

Part 1

By Gini Index:

Prediction on Train Data

	price	maintenance	capacity	airbag	profitable	Predicted
0	low	low	2	no	yes	yes
1	low	med	4	yes	no	no
2	low	high	4	no	no	no
3	med	med	4	no	no	no
4	med	med	4	yes	yes	yes
5	med	high	2	yes	no	no
6	high	med	4	yes	yes	yes
7	high	high	2	yes	no	no
8	high	high	5	yes	yes	yes

Accuracy = 100%

Prediction on Test Data

	price	maintenance	capacity	airbag	profitable	Predicted
0	med	high	5	no	yes	yes
1	low	low	4	no	yes	yes

Accuracy = 100%

By Information Gain:

Prediction on Train Data

	price	maintenance	capacity	airbag	profitable	Predicted
0	low	low	2	no	yes	yes
1	low	med	4	yes	no	no
2	low	high	4	no	no	no
3	med	med	4	no	no	no
4	med	med	4	yes	yes	yes
5	med	high	2	yes	no	no
6	high	med	4	yes	yes	yes
7	high	high	2	yes	no	no
8	high	high	5	yes	yes	yes

Accuracy = 100%

Prediction on Test Data

	price	maintenance	capacity	airbag	profitable	Predicted
0	med	high	5	no	yes	yes
1	low	low	4	no	yes	yes

Accuracy = 100%

DECISION TREE (GINI INDEX)

```
| maintenance = high
|   capacity = 2 : no
|   capacity = 5 : yes
|   capacity = 4 : no
| maintenance = med
|   price = med
|     airbag = yes : yes
|     airbag = no : no
|   price = low : no
|   price = high : yes
| maintenance = low : yes
```

DECISION TREE (INFORMATION GAIN)

```
| maintenance = high
|   capacity = 2 : no
|   capacity = 5 : yes
|   capacity = 4 : no
| maintenance = med
|   price = med
|     airbag = yes : yes
|     airbag = no : no
|   price = low : no
|   price = high : yes
| maintenance = low : yes
```

Without dummies:

1. Entropy of Root Node by my model:	0.9910760598382222	≈ 0.991
2. Information Gain of Root by My model:	0.18606356007860758	≈ 0.186
3. Gini Index of Root by My model:	0.49382716049382713	≈ 0.494
Gini of Split at Root by My model:	0.38888888888888884	≈ 0.388

With dummies:

1. Entropy of Root Node by Scikit-Learn:	0.991	≈ 0.991
2. Entropy of Root Node by my model:	0.9910760598382222	≈ 0.991
3. Information Gain of Root by My model:	0.14269027946047563	≈ 0.143
4. Information Gain of Root by Scikit-Learn:	0.143	≈ 0.143
5. Gini Index of Root by Scikit-Learn:	0.49382716049382713	≈ 0.494
6. Gini Index of Root by My model:	0.49382716049382713	≈ 0.494
Gini of Split at Root by My model:	0.4166666666666667	≈ 0.417

My model

DECISION TREE (GINI INDEX)

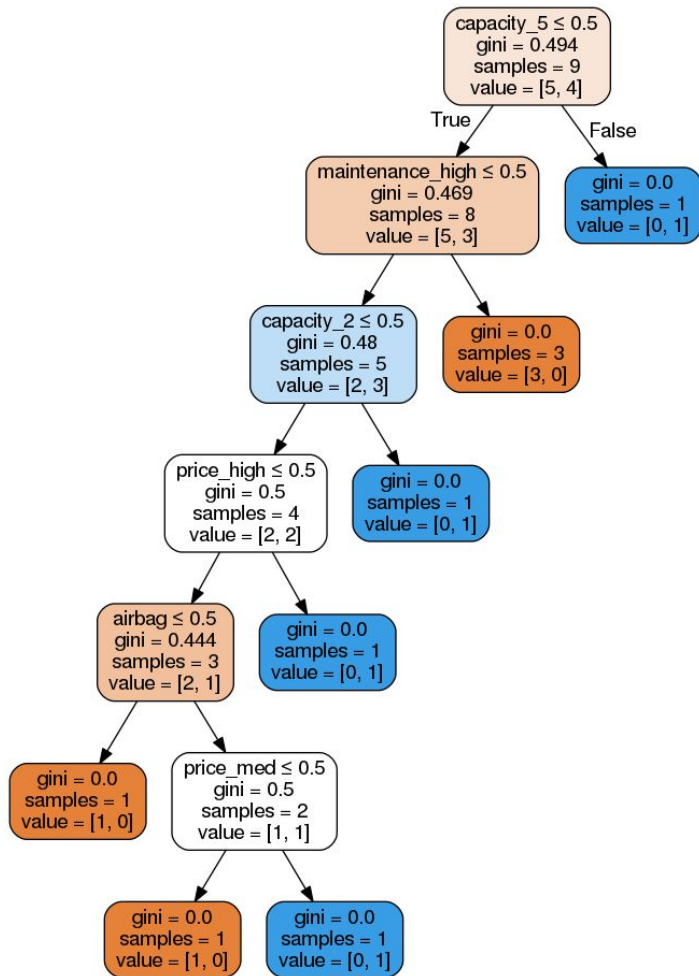
```
| capacity_5 = 0
|   maintenance_high = 0
|     capacity_4 = 1
|       price_low = 0
|         airbag = 1 : yes
|         airbag = 0 : no
|       price_low = 1 : no
|     capacity_4 = 0 : yes
|   maintenance_high = 1 : no
| capacity_5 = 1 : yes
```

