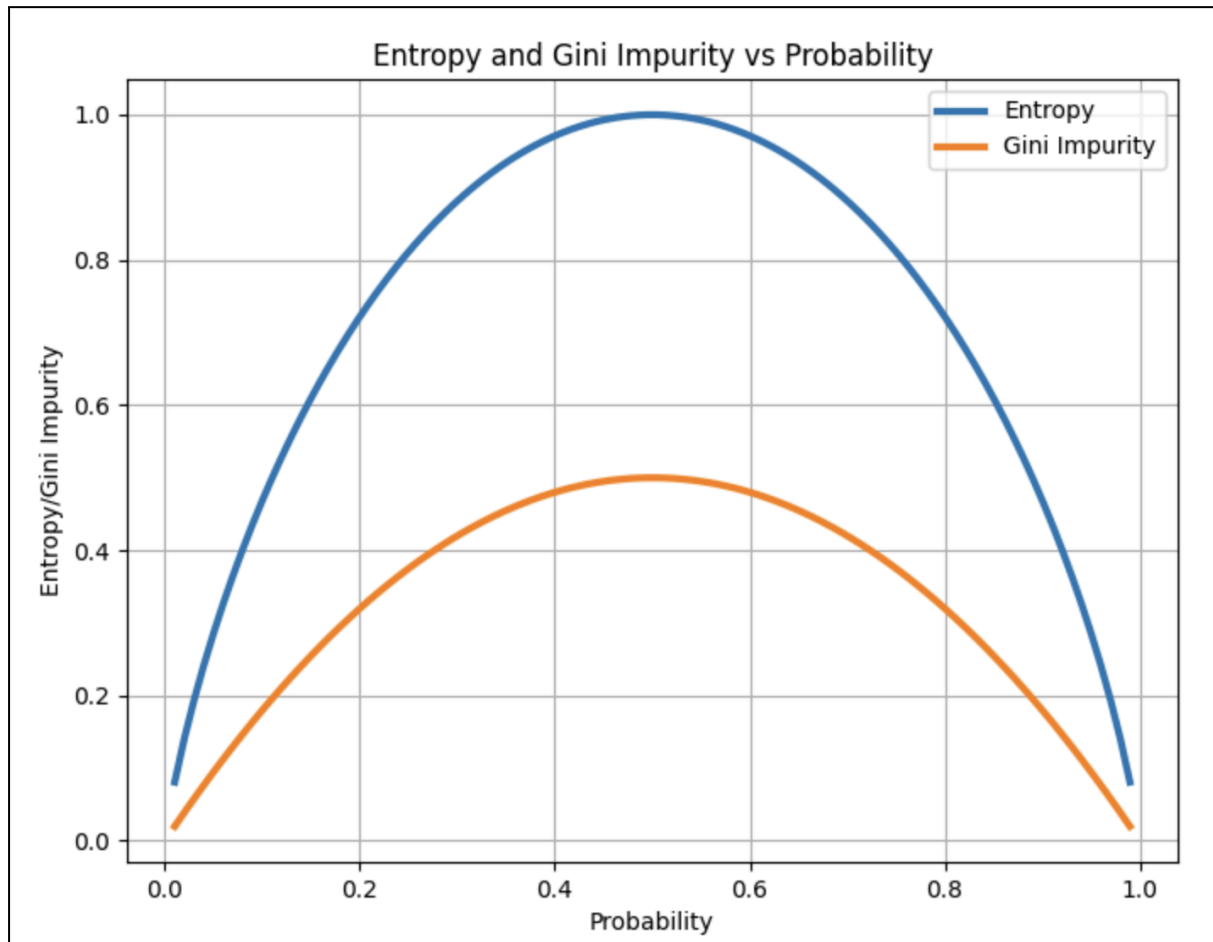


Homework 7 - CS 5805

Name : Jyothi Sevakula

Q1.



Q2.

