

CS 6375

ASSIGNMENT 2

Decision Trees

Names of students in your group:

1. Rutuja Kaushike (rnk170000)
2. Akhila Kancharana (axk180025)

CS 6375.502

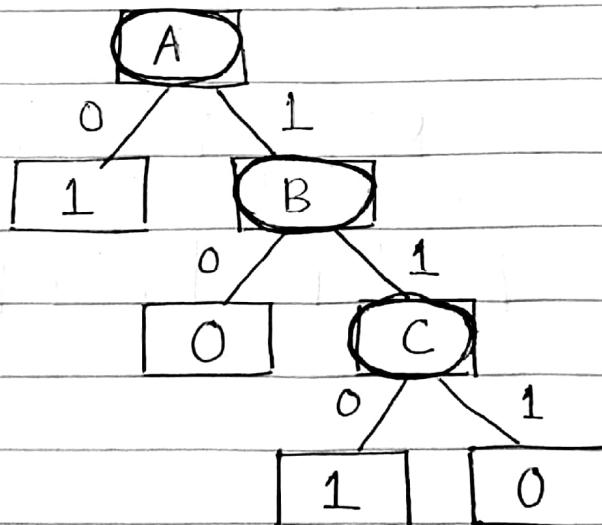
Number of free late days used: 0

Note: You are allowed a **total** of 4 free late days for the **entire semester**. You can use at most 2 for each assignment. After that, there will be a penalty of 10% for each late day.

Q1. Representing boolean functions:

(a) $Y = (\sim A \vee B) \wedge \sim(C \wedge A)$ where \sim represents NOT operator.

A	B	C	$\sim A$	$\sim A \vee B$	$C \wedge A$	$\sim(C \wedge A)$	$I \wedge III$
0	0	0	1	1	0	1	1
0	0	1	1	1	0	1	1
0	1	0	1	1	0	1	1
0	1	1	1	1	0	1	1
1	0	0	0	0	0	1	0
1	0	1	0	0	1	0	0
1	1	0	0	1	0	1	1
1	1	1	0	1	1	0	0



(b)

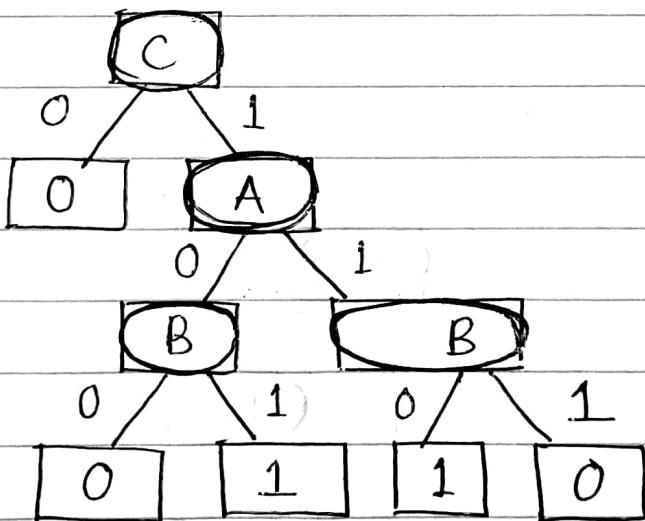
$$y = (A \oplus B) \wedge C$$

\oplus is XOR operator.

→

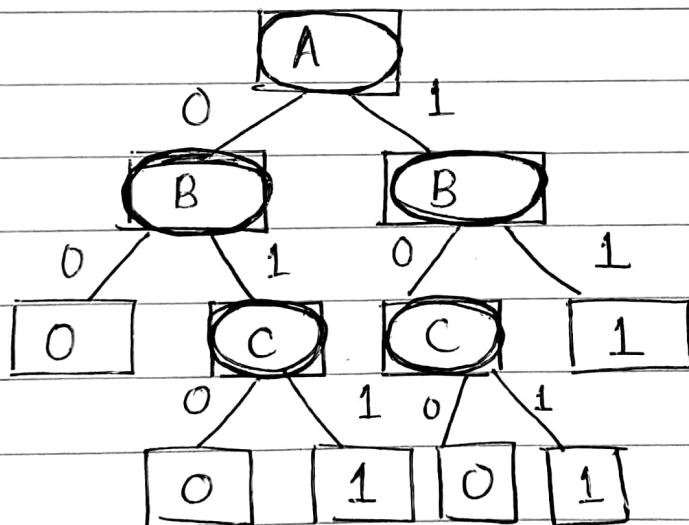
I

A	B	C	$A \oplus B$	$I \wedge C$
0	0	0	0	0
0	0	1	0	0
0	1	0	1	0
0	1	1	1	1
1	0	0	1	0
1	0	1	1	1
1	1	0	0	0
1	1	1	0	0



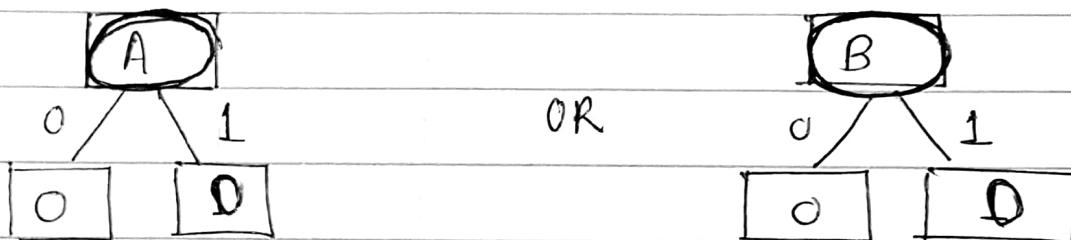
$$(c) \quad \psi = (A \vee B) \wedge (B \vee C) \wedge (A \vee C)$$

A	B	C	$A \vee B$	$B \vee C$	$A \vee C$	$I \wedge II \wedge III$
0	0	0	0	0	0	0
0	0	1	0	1	1	0
0	1	0	1	1	0	0
0	1	1	1	1	1	1
1	0	0	1	0	1	0
1	0	1	1	1	1	1
1	1	0	1	1	1	1
1	1	1	1	1	1	1



$$(a) Y = (A \vee B) \wedge \sim A \wedge \sim B.$$

		I		II		I \wedge II
A	B	$A \vee B$	$\sim A$	$\sim B$	$\sim A \wedge \sim B$	
0	0	0	1	1	1	0
0	1	1	1	0	0	0
1	0	1	0	1	0	0
1	1	1	0	0	0	0



Q2.