DECISION TREE (INFORMATION GAIN)

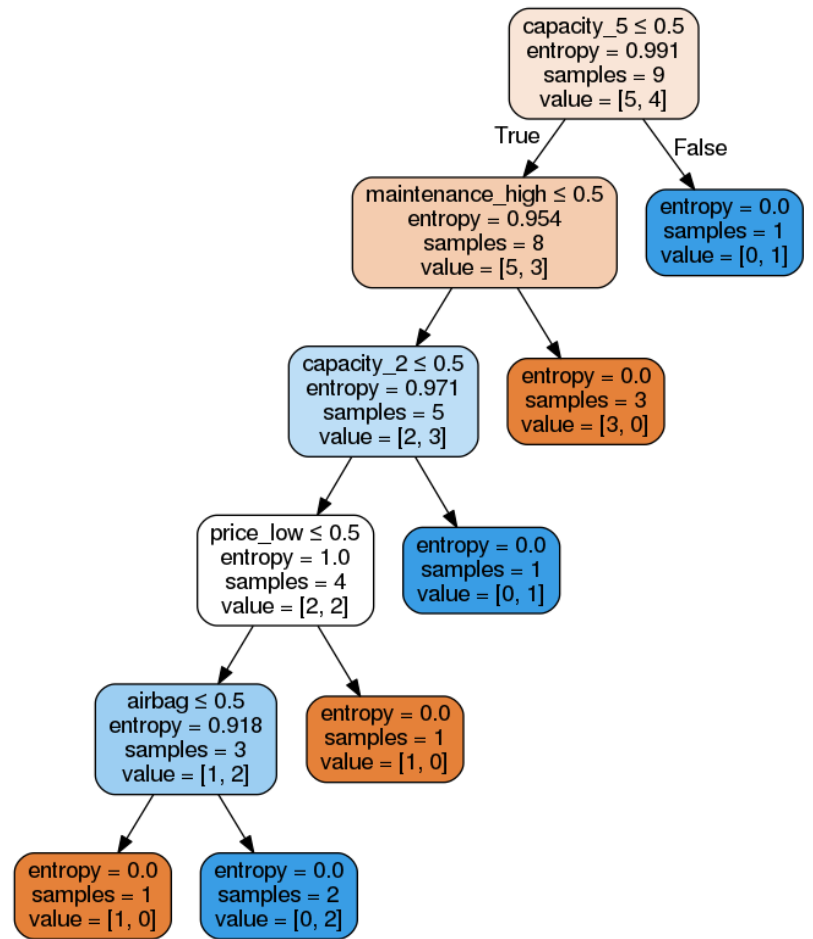
```
| capacity_5 = 0
|   maintenance_high = 0
|     capacity_4 = 1
|       price_low = 0
|         airbag = 1 : yes
|         airbag = 0 : no
|       price_low = 1 : no
|     capacity_4 = 0 : yes
|   maintenance_high = 1 : no
| capacity_5 = 1 : yes
```

Scikit Learn Model

DECISION TREE (GINI INDEX)



DECISION TREE (INFORMATION GAIN)



