

1. Classification: Playing Tennis (10 points) Provide necessary steps in the problem solving (e.g. list the steps and intermediate results to calculate entropy, information gain, and GainRatio, prior probabilities). Draw the final decision trees (e.g. use Powerpoint to draw and save as a picture or any other way you prefer)

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

- Build a decision tree using ID3 algorithm (Information Gain). (3 pts)
- Build a decision tree using CART algorithm (Gini index). (3 pts)
- Make predictions of (D15: Rain, Mild, Normal, Strong) using both trees. (2 pts)
- Use Naïve Bayes classifier to predict the result in (c). (2 pts)

Part A) Build a decision tree using ID3 algorithm (Information Gain).

$$\text{Entropy} = \sum_i -p_i \log_2 p_i \quad \begin{array}{l} \Delta = I(\text{before splitting}) - I(\text{after splitting}) \\ \Delta = I(\text{parent}) - \text{weighted_average}(I(\text{children})) \end{array}$$

First we need to calculate the Parent Entropy = $-\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.9403$

Now we have to calculate the entropy for each class (i.e. Outlook Temperature ...etc.) and subtract it from the parent. The one which has gain we choose it as root node.

For outlook we have three subclasses (sunny, overcast and rain) we will calculate the entropy for each subclass and average their entropy.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|---------|-------------|----------|--------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |

$$\text{Entropy (decision | outlook = Sunny)} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.9710$$

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D3 | Overcast | Hot | High | Weak | Yes |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |

Entropy (decision | outlook = Overcast) = 0 (pure set).

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|---------|-------------|----------|--------|-------------|
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$\text{Entropy (decision | outlook = Rain)} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.9710$$

Now we average them as follow:

$$p(\text{sunny}) * \text{Entropy}(\text{Sunny}) + p(\text{overcast}) * \text{Entropy}(\text{Overcast}) + p(\text{Rain}) * \text{Entropy}(\text{Rain})$$

$$= \frac{5}{14} * 0.97 + \frac{4}{14} * 0 + \frac{5}{14} * 0.97 = 0.6936$$

Now we calculate the information gain of outlook as follow.

$$\text{Information Gain (outlook)} = 0.9403 - 0.6936 = \mathbf{0.2467}$$

Now we do the same for temperature

For temperature we have three subclasses (hot, mild and cool) we will calculate the entropy for each subclass and average their entropy.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |

$$\text{Entropy (decision | Temperature = hot)} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D4 | Rain | Mild | High | Weak | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$\text{Entropy (decision | Temperature = mild)} = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9183$$

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D9 | Sunny | Cool | Normal | Weak | Yes |

$$\text{Entropy (decision | Temperature = cool)} = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113$$

Now we average them as follow:

$$p(\text{hot}) * \text{Entropy}(\text{hot}) + p(\text{mild}) * \text{Entropy}(\text{mild}) + p(\text{cool}) * \text{Entropy}(\text{cool})$$

$$= 4/14 * 1 + \frac{6}{14} * 0.9183 + \frac{4}{14} * 0.8113 = 0.9111$$

Now we calculate the information gain of outlook as follow: parent entropy- entropy of outlook.

$$\text{Information Gain (Temperature)} = 0.9403 - 0.9111 = \mathbf{0.0292}$$

Now we do the same for humidity

For humidity we have two subclasses (high and normal) we will calculate the entropy for each subclass and average their entropy.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D12 | Overcast | Mild | High | Strong | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$\text{Entropy (decision | Humidity= high)} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |

$$\text{Entropy (decision | Humidity= normal)} = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.5917$$

Now we average them as follow:

$$p(\text{high}) * \text{Entropy}(\text{high}) + p(\text{normal}) * \text{Entropy}(\text{normal})$$

$$= 7/14 * 0.9852 + \frac{7}{14} * 0.5917 = 0.7884$$

$$\text{Information Gain (Humidity)} = 0.9403 - 0.7884 = \mathbf{0.1519}$$

Now we do the same for wind

For humidity we have two subclasses (strong and weak) we will calculate the entropy for each subclass and average their entropy.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|--------|-------------|
| D2 | Sunny | Hot | High | Strong | No |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$\text{Entropy (decision | wind = strong)} = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|----------|-------------|----------|------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |

$$\text{Entropy (decision | wind = weak)} = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.8113$$

