

1. Create Decision tree using ID3 algorithm.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

The steps in ID3 algorithm are as follows:

1. Calculate entropy for the dataset.
2. For each attribute/feature.
 - 2.1. Calculate entropy for all its categorical values.
 - 2.2. Calculate information gain for the feature.
3. Find the feature with maximum information gain.
4. Repeat it until we get the desired tree.

Solution: **Complete entropy of dataset is:**

$$H(S) = - p(\text{yes}) * \log_2(p(\text{yes})) - p(\text{no}) * \log_2(p(\text{no}))$$

$$= - (9/14) * \log_2(9/14) - (5/14) * \log_2(5/14)$$

$$= - (-0.41) - (-0.53) = \mathbf{0.94}$$

First Attribute - Outlook

$$I(\text{Outlook}) = p(\text{sunny}) * H(\text{Outlook}=\text{sunny}) + p(\text{rain}) * H(\text{Outlook}=\text{rain}) + p(\text{overcast}) *$$

$$H(\text{Outlook}=\text{overcast})$$

$$= (5/14)*0.971 + (5/14)*0.971 + (4/14)*0$$

$$= 0.693$$

$$\text{Information Gain} = H(S) - I(\text{Outlook})$$

$$= 0.94 - 0.693$$

$$= 0.247$$

Second Attribute - Temperature:-

Categorical values - hot, mild, cool

$$H(\text{Temperature}=\text{hot}) = -(2/4)*\log(2/4)-(2/4)*\log(2/4) = 1$$

$$H(\text{Temperature}=\text{cool}) = -(3/4)*\log(3/4)-(1/4)*\log(1/4) = 0.811$$

$$H(\text{Temperature}=\text{mild}) = -(4/6)*\log(4/6)-(2/6)*\log(2/6) = 0.9179$$

Average Entropy Information for Temperature -

$$I(\text{Temperature}) = p(\text{hot})*H(\text{Temperature}=\text{hot}) + p(\text{mild})*H(\text{Temperature}=\text{mild}) + p(\text{cool})*H(\text{Temperature}=\text{cool})$$

$$= (4/14)*1 + (6/14)*0.9179 + (4/14)*0.811 = 0.9108$$

$$\text{Information Gain} = H(S) - I(\text{Temperature})$$

$$= 0.94 - 0.9108$$

$$= 0.0292$$

Third Attribute - Humidity

Categorical values - high, normal

$$H(\text{Humidity}=\text{high}) = -(3/7)*\log(3/7)-(4/7)*\log(4/7) = 0.983$$

$$H(\text{Humidity}=\text{normal}) = -(6/7)*\log(6/7)-(1/7)*\log(1/7) = 0.591$$

Average Entropy Information for Humidity -

$$I(\text{Humidity}) = p(\text{high})*H(\text{Humidity}=\text{high}) + p(\text{normal})*H(\text{Humidity}=\text{normal})$$

$$= (7/14)*0.983 + (7/14)*0.591$$

$$= 0.787$$

$$\text{Information Gain} = H(S) - I(\text{Humidity})$$

$$= 0.94 - 0.787 = 0.153$$

Fourth Attribute - Wind

Categorical values - weak, strong

$$H(\text{Wind}=\text{weak}) = -(6/8)*\log(6/8)-(2/8)*\log(2/8) = 0.811$$

$$H(\text{Wind}=\text{strong}) = -(3/6)*\log(3/6)-(3/6)*\log(3/6) = 1$$

Average Entropy Information for Wind -

$$I(\text{Wind}) = p(\text{weak})*H(\text{Wind}=\text{weak}) + p(\text{strong})*H(\text{Wind}=\text{strong})$$

