

Decision Tree

A tree is constructed to model the classification process. Once the tree is built, it is applied to each tuple in the database and results in the classification for that tuple.

Attribute Selection measure

It is selecting the best splitting criteria that separates the given data partition ' D ' of class labelled training tuples into individual classes. They determined how the tuples at the given node are to be split. There are 3 popular attribute selection measures -

- 1) Information Gain
- 2) Gain Ratio
- 3) Gini index

1) Information Gain

Let node N represent or hold the tuples of partition ' D '. The information to classify the tuple in ' D ' is given by -

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

where, p_i is the probability that an arbitrary tuple in D belongs to class C_i .

Info(D) is also known as entropy of D

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \cdot \text{Info}(D_j)$$

$\text{Info}_A(D)$ is the expected information required to classify a tuple from D based on the partitioning by A .

→ Information gain is defined as the ~~reqd~~ difference between the original information and the new requirement.

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

The attribute A with the highest information gain is chosen as the splitting attribute at node N .

Database name

All electronics customer database.

| Q10 | Rid | age | income | student | credit rating | class buys computer |
|-----|--------------|--------|--------|-----------|---------------|---------------------|
| 1 | <u>youth</u> | high | no | fair | no | |
| 2 | <u>youth</u> | high | no | excellent | no | |
| 3 | <u>mid</u> | high | no | fair | yes | |
| 4 | senior | medium | no | fair | yes | |
| 5 | senior | low | yes | fair | yes | |
| 6 | senior | low | yes | excellent | no | |
| 7 | <u>mid</u> | low | yes | excellent | yes | |
| 8 | <u>youth</u> | medium | no | fair | no | |
| 9 | <u>youth</u> | low | yes | fair | yes | |
| 10 | senior | medium | yes | fair | yes | |
| 11 | <u>youth</u> | medium | yes | excellent | yes | |
| 12 | <u>mid</u> | medium | no | excellent | yes | |
| 13 | <u>mid</u> | high | yes | fair | yes | |
| 14 | senior | medium | no | excellent | no | |

First determine classes.

class C1 buys-computer = "yes"

class C2 buys-computer = "no"

no. of tuples with class C1 ("yes") = 9

no. of tuples with class C2 ("no") = 5

$$\text{Info}(D) = - \frac{c_1}{c_1+c_2} \log_2 \left(\frac{c_1}{c_1+c_2} \right)$$

$$- \frac{c_2}{c_1+c_2} \log_2 \left(\frac{c_2}{c_1+c_2} \right)$$

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

$$= 0.940$$

1st step

(I)

start with attribute age

Age category youth there are 2 yes tuples and
3 no tuples.

$$I(2,3) = \frac{-2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$
$$= \underline{\underline{0.971}}$$

Age category middle aged there are 4 yes tuples
and 0 no tuples.

$$I(4,0) = \underline{\underline{0}}$$

Age category senior there are 3 yes tuples and
2 no tuples.

$$I(3,2) = \frac{-3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$
$$= \underline{\underline{0.971}}$$

Calculate Entropy

$$E(\text{age}) = \sum_{i=1}^n p_i n_i I(p_i, n_i)$$

$$= \frac{5}{14} \times I(2,3) + \frac{4}{14} \times I(4,0) + \frac{5}{14} \times I(3,2)$$
$$= \underline{\underline{0.694}}$$

Information gain of age is

$$\begin{aligned} \text{Gain}(\text{age}) &= \text{Info}(D) - E(\text{age}) \\ &= 0.940 - 0.694 \\ &= \underline{\underline{0.246}} \end{aligned}$$

start with income attribute

Income attribute there are 4 high, 2 yes tuples and 2 no tuples

$$I(2,2) = -\frac{2}{4} \log\left(\frac{2}{4}\right) - \frac{2}{4} \log\left(\frac{2}{4}\right)$$

$$= -\log_2\left(\frac{2}{4}\right) = \underline{\underline{1}}$$

Income category medium there are 4 yes tuples and 2 no tuples

$$I(4,2) = -\frac{4}{6} \log\left(\frac{4}{6}\right) - \frac{2}{6} \log\left(\frac{2}{6}\right)$$

$$= \underline{\underline{0.918}}$$

Income category low there are 3 yes tuples and 1 no tuple.