ID3 Algorithm

 \rightarrow

	X1	X2	X3	Class
	1	0	0	1
	0	1	0	1
	0	0	0	0
	1	0	1	0
	0	0	0	0
	1	1	0	1
	0	1	1	0
	1	0	0	1
	0	0	0	0
	1	0	0	1

Entropy of Class 4 :

$$\text{Entropy } [5+, 5-] = -\frac{5}{10} \log \frac{5}{10} - \frac{5}{10} \log \frac{5}{10}$$

$$= 1.$$

For X1 :

$$X_1 = 0$$

$$\text{Entropy } [1+, 4-] = -\frac{1}{5} \log \frac{1}{5} - \frac{4}{5} \log \frac{4}{5}$$

$$= -\left[\frac{1}{5} (-2.3219) + \frac{4}{5} (-0.3219) \right]$$

$$= (0.2)(2.3219) + (0.8)(0.3219)$$

$$= 0.46438 + 0.25752 \\ = 0.7219$$

$x_1 = 1$

$$\text{Entropy } [4+, 1-] = -\frac{4}{5} \log \frac{4}{5} - \frac{1}{5} \log \frac{1}{5}$$

Information Gain = Entropy(4) - weighted avg. of entropy of (x_1)

$$= 1 - \left(\frac{1}{2} \times 0.7219 + \frac{1}{2} \times 0.7219 \right)$$

$$= 1 - 0.7219$$

$$IG_{x_1} = 0.2781$$

for x_2 :

$x_2 = 0$

$$\text{Entropy } [3+, 4-] = -\frac{3}{7} \log \frac{3}{7} - \frac{4}{7} \log \frac{4}{7}$$

$$x_2 = 1 = 0.98527$$

$$\text{Entropy } [2+, 1-] = -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$= 0.9182$$

$$IG_{x_2} = \text{Entropy}(4) - \text{Weighted Avg}(x_2)$$

$$= 1 - \left(\frac{7}{10} \times 0.98527 + \frac{3}{10} \times 0.9182 \right)$$

$$= 1 - 0.965 = 0.035$$

for x_3

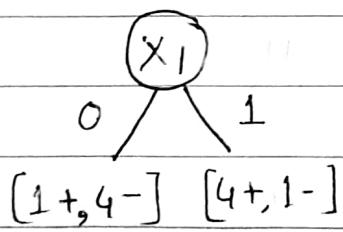
$$x_3 = 0$$

$$\text{Entropy } [5+, 3-] = -\frac{5}{8} \log \frac{5}{8} - \frac{3}{8} \log \frac{3}{8}$$
$$= 0.9543$$

$$\text{Entropy } [2+, 2-] = -\frac{2}{2} \log \left(\frac{2}{2}\right) = 0$$
$$x_3 = 1$$

$$\text{I.G.} = \text{Entropy}(4) - \text{Weighted Avg}(x_3)$$
$$x_3 = 1 - \left(\frac{8}{10} \times 0.9543 \right)$$
$$= 1 - 0.76348 = 0.23652$$

As IG from x_1 is the highest, root = x_1 .



consider node x_2

$$x_2 = 0 =$$
$$\underline{[0+, 3-]}$$

Now, let's classify further.

To classify again, we need to calculate Information gain for other attributes and select best classifier.

for splitting the tree.

Thus, for X_2 ,

$$X_2 = 0$$

$$\text{Entropy } [0+, 3-] = -0 \log_2 0 - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$X_2 = 1$$

$$\text{Entropy } [1+, 1-] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Information Gain} = 0.72192 - \frac{3}{5}(0) - \frac{2}{5}(1)$$

$$\therefore \text{Information Gain} = 0.32192$$

Now, let's solve for X_3 ,

$$\text{for } X_3 = 0$$

$$\begin{aligned} \text{Entropy } [0+, 3-] &= -\frac{1}{4} \log_2 \frac{1}{24} - \frac{3}{4} \log_2 \frac{3}{4} \\ &= 0.81127 \end{aligned}$$

$$\text{for } X_3 = 1$$

$$\text{Entropy } [0+, 1-] = -1 \cdot \log_2 1 \cdot = 0$$

(0.81127)

$$\text{Information Gain} = 0.72192 - \frac{4}{5} \cancel{\log_2 \frac{3}{4}}$$

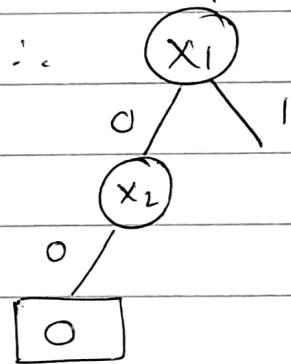
$$\therefore \text{Information Gain for } X_3 = 0.0729$$

Hence, now, we will split on $x_1=0$, we will consider x_2 as the next node.

for $x_1=0$ and $x_2=0$,
we have

x_1	x_2	x_3	y	i.e. for $x_1=0$ & $x_2=0$, $y=0$
0	1	0	1	$x_2=0$
0	0	0	0	Hence, our tree stops at
0	0	0	0	$x_2=0$
0	1	1	0	
0	0	0	0	

tree in progress will be:



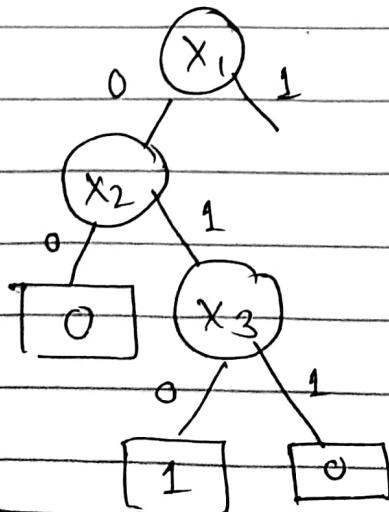
$x_1=0$, for $x_2=1$,
we have

x_1	x_2	x_3	y
0	1	0	1
0	1	1	0

i.e. When, $x_2=0$, $y=1$

\therefore our tree will be,

~~$x_2=1$, $y=0$~~



Now, for $x_1=1$,

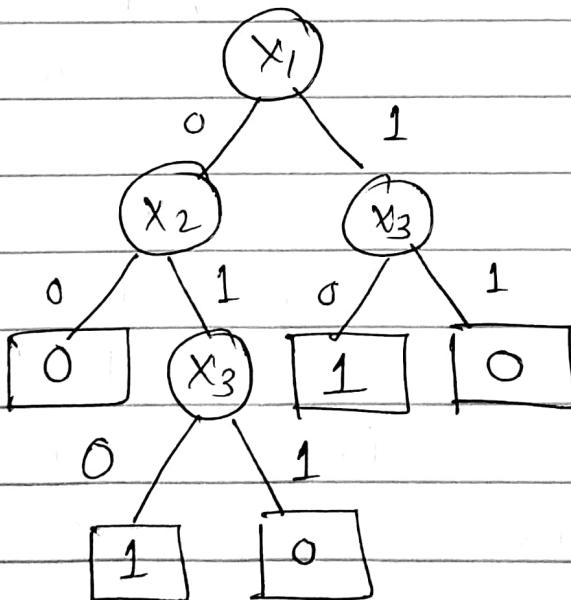
P70

we have

x_1	x_2	x_3	y
1	0	0	1
1	0	1	0
1	1	0	1
1	0	0	1
1	0	0	1

when, $x_1 = 1$, irrespective of what x_2 is
if $x_3 = 0, y = 1$
and $x_3 = 1, y = 0$

