

Question 1:

① Shannon entropy of target variable (edible) is calculated as follows:

$$H(\text{edible}) = - \left(\frac{16}{24} \times \log_2 \left(\frac{16}{24} \right) \right) - \left(\frac{8}{24} \times \log_2 \left(\frac{8}{24} \right) \right) = 0.92 \text{ bits}$$

② Each feature white, tall, shiny's entropies & information gains are calculated as follows:

a) White:

		True	False
white	0	9	5
	1	7	3

$$H(\text{white, edible}) = P(\text{white}=0) \times H(\text{white}=0) + P(\text{white}=1) \times H(\text{white}=1)$$

$$P(\text{white}=0) = \frac{14}{24}, E(9, 5) = - \left(\frac{9}{14} \times \log_2 \left(\frac{9}{14} \right) \right) - \left(\frac{5}{14} \times \log_2 \left(\frac{5}{14} \right) \right) = 0.94$$

$$P(\text{white}=1) = \frac{10}{24}, E(7, 3) = - \left(\frac{7}{10} \times \log_2 \left(\frac{7}{10} \right) \right) - \left(\frac{3}{10} \times \log_2 \left(\frac{3}{10} \right) \right) = 0.88$$

$$H(\text{white, edible}) = \frac{14}{24} \times 0.94 + \frac{10}{24} \times 0.88 = 0.915$$

$$\text{Information gain} = H(\text{edible}) - H(\text{white, edible}) = 0.92 - 0.915 = 0.005$$

b) Tall :

		edible	True	False
		Tall		
tall	0	6	4	
	1	10	14	

$tall = 0$

$$E(b, 4) = -\left(\frac{6}{10} \times \log_2 \frac{6}{10}\right) - \left(\frac{4}{10} \times \log_2 \frac{4}{10}\right) = 0.97$$

$$E(10, 4) = -\left(\frac{10}{14} \times \log_2 \frac{10}{14}\right) - \left(\frac{4}{14} \times \log_2 \frac{4}{14}\right) = 0.861$$

$$P(tall = 0) = 10/24$$

$$P(tall = 1) = 14/24$$

$$U(tall, edible) = 0.906$$

$$\text{Info gain} = U(\text{edible}) - U(tall, \text{edible}) = 0.92 - 0.906 = 0.014$$

c) Frilly :

		edible	True	False
		Frilly		
frilly	0	13	3	
	1	3	5	

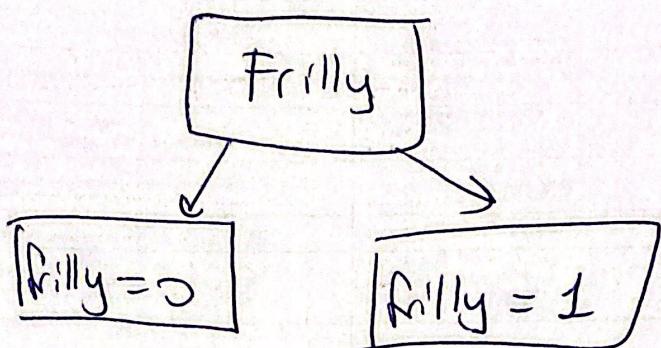
$$P(frilly = 0) = 16/24$$

$$P(frilly = 1) = 8/24$$

$$H(\text{frilly, edible}) = 0.78228 \rightarrow \text{Info gain} = 0.92 - 0.78228 = 0.1377$$

1 - 0.92

③ Comparing information gains feature 'Frilly' should be next Node with max info gain.



$\Rightarrow \text{Frilly} = 0$, consider attribute white's entropy.

