# Decision Tree Classification Algorithm:

## i) Information Gain:-

$$I(p,n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

where,

$p$ = no. of elements with class Yes.

$n$ = no. of elements with class No.

## ii) Entropy :-

$$E(A) = \sum_{i=1}^{v} \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

## iii) Gain

$$Gain(A) = I(p,n) - E(A)$$

# example 1:-

Suppose, we want to decide whether the weather is amenable to play "baseball". The target classification is "Should we play baseball?" which can be yes or no.

The weather attributes are outlook, humidity, temperature & wind speed. They have following values.

outlook = {sunny, overcast, rain}

temperature = {hot, mild, cool}

humidity = {high, normal}

wind = {weak, strong}.

| ay | outlook | Temp | humidity | wind | play baseball |
|---|---|---|---|---|---|
| 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 |

Table :- Training Data set.

Solution :-

## Step 1 :- CALCULATION OF INFORMATION GAIN.

class $P = ($play baseball $=$ "yes"$)$

class $N = ($play baseball $=$ "No"$)$

Total no. of records $= 14$

∴ no. of records with "yes" class $= 9$

& no. of records with "No" class $= 5$.

So,

Information Gain is,

$$I(P, N) = -\frac{P}{P+n} \log_2 \frac{P}{P+n} - \frac{n}{P+n} \log_2 \frac{n}{P+n}$$

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

$$= 0.940.$$

2:- COMPUTE ENTROPY ~~Gain~~ for ALL ATTRIBUTES

a) calculate Information gain for outlook.

For outlook = sunny.

$P_i$ = with "yes" class = 2

$n_i$ = with "No" class = 3

So,

$$I(P_i, n_i) = I(2,3) = -\tfrac{2}{5} \log_2(2/5) - (3/5) \log_2(3/5)$$

$$= 0.971.$$

similarly for all values of each outlook

$I(P_i, n_i)$ is calculated as,

| outlook | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---------|-------|-------|---------------|
| sunny | 2 | 3 | 0.971 |
| overcast | 4 | 0 | 0 |
| Rain | 3 | 2 | 0.971 |

So, Entropy is calculated as,

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

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

$$= 0.694$$

$\therefore \text{gain (s, outlook)} = I(P, N) - E(\text{outlook})$
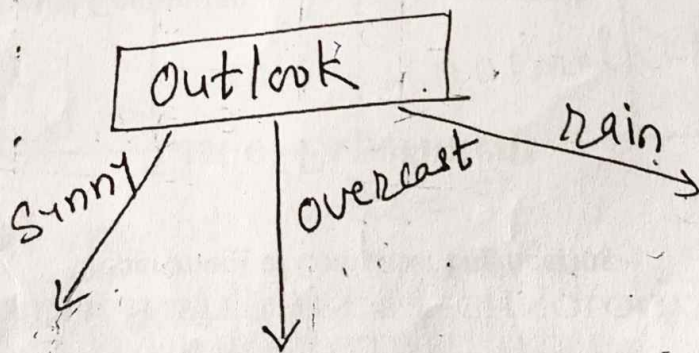
$$= 0.940 - 0.694$$

$$= \boxed{0.246}$$

Similarly,

$\quad \text{gain (s, Temparature)} = 0.029.$

$\quad \text{gain (s, Humidity)} = 0.151$

$\quad \text{gain (s, wind)} = 0.048$

As outlook attribute has highest gain, it is considered as root node with three branches (sunny, overeast, rain).

Ans 3 :-

As attribute outlook is considered as root node, we have to consider the remaining three attribute for sunny branch node.

So, consider outlook = sunny.

$$S_{sunny} = \{D1, D2, D8, D9, D11\}$$
$$= 5.$$

| Day | outlook | Temp | humidity | wind | play baseball |
|---|---|---|---|---|---|
| D1 | Sunny | Hot | high | weak | No |
| D2 | Sunny | Hot | high | strong | No |
| D8 | Sunny | mild | high | week | No |
| D9 | Sunny | cool | normal | weak | yes |
| D11 | Sunny | mild | normal | strong | yes. |

Total no. of tuples = 5

$P$ = no. of tuples with 'yes' = 2

$N$ = no. of tuples with 'No' = 3.

information gain is calculated as,

$$I(P, n) = I(2, 3)$$
$$= -(2/5) \log_2 (2/5) - (3/5) \log_2 (3/5)$$
$$= 0.971$$

- calculate gain for all values of Temparature

| Temparature | pi | ni | $I(P_i, n_i)$ |
|---|---|---|---|
| hot | 0 | 2 | 0 |
| mild | 1 | 1 | 1 |
| cold | 1 | 0 | 0 |

entropy for temparature,

$$E(temp) = \sum_{i=1}^{N} \frac{pi + ni}{P + n} I(P_i, n_i)$$

$$E(temparature) = 2/5 * \dot{I}(0, 2) + (2/5) * (1, 1) + (1/5) * I(1, 0)$$

$$= 0.4$$

∴ Gain (S sunny, temparature) $= I(P, n) - E(temp)$
$$= 0.971 - 0.4$$
$$= 0.571$$