Hence, our final tree is



PART-2(Programming Part):

Q1) The code has the following dependencies, which need to be installed before running this code:

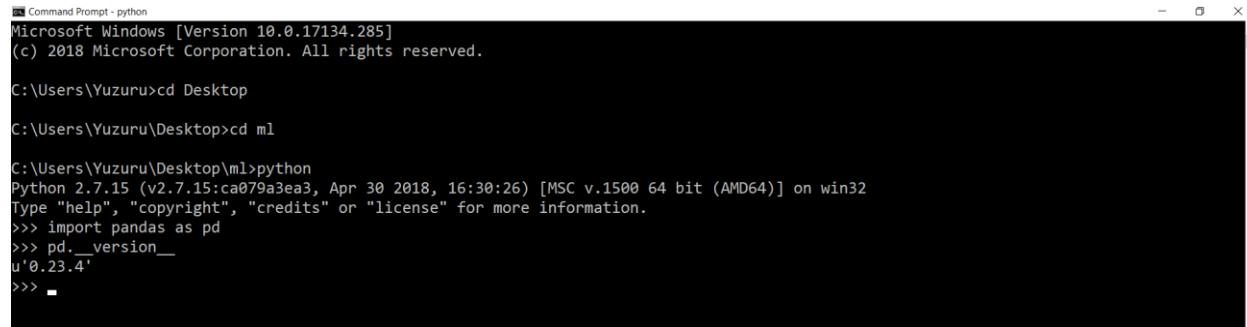
(a) Pandas. More details at: <https://pandas.pydata.org/>

(b) Scikit Learn for only one method in the driver code - train test split

You have to make changes only in the DecisionTree.py file to get everything in the driver.py to run.

A)

a) **Pandas(install):** Screen shot after installation and running pandas version

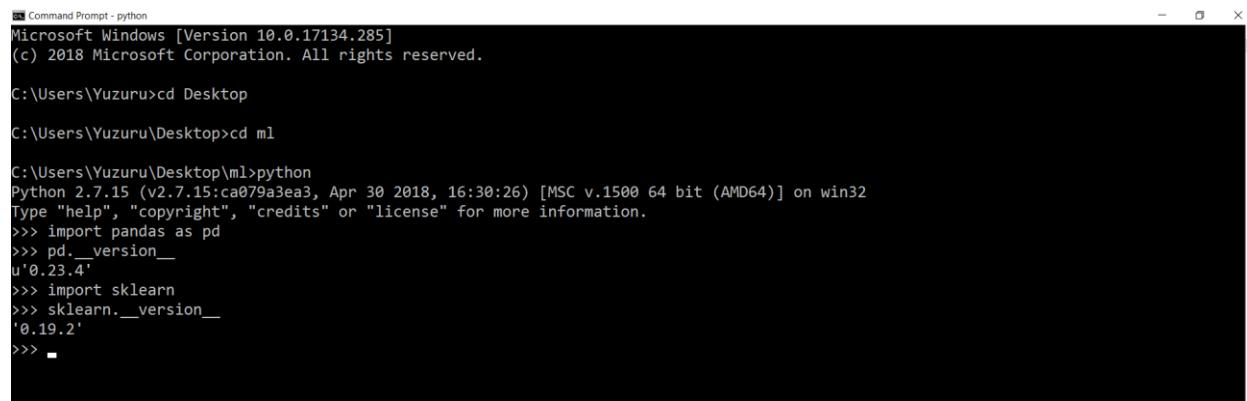


```
Command Prompt - python
Microsoft Windows [Version 10.0.17134.285]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Yuzuru>cd Desktop
C:\Users\Yuzuru\Desktop>cd ml

C:\Users\Yuzuru\Desktop\ml>python
Python 2.7.15 (v2.7.15:ca079a3ea3, Apr 30 2018, 16:30:26) [MSC v.1500 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import pandas as pd
>>> pd.__version__
u'0.23.4'
>>>
```

b) **Scikit (install):** Screen shot after installation and running scikit version.



```
Command Prompt - python
Microsoft Windows [Version 10.0.17134.285]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Yuzuru>cd Desktop
C:\Users\Yuzuru\Desktop>cd ml

C:\Users\Yuzuru\Desktop\ml>python
Python 2.7.15 (v2.7.15:ca079a3ea3, Apr 30 2018, 16:30:26) [MSC v.1500 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import pandas as pd
>>> pd.__version__
u'0.23.4'
>>> import sklearn
>>> sklearn.__version__
'0.19.2'
>>>
```

Q2) The code for each decision node in class Decision Node stores information about non-leaf nodes. Besides the information given in the code on line 278 of the original code, you need to add following properties for each node:

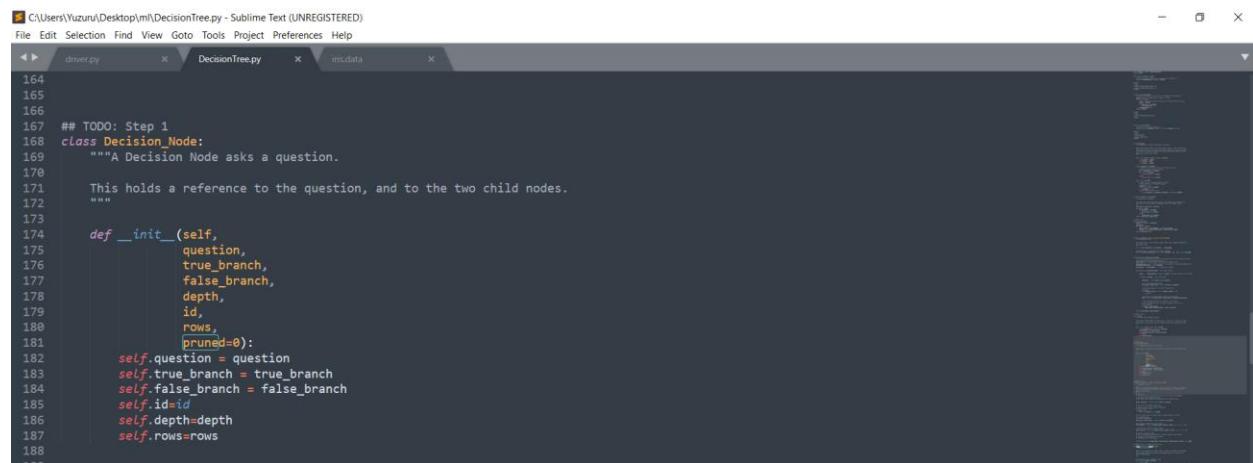
- depth of the node - starting from the root node, which is at depth 0
- id of the node, using the following schema:

The root node will have an id of 0. We assume that each node will split into two branches only - left and right. If a node has id of n, then its left child will have id of $2n + 1$ and right child will have id of $2n + 2$. This will ensure that there is no clash in node ids.