Now we average them as follow:

$$p(\text{strong}) * \text{Entropy}(\text{strong}) + p(\text{weak}) * \text{Entropy}(\text{weak})$$

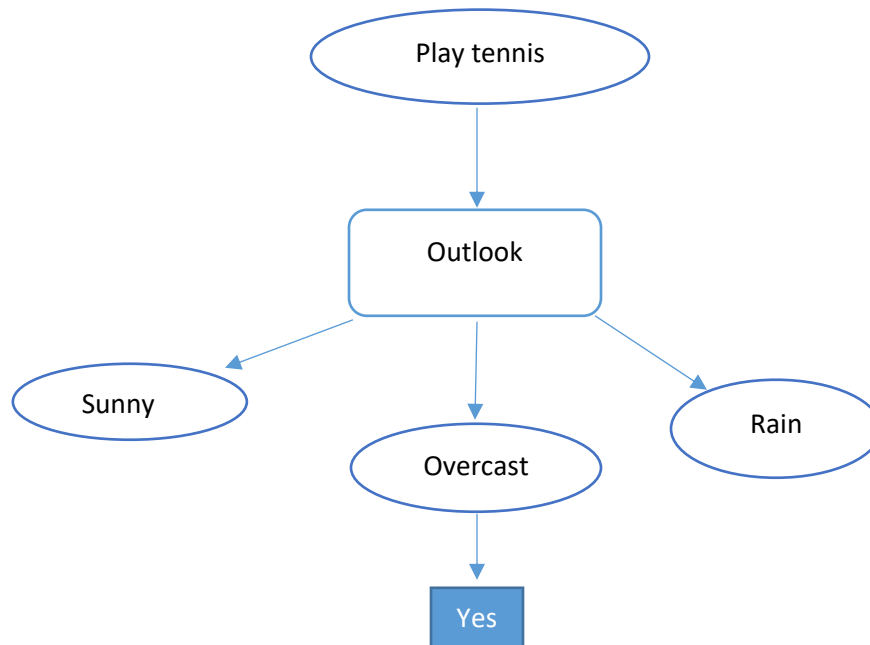
$$= 6/14 * 1 + \frac{8}{14} * 0.8113 = 0.8922$$

$$\text{Information Gain (wind)} = 0.9403 - 0.8922 = \mathbf{0.0481}$$

| | |
|---------------------------------------|---------------|
| Information Gain (outlook) | 0.2467 |
| Information Gain (Temperature) | 0.0292 |
| Information Gain (Humidity) | 0.1519 |
| Information Gain (wind) | 0.0481 |

According to ID3 algorithm we choose the node that has the highest information gain which is outlook.

Now the decision tree look like this



Now our new data is the sunny data. we will calculate parent node for it and calculate information gain for the other classes (temperature , humidity and wind) holding sunny as parent node.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|---------|-------------|----------|--------|-------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |

First we need to calculate the **Sunny Parent Entropy** = $-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = \mathbf{0.9710}$

We calculate entropy for **temperature** (i.e. average entropy for hot, mild, and cold)

$$E(\text{hot}) = -0 - \frac{2}{5} \log_2 \frac{2}{5} = 0.5288, E(\text{mild}) = 1, E(\text{cool}) = -\frac{1}{5} \log_2 \frac{1}{5} - 0 = 0.4644$$

Now we average them as follow

$$p(\text{hot}) * \text{Entropy}(\text{hot}) + p(\text{mild}) * \text{Entropy}(\text{mild}) + p(\text{cool}) * \text{Entropy}(\text{cool})$$

$$2/5 * 0.5288 + 2/5 * 1 + 1/5 * 0.4644 = 0.7044$$

$$\text{Information Gain @ sunny (temperature)} = 0.9710 - 0.7044 = \mathbf{0.2666}$$

We will do the same for humidity and wind

We calculate entropy for **humidity** (i.e. average entropy for high and normal)

$$E(\text{high}) = -0 - \frac{3}{5} \log_2 \frac{3}{5} = 0.4422 \quad E(\text{normal}) = -\frac{2}{5} \log_2 \frac{2}{5} - 0 = 0.5288$$

Now we average them as follow

$$p(\text{high}) * \text{Entropy}(\text{high}) + p(\text{normal}) * \text{Entropy}(\text{normal})$$

$$3/5 * 0.4422 + 2/5 * 0.5288 = 0.4768$$

$$\text{Information Gain @ sunny (humidity)} = 0.9710 - 0.4768 = \mathbf{0.4942}$$

We calculate entropy for **wind** (i.e. average entropy for strong and weak)

$$E(\text{strong}) = 1 \quad E(\text{weak}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.9183$$

Now we average them as follow

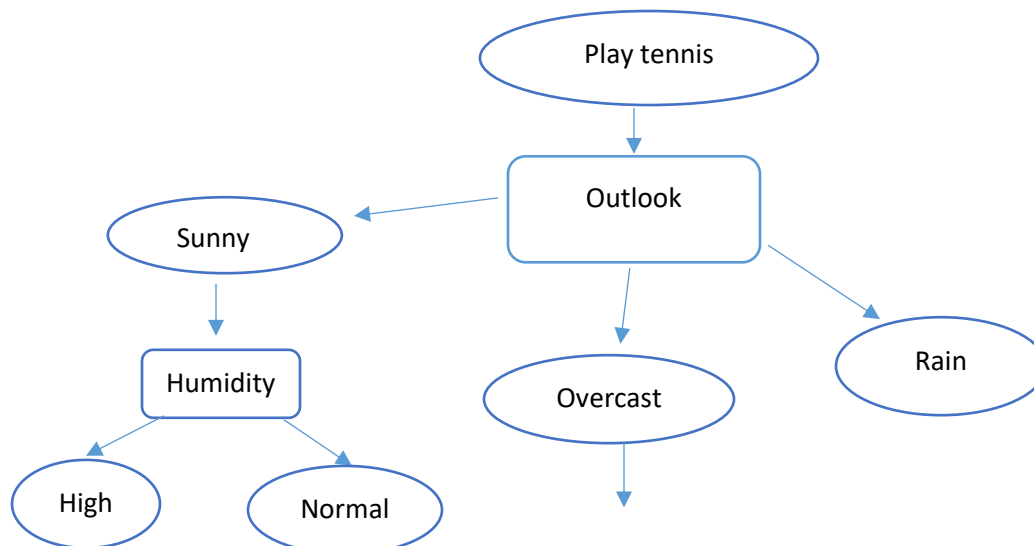
$$p(\text{strong}) * \text{Entropy}(\text{strong}) + p(\text{weak}) * \text{Entropy}(\text{weak})$$

$$2/5 * 1 + 3/5 * 0.9183 = 0.9510$$

$$\text{Information Gain @ sunny (wind)} = 0.9710 - 0.9510 = \mathbf{0.0200}$$

| | |
|---|---------------|
| Information Gain @ sunny (Temperature) | 0.2666 |
| Information Gain @ sunny (Humidity) | 0.4942 |
| Information Gain @ sunny (wind) | 0.0200 |

We will choose the one with high information gain, now our tree looks like this



Now, our new data is the rain data. We will calculate parent node for it and calculate information gain for the other class (temperature, humidity and wind) holding rain as parent node.

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|---------|-------------|----------|--------|-------------|
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

First we need to calculate the **Rain Parent Entropy** = $-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = \mathbf{0.8236}$

We calculate entropy for **temperature** (i.e. average entropy for hot, mild, and cold)

$$E(\text{hot}) = 0, \quad E(\text{mild}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9183, \quad E(\text{cool}) = 1$$

Now we average them as follow

$$p(\text{hot}) * \text{Entropy}(\text{hot}) + p(\text{mild}) * \text{Entropy}(\text{mild}) + p(\text{cool}) * \text{Entropy}(\text{cool})$$

$$0 + 3/5 * 0.9183 + 2/5 * 1 = 0.9510$$

$$\mathbf{\text{Information Gain @ Rain (temperature) = } 0.9710 - 0.9510 = \mathbf{0.0200}}$$

We calculate entropy for **humidity** (i.e. average entropy for high and normal)

$$E(\text{high}) = 1, \quad E(\text{normal}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9183$$