nilarly,

gain $(S_{sunny}, Humidity) = 0.971$

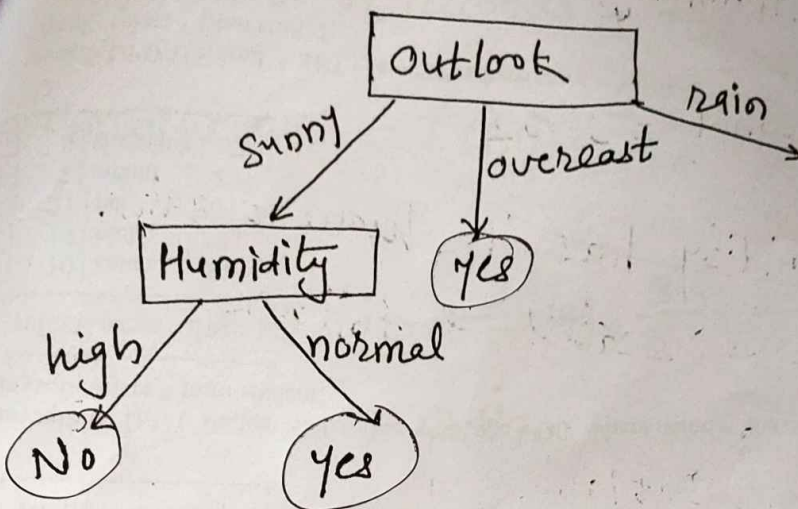gain $(S_{sunny}, Wind) = 0.02$

As Humidity has the highest gain, decision tree is created as follows,



step 4 :- Now consider temparature & wind for outlook = overcast.

| Day | outlook | Temp | humidity | wind | Play baseball |
|---|---|---|---|---|---|
| D3 | overcast | hot | high | weak | yes |
| D7 | overcast, | cool | normal | strong | yes |
| D12 | overcast | mild | high | strong | yes |

Since for attribute temparature & wind
assign class "yes" to overcast.



## step 5 :-

Now, consider temparature & wind for outlook=
rain.

$$Srain = \{D4, D5, D6, D10, D14\}$$
$$= 5$$

| Day | outlook | Temp | humidity | wind | play baseball |
|---|---|---|---|---|---|
| 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 | ~~Cool~~ mild | high | strong | No |

class P : play baseball = "yes" = 3

class N : play baseball = "No" = 2

Total No. of records = 5.

Information gain $= I(P,n) = \dfrac{-P}{P+n} \log_2 \dfrac{P}{P+n} - \dfrac{n}{P+n} \log_2 \dfrac{n}{P+n}$

$\therefore I(3,2) = -\dfrac{3}{5}\log_2 \dfrac{3}{5} - \dfrac{2}{5}\log_2 \dfrac{2}{5}$

$= 0.970$

Info. gain $I(P_i, n_i)$ for wind

| Wind | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|------|-------|-------|---------------|
| weak | 3 | 0 | 0 |
| strong | 0 | 2 | 0 |

Entropy for wind

$E(wind) = \dfrac{3}{5} I(3,0) + \dfrac{2}{5} I(0,2)$

$= 0$

$Gain(S_{gain}, wind) = I(P,n) - E(wind)$

$= 0.970 - 0$

$= 0.970$

info. gain $I(P_i, n_i)$ for temparature

| Temparature | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| Hot | 0 | 0 | 0 |
| mild | 2 | 1 | 0.918 |
| Cool | 1 | 1 | 1 |

Calculate Entropy,

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

$$= 0.951$$

∴ Gain (Srain, temparature) $= I(P,n) - E(\text{temp})$

$$= 0.970 - 0.951$$

$$= 0.019$$

As wind has highest gain, it will be placed below outlook = rain

The decision tree can be represented in rule format as,

If outlook = sunny and humidity = high then play baseball = yes.

If outlook = overcast then play baseball = yes

If outlook = rain and wind = strong then play baseball = no

if outlook = rain and wind = weak then play baseball = yes.

# Accuracy By Class :-

| TP rate | FP rate | Precision | Recall | F-measure | Roc Area | class |
|---|---|---|---|---|---|---|
| $\dfrac{TN}{N}$ | $\dfrac{FP}{P}$ | $\dfrac{TN}{TN+FN}$ | $\dfrac{TN}{N}$ | $\dfrac{(2*Precision*Recall)}{(Precision+recall)}$ | $\dfrac{(TP/P)}{(FP/FP+TN)}$ | No |
| $\dfrac{TP}{P}$ | $\dfrac{FP}{N}$ | $\dfrac{TP}{TP+FP}$ | $\dfrac{TP}{P}$ | $\dfrac{(2*Precision*recall)}{(Precision+recall)}$ | $\dfrac{(TP/P)}{(FP/FP+TN)}$ | Yes |

**Note :-**

Calculate everything in calculation and also solve the eg. using ID3 & match tree with your o/p: