**What is decision tree induction**

Ans:

A decision tree is a structure that includes a root node, branches & leaf nodes. It's a top down recursive divide & conquer algorithm.

Decision Tree to

→ Each internal node denotes a test on an attribute.

→ Each branch denotes the outcome of a test.

→ Each leaf node holds a class label.

→ The Topmost node in the tree is the root node.

**Define classification & it's step**

Ans:
Classification is a data mining function. It describes & distinguishes data classes.

It's a two step Process.

Step-1: Learning step → classification model, classifier is constructed. Using training dataset.

Step-2: Classification step → Classification model is used to Predict class lables for given data.

Decision Tree Algorithm এর ব্যবহৃত Common Symbol

| Symbol | Term |
|---|---|
| D | Training Data |
| $X_i$ | A data instance |
| X | A subset of instances |
| $A_j$ | A feature |
| $a_j^i$ | A feature's value |
| $C_i$ | A class label |
| DT | A decision tree |

⊞ $D = \{x_1, x_2, \ldots, x_n\}$ , D এর training data.

⊞ $X_i$ কে vector দিয়ে represent করা হয়, $x_i = \{x_{i1}, x_{i2}, \ldots x_{in}\}$ :

⊞ $D \Rightarrow \{A_1, A_2, A_3, \ldots A_n\}$ ; dataset ;D তে —
An সংখ্যক attribute থাকতে পারে,

⊞ প্রতিটি attribute -এর different ; attribute value
- থাকে, $\{A_{i1}, A_{ia}, A_{ik}\}$

⊞ Dataset এ class এর ১টা Column থাকে, $C = \{C_1, C_2, C_m\}$

⊞ প্রতিটি $x_i$ instance একটা Predefined class label
$C_i$ এর সাথে থাকে,

---

⊞ Decision Tree =এর প্রতিটি leaf decision দেয়,
⊞ Decision Tree এর little Prior knowledge
- দরকার হয়,
⊞ Decision Tree দ্রুত Class Predict করে.

Decision ☒ Write the advantages of decision tree.

__Ans:__

The decisions of

The advantages are

a. Simple to understand

b. Easy to implement

c. Requiring a little Prior knowledge.

d. Can handle numerical & categorical data.

e. ☒ Robust.

f. Dealing with large & noisy datasets.

☒ Decision Tree এর কয়েকটি Version আছে।

| Decision Tree Algorithm |
| --- |

Iterative Dichotomiser (ID3)

C 4.5 Decision Induction

C5 Commercial Version

**⊞ Write the working Process of Iterative Dicholomisens (ID3)**

Ans:

ID3 iteratively Partition the data into smaller Subsets until all the Subsets belong to a single class.

_ID3 এর সূত্র_

a. $Info(D) = -\sum_{i=1}^{n} P_i \log_2 P_i$  // Prior Probability Calculate করি

⤷ i class value / attribute value represent করে

⤷ আর D আলে Total data / Training data

b. $Info_A(D) = \sum_{j=1}^{n} \frac{|D_j|}{|D|} \times Info(D)$

c. $Gain(A) = Info(D) - Info_A(D)$

⇒ Info(D) এর - অন্য আরেকটা version হল Entropy. যেটার formula ; $Entropy = -\sum P(x) \log P(x)$,

**⊞ Define Entropy**

Ans. It is the measure of impurity, disorder, uncertainty in a bunch of examples.

**⊞ What is information gain.**

Ans: It measures how much information a feature gives us about a class.

প্রশ্নঃ Why Information Gain matters.

Ans:

It is the main key that is used by Decision Tree algorithms to construct a decision tree. It always tries to maximize information gain. An attribute with highest information gain splited first.

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

প্র Info gain dataset কে more appropriately ভাগ করে,

প্র যেই ID3 চার্টে X আসবে। ECE তে যা আমাদের যেই instance একটা Particular class এ belong করবে,

$$Info_A(D) = \sum \frac{|D_j|}{|D|} \times Info(D_j)$$

$$Gain(A) = Info(D) - Info_A(D)$$

Table: The playing tennis dataset

| Day | Outlook | Temperature | Humidity | Wind | Play |
|-----|---------|-------------|----------|------|------|
| $D_1$ | Sunny | Hot | High | Weak | No |
| $D_2$ | Sunny | Hot | High | Strong | No |
| $D_3$ | Overcast | Hot | High | Weak | Yes |
| $D_4$ | Rain | Mild | High | Weak | Yes |
| $D_5$ | Rain | Cool | Normal | Weak | Yes |
| $D_6$ | Rain | Cool | Normal | Strong | No |
| $D_7$ | Overcast | Cool | Normal | Strong | Yes |
| $D_8$ | Sunny | Mild | High | Weak | No |
| $D_9$ | Sunny | Cool | Normal | Weak | Yes |
| $D_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $D_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $D_{12}$ | Overcast | Mild | High | Strong | Yes |
| $D_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $D_{14}$ | Rain | Mild | High | Strong | No |