PREDICTIONS

My Model (Gini Index) (Test data)

	airbag	price_high	price_low	price_med	maintenance_high	maintenance_low	maintenance_med	capacity_2	capacity_4	capacity_5	profitable	Predicted
0	0	0	0	1	1	0	0	0	0	1	1	yes
1	0	0	1	0	0	1	0	0	1	0	1	no

Accuracy: 50%

My Model (Information Gain) (Test data)

	airbag	price_high	price_low	price_med	maintenance_high	maintenance_low	maintenance_med	capacity_2	capacity_4	capacity_5	profitable	Predicted
0	0	0	0	1	1	0	0	0	0	1	1	yes
1	0	0	1	0	0	1	0	0	1	0	1	no

Accuracy: 50%

Scikit Learn (Test data)

*****Prediction on TEST DATA*****

Model Prediction on Test Data (Scikit-Learn) (Gini Index) (dummy)

[1 0]

Model Prediction on Test Data (Scikit-Learn) (Information Gain) (dummy)

[1 0]

Therefore Accuracy is 50% in both cases.

My Model (Gini Index) (Train data)

	airbag	price_high	price_low	price_med	maintenance_high	maintenance_low	maintenance_med	capacity_2	capacity_4	capacity_5	profitable	Predicted
0	0	0	1	0	0	1	0	1	0	0	1	yes
1	1	0	1	0	0	0	1	0	1	0	0	no
2	0	0	1	0	1	0	0	0	1	0	0	no
3	0	0	0	1	0	0	1	0	1	0	0	no
4	1	0	0	1	0	0	1	0	1	0	1	yes
5	1	0	0	1	1	0	0	1	0	0	0	no
6	1	1	0	0	0	0	1	0	1	0	1	yes
7	1	1	0	0	1	0	0	1	0	0	0	no
8	1	1	0	0	1	0	0	0	0	1	1	yes

Accuracy: 100%

My Model (Information Gain) (Train data)

	airbag	price_high	price_low	price_med	maintenance_high	maintenance_low	maintenance_med	capacity_2	capacity_4	capacity_5	profitable	Predicted
0	0	0	1	0	0	1	0	1	0	0	1	yes
1	1	0	1	0	0	0	1	0	1	0	0	no
2	0	0	1	0	1	0	0	0	1	0	0	no
3	0	0	0	1	0	0	1	0	1	0	0	no
4	1	0	0	1	0	0	1	0	1	0	1	yes
5	1	0	0	1	1	0	0	1	0	0	0	no
6	1	1	0	0	0	0	1	0	1	0	1	yes
7	1	1	0	0	1	0	0	1	0	0	0	no
8	1	1	0	0	1	0	0	0	0	1	1	yes

Accuracy: 100%

Scikit Learn (Train data)

*****Prediction on TRAIN DATA*****

Model Prediction on Train Data (Scikit-Learn) (Gini Index) (dummy)

[1 0 0 0 1 0 1 0 1]

Model Prediction on Train Data (Scikit-Learn) (Information Gain) (dummy)

[1 0 0 0 1 0 1 0 1]

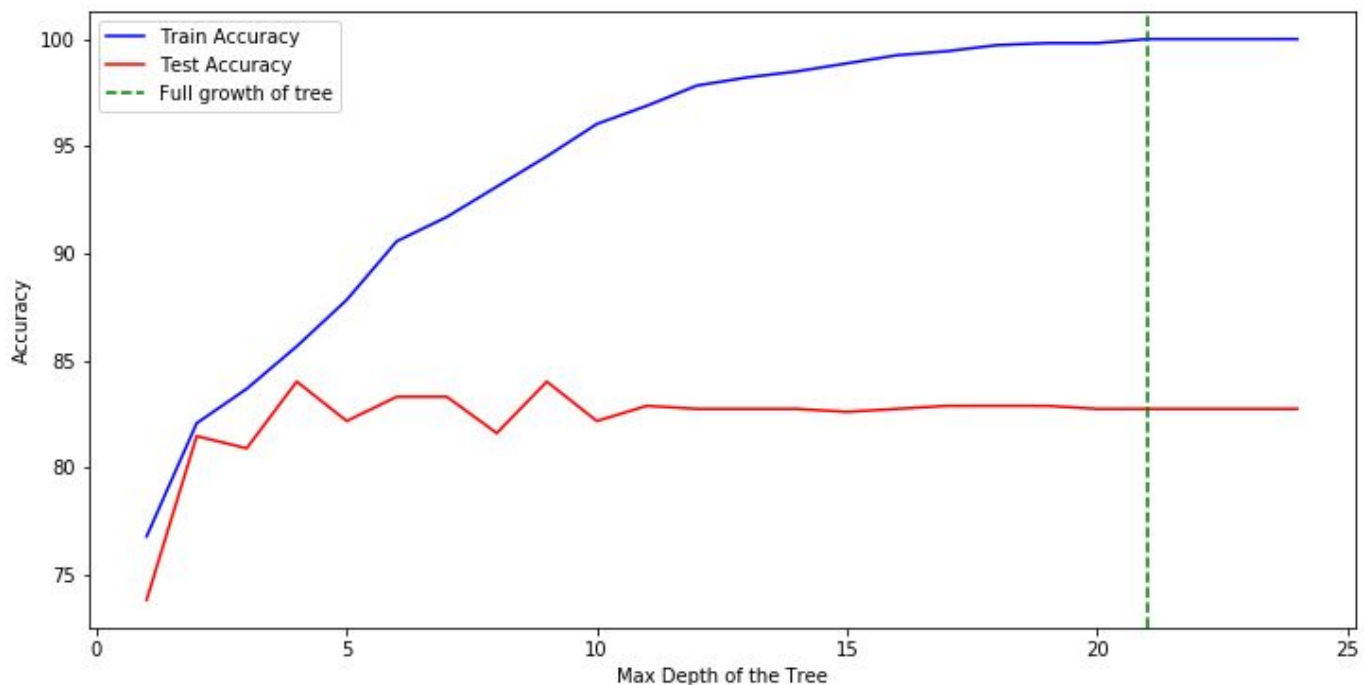
Therefore Accuracy is 100%

Part 2

All are Percentage Accuracy (in %)

Max_Depth	Train_Accuracy	Test_Accuracy
1	76.792453	73.833098
2	82.075472	81.471004
3	83.679245	80.905233
4	85.660377	84.016973
5	87.830189	82.178218
6	90.566038	83.309760
7	91.698113	83.309760
8	93.113208	81.612447
9	94.528302	84.016973
10	96.037736	82.178218
11	96.886792	82.885431
12	97.830189	82.743989
13	98.207547	82.743989
14	98.490566	82.743989
15	98.867925	82.602546
16	99.245283	82.743989
17	99.433962	82.885431
18	99.716981	82.885431
19	99.811321	82.885431
20	99.811321	82.743989
21	100.000000	82.743989
22	100.000000	82.743989
23	100.000000	82.743989
24	100.000000	82.743989

Tree of **maximum depth=4** and **maximum depth=9** achieved maximum Testing accuracy of **84.016973%**



Tree of **maximum depth=4** and **maximum depth=9** achieved maximum Testing accuracy of **84.016973%**

Yes, overfitting does occur. Since Accuracy on Training Data reached to 100% whereas Accuracy achieved on Testing Data is less than 85%. Also, Maximum testing accuracy was achieved at max_depth=4 and max_depth=9. At these depths, we can see that training error was around 85.66% and 94.52%. Now as the max_depth increases accuracy on only training data increases. No improvement has been recorded on Testing data.

Overfitting starts occurring from depth>=10.

Various Node asks question on one of the following words given below

If a document is related to **comp.graphics** newsgroup, then it is very likely that these documents will contain words like given below. These words resemble that if a document contains such words then it must belong to **comp.graphics** newsgroup.

- | | | |
|-------------------|--------------------|------------------|
| • graphics | • mac | • acm |
| • graphic | • password | • comp |
| • image | • program | • windows |
| • online | • slow | • port |
| • disk | • algorithm | • circle |
| • time | • format | |

Also, if a document is related to **alt.atheism** newsgroup, then it is very likely that these documents will contain words like

- | | | |
|--------------------|-------------------|-------------------|
| • god | • wrote | • thanks |
| • says | • dwyer | • addition |
| • don | • claiming | • book |
| • religious | • bible | • cheers |
| • face | • bill | • people |

In fact there are some nodes in this decision tree which can belong to **alt.atheism** and in general, are less seen in **comp.graphics** newsgroup. These nodes are:

- | | | |
|---------------|--------------|---------------|
| • your | • how | • am |
| • you | • he | • have |
| • what | • who | |