Q2)

outlook	Yes	No	
Sunny	2	3	5
overcast	4	0	4
rain	3	2	5
	9	5	14 ✓

$$E(\text{tennis}) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$= 0.407 + 0.53057 = \underline{0.937}$$

$$E(\text{Sunny}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= 0.5275 + 0.4405 = \underline{0.968}$$

$$E(\text{overcast}) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) - 0$$

$$= 0$$

$$E(\text{Rain}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.968$$

$$\rightarrow \text{weighted entropy}$$

$$= \frac{0.968(8) + 0 + 0.968(5)}{14}$$

$$= 0.691$$

$$\text{Information gain (outlook)}$$

$$= 0.937 - 0.691$$

$$\boxed{IG(\text{tennis}) = 0.246}$$

Temp	Yes	No	
Hot	2	2	4
mild	4	2	6
cool	3	1	4
	9	5	14

$$E(\text{Temp}) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$E(\text{hot}) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

$$E(\text{mild}) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) = 0.918$$

$$E(\text{cool}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.811$$

$$\rightarrow \text{weighted entropy}$$

$$= \frac{1 \times 4 + 0.918 \times 6 + 0.811 \times 4}{14}$$

$$= 0.91$$

$$IG(\text{Temp}) = 0.94 - 0.91$$

$$= 0.03$$

$$\boxed{IG(\text{Temp}) = 0.03}$$

Humidity	Yes	No	
High	3	4	7
Normal	6	1	7
	9	5	14

$$E(\text{high}) = -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) = 0.986$$

$$E(\text{normal}) = -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) = 0.5917$$

$$\rightarrow \text{weighted entropy}$$

$$= \frac{0.986 \times 7 + 0.5917 \times 7}{14}$$

$$= 0.78885$$

$$IG(\text{Humidity}) = 0.94 -$$

$$0.78885$$

$$= 0.1485$$

$$\boxed{IG(\text{Humidity}) = 0.1485}$$

wind	Yes No		
	Yes	No	
weak	6	2	8
strong	3	3	6
	9	5	14

$$E(\text{weak}) = -\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) = 0.811$$

$$E(\text{strong}) = -\frac{3}{6} \log_2 \left(\frac{3}{6} \right) \times 2 = 1$$

weighted entropy =

$$\frac{0.811 \times 8 + 1 \times 6}{14} = 0.892$$

$$I_4(\text{wind}) = 0.937 - 0.892 = 0.045$$

$$I_4(\text{wind}) = 0.045$$

~~E(wind)~~

On comparing all I_4 values, outlook has highest information gain so it will be the root node.

→ let's see outlook on sunny

Temp	Yes No		
	Yes	No	
Hot	0	2	2
mild	1	1	2
cold	1	0	1
	2	3	5

$$E(\text{Sunny} | \text{Temp} = \text{hot}) = 0 - \frac{2}{2} \log_2(1) = 0$$

$$E(\text{Sunny} | \text{Temp} = \text{mild}) = -\frac{1}{2} \log_2 \left(\frac{1}{2} \right) \times 2 = 1$$

$$E(\text{Sunny} | \text{Temp} = \text{cold}) = -\frac{1}{1} \log_2 \left(\frac{1}{1} \right) = 0$$

$$E(\text{Sunny} | \text{Temp}) = 0 \times \frac{2}{5} + 1 \times \frac{2}{5} + 0 \times \frac{1}{5} = 0.4$$

$$I_4(\text{Sunny} | \text{Temp}) = 0.937 - 0.4 = 0.537$$

$$I_4(\text{Sunny} | \text{Temp}) = 0.537$$

Humidity (outlook → sunny)

	Yes No		
	Yes	No	
high	0	3	3
normal	2	0	2
	2	3	5

$$E(\text{Sunny} | \text{Hum} = \text{high}) = -\frac{3}{3} \log_2(1) = 0$$

$$E(\text{Sunny} | \text{Hum} = \text{normal}) = -\frac{2}{2} \log_2(1) = 0$$

$$E(\text{Sunny} | \text{Hum}) = 0$$

$$I_4(\text{Sunny} | \text{Humidity}) = 0.937 - 0 = 0.937$$

	I_4	
outlook	0.246	→ we pick this as root node.
Temp	0.03	
humidity	0.1485	
wind	0.045	

wind (outlook → sunny)

	Yes No		
	Yes	No	
strong	1	1	2
weak	1	2	3
	2	3	5

$$E(\text{Sunny} | \text{wind} = \text{strong}) = -\frac{1}{2} \log_2 \left(\frac{1}{2} \right) \times 2 = 1$$

$$E(\text{Sunny} | \text{wind} = \text{weak}) = -\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) = 0.918$$

~~E(wind)~~

$$I_4(\text{Sunny} | \text{wind}) = 0.937 - \left(\frac{1}{5} \times 1 + \frac{2}{5} \times 0.918 \right) = 0.420$$

on comparing all I_4 above we can see humidity is the highest info gain for outlook = sunny.