Entropy ($\text{Frilly} = 0$) = 0,696205 \rightarrow calculated below.

when $\text{frilly} = 0$ & $\text{white} = 0$, there are 6 True's & 3 False's.

$$\text{entropy}(\text{white} = 0) = -\frac{6}{9} \log_2\left(\frac{6}{9}\right) - \frac{3}{9} \log_2\left(\frac{3}{9}\right) = 0,92$$

entropy ($\text{white} = 1$) = 0 since all target values are True's.

$$\begin{aligned} \text{Information gain} (\text{frilly} = 0, \text{white}) &= 0,696205 - \underbrace{\frac{9}{16} \times 0,92}_{\substack{\text{entropy} \\ \text{white} \\ = 0}} \\ \text{Information gain} &= 0,696205 - 0,5175 \\ &= 0,178 \end{aligned}$$

$$\text{Information gain} (\text{frilly} = 0, \text{white}) = 0,178$$

$\Rightarrow \text{Frilly} = 0$, consider attribute tall's entropy.

$$\text{Entropy}(\text{Frilly} = 0) = 0.696205$$

when $\text{Frilly} = 0 \& \text{tall} = 0$ there are
→ 4 True
→ 3 False

when $\text{Frilly} = 0 \& \text{tall} = 1$ there are
→ 9 True
→ 0 False

$$\text{Entropy}(\text{tall} = 0) = -\frac{4}{7} \left(\log_2 \left(\frac{4}{7} \right) \right) - \frac{3}{7} \left(\log_2 \left(\frac{3}{7} \right) \right) = 0.981$$

$$\text{Entropy}(\text{tall} = 1) = 0$$

Information gain ($\text{Frilly} = 0, \text{tall}$) = ~~P(tall)~~

$$= 0.696205 - \underbrace{\left(P(\text{tall} = 0) \cdot E(\text{tall} = 0) \right)}_{\frac{7}{16} \quad 0.981} - 0$$

0.429

④ Therefore, $\text{Gain}(\text{white}) = 0.178$) for $\text{Frilly} = 0$
 $\text{Gain}(\text{tall}) = 0.2672$

⑤ Now consider the branch $\text{Frilly} = 1$

	<u>white</u>	<u>tall</u>	<u>frilly</u>
0	0	0	T
0	1	0	F
0	1	1	F
1	0	0	F
1	1	1	F
1	1	1	F
0	1	1	F

$$\text{entropy} (\text{frilly} = 1) = 0,9544$$

a) Consider white:

(white=0) & frilly=1 there are $\rightarrow 2T$
 $\rightarrow 3F$

$$\text{entropy} (\text{white} = 0) = -\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \cdot \log_2 \left(\frac{3}{5} \right) = -0,4 \times -1,32$$

$$-0,6 \times -0,74 \\ = 0,97 \leftarrow$$

$$\text{entropy} (\text{white} = 1) = 0$$

$$\text{Info gain} (\text{frilly} = 1, \text{white} = 0) = 0,9544 - \underbrace{\left(P(\text{white} = 0), E(\text{white} = 0) \right)}_{0,97} \underbrace{-}_{0,97}$$

$$= 0,35$$

b) Consider tall:

tall=0 & frilly=1 there are $\rightarrow 2T$
 $\rightarrow 1F$

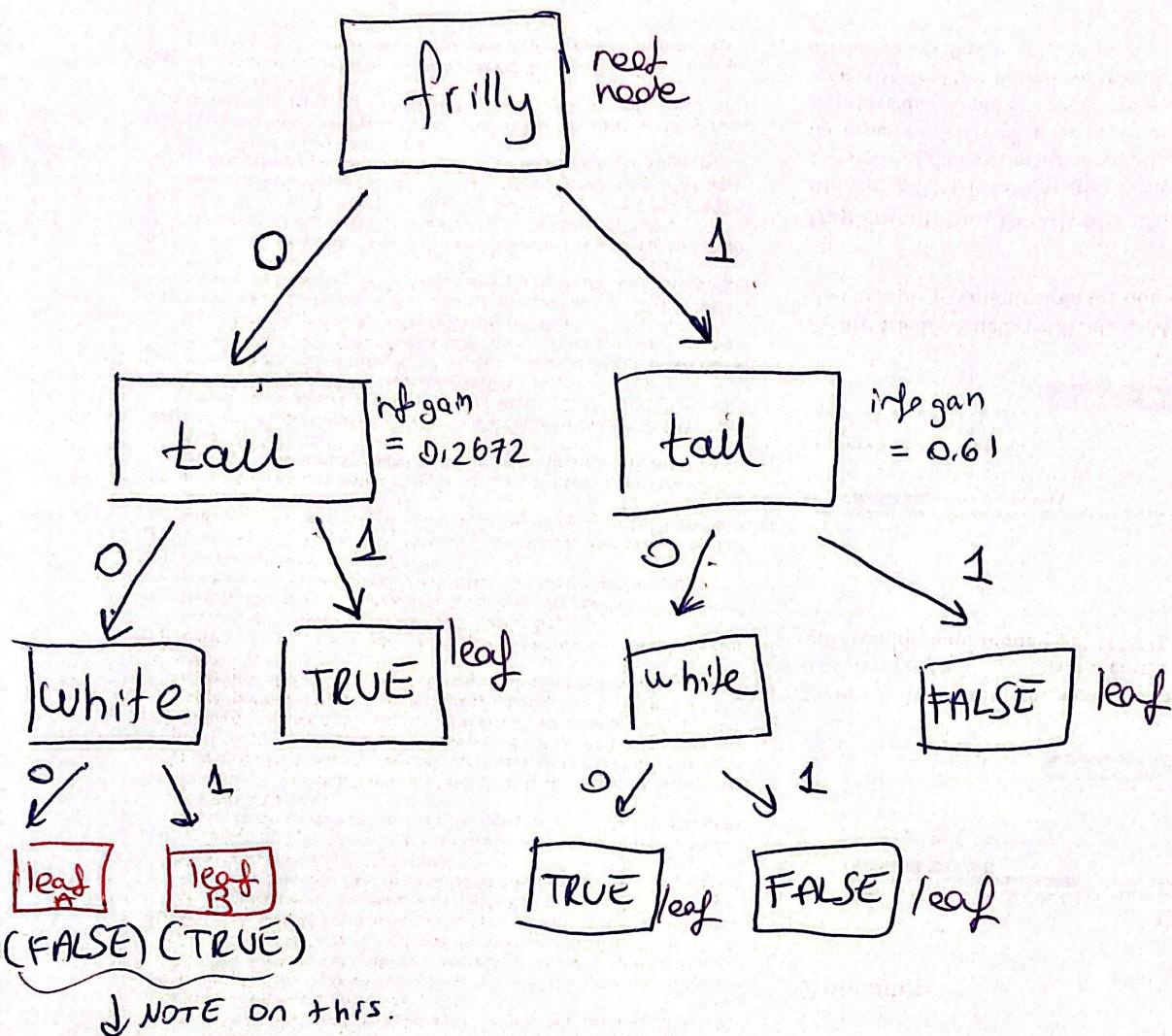
tall=1 & frilly=1 there are $\rightarrow 5F$
 $\rightarrow 0T$

$$\text{entropy} (\text{tall} = 0) = -\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) = 0,9182 \quad 0,344$$

$$\text{entropy} (\text{tall} = 1) = 0 \quad / \quad \text{inf gain} (\text{frilly} = 1, \text{tall}) = 0,9544 - 0,9182 \times \frac{3}{8}$$

$$= 0,61 \quad \cancel{0,344}$$

Therefore the tree structure is?



① When $\text{frilly} = 0, \text{tall} = 0, \text{white} = 0$ there are → 1T

② " " " , $\text{white} = 1$ " " → 3F → 3T → OF

In the raw data table case ① some entries point to different classifications. In this case my opinion is to classify $\text{frilly} = 0, \text{tall} = 0, \text{white} = 0$ (leaf A) as false and leaf B as True. So we are ignoring the false entries in case ①.

Question 2:

The tree looks same with Question 1.