$$= (8/14)*0.811 + (6/14)*1$$

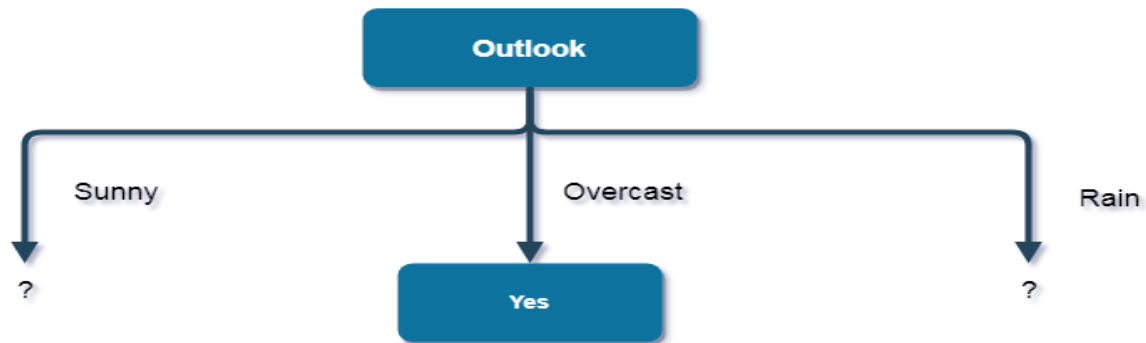
$$= 0.892$$

$$\text{Information Gain} = H(S) - I(\text{Wind})$$

$$= 0.94 - 0.892$$

$$= 0.048$$

Here, the attribute with maximum information gain is Outlook. So, the decision tree built so far -



Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

$E = .97$

Outlook	Temperature	Humidity	Windy	PlayTennis
Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes

$E = 0$

Outlook	Temperature	Humidity	Windy	PlayTennis
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Rainy	Mild	Normal	False	Yes
Rainy	Mild	High	True	No

$E = .97$

Average Entropy = .64 (weighted .69)

Here, when Outlook == overcast, it is of pure class(Yes).

Now, we have to repeat same procedure for the data with rows consist of Outlook value as Sunny and then for Outlook value as Rain.

Categorical values - hot, mild, cool

$$H(\text{Sunny, Temperature=hot}) = -0-(2/2)*\log(2/2) = 0$$

$$H(\text{Sunny, Temperature=cool}) = -(1)*\log(1)- 0 = 0$$

$$H(\text{Sunny, Temperature=mild}) = -(1/2)*\log(1/2)-(1/2)*\log(1/2) = 1$$

Average Entropy Information for Temperature -

$$I(\text{Sunny, Temperature}) = p(\text{Sunny, hot})*H(\text{Sunny, Temperature=hot}) + p(\text{Sunny, mild})*H(\text{Sunny, Temperature=mild}) + p(\text{Sunny, cool})*H(\text{Sunny, Temperature=cool})$$

$$= (2/5)*0 + (1/5)*0 + (2/5)*1$$

$$= 0.4$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny, Temperature})$$

$$= 0.971 - 0.4$$

$$= 0.571$$

Second Attribute - Humidity

Categorical values - high, normal

$$H(\text{Sunny, Humidity=high}) = - 0 - (3/3)*\log(3/3) = 0$$

$$H(\text{Sunny}, \text{Humidity}=\text{normal}) = -(2/2) * \log(2/2) - 0 = 0$$

Average Entropy Information for Humidity -

$$I(\text{Sunny}, \text{Humidity}) = p(\text{Sunny}, \text{high}) * H(\text{Sunny}, \text{Humidity}=\text{high}) + p(\text{Sunny}, \text{normal}) * H(\text{Sunny}, \text{Humidity}=\text{normal})$$

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

$$= 0$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny}, \text{Humidity})$$

$$= 0.971 - 0$$

$$= 0.971$$

Third Attribute - Wind

Categorical values - weak, strong

$$H(\text{Sunny}, \text{Wind}=\text{weak}) = -(1/3) * \log(1/3) - (2/3) * \log(2/3) = 0.918$$

$$H(\text{Sunny}, \text{Wind}=\text{strong}) = -(1/2) * \log(1/2) - (1/2) * \log(1/2) = 1$$

Average Entropy Information for Wind -

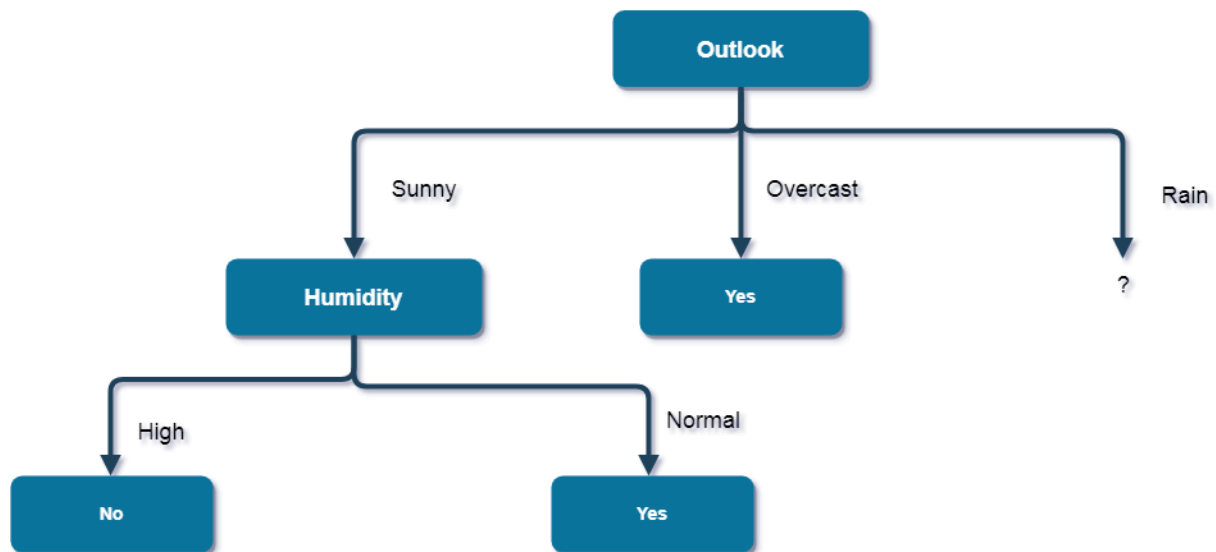
$$I(\text{Sunny, Wind}) = p(\text{Sunny, weak}) * H(\text{Sunny, Wind=weak}) + p(\text{Sunny, strong}) * H(\text{Sunny, Wind=strong})$$

$$= (3/5) * 0.918 + (2/5) * 1$$

$$= 0.9508$$

$$\text{Information Gain} = H(\text{Sunny}) - I(\text{Sunny, Wind}) = 0.971 - 0.9508 = 0.0202$$

Here, the attribute with maximum information gain is Humidity. So, the decision tree built so far -



Here, when Outlook = Sunny and Humidity = High, it is a pure class of category "no". And When Outlook = Sunny and Humidity = Normal, it is again a pure class of category "yes". Therefore, we don't need to do further calculations.

Now, finding the best attribute for splitting the data with Outlook=Sunny values { Dataset rows = [4, 5, 6, 10, 14]}.

Complete entropy of Rain is -

$$H(S) = - p(\text{yes}) * \log_2(p(\text{yes})) - p(\text{no}) * \log_2(p(\text{no}))$$

$$= - (3/5) * \log_2(3/5) - (2/5) * \log_2(2/5)$$

$$= 0.971$$

First Attribute - Temperature

Categorical values - mild, cool

$$H(\text{Rain}, \text{Temperature}=\text{cool}) = -(1/2)*\log_2(1/2) - (1/2)*\log_2(1/2) = 1$$

$$H(\text{Rain}, \text{Temperature}=\text{mild}) = -(2/3)*\log_2(2/3) - (1/3)*\log_2(1/3) = 0.918$$