	I_4	
Temp	0.537	
Hum	0.937	→ we pick this child
wind	0.420	

For outlook=sunny we can say it will first branch/divide to humidity

Humidity (outlook→sunny)

high → No
Normal → yes } direct labels

→ let's see outlook on overcast

$$E(\text{overcast}) = 0$$

↓
All labels are yes

→ let's see outlook on rain

wind	yes	no	
weak	3	0	3
strong	0	2	2
	3	2	5

$$E(\text{Rain} | \text{wind} = \text{weak})$$

$$= -\frac{3}{3} \log_2(1) = 0$$

$$E(\text{Rain} | \text{wind} = \text{strong})$$

$$= -\frac{2}{2} \log_2(1) = 0$$

$$E(\text{Rain} | \text{wind}) = \frac{3}{5} \times 0 + \frac{2}{5} \times 0$$

$$IG(\text{Rain} | \text{wind}) = 0.968 - 0 = 0.968$$

Temp	yes	no	
mild	2	1	3
cool	1	1	2
	3	2	5

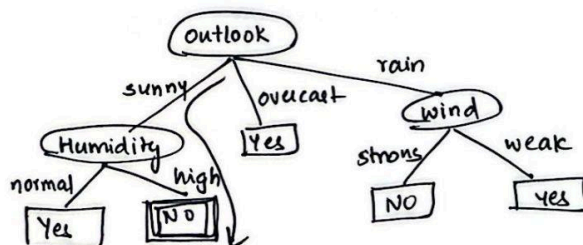
$$E(\text{rain} | \text{Temp} = \text{mild}) = 0.918$$

$$E(\text{rain} | \text{Temp} = \text{cool}) = 1$$

$$E(\text{rain} | \text{temp}) = 0.95$$

$$IG = 0.968 - 0.95 = 0.018$$

∴ we choose wind as it has highest IG. for outlook=rain



∴ For given record, outlook=sunny, Temp=cool, humidity=high, wind=strong

∴ It results in "No" therefore the player does not play Tennis

Q3.

Q3) Gini Impurity approach:-

outlook	yes	no	
sunny	2	3	5
overcast	4	0	4
rain	3	2	5
	9	5	14

Gini(parent | outlook=sunny)

$$Gini(sunny) = 1 - \left(\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right) = 0.48$$

Gini(parent | outlook=overcast)

$$= 1 - \left(\frac{4}{4} \right)^2 = 0$$

$$Gini(parent | outlook=rain) = 1 - \left(\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right)$$

$$= 0.48$$

Gini(parent | outlook)

$$= \frac{0.48 \times 5 + 0 \times 4 + 0.48 \times 5}{14} = 0.342$$

$$\boxed{Gini(parent | outlook) = 0.342}$$

Temp	yes	no	
Hot	2	2	4
mild	4	2	6
cool	3	1	4
	9	5	14

$$Gini(parent | Temp=hot) = 1 - \left(\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right) = 0.5$$

$$Gini(parent | Temp=mild) = 1 - \left(\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right) = 0.445$$

$$Gini(parent | Temp=cool) = 1 - \left(\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = 0.375$$

$$Gini(parent | Temp) = \frac{0.5 \times 4 + 0.445 \times 6 + 0.375 \times 4}{14}$$

$$\boxed{Gini(parent | Temp) = 0.441}$$

Humidity	yes	no	
High	3	4	7
Normal	6	1	7
	9	5	14

$$Gini(parent | hum=high) = 1 - \left(\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right)$$

$$= 0.49$$

$$Gini(parent | hum=normal) = 1 - \left(\left(\frac{6}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right)$$

$$= 0.245$$

$$Gini(parent | humidity) = \frac{0.49 \times 7 + 0.245 \times 7}{14}$$

$$\boxed{Gini(parent | humidity) = 0.3675}$$

wind	yes	no	
weak	6	2	8
strong	3	3	6
	9	5	14

$$Gini(parent | wind=weak) = 1 - \left(\left(\frac{6}{8} \right)^2 + \left(\frac{2}{8} \right)^2 \right)$$

$$= 0.375$$

$$Gini(parent | wind=strong) = 1 - \left(\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right)$$