Now we average them as follow: $p(\text{high}) * \text{Entropy}(\text{high}) + p(\text{normal}) * \text{Entropy}(\text{normal})$

$$2/5 * 1 + 3/5 * 0.9183 = 0.9510$$

$$\mathbf{\text{Information Gain @ Rain (humidity) = } 0.9710 - 0.9510 = \mathbf{0.0200}}$$

We calculate entropy for **wind** (i.e. average entropy for strong and weak)

$$E(\text{strong}) = -0 - \frac{2}{5} \log_2 \frac{2}{5} = 0.5288, \quad E(\text{weak}) = -\frac{3}{5} \log_2 \frac{3}{5} - 0 = 0.4422$$

Now we average them as follow: $p(\text{strong}) * \text{Entropy}(\text{strong}) + p(\text{weak}) * \text{Entropy}(\text{weak})$

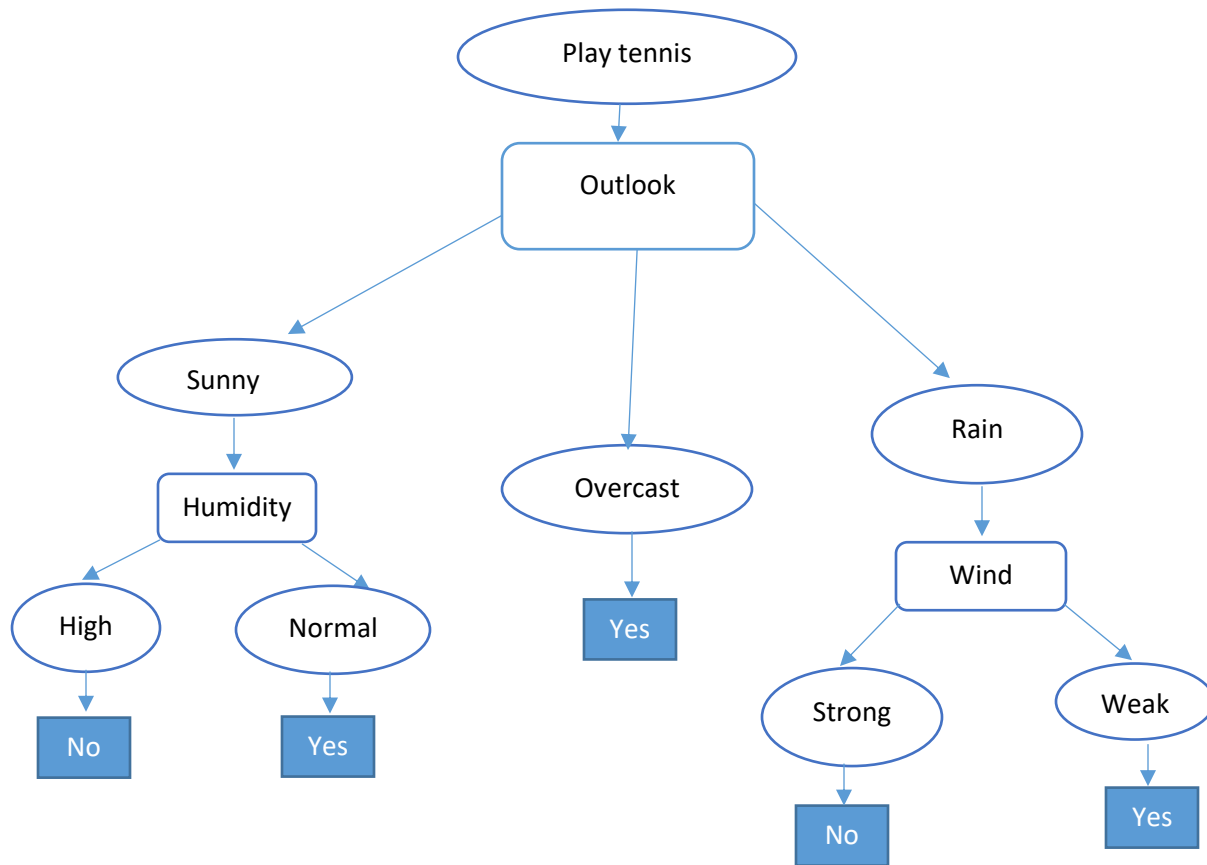
$$2/5 * 0.5288 + 3/5 * 0.4422 = 0.4768$$

$$\mathbf{\text{Information Gain @ sunny (wind) = } 0.9710 - 0.9510 = \mathbf{0.4942}}$$

| | |
|---|---------------|
| Information Gain @ sunny (Temperature) | 0.0200 |
| Information Gain @ sunny (Humidity) | 0.0200 |
| Information Gain @ sunny (wind) | 0.4942 |

We will choose the one with high information gain, now our tree looks like this

Now our final decision tree is as follow



Part b) Build a decision tree using CART algorithm (Gini index).