Average Entropy Information for Temperature -

$$I(\text{Rain}, \text{Temperature}) = p(\text{Rain}, \text{mild}) * H(\text{Rain}, \text{Temperature}=\text{mild}) + p(\text{Rain}, \text{cool}) * H(\text{Rain}, \text{Temperature}=\text{cool})$$

$$= (2/5)*1 + (3/5)*0.918$$

$$= 0.9508$$

$$\text{Information Gain} = H(\text{Rain}) - I(\text{Rain}, \text{Temperature})$$

$$= 0.971 - 0.9508$$

$$= 0.0202$$

Second Attribute - Wind

Categorical values - weak, strong

$$H(\text{Wind}=\text{weak}) = -(3/3)*\log(3/3) - 0 = 0$$

$$H(\text{Wind}=\text{strong}) = 0 - (2/2)*\log(2/2) = 0$$

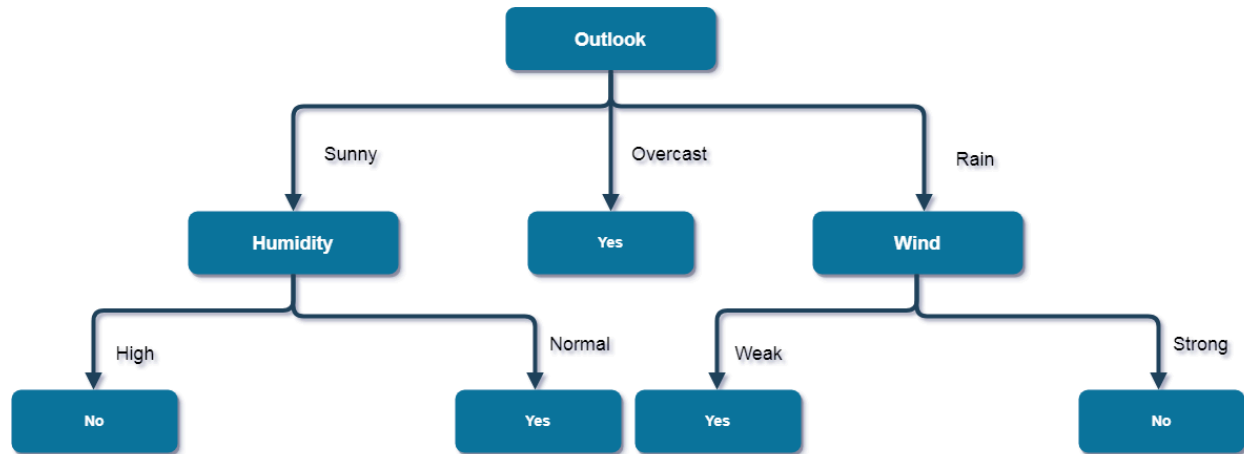
Average Entropy Information for Wind -

$$I(\text{Wind}) = p(\text{Rain}, \text{weak}) * H(\text{Rain}, \text{Wind}=\text{weak}) + p(\text{Rain}, \text{strong}) * H(\text{Rain}, \text{Wind}=\text{strong})$$

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

$$\text{Information Gain} = H(\text{Rain}) - I(\text{Rain}, \text{Wind}) = 0.971 - 0 = 0.971$$

Here, the attribute with maximum information gain is Wind. So, the decision tree built so far -



2. Build a **decision tree** using ID3 algorithm for the given training data in the table (Buy Computer data), and predict the class of the following new example: **age** \leq **30**, **income**=medium, **student**=yes, **credit-rating**=fair

age	income	student	Credit rating	Buys computer
≤ 30	high	no	fair	no
≤ 30	high	no	excellent	no
31...40	high	no	fair	yes
> 40	medium	no	fair	yes
> 40	low	yes	fair	yes
> 40	low	yes	excellent	no
31...40	low	yes	excellent	yes
≤ 30	medium	no	fair	no
≤ 30	low	yes	fair	yes
> 40	medium	yes	fair	yes
≤ 30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
> 40	medium	no	excellent	no