$$= 0.5$$

$$Gini(parent | wind) = 0.428$$

parent

	Gini
outlook	0.342
Temp	0.441
hum	0.3675
wind	0.428

less gini index so we choose "outlook" as root node

∴ let's see outlook on sunny,

Temp	yes	no	
hot	0	2	2
mild	1	1	2
cold	1	0	1
	2	3	5

$$Gini(Sunny|Temp=hot) = 1 - \left(\frac{2}{5}\right)^2 = 0$$

$$Gini(Sunny|Temp=mild) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.5$$

$$Gini(Sunny|Temp=cold) = 1 - \left(\frac{1}{5}\right)^2 = 0$$

$$Gini(Sunny|Temp) = \frac{0.5 \times 2}{5} = 0.2$$

$$\boxed{Gini(Sunny|Temp) = 0.2}$$

Humidity	yes	no	
high	0	3	3
normal	2	0	2
	2	3	5

$$Gini(Sunny|humidity=high) = 1 - 1 = 0$$

$$Gini(Sunny|humidity=normal) = 1 - \left(\frac{2}{5}\right)^2 = 0$$

$$\boxed{Gini(Sunny|Humidity) = 0}$$

wind	yes	no	
strong	1	1	2
weak	1	2	3
	2	3	5

$$Gini(Sunny|wind=strong) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.5$$

$$Gini(Sunny|wind=weak) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.44$$

$$\boxed{Gini(Sunny|wind) = 0.466}$$

outlook=Sunny	Gini
Temp	0.2
hum	0
wind	0.466

we choose humidity as child of outlook on sunny.

For humidity high → ~~yes~~ no
normal → yes } direct labels.

let's see outlook on overcast,

all records has value yes } direct label.

let's see outlook on ~~wind~~ rain

wind	yes	no	
weak	3	0	3
strong	0	2	2
	3	2	5

$$Gini(Rain|wind=weak) = 1 - \left(\frac{3}{5}\right)^2 = 0$$

$$Gini(Rain|wind=strong) = 1 - \left(\frac{2}{5}\right)^2 = 0$$

$$\boxed{Gini(Rain|wind) = 0}$$

Temp	yes	no	
mild	2	1	3
cold	1	1	2
	3	2	5

$$Gini(Rain|Temp=mild) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.445$$

$$Gini(Rain|Temp=cold) = 1 - \left(\frac{1}{5}\right)^2 = 0.5$$

$$Gini(Rain|Temp) = \frac{0.445 \times 3 + 0.5 \times 2}{5}$$

$$\boxed{Gini(Rain|Temp) = 0.467}$$

outlook=rain

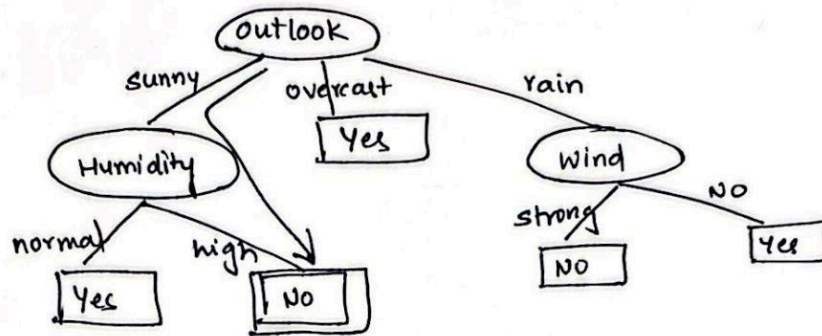
wind
Temp

Gini

0

0.467

we choose ^{'wind'} ~~rain~~ as child of outlook on rain



For the record-test,

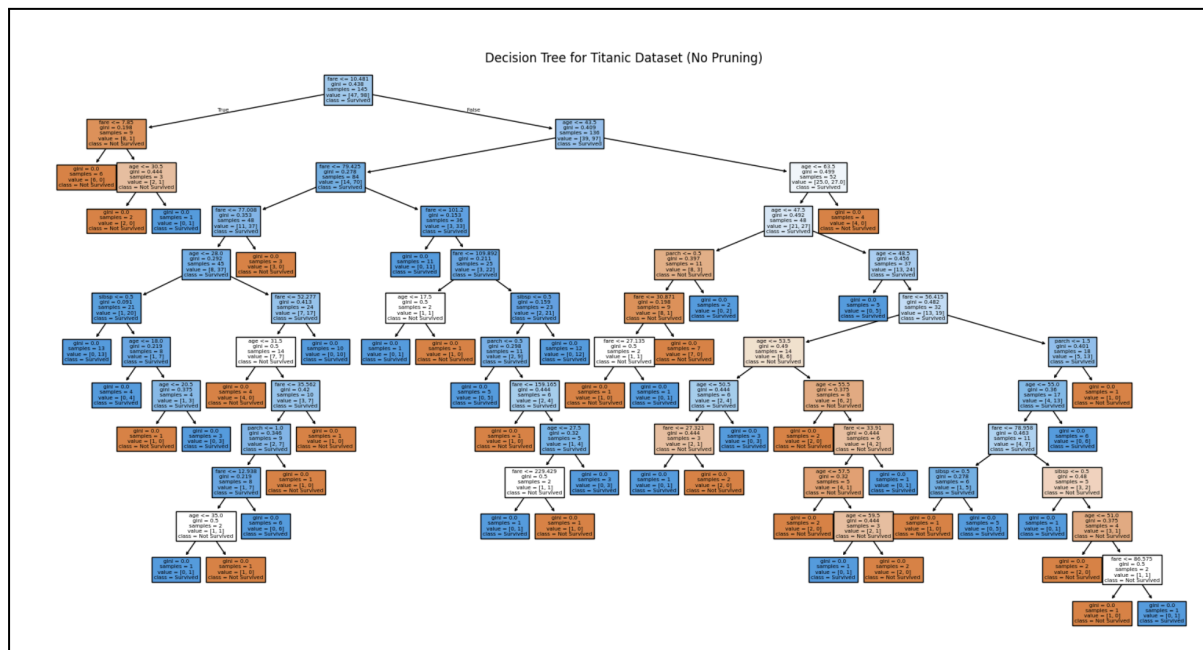
outlook=sunny, temp=cool, humidity = high & wind=strong.

∴ Based on above decision tree we get the player does not play tennis i.e. "No"

Q4.

```
Training Accuracy (No Pruning): 1.00  
Test Accuracy (No Pruning): 0.62  
Decision Tree Parameters: {'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': None, 'max_leaf_nodes': None, 'min_impur
```

Decision Tree Parameters: {'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': None, 'max_leaf_nodes': None, 'min_impurity_decrease': 0.0, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'monotonic_cst': None, 'random_state': 5805, 'splitter': 'best'}



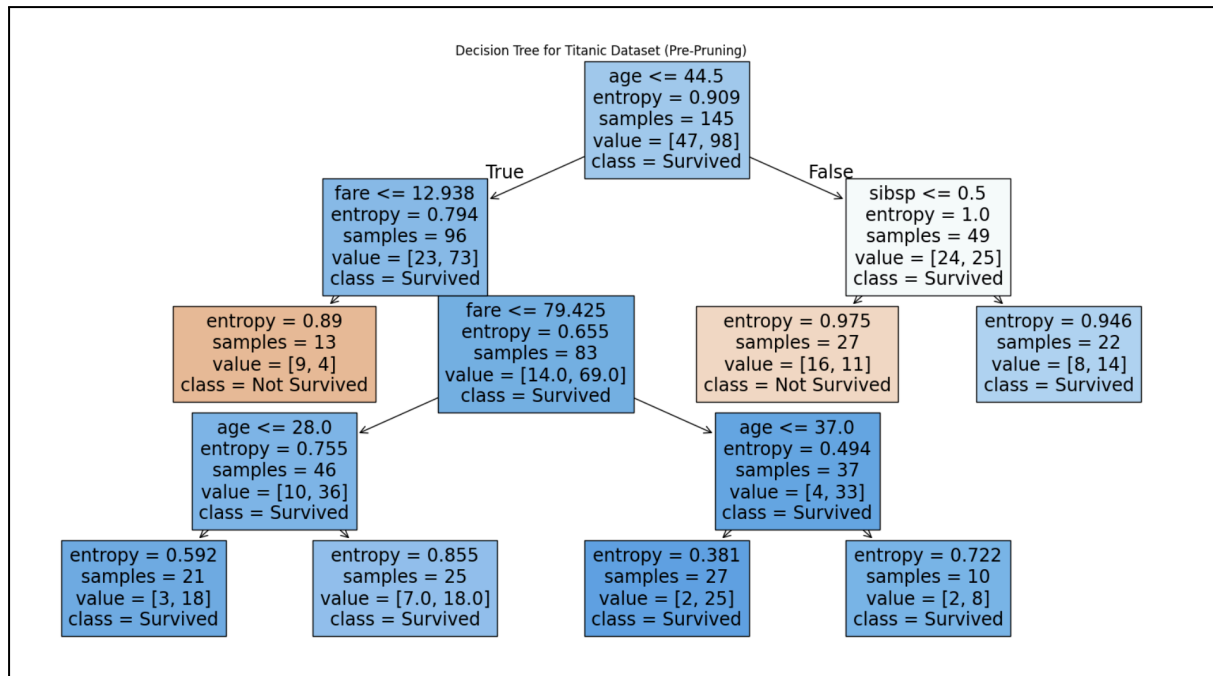
Comments → Overfitting the training data

