**Decision tree** :- A classifier (Tree Structured), Decision Node ① (Test) /feature /Attribute

① ✓ Leaf node - (classification /value)
✓ edge — value
→ Also used in Regression

②
Ⓓ
(Employed)

Ⓓ₂ NO / Yes Ⓓ₁

(Credit Score)   (income)

H / L      H / L

Ⓐ    Ⓡ    Ⓐ         Ⓡ

③ Training Data — Algo → (m/c)

generate

tuble (Test)

→ Loan or not prediction

→ Test is performed on the feature / Attribute ) Train. DS

Ⓓ ④

| Employed | Credit Score | income | Status — Target Leaf node |
|---|---|---|---|
| Y | H | H | A |
| Y | H | H ) D₁ | A |
| Y | H | H | A |
| H | L | L ) D₂ | R |
| H | L | L | R |

funter do the splitting / leaf node having class or value.

Qustion:1 - for the following Medical diagnostic data ..... DT.

② 

$$T \cdot G \underline{\quad} - \frac{P}{P+N}$$

## Information Gain :- Measure of how much information the Answer to a specific Question provides.

## Entropy :

measure of how much uncertainty in the dataset / information Gain / Information.

Info gain ↑ = entropy ↓

---

Information Gain :-

$$I(P,n) = -\frac{P}{S} \log_2 \frac{P}{S} - \frac{n}{S} \log_2 \frac{n}{S} \qquad \text{---(I)}$$

$S \rightarrow$ total Sample space

$S = P+N$

$$E(A) = \sum_{i=1}^{V} \frac{P_i + n_i}{P+n} (I(P_i, n_i) \qquad \text{---(II)} \cdot$$

$$Gain(A) = I(P,n) - E(A) \qquad \text{---(III)}$$

---

$$\log_2 x = \frac{\log_{10} x}{\log_{10} 2}$$

Question :1 :- for the following Medical diagnosis data, cocute DT- ③

| | Sore Throat | fever | Swallen Glands | Congestion | Headache | Diagnosis |
|---|---|---|---|---|---|---|
| 1 | Yes | Yes | Yes | Yes | Yes | Spep throat |
| 2 | No | No | No | Yes | Yes | Allergy |
| 3 | Yes | Yes | No | Yes | No | Cold |
| 4 | Yes | No | Yes | No | No | S.T |
| 5 | No | Yes | No | Yes | No | Cold |
| 6 | No | No | No | Yes | No | Allergy |
| 7 | No | No | Yes | No | No | S.T |
| 8 | Yes | No | No | Yes | Yes | Allergy |
| 9 | No | Yes | No | Yes | Yes | Cold |
| 10 | Yes | Yes | No | Yes | Yes | cold |

Sample space : -  S.T + All + cold = 10
                   ↓       ↓       ↓
                   3       3       4

$I(P, n)$

$$I(ST, All, cold) = -\left[\left(\frac{3}{10}\right)\log_2\left(\frac{3}{10}\right) + \frac{3}{10}\log_2\left(\frac{3}{10}\right) + \left(\frac{4}{10}\right)\log_2\left(\frac{4}{10}\right)\right]$$

$$= \underline{1.562}$$

finding
Splitting Attribute! . [ Select Attribute with Highest gain ]

Sore Throat : ―

| | ST | A | C | |
|---|---|---|---|---|
| Yes | 2 | 1 | 2 | (Infogain) × P ⎤ + E(A) |
| No | 1 | 2 | 2 | (Info gain) × P ⎦ |

/ 7        / 4

E ( Sore Throat )

$$I(Yes) = -\left[\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \left(\frac{1}{5}\right) \log_2\left(\frac{1}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right]$$

I (Yes) = 1·52

I (No) = 1·52

$$E ( Sore throat ) = \frac{5}{10} \times 1·52 + \frac{5}{10} \times 1·52 = 1·52$$

Gain ( Sore throat ) = 1·562 - 1·52 = 0·05

---

Sore Throat — 0·05
fever → ⟨0·72⟩

Swollen glands → 0·88 ].

Congestion — 0·45

Headache → 0·05

Decision Tree

Create a Decision tree for given following Data:- ⑤

| Day | Outlook | Tem | Humidity | Wind | Play |
|-----|---------|-----|----------|------|------|
| 1 | Sunny | Hot | High | weak | No |
| 2 | Sunny | " | " | Strong | no |
| 3 | Overcast | " | " | weak | Yes |
| 4 | Rain | mild | " | " | " |
| 5 | Rain | cold | Normal | " | " |
| 6 | Rain | " | " | Strong | No |
| 7 | overcast | " | " | " | Yes |
| 8 | Sunny | mild | high | weak | No |
| 9 | " | cold | Normal | " | Yes |
| 10 | Rain | mild | " | " | Yes |
| 11 | Sunny | " | " | Strong | yes |
| 12 | overcast | " | High | " | " |
| 13 | " | Hot | Normal | weak | " |
| 14 | Rain | mild | High | Strong | No |

CART $\rightarrow$ , Gini(S) $= 1 -$

entropy
$\downarrow$
Gini

$$Gini(E) = 1 - \sum_{J=1}^{C} P_j^2$$

$$= 1 -$$

entropy / IG

ID3 — Iterative Dichotomiser 3

CART — Classification & Regression tree.

Gini Index

CART ( Classification & Regression trees) — <u>Decision tree</u>
<u>Gini Index</u> :- Probability of each class :-

Sum of squared probability of each class. we can formulate it as below :-

$$Gini = 1 - \sum (P_i)^2 \text{ for } i=1 \text{ to number of classes.}$$

<u>outlook</u>: outlook is nominal feature

Sunny ← ↓ → rain
Overcast

| Outlook | Yes | No | no. of instances |
|---------|-----|-----|------------------|
| Sunny | 2 | 3 | 5 |
| Overcast | 4 | 0 | 4 |
| Rain | 3 | 2 | 5 |

$$Gini (Outlook = Sunny) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2$$
$$= 1 - 0.16 - 0.36 = 0.48$$

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

$$Gini (Outlook = Rain) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$
$$= 1 - 0.36 - 0.16 = 0.48$$

$$Gini (Outlook) = \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48$$

Gini(outlook) = 0.342

Gini(temp) = 0.439

Gini(Humidity) = 0.367

Gini(wind) = 0.428

✓

lowest cost

8



Outlook

Sunny

Overcast

Yes

Rain

| 1 | Sunny | Hot | High | No | weak |
|---|---|---|---|---|---|
| 2 | S | Hot | High | No | Strong |
| 8 | S | mild | High | No | weak |
| 9 | S | Cold | Normal | Yes | weak |
| 11 | S | mild | Normal | Yes | Strong |

| 4 | Rain | mild | High | weak | Y |
|---|---|---|---|---|---|
| 5 | " | cold | Nor | " | ✗ |
| 6 | " | " | " | Strong | N |
| 10 | " | mild | " | weak | ✓ |
| 14 | " | mild | High | Strong | N |

Humidity

High → No

Normal → Yes

Wind

Strong → No

weak → Yes