If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

If a data set D is split on A into two subsets D1 and D2, the gini index $gini(D)$ is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

§ Reduction in Impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

First we compute gini index for dataset as follow:

$$1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.4592$$

Now we will calculate gini index for each class attributes as follow

First outlook:

$$Gini(\text{PlayTennis} | \text{Outlook} = \text{Sunny}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.4800$$

$$Gini(\text{PlayTennis} | \text{Outlook} = \text{Overcast}) = 1 - \left(\frac{4}{4}\right)^2 - 0^2 = 0$$

$$Gini(\text{PlayTennis} | \text{Outlook} = \text{Rain}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.4800$$

Therefore, the Gini index after the outlook test is applied is

$$Gini(\text{outlook}) = 5/14 * 0.4800 + 4/14 * 0 + 5/14 * 0.4800 = 0.3429$$

$$\Delta gini(\text{outlook}) = gini(\text{play_tennis}) - gini(\text{outlook}) = 0.4592 - 0.3429 = \mathbf{0.1163}$$

Now we do the same for temperature

$$\text{Gini (PlayTennis| Temperature =Hot)} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5000$$

$$\text{Gini (PlayTennis| Temperature =Mild)} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.3056$$

$$\text{Gini (PlayTennis| Temperature =Cool)} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.3750$$

Therefore, the Gini index after the Temperature test is applied is

$$4/14 * 0.5000 + 6/14 * 0.3056 + 4/14 * 0.3750 = 0.3810$$

$$\text{Delta gini(temperature)} = \text{gini(play_tennis)} - \text{gini (temperature)} = 0.4592 - 0.3810 = \mathbf{0.0782}$$

Now we do the same for humidity

$$\text{Gini (PlayTennis|Humidity=High)} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.4898$$

$$\text{Gini (PlayTennis|Humidity=Normal)} = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.2449$$

Therefore, the Gini index after the humidity test is applied is

$$7/14 * 0.4898 + 7/14 * 0.2449 = 0.3674$$

$$\text{Delta gini(Humidity)} = \text{gini(play_tennis)} - \text{gini (humidity)} = 0.4592 - 0.3674 = \mathbf{0.0918}$$

Now we do the same for wind

$$\text{Gini (PlayTennis|Wind=strong)} = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.5000$$

$$\text{Gini (PlayTennis|Wind=weak)} = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.3750$$

Therefore, the Gini index after the Wind test is applied is

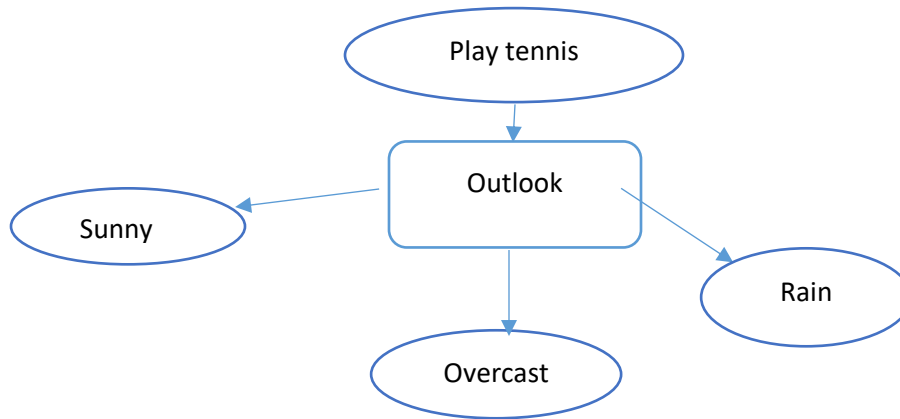
$$6/14 * 0.5 + 8/14 * 0.3750 = 0.4286$$

$$\text{Delta gini(wind)} = \text{gini(play_tennis)} - \text{gini (wind)} = 0.4592 - 0.4286 = \mathbf{0.0306}$$

| | |
|---------------------------|---------------|
| Delta gini(outlook)= | 0.1163 |
| Delta gini(temperature)= | 0.0782 |
| Delta gini(Humidity)= | 0.0918 |
| Delta gini(wind)= | 0.0306 |

We will choose the largest delta gini index which is outlook as node.