The training accuracy is 1.00 (100%), indicating that the model has perfectly fit the training data. However, the test accuracy is significantly lower at 0.62 (62%). This suggests that the model has **overfitted the training data**, capturing noise and details specific to the training data that do not generalize well to unseen data. The model has likely created many branches that lead to an overly complex decision boundary, which does not align well with the real underlying patterns in the data.

Q5.

```
Best Parameters: {'criterion': 'entropy', 'max_depth': 4, 'max_features': 'sqrt', 'min_samples_leaf': 10, 'min_samples_split': 30, 'splitter': 'best'}  
Training Accuracy (Pre-pruned Tree): 0.74  
Test Accuracy (Pre-pruned Tree): 0.73
```

Best Parameters: {'criterion': 'entropy', 'max_depth': 4, 'max_features': 'sqrt', 'min_samples_leaf': 10, 'min_samples_split': 30, 'splitter': 'best'}



Comments → Improved model performance (reduced overfitting)

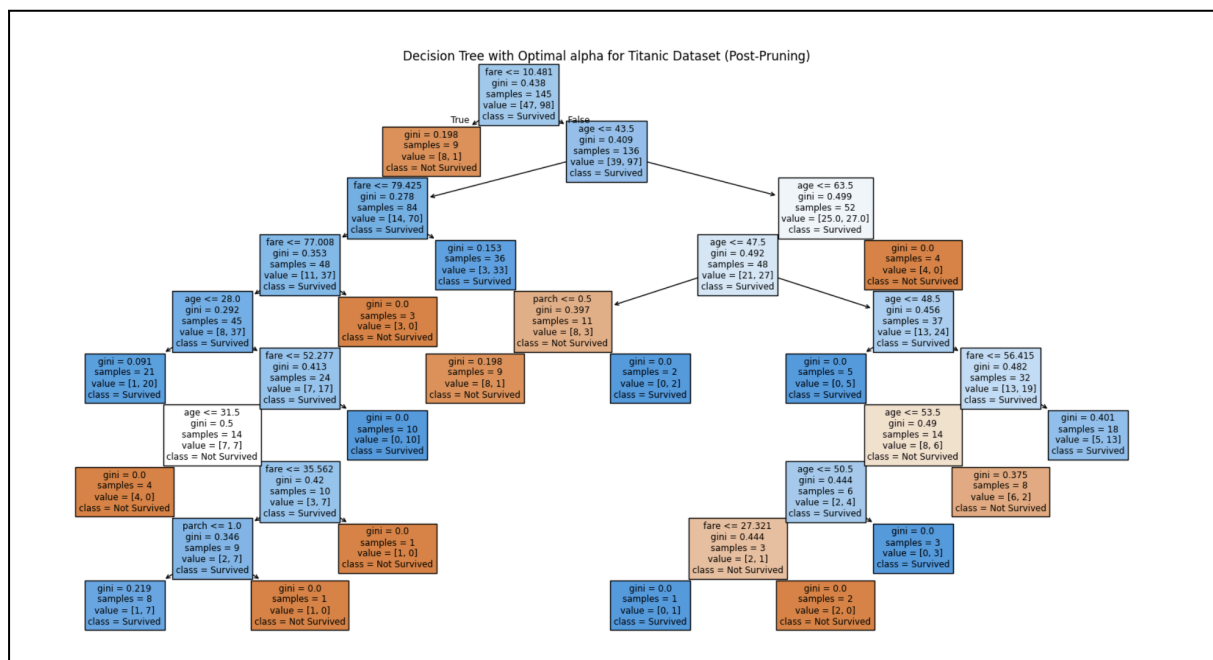
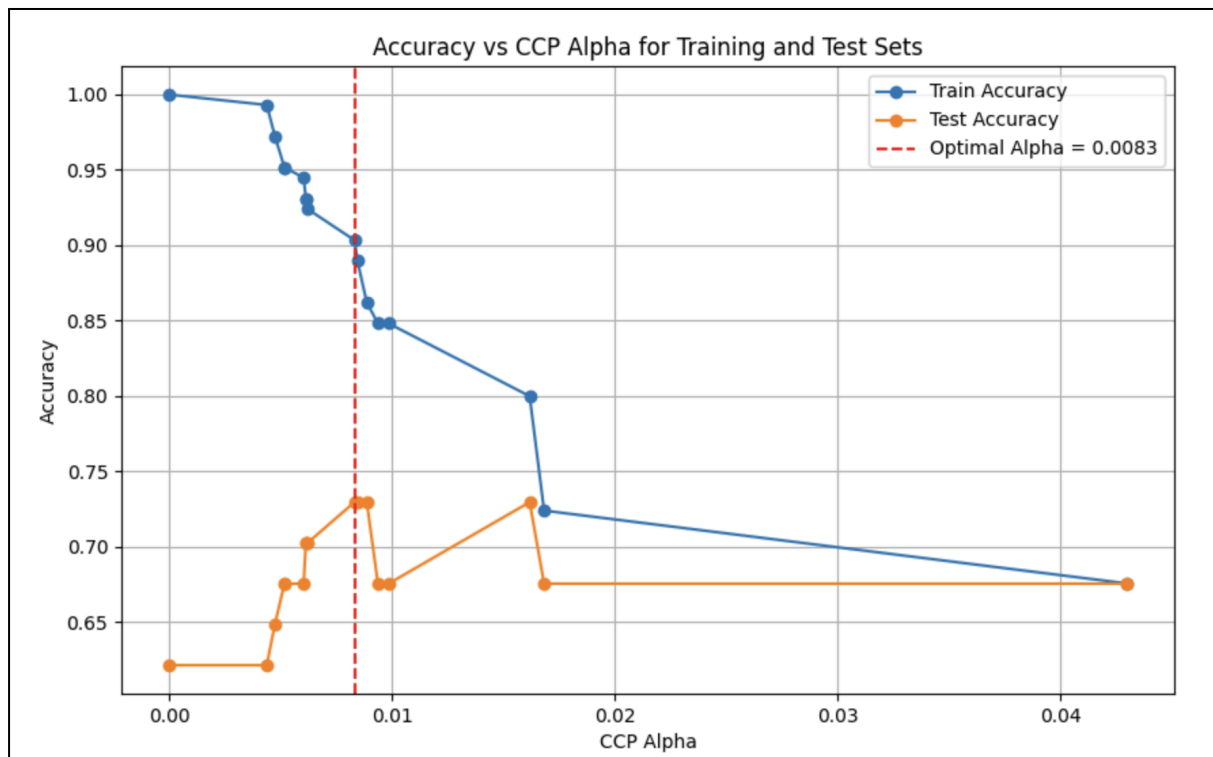
After applying pre-pruning using Grid Search with cross-validation. The training accuracy (0.74) and test accuracy (0.73) are very close. This suggests that the **model is not overfitting**, as it performs similarly on both the training and test data. **Pre-pruning improved the model's performance** by reducing overfitting, leading to a more generalizable model. Although the training accuracy is slightly lower than in the no-pruning scenario, the consistency between training and test accuracies (0.74 and 0.73) indicates that the model is better suited for making reliable predictions on unseen data.

Q6.

Optimal CCP Alpha: 0.00831417624521073

Training Accuracy (Post-pruned Tree): 0.90

Test Accuracy (Post-pruned Tree): 0.73



Comments → Reduced overfitting when compared to no-pruned but less effective when compared to pre-pruned.

After applying post-pruning using optimum alpha in the cost complexity function. The training accuracy (0.9) and test accuracy (0.73). In short, ***Post-pruning reduced overfitting and improved test performance compared to the no-pruned model.*** However, ***it was slightly less effective at controlling overfitting compared to the pre-pruned model***, as indicated by the larger gap between training and test accuracies.