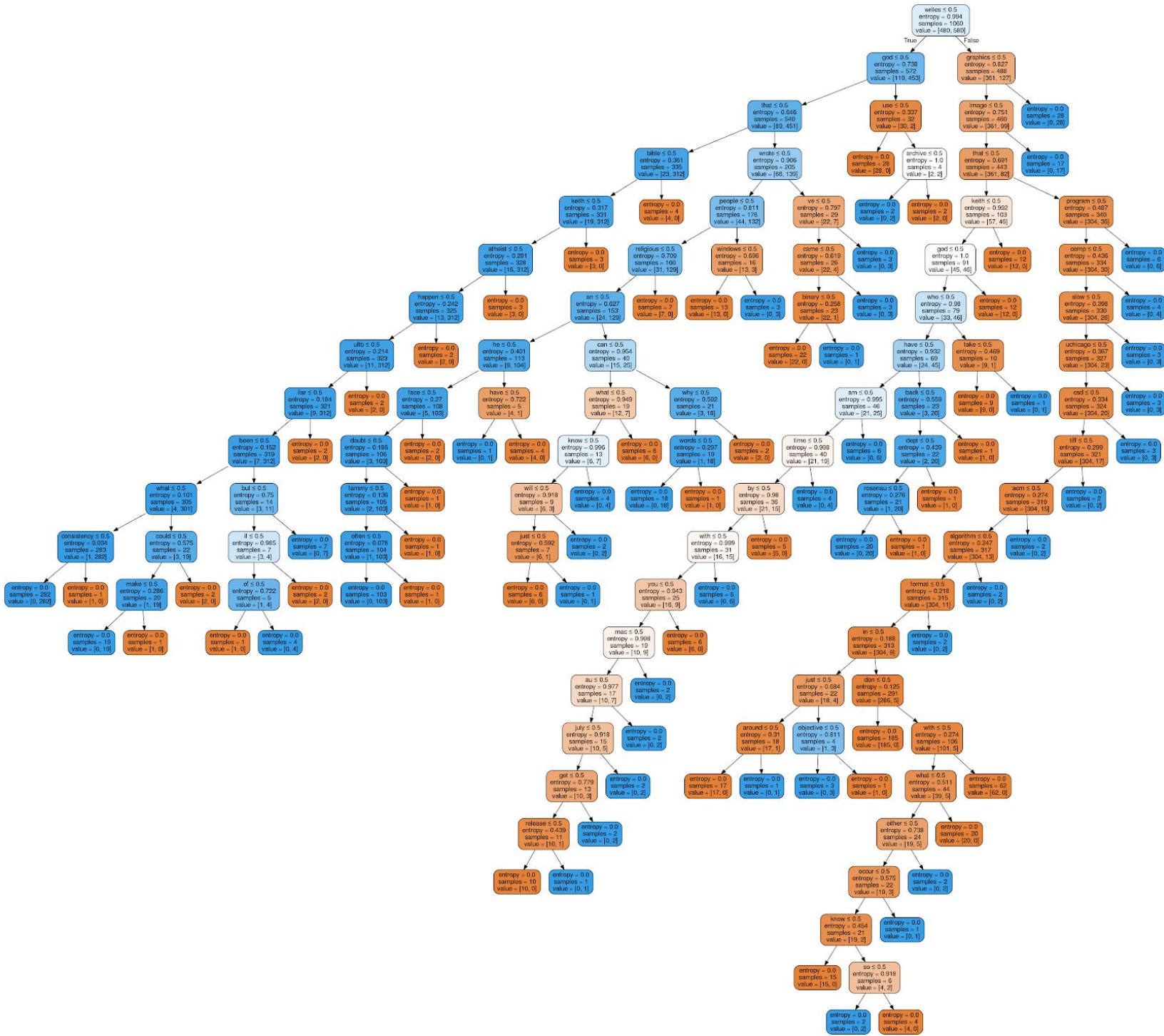

```

50 | ve = 0
51 | online = 0 : no
52 | online = 1 : yes
53 | ve = 1 : yes
54 | came = 1 : yes
55 | god = 1
56 | use = 0 : no
57 | use = 1
58 | will = 1 : no
59 | will = 0 : yes
60 | writes = 1
61 | graphics = 0
62 | image = 0
63 | that = 1
64 | program = 0
65 | comp = 0
66 | csd = 0
67 | slow = 0
68 | uchicago = 0 : no
69 | uchicago = 1 : yes
70 | slow = 1 : yes
71 | csd = 1 : yes
72 | comp = 1 : yes
73 | program = 1 : yes
74 | that = 0
75 | keith = 0
76 | god = 0
77 | who = 0
78 | have = 0
79 | am = 0 : no
80 | am = 1 : yes
81 | have = 1
82 | rosenau = 0 : yes
83 | rosenau = 1 : no
84 | who = 1
85 | disk = 0 : no
86 | disk = 1 : yes
87 | god = 1 : no
88 | keith = 1 : no
89 | image = 1 : yes
90 | graphics = 1 : yes

```

Please look at “**Max_Accuracy_Tree.txt**” file for getting a clear view.

