Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

# Team Project Computer Science 1001

## I.    Topic

This project objection is to make a program with the decision as the baseline. We are allowed to use any possible ways in order to meet the requirements of the dataset. For the dataset, there are two different datasets that we can choose either the Google Play Store or the Red Wine Quality. However, every dataset has their own tasks which are classification and regression. In this project, the aim of the classification method  is to predict whether the user rating of the app is above 4.5 or less. If we choose the classification, we also need to give the accuracy report. In contrast,  the regression method is used to predict the user rating of the app. If we choose the regression, we need to provide the report about the mean square error. However, for the classification, knowing the real value of the prediction is not necessary whereas for the regression is required.

In our project on the decision tree, we are interested in configuring the code of the second dataset which is the Red Wine Quality. Then, we also choose classification which is used to predict whether the user quality of the wine is above 6 or less. We also provide a report about the accuracy of the program. For us, this topic is seen as the most interesting to solve.

## II.    Approaches

Before we started to make the project, our team did some research in order to make some consideration about the method that we are going to use for the  machine learning program. We searched the basic concept and logic of machine learning and these are what we found out; information gain, Entropy, Classification and Regression Tree (CART), GINI index, and decision tree to build the program.

The first step that we do is to find the class of the root. The purpose of this is to make the tree as well as the branches. As we find it, we are able to create the information gain which also allows us to create the tree and the branches. In order to construct a decision tree, information gain is the main point which has an aim to choose the most compatible class for the first branch of the root. It also decreases the 'uncertainty' of the result. In the matter of calculating the

Evan Justin                  Ryan Christopher                  Tantama Trisan
119010521                    119010507                         119010543

information gain, we first tried to use the entropy. Later then, we started to realize that entropy is not appropriate or not suitable for this task since it uses some complex mathematical functions.

Then we figure out another possible way to calculate the information gain. We then think about using other methods which later seem ineffective and inefficient. They also do not work well. However, we keep on trying to use another method and after some time, we finally found out that Gini index is the perfect solution to calculate the information gain.

## III.    Final Approach

The aim of our final program is to calculate the information gain using the Gini index. Gini index is the fastest method in the matter of computing the information gain  since it has no complex mathematical functions. Despite the fast calculation, Gini index also provides users the most accurate amount and less than the entropy index. Moreover, the Gini index can also calculate about how often an element or data will be incorrectly labeled when it was chosen at random and randomly labeled according to the distribution of labels in the subset. So that Gini index is the best method among others. Thus, we use this method in our program. By using the Gini index, the information gain can be calculated in order to find the best class to be placed as the first root and the next one.

## IV.    Working Principle and Functions

Our program consisted of 3 classes which include Question, Leaf and Decision. It also has 14 functions.

- load_csv(filename)

First our program will run the dataset with the build up function "reader" inside function "load_csv(filename)". Then our program changes all the data into float with 2D array.

```
4    def load_csv(filename):
5        file = open(filename,"rt")
6        lines = reader(file)
7        dataset = list(lines)
8        del dataset[0]
9        return dataset
10
11   traindatadirectory = 'train.csv'
12   testdatadirectory = 'test.csv'
13   training_data = [[float(y) for y in x] for x in load_csv(traindatadirectory)]
14   testing_data = [[float(y) for y in x] for x in load_csv(testdatadirectory)]
```

Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

- quality_classifier(dataset)

In the function "quality_classifier(dataset)", we classify the quality of each data from training data to ease the prediction we separated the quality of red wine from below 6 and above 6 and represent both with number 0 and 1

```python
17    def quality_classifier(dataset):
18        for i in range (len(dataset)):
19
20            temp = dataset[i]
21
22            for k in range(len(temp)):
23                temp[k] = float(temp[k])
24
25            if temp[-1] > 6:
26                temp[-1] = 1
27            else:
28                temp[-1] = 0
```

- data_count(training_data)

Then "data_count(training_data)" used to store the dictionary for quality with key 1 and 0 and increase the value by 1 if it meets the condition

```python
33    def data_count(training_data):
34        qualitydirectory = {0:0, 1:0}
35        for k in training_data:
36            if k[-1] == 0:
37                qualitydirectory[0] += 1
38            else:
39                qualitydirectory[1] += 1
40        return qualitydirectory
```

- values(rows,col) and num_check(value)

Then function "value(rows, col)" and "num_check(value)" are used to check each row and column to be int and float

```python
43    def values(rows, col):
44        return set([row[col] for row in rows])
45
46
47    def num_check(value):
48        return isinstance(value, int) or isinstance(value, float)
```

Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

- Class Question

Class "Question" is used to build a condition for the tree for each value in each column and row the question is based on the header chosen. Each column number and value stored are going to be compared in function match

```
50    class Question:
51
52        def __init__(self, column, value):
53            self.column = column
54            self.value = value
55
56        def match(self, example):
57            val = example[self.column]
58            if num_check(val):
59                return val >= self.value
60            else:
61                return val == self.value
62
63        def __repr__(self):
64            condition = "=="
65            if num_check(self.value):
66                condition = ">="
67            return "Is %s %s %s?" % (
68                header[self.column], condition, str(self.value))
```

- partition(rows, question)

Function "partition(rows, partition)" is used to check if each row match the question and separated into 2 group "Trow" and "Frow"

```
70    def partition(rows, question):
71        Trows, Frows = [], []
72        for row in rows:
73            if question.match(row):
74                Trows.append(row)
75            else:
76                Frows.append(row)
77        return Trows, Frows
```

- gini_index(rows)

Evan Justin                     Ryan Christopher                     Tantama Trisan
119010521                       119010507                            119010543

Function "gini_Index(rows)" is the most crucial part for this program. This function is used to calculate the impurity inside the dataset given. It started by calculating the probability of labels.

```python
80    def gini_index(rows):
81        counts = data_count(rows)
82        impurity = 1
83        for lbl in counts:
84            prob_of_lbl = counts[lbl] / float(len(rows))
85            impurity -= prob_of_lbl**2
86        return impurity
```

- info_gain(left, right, current_uncertainty)

Function "info_gain(left, right, current_uncertainty)" is used to calculate the uncertainty node by using the result from "gini_index(rows)" calculation

```python
88    def info_gain(left, right, current_uncertainty):
89
90        p = float(len(left)) / (len(left) + len(right))
91        return current_uncertainty - p * gini_index(left) - (1 - p) * gini_index(right)
```

- best_split(rows)

Function "best_split(rows)" is used to find the question by iterating the feature and value. Then we split the dataset from the question and then calculate the new information gain after the split

Evan Justin
119010521

Ryan Christopher
119010507

Tantama Trisan
119010543

```
93    def best_split(rows):
94        gain = 0
95        question_limit = None
96        current_uncertainty = gini_index(rows)
97        n_features = len(rows[0]) - 1
98
99        for col in range(n_features):
100           values = set([row[col] for row in rows])
101           for val in values:
102               question = Question(col, val)
103               Trows, Frows = partition(rows, question)
104               if len(Trows) == 0 or len(Frows) == 0:
105                   continue
106               gain = info_gain(Trows, Frows, current_uncertainty)
107               if gain >= gain:
108                   gain, question_limit = gain, question
109
110       return gain, question_limit
```

- Class Leaf

"Class Leaf" is to determine the prediction from function data count with the help of initializer function.

```
113    class Leaf:
114
115        def __init__(self, rows):
116            self.predictions = list(data_count(rows))[0]
```

- Class Decision_Node

"Class Decision Node" is used to generate question for nodes, false branch and true branch each represented by "Fbranch" and "TBranch"

Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

```
119    class Decision_Node:
120
121        def __init__(self,
122                        question,
123                        Tbranch,
124                        Fbranch):
125            self.question = question
126            self.Tbranch = Tbranch
127            self.Fbranch = Fbranch
```

- tree_function(rows)

Function "tree_function(rows)" is used to grow the decision tree from nodes given by the dataset header. It return the leaf of the tree and uses local variable named as "Trow" and "Frows".

```
130    def tree_function(rows):
131
132        gain, question = best_split(rows)
133
134        if gain == 0:
135            return Leaf(rows)
136        Trows, Frows = partition(rows, question)
137        Tbranch = tree_function(Trows)
138
139        Fbranch = tree_function(Frows)
140
141        return Decision_Node(question, Tbranch, Fbranch)
```

- print_tree(node, spacing = "")

Function "print_tree(node, spacing = "")" is used to print the result of the tree from root to nodes until the leaf which is the prediction from the actual quality.

Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

```
144    def print_tree(node, spacing=""):
145
146        if isinstance(node, Leaf):
147            print (spacing + "Predict", node.predictions)
148            return
149
150        print (spacing + str(node.question))
151        print (spacing + '--> True:')
152        print_tree(node.Tbranch, spacing + "  ")
153        print (spacing + '--> False:')
154        print_tree(node.Fbranch, spacing + "  ")
```

- classify(row, node)

Function "classify(row, node)" is a recursion function to return the classification from each row and nodes which represent by "node.Tbranch" or "node.Fbranch".

```
157    def classify(row, node):
158
159        if isinstance(node, Leaf):
160            return node.predictions
161
162        if node.question.match(row):
163            return classify(row, node.Tbranch)
164        else:
165            return classify(row, node.Fbranch)
```

- main function

Main function is a function to set the test sample of the program, calculate the prediction of the program. And lastly print the result of the tree from root to leaf, print total test sample, print correct prediction and calculate the accuracy. It also gave a comparison of the actual quality with the predicted quality with 0 and 1.

Evan Justin
119010521

Ryan Christopher
119010507

Tantama Trisan
119010543

```
167    if __name__ == '__main__':
168
169        my_tree = tree_function(training_data)
170        print_tree(my_tree)
171        test_sample = 0
172        correct_prediction=0
173
174 ∨     for row in testing_data:
175            prediction = classify(row,my_tree)
176            real = (row[-1])
177            print('Real', real, '|', prediction,'Predict')
178            if prediction == real:
179                test_sample+=1
180                correct_prediction+=1
181            else:
182                test_sample+=1
183        print("Results:")
184        print("Total sample data: " ,test_sample)
185        print("Total correct prediction: " ,correct_prediction)
186        accuracy = (correct_prediction/test_sample)*100
187        print("Accuracy: ",accuracy)
```

## V.    Summary of the result

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL
Is alcohol >= 11.5?
--> True:
  Is sulphates >= 0.7?
  --> True:
    Is sulphates >= 0.84?
    --> True:
      Predict 0
    --> False:
      Is sulphates >= 0.71?
      --> True:
        Is sulphates >= 0.83?
        --> True:
          Predict 0
        --> False:
          Is sulphates >= 0.72?
          --> True:
            Is sulphates >= 0.77?
            --> True:
              Predict 0
            --> False:
              Is pH >= 3.44?
              --> True:
                Predict 0
              --> False:
                Predict 0
          --> False:
            Predict 0
      --> False:
        Predict 0
  --> False:
    Is sulphates >= 0.59?
    --> True:
      Predict 0
    --> False:
      Is sulphates >= 0.55?
      --> True:
        Predict 0
```

Tree from root to leaf

Evan Justin                    Ryan Christopher                    Tantama Trisan
119010521                      119010507                          119010543

Comparison for actual quality and the code's prediction



The result from the program's code

## VI.    Conclusion

It can be seen clearly that our program has 86,667 on the accuracy based on the data that we provide. In the matter to predict whether the user quality of the wine is above 6 or less, our program is accurate enough. Thus, we can state that our program is a successful program.

## VII.    Workload

Evan Justin - 119010521 ---------- 33,3%

Ryan Christopher - 119010507 ---------- 33,4%

| Evan Justin | Ryan Christopher | Tantama Trisan |
|-------------|------------------|----------------|
| 119010521   | 119010507        | 119010543      |

Tantama - 119010543 ---------- 33,3%