Comparison with Pre-Pruning:

- The pre-pruned model (training: 0.74, test: 0.73) achieved a better balance between training and test performance than the post-pruned model. With pre-pruning, the tree was constrained from growing too deep, leading to a simpler model that generalized better.
- In contrast, the post-pruning approach allows the tree to grow fully before pruning back some branches. Although post-pruning improves over the no-pruning model, it appears to be slightly less effective at controlling overfitting than pre-pruning in this case.

Comparison with No Pruning:

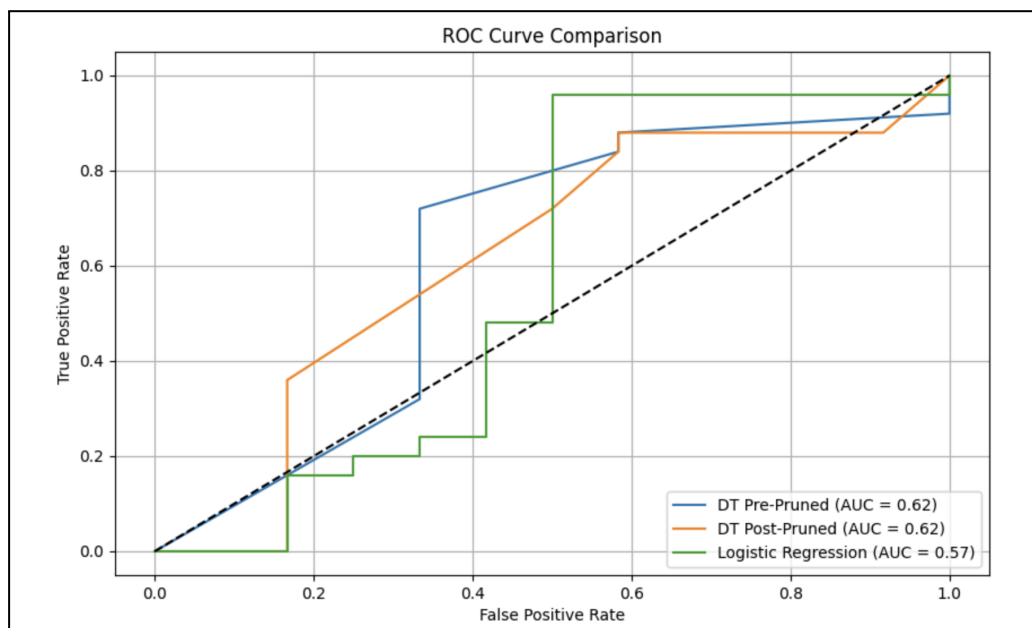
- Compared to the no-pruned model (training: 1.00, test: 0.62), post-pruning significantly reduces overfitting, improving the test accuracy from 0.62 to 0.73.
- The training accuracy has dropped from 1.00 to 0.90, showing that post-pruning effectively removes branches that capture noise or overly specific patterns in the training data.

Q7.

```
Training Accuracy (Logistic Regression): 0.73
Test Accuracy (Logistic Regression): 0.78
```

Q8.

Index	Model	Train Accuracy	Test Accuracy	Recall	AUC	Confusion Matrix
0	Pre-Pruning	0.74	0.73	0.88	0.62	$\begin{bmatrix} 5 & 7 \\ 3 & 22 \end{bmatrix}$
1	Post-Pruning	0.90	0.73	0.88	0.62	$\begin{bmatrix} 5 & 7 \\ 3 & 22 \end{bmatrix}$
2	Logistic Regression	0.73	0.78	0.96	0.57	$\begin{bmatrix} 5 & 7 \\ 1 & 24 \end{bmatrix}$



Comment: → Pre-pruning is the best classifier

Among DT Pre-pruned, DT Post-pruned and logistic regression comparing based on AUC, both DT Post-pruned and DT Pre-pruned have highest which is 0.62. Among post and pre pruned model, post-pruning model is little over-fitted on the training dataset hence the gap between train and test accuracy for post-pruning is higher. Hence, I would say Pre-pruning model classifier works best among all three model as the difference in test and train accuracy is very less when compared to others.(i.e 0.01) and also AUC is high.