
467 # Aideen McLoughlin - 17346123
468 # Calculate the entropy of the passed data set
469 def entropy(dataset):
470
471     # Initialise entropy to zero value
472     entropy = 0
473
474     # Get the ale, lager and stout subset DataFrames from the passed dataset
475     subsets = split_data_styles(dataset)
476
477     # For each subset
478     for index in range(len(subsets)):
479
480         # Get the percentage of the dataset which is in the subset
481         probability = len(subsets[index]) / len(dataset)
482
483         # If the probability is not zero,
484         # Subtract plog2(p) from the entropy value where p is probability
485         if probability != 0:
486             entropy = entropy - (probability)*log2(probability)
487
488     # Return the entropy value
489     return entropy
490
491 # Louise Kilheeney - 16100463
492 # function to calculate the information gain
493 def information_gain(train_target, subsets):
494     #getting the entropy value of the train_target
495     entropyTarget = entropy(train_target)
496     total = len(train_target)
497
498     Gain = entropyTarget
499
500     #for each subset
501     for i in range(0, len(subsets)):
502         #length of each subset
503         numBeer = len(subsets[i])
504         #Gain = Ent(S) - |S beer =ale|/|S|*( Ent( S beer=ale )) - |S beer=stout
505         #|/|S|*( Ent( S beer=stout )) - |S beer=lager|/|S|*( Ent( S beer=lager )) - ....
506         firstPart = numBeer/total
507         secondPart = entropy(subsets[i])
508         Gain -= (firstPart*secondPart)
509     # Return the information gain value
510     return Gain

```

```

511
512 # Aideen McLoughlin - 17346123
513 # Test the built tree using the testing data
514 def test_tree(root_node, testing_data, split):
515
516     # Get the style as the target values, and then drop them from the testing dataset
517     test_target = testing_data['style'].values
518     testing_data = testing_data.drop(columns=['style'])
519
520     # get the results of the tree predictions for all the testing data values
521     test_results = test_data(testing_data, root_node, [])
522
523     # Initialise the number of correct entries to 0
524     correct = 0
525
526     accurate = []
527     # For each test result, check if it is accurate using the style values we removed
528     # from the Dataframe earlier
529     # Keep a count of the number of correct predictions
530     # If the prediction is wrong, print the incorrect prediction
531     for index in range(0, len(test_results)):
532         if test_results[index] == test_target[index]:
533             correct = correct + 1
534             accurate.append(1)
535         else:
536             accurate.append(0)
537
538     # Calculate the accuracy to 2 decimal places, and return it
539     accuracy = round(correct/len(test_results), 2)
540
541     df = pd.DataFrame()
542     df['Actual'] = test_target
543     df['Predicted'] = test_results
544     df['Accuracy'] = accurate
545
546     filename =
547     "results/python-results-"+str(round(split, 2))+ "-" + str(round(accuracy, 2)) + ".csv"
548     df.to_csv(filename, index=False, header=True)
549
550     return accuracy
551
552 # Aideen McLoughlin - 17346123
553 # Using the root node of the constructed tree, predict the output of all the test data
554 # inputs
555 def test_data(data, node, test_results):
556
557     # For each data value, get the predicted result
558     for item in range(0, len(data)):
559         test_results.append(test_lr(node, data.iloc[item]))
560
561     # return the set of all predicted results
562     return test_results
563
564 # Louise Kilheeneey - 16100463
565 # Get the final leaf node destination for a data row
566 def test_lr(node, row):
567
568     # Decide which child path of a node to proceed into, based on the input value.
569     # This function will call itself recursively until it reaches a leaf node
570     # That leaf node will be returned to the function which called it
571     if node.isLeaf:
572         return node.label
573     else:
574         if row[node.label] <= node.divisor:
575             return test_lr(node.children[0], row)

```

```
575         else:
576             return test_lr(node.children[1], row)
577
578
579 # Louise Kilheeney - 16100463
580 # Node class
581 class Node:
582     def __init__(self, isLeaf, label, divisor):
583         self.label = label
584         self.isLeaf = isLeaf
585         self.divisor = divisor
586         self.children = []
587
588
589 # Louise Kilheeney - 16100463
590 # Set the main function to run when the file is run
591 if __name__ == '__main__':
592     main()
593
```

```

1  # weka_implementation.py
2
3  from weka.classifiers import Classifier
4  from weka.core.converters import Loader
5  import weka.core.jvm as jvm
6  from weka.core.dataset import create_instances_from_lists
7  from weka.filters import Filter
8  from weka.core.classes import Random
9  from weka.classifiers import Classifier, Evaluation, PredictionOutput
10 import weka.plot.graph as graph
11 import graphviz
12 import time
13 import pandas as pd
14 import numpy as np
15 from graphviz import Digraph
16
17 # Louise Kilheeney - 16100463
18 # Taking in the data file location, and the train/test split proportion
19 # Build a weka C4.5 implementation using the Python Weka Wrapper API
20 def build_weka_tree(split):
21     jvm.start()
22
23     accuracy, time_to_build = build_weka(split)
24
25     jvm.stop()
26     return accuracy, time_to_build
27
28 # Louise Kilheeney - 16100463
29 def build_weka(split):
30     # Load the data file
31     loader = Loader(classname="weka.core.converters.CSVLoader")
32     train = loader.load_file('train_data_generated.csv')
33     test = loader.load_file('test_data_generated.csv')
34
35     # Store the target values for the test data
36     # so that the accuracy of the formula can be checked
37     test_target = pd.read_csv('test_data_generated.csv')['style'].values
38     # Get the dataset used to train the model,
39     # so that we can identify what the key values for the class are.
40     # As data is split randomly, we cannot assume it is in [ale, lager, stout] order
41     train_classes = pd.read_csv('train_data_generated.csv')['style'].values
42
43     # Set the class to be column 3 - the style column
44     train.class_index = 3
45
46     # Set the class to be column 3 - the style column
47     test.class_index = 3
48
49     # Store the time before starting to build the tree
50     starttime = time.time()
51
52     # initialise the time_to_run and accuracy to 0
53     time_to_run = 0
54     accuracy = 0
55
56     # Build and Train the weka tree
57     cls = Classifier(classname="weka.classifiers.trees.J48")
58
59     # Check that the data is valid
60     # If so, Build and Train the weka tree
61     if len(list(np.unique(train_classes))) != 1:
62         cls.build_classifier(train)
63         graph = cls.graph
64         # Store the time once the tree has been built
65         endtime = time.time()
66
67         # Create a list to store the predicted values in

```



```

68     pred = []
69     accurate = []
70     # Get the class labels in the order that they were allocated when training the
    model
71     classes = pd.Series(train_classes).drop_duplicates().tolist()
72
73     correct = 0
74     total = 0
75     # loop through test dataset, incrementing total every time
76     # and incrementing count if the predicted value was correct
77     for index, inst in enumerate(test):
78         total = total + 1
79         predicted = classes[int(cls.classify_instance(inst))]
80         pred.append(predicted)
81         act = test_target[index]
82         if predicted == test_target[index]:
83             correct = correct + 1
84             accurate.append(1)
85         else:
86             accurate.append(0)
87
88     # Get the accuracy of the weka implementation
89     accuracy = (correct/total)
90
91     # store the results in a csv file - the predicted class and the actual class
92     df = pd.DataFrame()
93     df['Actual'] = test_target
94     df['Predicted'] = pred
95     df['Accuracy'] = accurate
96     filename =
    "results/weka-results-"+str(round(split,2))+ "-" + str(round(accuracy,2)) + ".csv"
97     df.to_csv(filename, index=False, header=True)
98
99     time_to_run = endtime-starttime
100    # If the data is invalid, create a node to indicate failure
101    elif len(train) == 0:
102        graph = Digraph('python_tree_implementation')
103        graph.node(name='A', label="Fail", shape='box', style='filled')
104    else:
105        graph = Digraph('python_tree_implementation')
106        graph.node(name='A', label=train_classes[0], shape='box', style='filled')
107
108    # Render a png image of the weka tree to display in the PySimpleGUI popup
109    g = graphviz.Source(graph)
110    g.format = "png"
111    g.render('weka-test.gv', view=False)
112    return round(accuracy*100, 2), time_to_run

```

```
1  # 10tests.py
2
3  # Aideen McLoughlin - 17346123
4
5  import pandas as pd
6  import numpy as np
7  from implementation_from_scratch import *
8  from weka_implementation import *
9  import weka.core.jvm as jvm
10 # Define a range of train/test split proportions
11 proportions = [0.05, 0.15, 0.25, 0.35, 0.45, 0.55, 0.65, 0.75, 0.85, 0.95]
12
13 jvm.start()
14
15 # For each proportion
16 for split in proportions:
17     # Read in the data
18     data = pd.read_csv('beer.csv')
19
20     # declare Queue for Tree output storage
21     Q = Queue()
22     # Create tree
23     createTree(data, split, mp.Event(), Q)
24
25     # Get the tree output
26     queue_data = Q.get()
27     root_node = queue_data[0]
28     test_data = queue_data[1]
29
30     # Get the accuracy of the tree and test it
31     python_accuracy = test_tree(root_node, test_data, split)
32
33     # Build and test the weka tree
34     weka_accuracy, weka_time_to_build = build_weka(split)
35
36 jvm.stop()
```