$$I(3,1) = -\frac{3}{4} \log\left(\frac{3}{4}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right)$$

$$= \underline{\underline{0.811}}$$

Calculate Entropy

$$E(\text{income}) = \sum_{i=1}^V p_i n_i I(p_i, n_i)$$

$$= \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1)$$

$$= \underline{\underline{0.911}}$$

$$\text{Gain}(\text{income}) = \text{Info}(D) - E(\text{income})$$

$$= 0.940 - 0.911$$

$$= \underline{\underline{0.029}}$$

start with student attribute.

student category yes there are 6 yes tuples
and 1 no tuple

$$I(6,1) = \frac{6}{7} \log_2 \left(\frac{6}{7} \right) + \frac{1}{7} \log_2 \left(\frac{1}{7} \right)$$
$$= \underline{\underline{0.592}}$$

student category no there are 3 yes tuples
and 4 no tuples.

$$I(3,4) = \frac{3}{7} \log_2 \left(\frac{3}{7} \right) + \frac{4}{7} \log_2 \left(\frac{4}{7} \right)$$
$$= \underline{\underline{0.985}}$$

calculate entropy

$$E(\text{student}) = \sum_{j=1}^7 \frac{p_j n_j}{p+n} I(p_j, n_j)$$
$$= \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4)$$
$$= \underline{\underline{0.7885}}$$

$$\text{Gain}(\text{student}) = \text{Info}(D) - E(\text{student})$$
$$= 0.940 - \underline{\underline{0.7885}}$$
$$= \underline{\underline{0.1515}}$$

start with credit rating attribute.

credit rating category fair there are 6 yes
tuples and 2 no tuples.

$$I(6,2) = \frac{6}{8} \log_2 \left(\frac{6}{8} \right) + \frac{2}{8} \log_2 \left(\frac{2}{8} \right)$$
$$= \underline{\underline{0.811}}$$

coedit rating category excellent there are 3 yes tuples and 3 no tuples

$$I(3,3) = \frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \frac{3}{6} \log_2 \left(\frac{3}{6} \right)$$

$$= \underline{\underline{1}}$$

Calculate entropy,

$$E(\text{credit rating}) = \sum_{i=1}^8 \frac{P_i + N_i}{P+N} I(P_i, N_i)$$

$$= \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$

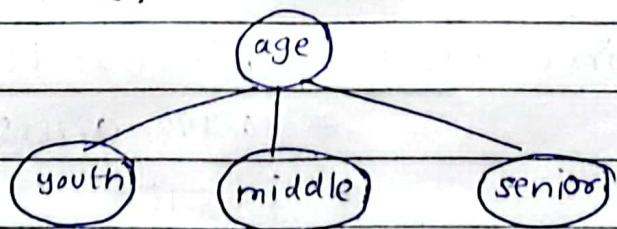
$$= \underline{\underline{0.892}}$$

$$\text{Gain}(\text{credit rating}) = \text{Info}(D) - E(\text{credit rating})$$

$$= 0.940 - 0.892$$

$$= \underline{\underline{0.048}}$$

Since, Age attribute has the highest information gain among the attributes. It is selected as the next attribute.



^{2nd iteration} (II) Now age = youth total tuples = 5

There are 2 yes and 3 no,

$$\text{Info}(D) I(2,3) = 0.971$$

Income category high \rightarrow 0 yes and 2 no

$$I(0,2) = 0$$

Income = medium \rightarrow 1 yes & 1 no

$$I(1,1) = \underline{\underline{1}}$$

Income = low \rightarrow 1 yes & 0 no

$$I(1,0) = \underline{\underline{0}}$$

$$E(\text{income}) = \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0)$$
$$= \underline{\underline{0.4}}$$

$$\text{Gain}(\text{income}) = \text{Info}(D) - E(\text{income}) = \underline{\underline{0.571}}$$

now, student = yes \rightarrow 2 yes and 0 no

$$I(2,0) = \underline{\underline{0}}$$

student = no \rightarrow 0 yes and 3 no

$$I(0,3) = \underline{\underline{0}}$$

$$E(\text{student}) = \underline{\underline{0}}$$

$$\text{Gain}(\text{student}) = \text{Info}(D) - E(\text{student})$$

$$= \underline{\underline{0.971}}$$

now, credit rating = fair \rightarrow 1 yes and 2 no

$$I(1,2) = \frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) = \underline{\underline{0.918}}$$

credit rating = excellent \rightarrow 1 yes and 1 no

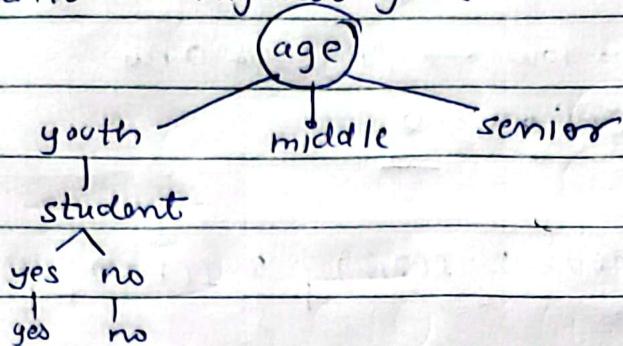
$$I(1,1) = \underline{\underline{1}}$$

$$E(\text{credit rating}) = \frac{2}{5} I(1,1) + \frac{3}{5} I(1,2) = \underline{\underline{0.951}}$$

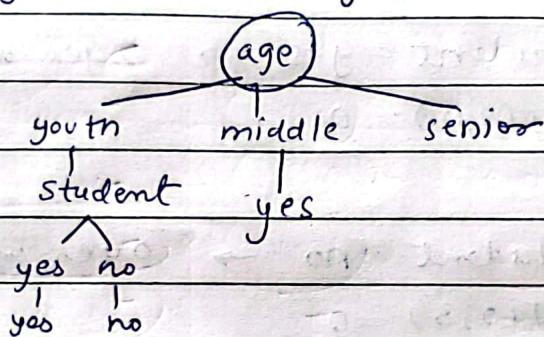
$$\text{Gain (credit rating)} = \text{Info}(D) - E(\text{credit rating})$$

$$= 0.02$$

student has highest gain



~~Iter 3~~ (III) The no. of tuples for category middle age is 4 since the attribute's income & credit rating - The class : buys computer is yes. so you can assign the class yes to middle-aged.



~~Iter 4~~ (IV) senior age category age = senior tuples = 5
3 yes 2 no

$$\text{Info}(D) = I(3,2) = 0.971$$

$$\text{Income} = \text{high} \rightarrow 0 \text{ yes} \& 0 \text{ no } I(0,0) = 0$$

$$\text{Income} = \text{medium} \rightarrow 2 \text{ yes} \& 1 \text{ no } I(2,1) = 0.918$$

$$\text{Income} = \text{low} \rightarrow 1 \text{ yes} \& 1 \text{ no } I(1,1) = 1$$

$$E(\text{income}) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) + \frac{0}{5} I(0,0)$$

$$= 0.951$$

$$\text{Gain (Income)} = \text{Info}(D) - E(\text{income})$$

$$= 0.971 - 0.951 = 0.02$$

credit rating = fair \rightarrow 3 Yes & 0 No

$$I(3,0) = 0$$

credit rating = excellent \rightarrow 0 Yes & 2 No

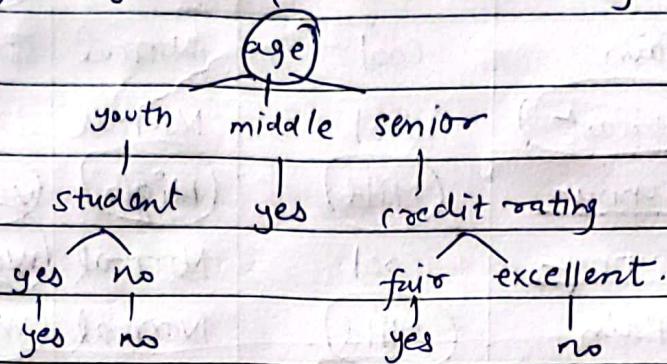
$$I(0,2) = 0$$

$$E(\text{credit rating}) = 0$$

$$\text{Gain}(\text{credit rating}) = \text{Info}(D) - E(\text{credit rating})$$

$$= 0.971 - 0 = 0.971$$

Highest gain is from credit rating



Classification Rules

\rightarrow If age = youth & student = yes then
buys_computer = yes

\rightarrow If age = youth & student = no then
buys_computer = no

\rightarrow If age = middle-aged then buys_computer = yes

\rightarrow If age = senior and credit_rating = fair then
buys_computer = yes

\rightarrow If age = senior and credit_rating = excellent then
buys_computer = no.