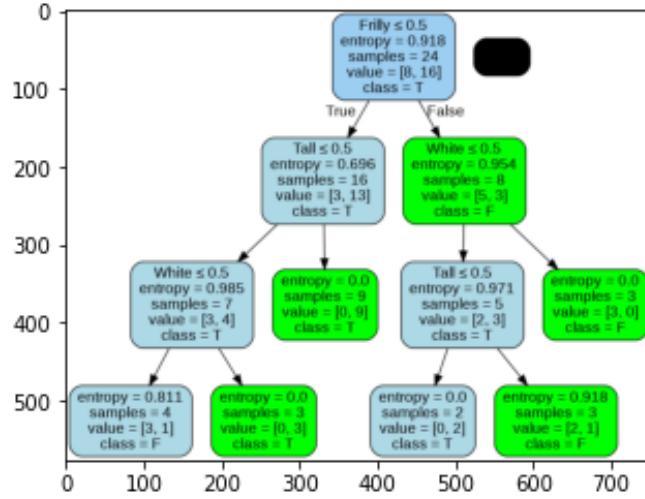


Figure 1: Decision tree

Question 3:

Tree built with entropy criteria:

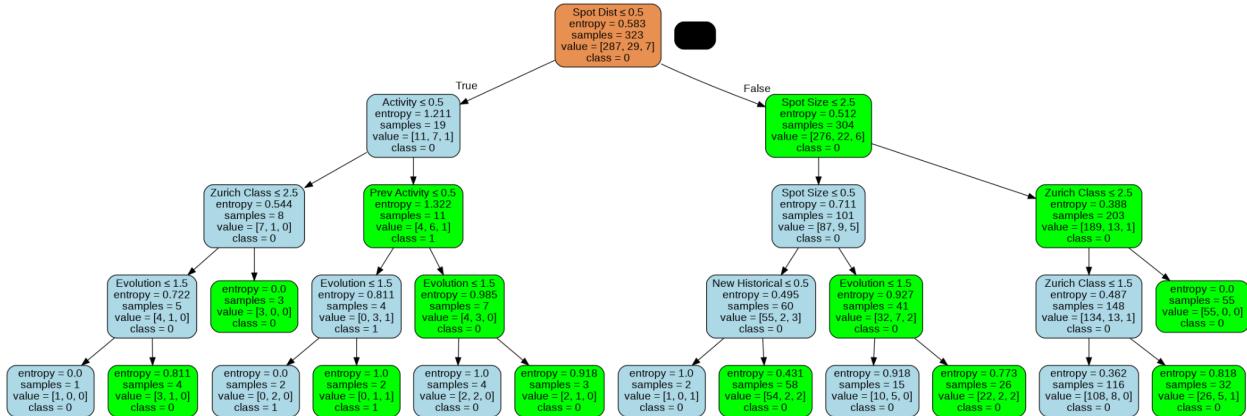


Figure 2: Entropy based built tree

Tree built with gini index:

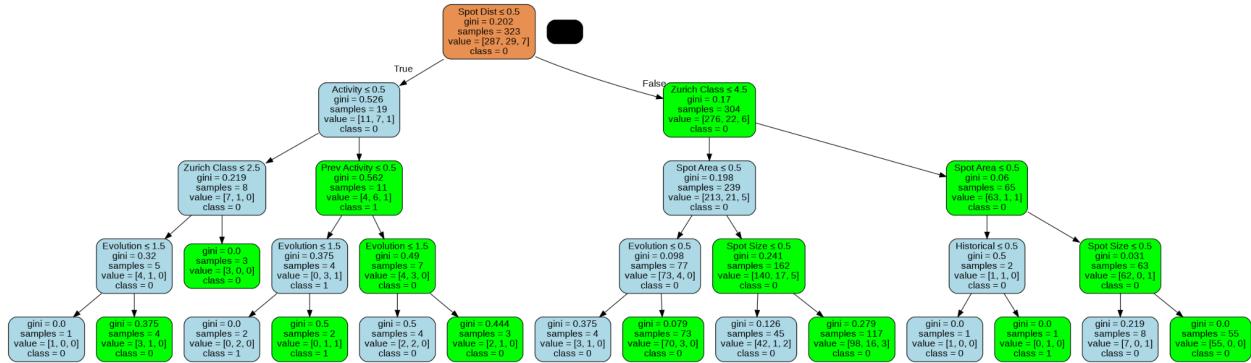


Figure 3: Gini Index based tree

Gini index seems to work better.

The codes for both questions are given at the Appendix.

Appendix:

Code:

Question 2:

```

import pandas as pd
df = pd.read_excel('AlienMushrooms.xlsx')
df.head(8)
X = df.drop(['Edible'],axis = 1)
Y = df.Edible
from sklearn import tree
model = tree.DecisionTreeClassifier(criterion='entropy')
model = model.fit(X,Y)
import pydotplus
import collections

# for a two-class tree, call this function like this:
# writegraphToFile(clf, ('F', 'T'), dirname+graphfilename)
def writegraphToFile(classifier, classnames, pathname):
    dot_data = tree.export_graphviz(model, out_file=None,      # merely to
write the tree out
                                    feature_names=X.columns,
                                    class_names=classnames,
                                    filled=True, rounded=True,
                                    special_characters=True)
    graph = pydotplus.graph_from_dot_data(dot_data)
    colors = ('lightblue', 'green')
    edges = collections.defaultdict(list)
    for edge in graph.get_edge_list():
        edges[edge.get_source()].append(int(edge.get_destination()))
    for edge in edges:
        edges[edge].sort()
        for i in range(2):
            dest = graph.get_node(str(edges[edge][i]))[0]
            dest.set_fillcolor(colors[i])
    graph.write_png(pathname)

writegraphToFile(model, ('F', 'T'), 'deno.png')

```

Question 3:

```

df = pd.read_excel('FlareData.xlsx')
df.head()
input = df.drop(['C class', 'M class', 'X class'], axis='columns')
input.head()
target = df['C class']
target.head()
from sklearn.preprocessing import LabelEncoder
le_model = LabelEncoder()
input['Zurich Class'] = le_model.fit_transform(input['Zurich Class'])
input['Spot Size'] = le_model.fit_transform(input['Spot Size'])
input['Spot Dist'] = le_model.fit_transform(input['Spot Dist'])
input['Activity'] = le_model.fit_transform(input['Activity'])
input['Evolution'] = le_model.fit_transform(input['Evolution'])
input['Prev Activity'] = le_model.fit_transform(input['Prev Activity'])
input['Historical'] = le_model.fit_transform(input['Historical'])
input['New Historical'] = le_model.fit_transform(input['New Historical'])
input['Area'] = le_model.fit_transform(input['Area'])
input['Spot Area'] = le_model.fit_transform(input['Spot Area'])
input.head(10)

model2 = tree.DecisionTreeClassifier(criterion='entropy', max_depth=4)
model2 = model2.fit(input, target)
featurelabels= ['Zurich Class', 'Spot Size', 'Spot
Dist', 'Activity', 'Evolution', 'Prev Activity', 'Historical', 'New
Historical', 'Area', 'Spot Area' ]
def writegraphToFile2(classifier, classnames, pathname):
    dot_data = tree.export_graphviz(model2, out_file=None,      # merely to
write the tree out
                                    feature_names=featurelabels,
                                    class_names=classnames,
                                    filled=True, rounded=True,
                                    special_characters=True)
    graph = pydotplus.graph_from_dot_data(dot_data)
    colors = ('lightblue', 'green')
    edges = collections.defaultdict(list)
    for edge in graph.get_edge_list():
        edges[edge.get_source()].append(int(edge.get_destination()))
    for edge in edges:
        edges[edge].sort()
        for i in range(2):

```

```
dest = graph.get_node(str(edges[edge][i]))[0]
dest.set_fillcolor(colors[i])
graph.write_png(pathname)
writegraphtofile2(model2, ('0', '1', '2'), 'ent.png')
# Read Images
img = mpimg.imread('ent.png')

# Output Images
plt.imshow(img)

model2 = tree.DecisionTreeClassifier(max_depth=4)
model2 = model2.fit(input,target)

writegraphtofile2(model2, ('0', '1', '2'), 'ent2.png')

# Read Images
img = mpimg.imread('ent2.png')

# Output Images
plt.imshow(img)
```

Deniz Aytemiz