Tree Build by Scikit Learn



For details please look at “output_entropy_3.png” file for details.

Scikit Learn Model:

Number of misclassification on Training Data: 0 out of 1060

(100% Accuracy)

Number of misclassification on Testing Data: 162 out of 707

(77.086% Accuracy)

Fully Grown Tree Build by My Model

[illegible]

```

50 | doubt = 1 : no
51 | claiming = 1 : no
52 | tammy = 1 : no
53 | face = 1 : no
54 | he = 1
55 | graphic = 0 : no
56 | graphic = 1 : yes
57 | an = 1
58 | can = 1
59 | dwyer = 0
60 | coming = 0 : yes
61 | coming = 1 : no
62 | dwyer = 1 : no
63 | can = 0
64 | how = 0
65 | know = 0
66 | will = 0
67 | port = 0 : no
68 | port = 1 : yes
69 | will = 1 : yes
70 | know = 1 : yes
71 | how = 1 : no
72 | religious = 1 : no
73 | people = 1
74 | windows = 0 : no
75 | windows = 1 : yes
76 | wrote = 1
77 | came = 0
78 | ve = 0
79 | online = 0 : no
80 | online = 1 : yes
81 | ve = 1 : yes
82 | came = 1 : yes
83 | god = 1
84 | use = 0 : no
85 | use = 1
86 | will = 1 : no
87 | will = 0 : yes
88 | writes = 1
89 | graphics = 0
90 | image = 0
91 | that = 1
92 | program = 0
93 | comp = 0
94 | csd = 0
95 | slow = 0
96 | uchiago = 0
97 | acm = 0
98 | tiff = 0

```

```

99      | algorithm = 0
100      | format = 0
101      | in = 1
102      | don = 0 : no
103      | don = 1
104      |   with = 1 : no
105      |   with = 0
106      |   what = 0
107      |   either = 0
108      |   rahul = 0
109      |   here = 0 : no
110      |   here = 1
111      |   | atheists = 1 : no
112      |   | atheists = 0 : yes
113      |   rahul = 1 : yes
114      |   either = 1 : yes
115      |   what = 1 : no
116      | in = 0
117      | just = 0
118      |   tracing = 0 : no
119      |   tracing = 1 : yes
120      | just = 1
121      |   says = 0 : yes
122      |   says = 1 : no
123      | format = 1 : yes
124      | algorithm = 1 : yes
125      | tiff = 1 : yes
126      | acm = 1 : yes
127      | uchicago = 1 : yes
128      | slow = 1 : yes
129      | csd = 1 : yes
130      | comp = 1 : yes
131      | program = 1 : yes
132      | that = 0
133      | keith = 0
134      | god = 0
135      |   who = 0
136      |   have = 0
137      |   am = 0
138      |   time = 0
139      |   by = 0
140      |   with = 0
141      |   you = 0
142      |   july = 0
143      |   mac = 0
144      |   au = 0
145      |   get = 0
146      |   password = 0 : no
147      |   password = 1 : yes

```

```

148      |   get = 1 : yes
149      |   au = 1 : yes
150      |   mac = 1 : yes
151      |   july = 1 : yes
152      |   you = 1 : no
153      |   with = 1 : yes
154      |   by = 1 : no
155      |   time = 1 : yes
156      |   am = 1 : yes
157      | have = 1
158      |   rosenau = 0
159      |   your = 0 : yes
160      |   your = 1
161      |   | or = 0 : yes
162      |   | or = 1 : no
163      |   rosenau = 1 : no
164      |   who = 1
165      |   disk = 0 : no
166      |   disk = 1 : yes
167      |   god = 1 : no
168      |   keith = 1 : no
169      | image = 1 : yes
170      | graphics = 1 : yes

```

Please look at “Full_Grown_Tree.txt” file for getting a clear view.

***** THE END *****

By Vaibhav Poddar
16CS10051