- the rows of data that it contains

A2) Printed id's and depth of nodes.

Code:



```
164
165
166
167 ## TODO: Step 1
168 class Decision_Node:
169     """A Decision Node asks a question.
170
171     This holds a reference to the question, and to the two child nodes.
172     """
173
174     def __init__(self,
175                  question,
176                  true_branch,
177                  false_branch,
178                  depth,
179                  id,
180                  rows,
181                  pruned=0):
182         self.question = question
183         self.true_branch = true_branch
184         self.false_branch = false_branch
185         self.id=id
186         self.depth=depth
187         self.rows=rows
188
189
```



```
271 ## TODO: Step 4
272 def print_tree(node, spacing=""):
273     """World's most elegant tree printing function."""
274
275     # Base case: we've reached a leaf
276     if isinstance(node, Leaf):
277         print (spacing + "Predict", node.predictions)
278         print (spacing+"Max Class",node.predictions_cls)
279         print (spacing+"<{id:"+str(node.id)+"}"+depth:"+str(node.depth)+"}")
280         return
281
282     # Print the question at this node
283     print (spacing + str(node.question))
284     print ('{id:' + str(node.id) + ", depth:" + str(node.depth) + "}")
285     # Call this function recursively on the true branch
286     print (spacing + "-> True:")
287     print_tree(node.true_branch, spacing + " ")
288
289     # Call this function recursively on the false branch
290     print (spacing + "-> False:")
291     print_tree(node.false_branch, spacing + " ")
```

Output:

```
Command Prompt
C:\Users\Yuzuru\Desktop\ml>python driver.py
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
  Is PetalW >= 1.8?
{id:2, depth:1}
--> True:
  Is Petall >= 4.9?
{id:6, depth:2}
--> True:
  Predict {'Iris-virginica': 43}
  Max Class Iris-virginica
  {id:14,depth:3}
--> False:
  Is SepalW >= 3.2?
{id:13, depth:3}
--> True:
  Predict {'Iris-versicolor': 1}
  Max Class Iris-versicolor
  {id:28,depth:4}
--> False:
  Predict {'Iris-virginica': 2}
  Max Class Iris-virginica
  {id:57,depth:4}
--> False:
  Is Petall >= 5.0?
{id:5, depth:2}
--> True:
  Is PetalW >= 1.6?
{id:12, depth:3}
--> True:
  Is Petall >= 5.8?
{id:26, depth:4}
--> True:
  Predict {'Iris-virginica': 1}
  Max Class Iris-virginica
  {id:54,depth:5}
--> False:
  Predict {'Iris-versicolor': 2}
  Max Class Iris-versicolor
  {id:53,depth:5}
--> False:
953 PM
9/21/2018
```

```
Command Prompt
id = 26 depth =4
id = 11 depth =3
*****Tree before pruning*****
Accuracy on test = 0.93
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
  Is Petall >= 4.8?
{id:2, depth:1}
--> True:
  Is Petall >= 5.0?
{id:6, depth:2}
--> True:
  Predict {'Iris-virginica': 33}
  Max Class Iris-virginica
  {id:14,depth:3}
--> False:
  Is PetalW >= 1.8?
{id:13, depth:3}
--> True:
  Is SepalW >= 3.2?
{id:28, depth:4}
--> True:
  Predict {'Iris-versicolor': 1}
  Max Class Iris-versicolor
  {id:58,depth:5}
--> False:
  Predict {'Iris-virginica': 4}
  Max Class Iris-virginica
  {id:57,depth:5}
--> False:
  Predict {'Iris-versicolor': 2}
  Max Class Iris-versicolor
  {id:27,depth:4}
--> False:
  Is PetalW >= 1.7?
{id:5, depth:2}
--> True:
  Predict {'Iris-virginica': 1}
  Max Class Iris-virginica
  {id:12,depth:3}
--> False:
952 PM
9/21/2018
```

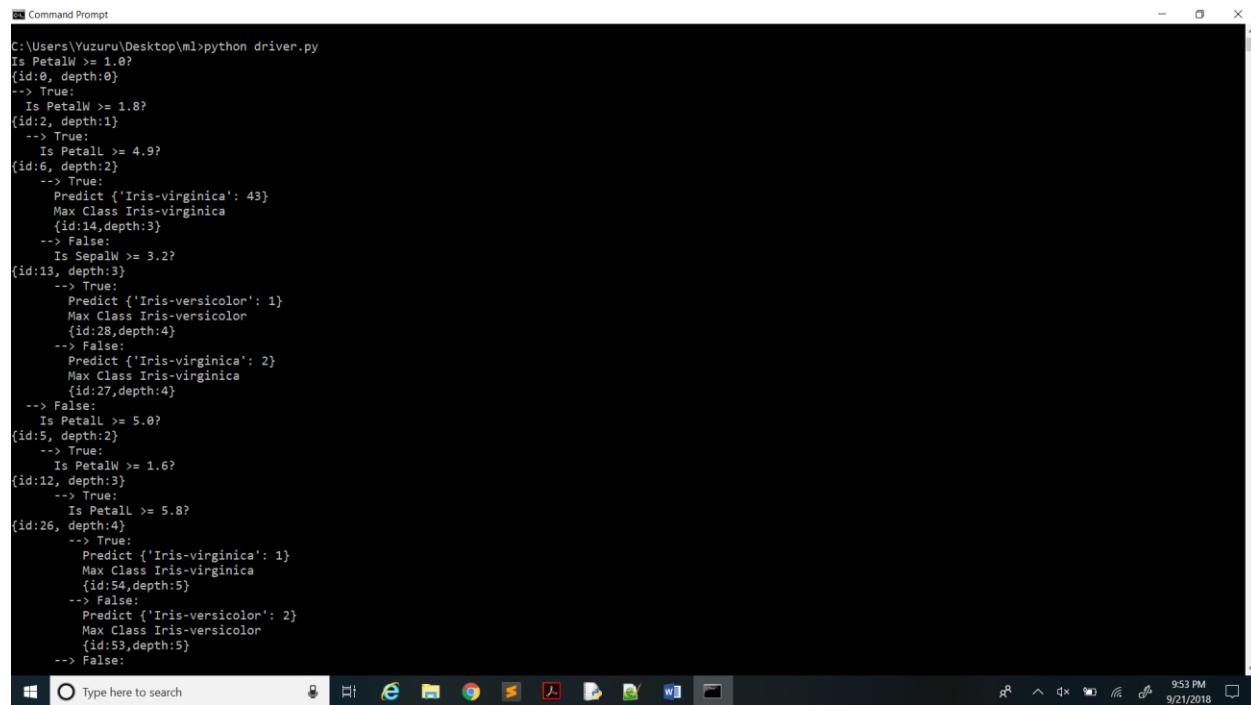
Q3) We will modify the definition of the leaf node in line 267 so that it stores _the following properties besides the ones given:

- predicted class label. Note that the predictions is a dictionary that contains the class counts. You will need to extract the majority class in it and output that.
- id of the leaf node
- depth of the leaf node

Code:

```
150 ## TODO: Step 2
151 class_Leaf:
152     """A Leaf node classifies data.
153
154     This holds a dictionary of class (e.g., "Apple") -> number of times
155     it appears in the rows from the training data that reach this leaf.
156     """
157
158     def __init__(self, rows, id, depth):
159         self.predictions = class_counts(rows)
160         count_temp=self.predictions
161         self.predictions_cls=max(count_temp)
162         self.id=id
163         self.depth=depth
164
```

Output:



```
C:\Users\Yuzuru\Desktop\ml>python driver.py
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
    Is PetalW >= 1.8?
{id:2, depth:1}
--> True:
    Is Petall >= 4.9?
{id:6, depth:2}
--> True:
    Predict {'Iris-virginica': 43}
    Max Class Iris-virginica
    {id:14,depth:3}
--> False:
    Is SepalW >= 3.2?
{id:13, depth:3}
--> True:
    Predict {'Iris-versicolor': 1}
    Max Class Iris-versicolor
    {id:28,depth:4}
--> False:
    Predict {'Iris-virginica': 2}
    Max Class Iris-virginica
    {id:27,depth:4}
--> False:
    Is Petall >= 5.0?
{id:5, depth:2}
--> True:
    Is PetalW >= 1.6?
{id:12, depth:3}
--> True:
    Is Petall >= 5.8?
{id:26, depth:4}
--> True:
    Predict {'Iris-virginica': 1}
    Max Class Iris-virginica
    {id:54,depth:5}
--> False:
    Predict {'Iris-versicolor': 2}
    Max Class Iris-versicolor
    {id:53,depth:5}
--> False:
```

```

Command Prompt
id = 26 depth =4
id = 11 depth =3
*****Tree before pruning*****
Accuracy on test = 0.93
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
| Is PetalL >= 4.8?
{id:2, depth:1}
--> True:
| Is PetalL >= 5.0?
{id:6, depth:2}
--> True:
| Predict {'Iris-virginica': 33}
| Max Class Iris-virginica
| {id:14,depth:3}
--> False:
| Is PetalW >= 1.8?
{id:13, depth:3}
--> True:
| Is SepalW >= 3.2?
{id:28, depth:4}
--> True:
| Predict {'Iris-versicolor': 1}
| Max Class Iris-versicolor
| {id:58,depth:5}
--> False:
| Predict {'Iris-virginica': 4}
| Max Class Iris-virginica
| {id:57,depth:5}
--> False:
| Predict {'Iris-versicolor': 2}
| Max Class Iris-versicolor
| {id:27,depth:4}
--> False:
| Is PetalW >= 1.7?
{id:5, depth:2}
--> True:
| Predict {'Iris-virginica': 1}
| Max Class Iris-virginica
| {id:12,depth:3}
--> False:

```

Q4) Next, you need to modify the build tree function. This is the main function that builds the tree from top (root mode) to bottom (leaf nodes) recursively. You need to modify the tree definition based on the conditions of parts 1 and 2 above. To construct each node, you need to pass besides the rows, other parameters such as depth, id, and header information. The header information is necessary because we will not hard code this information, like it is done in line 27 of the original code. The new definition of the build tree function would be as follows:

Make sure you finish and test this function before going ahead. Also, the example finds the best split using the find best split function that internally uses the function info gain function, which uses the Gini index for impurity (line 138 of original code). In class, we have studied Entropy measure. You would need to replace the gini function with an entropy function.

A4) Build Tree:

```
192     ## TODO: Step 3
193     def build_tree(rows, header, depth=0, id=0):
194         """Builds the tree.
195
196         Rules of recursion: 1) Believe that it works. 2) Start by checking
197         for the base case (no further information gain). 3) Prepare for
198         giant stack traces.
199         """
200
201         # depth = 0
202         # Try partitioning the dataset on each of the unique attribute,
203         # calculate the information gain,
204         # and return the question that produces the highest gain.
205
206         gain, question = find_best_split(rows, header)
207
208         # Base case: no further info gain
209         # Since we can ask no further questions,
210         # we'll return a leaf.
211         if gain == 0:
212             return Leaf(rows, id, depth)
213
214         # If we reach here, we have found a useful feature / value
215         # to partition on.
216         # nodelist.append(id)
217         true_rows, false_rows = partition(rows, question)
218
219         # Recursively build the true branch.
220         true_branch = build_tree(true_rows, header, depth + 1, 2 * id + 2 )
221
222         # Recursively build the false branch.
223         false_branch = build_tree(false_rows, header, depth + 1, 2 * id + 1)
224
225         # Return a Question node.
226         # This records the best feature / value to ask at this point,
227         # as well as the branches to follow
228         # depending on the answer.
229
230     return Decision_Node(question, true_branch, false_branch, depth, id, rows)
```

Entropy:

```
92     ## TODO: Step 3
93     def entropy(rows):
94         counts = class_counts(rows)
95         impurity=0
96         for lbl in counts:
97             prob_of_lbl = counts[lbl] / float(len(rows))
98             impurity += (prob_of_lbl*(math.log(prob_of_lbl,2)))
99
100     return (impurity**-1)
```

Q5) Now, we will modify the print tree function to print more details from the tree. Besides the current details that are outputted, we would also like to output the depth, node id of the decision nodes. For the leaf label, we need to see the predicted class label. You can output the predicted dictionary also, but that is optional.