Classification

(I) Decision Tree

Q1. Table.

| Day | outlook | Temperature | Humidity | Wind | Class : play Football |
|-----|----------|-------------|----------|--------|-----------------------|
| 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 |

class C1 : plays football = YES

class C2 : plays football = NO.

Iteration 1

Total number of records = 14

Number of Tuples with class C1 = 9

Number of Tuples with class C2 = 5

$$\text{Info}(D) = \frac{c_1}{c_1 + c_2} \log_2 \left(\frac{c_1}{c_1 + c_2} \right)$$

$$- \left(\frac{c_2}{c_1 + c_2} \right) \log_2 \left(\frac{c_2}{c_1 + c_2} \right)$$

$$= 0.940$$

$$= -\left(\frac{9}{14}\right) \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right).$$

$$= 0.940$$

For Outlook Property Attribute.

outlook category sunny there are 2 YES tuples and 3 NO tuples.

$$I(2,3) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= 0.971$$

outlook category overcast there are 4 YES tuples and 0 NO tuples.

$$I(4,0) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) - \frac{0}{4} \log_2\left(\frac{0}{4}\right)$$

$$= 0$$

outlook category Rain there are 3 YES tuples and 2 NO tuples.

$$I(3,2) = -\frac{3}{5} \log_2\left(\frac{2}{5}\right) - \frac{2}{5} \log_2\left(\frac{5}{5}\right)$$

$$= 0.971$$

Calculate entropy

$$E(\text{outlook}) = \sum_{i=1}^v p_i + n_i I(p_i, n_i)$$

$$= \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$= 0.6935.$$

Calculate Info Gain

$$\begin{aligned}\text{Info Gain} &= \text{Info}(D) - E(\text{outlook}) \\ &= 0.940 - 0.6935 \\ &\underline{= 0.2465}\end{aligned}$$

For Temperature Attribute,

Temp category Hot there are 2 YES tuples and
2 NO tuples.

$$\begin{aligned}I(2,2) &= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \\ &= \underline{\underline{0}}\end{aligned}$$

Temp category Mild there are 4 YES tuples and
2 NO tuples.

$$\begin{aligned}I(4,2) &= -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) \\ &= \underline{\underline{0.9182}}\end{aligned}$$

Temp category Cool there are 3 YES tuples and
1 NO tuple

$$\begin{aligned}I(3,1) &= -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \\ &= \underline{\underline{0.8112}}\end{aligned}$$

$$E(\text{temperature}) = \sum_{i=1}^k \frac{P_i + n_i}{P+n} I(P_i, n_i)$$

$$\begin{aligned}&= \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1) \\ &= \underline{\underline{0.911}}\end{aligned}$$

$$\begin{aligned}\text{Info Gain} &= \text{Info}(D) - E(\text{temperature}) \\ &= 0.940 - 0.911 \\ &= \underline{\underline{0.029}}\end{aligned}$$

For Humidity Attribute,

High \rightarrow 3 YES and 4 NO tuples.

$$\begin{aligned}I(3,4) &= -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) \\ &= \underline{\underline{0.9852}}\end{aligned}$$

Normal \rightarrow 6 YES and 1 NO tuples.

$$\begin{aligned}I(6,1) &= -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) \\ &= \underline{\underline{0.5916}}\end{aligned}$$

$$\begin{aligned}E(\text{humidity}) &= \sum_{i=1}^v \frac{P_i + n_i}{P+n} I(P_i, n_i) \\ &= \frac{7}{14} I(3,4) + \frac{7}{14} I(6,1) \\ &= \underline{\underline{0.7884}}\end{aligned}$$

$$\begin{aligned}\text{Info Gain} &= \text{Info}(D) - E(\text{humidity}) \\ &= 0.940 - 0.7884 \\ &= \underline{\underline{0.1516}}\end{aligned}$$

For Wind attribute,

weak \rightarrow 6 YES and 2 NO tuples

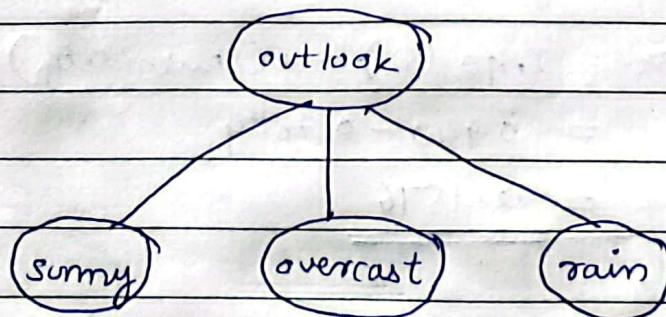
$$I(6,2) = \frac{6}{8} \log_2\left(\frac{6}{8}\right) + \frac{2}{8} \log_2\left(\frac{2}{8}\right)$$
$$= \underline{\underline{0.8112}}$$

strong \rightarrow 3 YES and 3 NO tuples

$$I(3,3) = \frac{3}{6} \log_2\left(\frac{3}{6}\right) + \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$
$$= \underline{\underline{1}}$$

$$E(\text{wind}) = \sum_{j=1}^V \frac{p_i + n_j}{P+N} I(p_i, n_j)$$
$$= \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$
$$= \underline{\underline{0.8921}}$$

$$\text{Info Gain} = \text{Info}(D) - E(\text{wind})$$
$$= 0.940 - 0.8921$$
$$= \underline{\underline{0.0479}}$$



Iteration 2 - outlook = sunny.

For attribute outlook = sunny

Tuples = 5

There are 2 YES and 3 NO

$$I(2,3) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$\text{Info}(D) = \underline{0.9709}$$

Check for attribute temperature under outlook =
sunny.

hot

temperature = high \rightarrow 0 YES and 2 NO.

$$I(0,2) = -\frac{0}{2} \log_2 \left(\frac{0}{2}\right) - \frac{2}{2} \log_2 \left(\frac{2}{2}\right)$$

$$= 0$$

temperature = mild \rightarrow 1 YES and 1 NO

$$I(1,1) = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)$$

$$= 0$$

temperature = cool \rightarrow 1 YES and 0 NO

$$I(1,0) = -1 \log_2 \left(\frac{1}{1}\right) - 0 \log_2 \left(\frac{0}{1}\right)$$

$$= 0$$

$$E(\text{temperature}) = \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0)$$

$$= \underline{0.4}$$

$$\text{Gain}(\text{temperature}) = \text{Info}(D) - E(\text{temp})$$

$$= 0.9709 - 0.4 = \underline{0.5709}$$

check for attribute humidity under outlook = sunny

humidity = high \rightarrow 0 Yes and 3 No
 $I(0,3) = 0$

humidity = normal \rightarrow 2 Yes and 0 No
 $I(2,0) = 0$

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

$$\begin{aligned} \text{Gain(humidity)} &= \text{Info}(D) - E(\text{humidity}) \\ &= 0.971 - 0 = \underline{\underline{0.971}} \end{aligned}$$

check for attribute wind under outlook = sunny

wind = weak \rightarrow 1 Yes and 2 No

$$I(1,2) = \frac{-1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = \underline{\underline{0.918}}$$

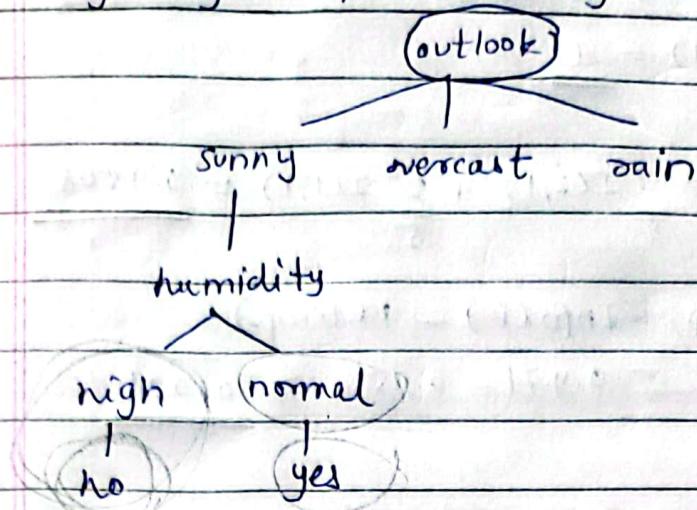
wind = strong \rightarrow 1 Yes and 1 No

$$I(1,1) = \frac{-1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) = \underline{\underline{1}}$$

$$\begin{aligned} E(\text{wind}) &= \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1) \\ &= \underline{\underline{0.9508}} \end{aligned}$$

$$\begin{aligned} \text{Gain(wind)} &= 0.971 - 0.9508 \\ &= \underline{\underline{0.021}} \end{aligned}$$

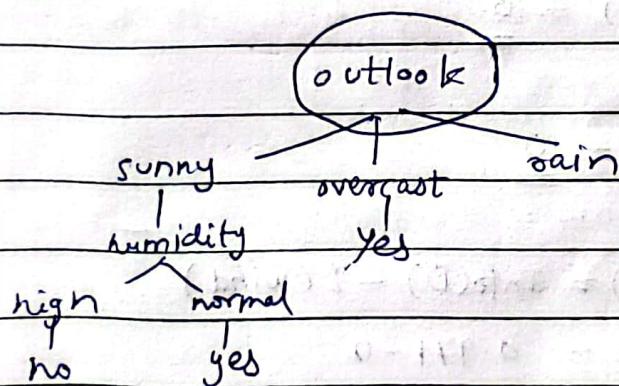
Highest gain is from humidity



Iteration 3 - outlook = overcast

Total tuples = 4

there are 4 Yes & 0 No hence outcome will always be Yes.



3 Yes & 2 No

Iteration 4 - outlook = rain Total tuples = 5

$$Info(D) = I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

0.971

check for attribute temperature under outlook =

temp = hot \rightarrow 0 Yes & 0 No

$$I(0,0) = 0$$

temp = mild \rightarrow 2 Yes & 1 No

$$I(2,1) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.918$$

$\text{temp} = \text{cool} \rightarrow 1 \text{ Yes} \& 1 \text{ No}$

$$I(1,1) = \underline{\underline{1}}$$

$$E(\text{temp}) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = \underline{\underline{0.9508}}$$

$$\begin{aligned}\text{Gain}(\text{temp}) &= \text{Info}(D) - E(\text{temp}) \\ &= 0.971 - 0.9508 = \underline{\underline{0.021}}\end{aligned}$$

check for attribute wind under outlook = rain

wind = weak \rightarrow 3 Yes & 0 No

$$I(3,0) = \underline{\underline{0}}$$

wind = strong \rightarrow 0 Yes & 2 No

$$I(0,2) = \underline{\underline{0}}$$

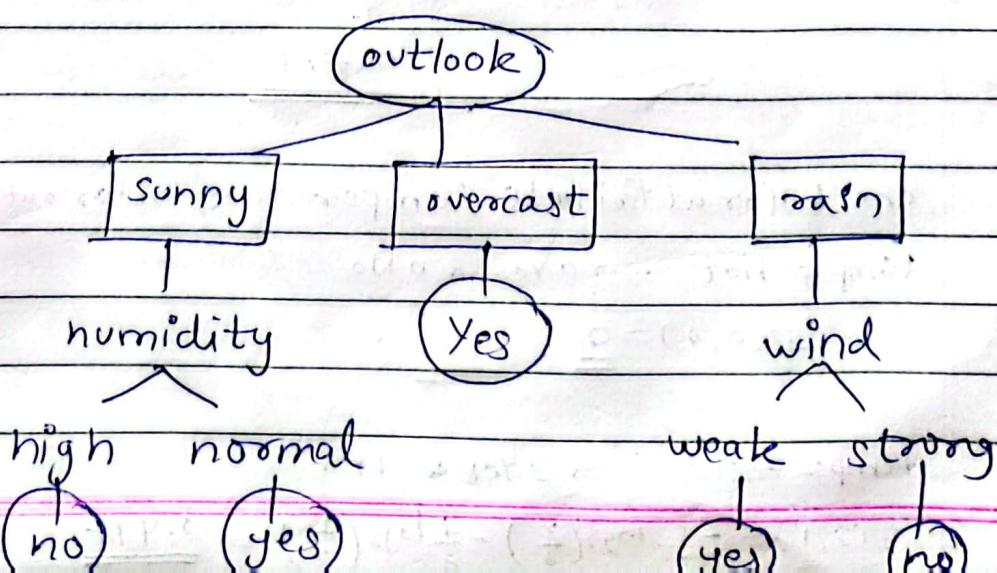
$$E(\text{wind}) = \underline{\underline{0}}$$

Gain(wind) = Info(D) - E(wind)

$$= 0.971 - 0$$

$$= \underline{\underline{0.971}}$$

Since, wind has highest gain consider wind



classification rules

→ If outlook = sunny and humidity = high then
play-football = yes'

→ If outlook = sunny and humidity = normal then
play-football = no

→ If outlook = overcast then play-football = yes

→ If outlook = rain and wind = weak then
play-football = yes

→ If outlook = rain and wind = strong then
play-football = no