So our decision tree will be like this



Now we will calculate the gini index for sunny as follow

$$1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.4800$$

Now we do the same for temperature

$$\text{Gini (sunny| Temperature =Hot)} = 1 - 0^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini (sunny| Temperature =Mild)} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5000$$

$$\text{Gini (sunny| Temperature =Cool)} = 1 - \left(\frac{1}{1}\right)^2 - 0^2 = 0$$

Therefore, the Gini index after the Temperature test is applied is

$$2/5 * 0 + 2/5 * 0.5000 + 1/5 * 0 = 0.2000$$

$$\text{Delta gini(temperature @sunny)} = \text{gini(sunny)} - \text{gini (temperature)} = 0.4800 - 0.2000 = \mathbf{0.2800}$$

Now we do the same for humidity

$$\text{Gini (sunny| Humidity=High)} = 1 - 0^2 - \left(\frac{3}{3}\right)^2 = 0$$

$$\text{Gini (sunny| Humidity =Normal)} = 1 - \left(\frac{2}{2}\right)^2 - 0^2 = 0$$

Therefore, the Gini index after the Humidity test is applied is

$$3/5 * 0 + 2/5 * 0 = 0$$

$$\text{Delta gini(humidity @sunny)} = \text{gini(sunny)} - \text{gini (humidity)} = 0.4800 - 0 = \mathbf{0.4800}$$

Now we do the same for wind

$$\text{Gini (sunny| Wind=strong)} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5000$$

$$\text{Gini (sunny| Wind=weak)} = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.4444$$

Therefore, the Gini index after the Wind test is applied is

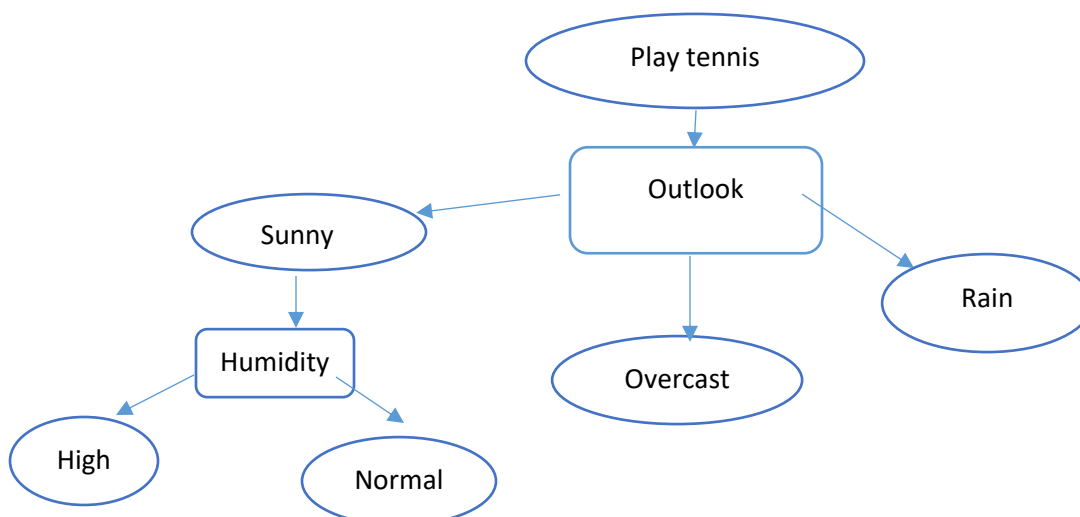
$$2/5 * 0.5000 + 3/5 * 0.4444 = 0.$$

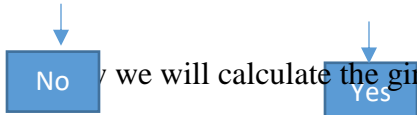
$$\text{Delta gini(wind @sunny)} = \text{gini(sunny)} - \text{gini (wind)} = 0.4800 - 0.4666 = \mathbf{0.0134}$$

| | |
|----------------------------------|---------------|
| Delta gini(temperature @sunny)= | 0.2800 |
| Delta gini(Humidity @sunny)= | 0.4800 |
| Delta gini(wind @sunny)= | 0.0134 |

We will choose the largest delta gini index which is humidity as node

So our decision tree will be like this




 we will calculate the gini index for rain as follow

$$1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.4800$$

Now we do the same for temperature

$$\text{Gini (rain| Temperature =Hot)} = 0$$

$$\text{Gini (rain| Temperature =Mild)} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.4444$$

$$\text{Gini (rain| Temperature =Cool)} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5000$$

Therefore, the Gini index after the Temperature test is applied is

$$0 + \frac{3}{5} * 0.4444 + \frac{2}{5} * 0.5000 = 0.4666$$

$$\text{Delta gini(temperature @rain)} = \text{gini(rain)} - \text{gini (temperature)} = 0.4800 - 0.4666 = \mathbf{0.0134}$$

Now we do the same for humidity

$$\text{Gini (rain| Humidity=High)} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5000$$

$$\text{Gini (rain| Humidity =Normal)} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.4444$$

Therefore, the Gini index after the Humidity test is applied is

$$\frac{2}{5} * 0.5000 + \frac{3}{5} * 0.4444 = 0.4666$$

$$\text{Delta gini(humidity @rain)} = \text{gini(rain)} - \text{gini (humidity)} = 0.4800 - 0.4666 = \mathbf{0.0134}$$

Now we do the same for wind

$$\text{Gini (sunny| Wind=strong)} = 1 - 0^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini (sunny| Wind=weak)} = 1 - \left(\frac{3}{3}\right)^2 - 0^2 = 0$$

Therefore, the Gini index after the Wind test is applied is

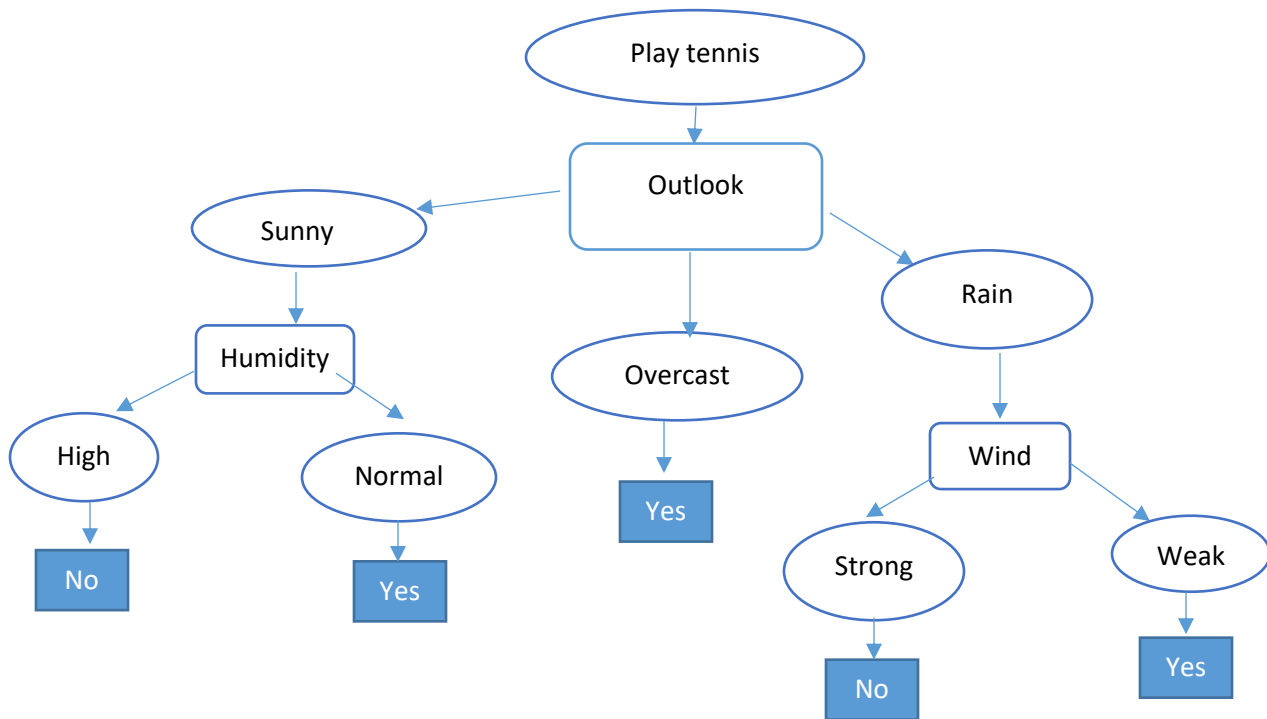
$$\frac{2}{5} * 0 + \frac{3}{5} * 0 = 0$$