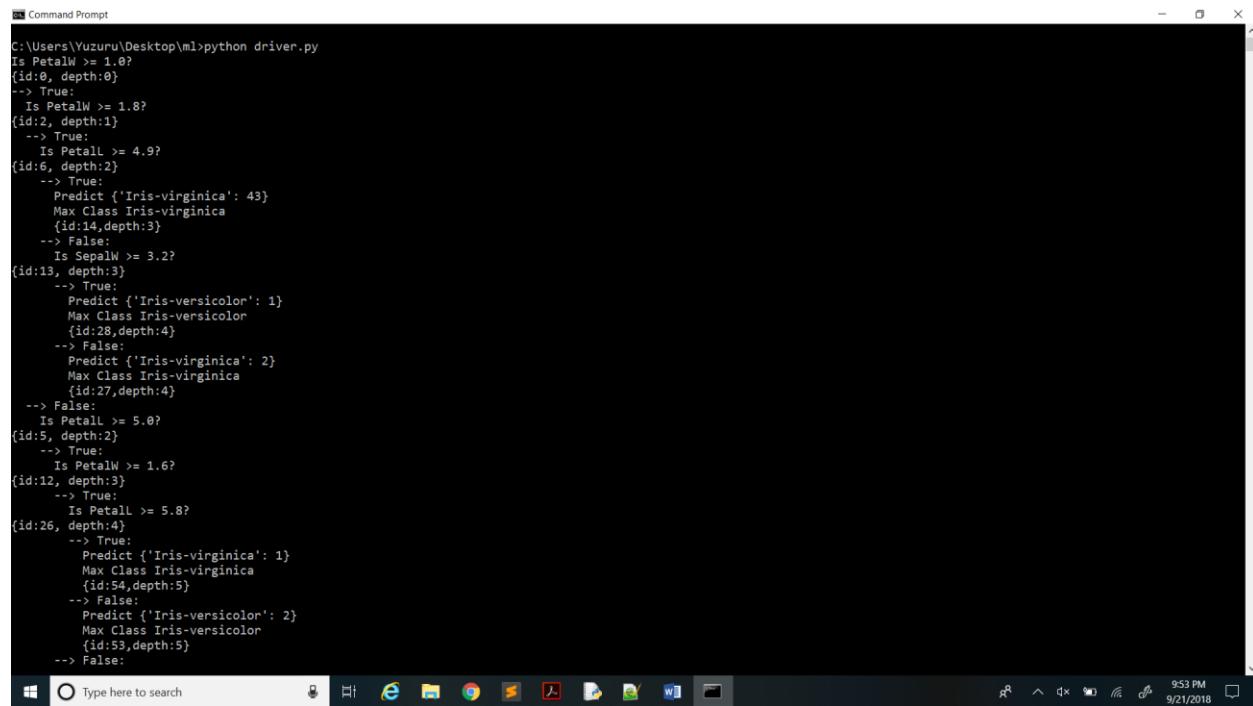
A5)

Code: Max Prediction and printing

```
150 ## TODO: Step 2
151 class Leaf:
152     """A Leaf node classifies data.
153
154     This holds a dictionary of class (e.g., "Apple") -> number of times
155     it appears in the rows from the training data that reach this leaf.
156     """
157
158     def __init__(self, rows, id, depth):
159         self.predictions = class_counts(rows)
160         count_temp=self.predictions
161         self.predictions_cls=max(count_temp)
162         self.id=id
163         self.depth=depth
164
```

```
270
271 ## TODO: Step 4
272 def print_tree(node, spacing=""):
273     """World's most elegant tree printing function."""
274
275     # Base case: we've reached a leaf
276     if isinstance(node, Leaf):
277         print (spacing + "Predict", node.predictions)
278         print (spacing+"Max Class",node.predictions_cls)
279         print (spacing+"{id:"+str(node.id)+" , "+ "depth:"+str(node.depth)+"}")
280         return
281     # Print the question at this node
282     print (spacing + str(node.question))
283     print ("{id:" + str(node.id)+", depth:"+str(node.depth)+"}")
284     # Call this function recursively on the true branch
285     print (spacing + '--> True:')
286     print_tree(node.true_branch, spacing + " ")
287
288     # Call this function recursively on the false branch
289     print (spacing + '--> False:')
290     print_tree(node.false_branch, spacing + " ")
```

Output: Max Class Prediction at Leaf nodes.



```
C:\Users\Yuzuru\Desktop\ml>python driver.py
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
  Is PetalW >= 1.8?
  {id:2, depth:1}
  --> True:
    Is Petall >= 4.9?
    {id:6, depth:2}
    --> True:
      Predict {'Iris-virginica': 43}
      Max Class Iris-virginica
      {id:14,depth:3}
    --> False:
      Is SepalW >= 3.2?
      {id:13, depth:3}
      --> True:
        Predict {'Iris-versicolor': 1}
        Max Class Iris-versicolor
        {id:28,depth:4}
      --> False:
        Predict {'Iris-virginica': 2}
        Max Class Iris-virginica
        {id:27,depth:4}
      --> False:
        Is Petall >= 5.0?
        {id:5, depth:2}
        --> True:
          Is PetalW >= 1.6?
          {id:12, depth:3}
          --> True:
            Is Petall >= 5.8?
            {id:26, depth:4}
            --> True:
              Predict {'Iris-virginica': 1}
              Max Class Iris-virginica
              {id:54,depth:5}
            --> False:
              Predict {'Iris-versicolor': 2}
              Max Class Iris-versicolor
              {id:53,depth:5}
            --> False:
```

Q6) Next, you will write two functions - one that returns a list of all inner nodes (i.e. those nodes that are not leaf nodes) and another function that returns a list of leaf nodes. Note that you have to return the nodes, not just their ids. You can call these functions getInnerNodes and getLeafNodes. Both of these should have one input parameter - the root node of the tree. You can see an example of this in the print tree function - without the spacing parameter, of course.

A6) Code:

```
302 ## TODO: Step 5
303 def getLeafNodes(node, leafNodes =[]):
304
305     # Returns a list of all leaf nodes of a tree
306     if isinstance(node, Leaf):
307         leafNodes.append(node)
308         return
309
310     getLeafNodes(node.true_branch)
311     getLeafNodes(node.false_branch)
312
313     return leafNodes
314
315
316 def getInnerNodes(node, innerNodes =[]):
317
318     # Returns a list of all non-leaf nodes of a tree
319     # Returns a list of all leaf nodes of a tree
320     if isinstance(node, Leaf):
321
322         return
323
324     innerNodes.append(node)
325     getInnerNodes(node.true_branch)
326     getInnerNodes(node.false_branch)
327
328     return innerNodes
329
```

Output:

```
{id:1,depth:1}
***** Leaf nodes *****
id = 14 depth =3
id = 28 depth =4
id = 27 depth =4
id = 54 depth =5
id = 53 depth =5
id = 25 depth =4
id = 24 depth =4
id = 23 depth =4
id = 1 depth =1
***** Non-leaf nodes *****
id = 0 depth =0
id = 2 depth =1
id = 6 depth =2
id = 13 depth =3
id = 5 depth =2
id = 12 depth =3
id = 26 depth =4
id = 11 depth =3
*****Tree before pruning*****
```

Q7) We will modify the classify function so that it outputs just the predicted class label, instead of class counts. The classify function works on any 1 row of data. We will also create a function called computeAccuracy that takes in a set of rows and returns their accuracy.

Below is the signature of that function:

A7) Code:

Classify:

```
254 ## TODO: Step 6
255 def classify(row, node):
256     """See the 'rules of recursion' above."""
257
258     # Base case: we've reached a leaf
259     if isinstance(node, Leaf):
260         return node.predictions
261
262     # Decide whether to follow the true-branch or the false-branch.
263     # Compare the feature / value stored in the node,
264     # to the example we're considering.
265     if node.question.match(row):
266         return classify(row, node.true_branch)
267     else:
268         return classify(row, node.false_branch)
269
270
```



Accuracy:

```
330
331 ## TODO: Step 6
332 def computeAccuracy(rows, node):
333     totalRows = len(rows)
334     numAccurate = 0
335     for row in rows:
336         pred_label=classify(row,node)
337         if row[-1] in pred_label.keys():
338             numAccurate = numAccurate + 1
339
340     return round(numAccurate/float(totalRows), 2)
341
342
```



Output: Accuracy after pruning is slightly greater.

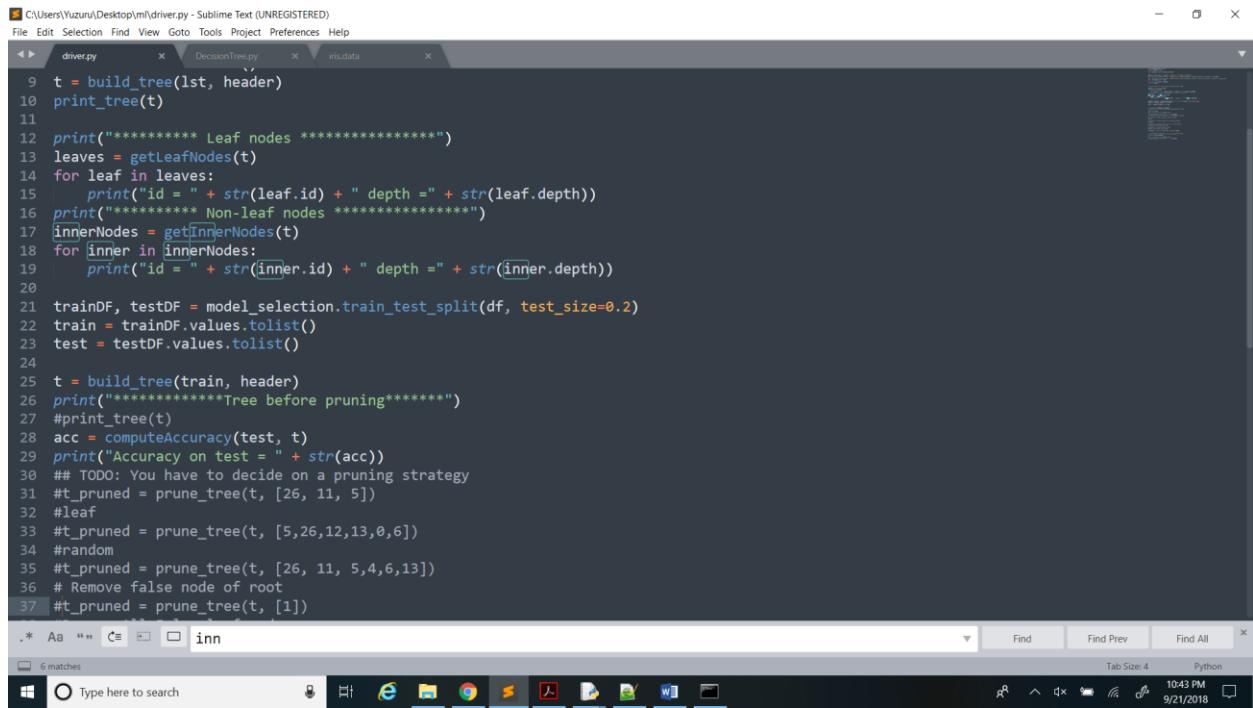
```
Accuracy on test = 0.93
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
| Is PetalW >= 1.8?
{id:2, depth:1}
--> True:
| Is Petall. >= 4.9?
{id:6, depth:2}
--> True:
| Predict ('Iris-virginica': 34)
| Max Class Iris-virginica
| {id:14, depth:3}
--> False:
```



```
--> False:
Predict {'Iris-virginica': 1, 'Iris-versicolor': 38}
Max Class Iris-virginica
{id:11,depth:3}
--> False:
Predict {'Iris-setosa': 43}
Max Class Iris-setosa
{id:1,depth:1}
Accuracy on test = 0.97
```

Q8) For testing our model, we will use Scikit Learn's train test split method. This will allow you to choose a random sample of the dataset for training and remaining for test.

A8) Code:



The screenshot shows a Sublime Text editor window with three tabs: 'driver.py', 'DecisionTree.py', and 'iris.data'. The 'driver.py' tab contains the following code:

```
9 t = build_tree(lst, header)
10 print_tree(t)
11
12 print("***** Leaf nodes *****")
13 leaves = getLeafNodes(t)
14 for leaf in leaves:
15     print("id = " + str(leaf.id) + " depth =" + str(leaf.depth))
16 print("***** Non-leaf nodes *****")
17 innerNodes = getInnerNodes(t)
18 for inner in innerNodes:
19     print("id = " + str(inner.id) + " depth =" + str(inner.depth))
20
21 trainDF, testDF = model_selection.train_test_split(df, test_size=0.2)
22 train = trainDF.values.tolist()
23 test = testDF.values.tolist()
24
25 t = build_tree(train, header)
26 print("*****Tree before pruning*****")
27 #print_tree(t)
28 acc = computeAccuracy(test, t)
29 print("Accuracy on test = " + str(acc))
30 ## TODO: You have to decide on a pruning strategy
31 #t_pruned = prune_tree(t, [26, 11, 5])
32 #leaf
33 #t_pruned = prune_tree(t, [5,26,12,13,0,6])
34 #random
35 #t_pruned = prune_tree(t, [26, 11, 5,4,6,13])
36 # Remove false node of root
37 #t_pruned = prune_tree(t, [1])
```

The status bar at the bottom indicates 'Tab Size: 4' and the date/time '10:43 PM 9/21/2018'.

Q9) You would need to decide a pruning strategy?

A9)

(a) Prune each node that is 1 level above the leaf e.g. node 2 in the figure above and see if it helps improve the test error

Code:

```
32 #leaf
33 t_pruned = prune_tree(t, [5,26,12,13,0,6])
```

Output: Getting an over fit of data, therefore accuracy is same.

```
*****Tree before pruning*****
Accuracy on test = 0.97
*****Tree after pruning*****
Predict {'Iris-virginica': 37, 'Iris-setosa': 40, 'Iris-versicolor': 43}
Max Class Iris-virginica
{id:0,depth:0}
Accuracy on test = 0.97

C:\Users\Yuzuru\Desktop\ml>
```

(b) Randomly select n nodes and prune them as a batch, and see if it improves accuracy

Code: `t_pruned = prune_tree (t, [26, 11, 5,4,6,13])`

Output: After pruning the accuracy is greater.

```
*****Tree before pruning*****
Accuracy on test = 0.93
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
  Is PetalW >= 1.8?
  {id:2, depth:1}
  --> True:
    Predict {'Iris-virginica': 39, 'Iris-versicolor': 1}
    Max Class Iris-virginica
    {id:6,depth:2}
  --> False:
    Predict {'Iris-virginica': 3, 'Iris-versicolor': 36}
    Max Class Iris-virginica
    {id:5,depth:2}
--> False:
  Predict {'Iris-setosa': 41}
  Max Class Iris-setosa
  {id:1,depth:1}
Accuracy on test = 1.0

C:\Users\Yuzuru\Desktop\ml>
```

(c) Anything else that you can think of

Code:

a) Remove false node of root

`t_pruned = prune_tree (t, [1])`

Output: Overfitting happens when you remove the false node of root.

```
Command Prompt
id = 13 depth =3
id = 5 depth =2
id = 12 depth =3
id = 26 depth =4
id = 11 depth =3
*****Tree before pruning*****
Accuracy on test = 0.93
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
```

```
Max Class Iris-versicolor
{id:5,depth:2}
--> False:
Predict {'Iris-setosa': 41}
Max Class Iris-setosa
{id:1,depth:1}
Accuracy on test = 0.93
C:\Users\Yuzuru\Desktop\ml>
```

b) Remove All False leaf nodes

`t_pruned = prune_tree (t, [1, 11,25,53,27])`

output: accuracy is greater after pruning all false nodes.

```
Command Prompt
id = 13 depth =3
id = 5 depth =2
id = 12 depth =3
id = 26 depth =4
id = 11 depth =3
*****Tree before pruning*****
Accuracy on test = 0.87
*****Tree after pruning*****
Is PetalW >= 1.0?
{id:0, depth:0}
--> True:
  Is PetalL >= 4.8?
{id:2, depth:1}
```

```
Max Class Iris-virginica
{id:12,depth:3}
--> False:
  Predict {'Iris-versicolor': 37}
  Max Class Iris-versicolor
  {id:11,depth:3}
--> False:
  Predict {'Iris-setosa': 39}
  Max Class Iris-setosa
  {id:1,depth:1}
Accuracy on test = 0.97
C:\Users\Yuzuru\Desktop\ml>
```

Q10) After making the above changes, you have to get everything in the driver.py file to run without error. and then choose another dataset from UCI Machine Learning dataset that is suitable for classification tasks and get it to work using the code that you created earlier.

A10) Car Data Set:

Output: As Data Set is too big we are just adding imp parts of output. Accuracy after pruning (0.9) is greater than (0.88)

```
C:\Users\Yuzuru\Desktop\ml>python driver.py
Is Class == low?
{id:0, depth:0}
--> True:
  Predict {'unacc': 576}
  Max Class unacc
  {id:2,depth:1}
--> False:
  Is Petall == 2?
{id:1, depth:1}
--> True:
  Predict {'unacc': 384}
  Max Class unacc
  {id:4,depth:2}
--> False:
  Is Sepall == vhigh?
{id:3, depth:2}
--> True:
  Is PetalW == small?
{id:8, depth:3}
--> True:
  Is Class == med?
{id:18, depth:4}
--> True:
  Predict {'unacc': 32}
  Max Class unacc
  {id:38,depth:5}
--> False:
  Is SepalW == 2?
{id:37, depth:5}
--> True:
  Is Petall == 2?
{id:76, depth:6}
--> True:
  Predict {'unacc': 4}
  Max Class unacc
  {id:154,depth:7}
--> False:
  Predict {'acc': 2, 'unacc': 2}
  Max Class unacc
  {id:153,depth:7}
--> False:
```

```
■ Select Command Prompt
Accuracy on test = 0.88
*****Tree after pruning*****
Is Petall == 2?
{id:0, depth:0}
--> True:
  Predict {'unacc': 464}
  Max Class unacc
  {id:2,depth:1}
--> False:
  Is Class == low?
{id:1, depth:1}
--> True:
  Predict {'unacc': 308}
  Max Class unacc
  {id:4,depth:2}
--> False:
  Is Sepall == vhigh?
{id:3, depth:2}
--> True:
  Predict {'acc': 54, 'unacc': 103}
  Max Class unacc
  {id:8,depth:3}
--> False:
  Is Class == med?
{id:7, depth:3}
--> True:
  Is PetalW == small?
{id:16, depth:4}
--> True:
  Is Sepall == high?
{id:34, depth:5}
--> True:
  Is SepalW == 4?
{id:70, depth:6}
--> True:
  Is PetalL == more?
{id:142, depth:7}
--> True:
  Predict {'acc': 1, 'unacc': 1}
  Max Class unacc
  {id:286,depth:8}
--> False:
  Predict {'acc': 1, 'unacc': 2}

10:53 PM 9/21/2018
```

```
■ Select Command Prompt
Is Sepall == low?
{id:2051, depth:11}
--> True:
  Predict {'acc': 1, 'vgood': 2}
  Max Class vgood
  {id:4104,depth:12}
--> False:
  Predict {'acc': 2, 'vgood': 2}
  Max Class vgood
  {id:4103,depth:12}
--> False:
  Is PetalW == med?
{id:511, depth:9}
--> True:
  Predict {'acc': 6, 'vgood': 8}
  Max Class vgood
  {id:1024,depth:10}
--> False:
  Is Sepall == low?
{id:1023, depth:10}
--> True:
  Is Petall == 4?
{id:2048, depth:11}
--> True:
  Predict {'acc': 1, 'vgood': 2}
  Max Class vgood
  {id:4098,depth:12}
--> False:
  Predict {'acc': 2, 'vgood': 2}
  Max Class vgood
  {id:4097,depth:12}
--> False:
  Is Petall == more?
{id:2047, depth:11}
--> True:
  Predict {'vgood': 2}
  Max Class vgood
  {id:4096,depth:12}
--> False:
  Predict {'acc': 1, 'vgood': 1}
  Max Class vgood
  {id:4095,depth:12}

Accuracy on test = 0.9
10:53 PM 9/21/2018
```