$$\text{Delta gini(wind @rain)} = \text{gini(rain)} - \text{gini (wind)} = 0.4800 - 0 = \mathbf{0.4800}$$

| | |
|---------------------------------|---------------|
| Delta gini(temperature @rain)= | 0.0134 |
| Delta gini(Humidity @rain)= | 0.0134 |
| Delta gini(wind @rain)= | 0.4800 |

We will choose the largest delta gini index which is wind as node

So our final decision tree will be like this



Part c) Make predictions of (D15: Rain, Mild, Normal, Strong) using both trees. (2 pts)

Based on both trees the player will not play on D15.

Part d) Use Naïve Bayes classifier to predict the result in (c). (2 pts)

Bayes theorem:
$$P(C | A) = \frac{P(A | C)P(C)}{P(A)}$$

The diagram shows the formula $P(C | \mathbf{x}) = \frac{P(C)p(\mathbf{x} | C)}{p(\mathbf{x})}$. Arrows point from the labels 'prior' to $P(C)$, 'likelihood' to $p(\mathbf{x} | C)$, 'evidence' to $p(\mathbf{x})$, and 'posterior' to $P(C | \mathbf{x})$.

Class: C1:play_tennis = yes , C2: play_tennis = no

Features:X = (outlook= sunny, overcast and rain, Temperature = hot, mild and cool, Humidity= high and normal, Wind= strong and weak)

Now we calculate probabilities

$P(C_i): P(\text{play_tennis} = \text{yes}) = 9/14 = 0.6428$ and $P(\text{play_tennis} = \text{no}) = 5/14 = 0.3571$

Compute $P(X|C_i)$ for each class

D15: Rain, Mild, Normal, Strong

$P(\text{Outlook} = \text{rain} | \text{play_tennis} = \text{yes}) = 3/5 = 0.6$

$P(\text{Outlook} = \text{rain} | \text{play_tennis} = \text{no}) = 2/5 = 0.4$

$P(\text{Temperature} = \text{mild} | \text{play_tennis} = \text{yes}) = 4/6 = 0.67$

$P(\text{Temperature} = \text{mild} | \text{play_tennis} = \text{no}) = 2/6 = 0.33$

$P(\text{Humidity} = \text{normal} | \text{play_tennis} = \text{yes}) = 6/7 = 0.85$

$P(\text{Humidity} = \text{normal} | \text{play_tennis} = \text{no}) = 1/7 = 0.14$

$P(\text{Wind} = \text{strong} | \text{play_tennis} = \text{yes}) = 3/6 = 0.5$

$P(\text{Wind} = \text{strong} | \text{play_tennis} = \text{no}) = 3/6 = 0.5$

$P(X|C_i) : P(X | \text{play_tennis} = \text{yes}) = 0.6 * 0.67 * 0.85 * 0.5 = 0.1709$

$P(X | \text{play_tennis} = \text{no}) = 0.4 * 0.33 * 0.14 * 0.5 = 0.0092$

$P(X|C_i) * P(C_i) : P(X | \text{play_tennis} = \text{yes}) * P(\text{play_tennis} = \text{yes}) = 0.1709 * 0.6428 = 0.1099$

$P(X | \text{play_tennis} = \text{no}) * P(\text{play_tennis} = \text{no}) = 0.0092 * 0.3571 = 0.0033$

Using naïve bayes the player will play tennis on D15.