**Find the root using ID3 algorithm from the given dataset**

$Info\ (D) = -\dfrac{9}{14}\ \log_2\ (9/14) - 5/14\ \log_2\ (5/14)$

$= -\ 0/64 \times (-0.191) - 0.35 \times (-0.45)$

$\approx 0.122 + 0.15$

$Info\ (D) = -\dfrac{9}{14}\ \log_2\ (9/14) - \dfrac{5}{14}\ \log_2\ (5/14)$

$= -\ 0.64 \times (-0.64) - 0.35 \times (-1.51)$

$=\ 0.40 + 0.52$

$=\ 0.92$

Attribute → Outlook

$\underset{Outlook}{Info}\ (D) = \dfrac{5}{14} \times \left(-\dfrac{2}{5}\ \log_2\ \dfrac{2}{5}\ -\ \dfrac{3}{5}\ \log_2 \dfrac{3}{5}\right)\ +$

Sunny

$\dfrac{4}{14} \times \left(-\dfrac{4}{4}\ \log_2\ \dfrac{4}{4}\ -\ \dfrac{0}{4}\ \log_2 \dfrac{0}{4}\right)\ +$

Overcast

$\dfrac{5}{14} \times \left(-\dfrac{3}{5}\ \log\ \dfrac{3}{5}\ -\ \dfrac{2}{5}\ \log\ \dfrac{2}{5}\right)$

Rain

$= 0.35 \times \left\{-0.4 \times (-1.32) - 0.6 \times (-0.73)\right\} +$

$0.28 \times \left\{-1 \times (0) - 0 \times \log_2 (0/4)\right\} +$

$0.35 \times \left\{-0.6 \times (-0.73) - 0.4 \times (-1.32)\right\}$

$= 0.35 \times 0.966 + 0.28 \times 0 + 0.35 \times 0.966$

$= 0.3381 + 0.3381 = 0.6762$

Gain (Outlook) $= Info(D) - Info_{Outlook}(D)$

$= 0.92 - 0.6762$

$= 0.2438$

Attribute → Temperature

$Info_{Temperature}(D) = \frac{4}{14} \times \left(-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}\right) +$

<u>Hot</u>

$\frac{6}{14} \times \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6}\right)$

Mild

$+$

$\frac{4}{14} \times \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}\right)$

Cool

$+$

$= 0.28 \times \{-0.5 \times (-1) - 0.5(-1)\} + 0.42 \times \{-0.66 \times (-0.599)$

$- 0.93 \times (-1.59)\} + 0.28 \times \{-0.75 \times (-0.41) - 0.25 \times (-2)$

$= 0.28 (0.5$

$= 0.28 \times 1 + 0.42 \times 0.92004 + 0.28 \times 0.8075$

$= 0.28 + 0.386568 + 0.2261$

$= 0.892$

Gain (Temperature) $= Info(D) - Info_{Temperature}(D) = 0.92 - 0.892$

$= 0.028$

Attribute → Humidity

$$Info_{Humidity}^{(D)} = \underline{\frac{7}{14} \times \left(- \frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log \frac{4}{7}\right)}_{\text{High}} +$$

$$\underline{\frac{7}{14} \times \left(- \frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7}\right)}_{\text{Normal}}$$

$$= 0.5 \times \{-0.42 \times (-1.25) - 0.57 \times (-0.81)\} +$$

$$0.5 \times \{-0.85 \times (-0.23) - 0.14 \times (-2.83)\}$$

$$= 0.5 \times 0.98 + 0.5 \times 0.5917$$

$$= 0.49 + 0.29585$$

$$= 0.78$$

Gain (Humidity) $= Info(D) - Info_{Humidity}^{(D)}$

$$= 0.92 - 0.785$$

$$= 0.135$$

## Attribute ⇒ Wind

$$Info_{Wind}^{(D)} = \frac{8}{14}\left(-\frac{6}{8}\log_2\frac{6}{8} - \frac{2}{8}\log_2\frac{2}{8}\right) +$$

<div align="center">Weak</div>

$$\frac{6}{14}\left(-\frac{3}{6}\log_2\frac{3}{6} - \frac{3}{6}\log_2\frac{3}{6}\right) \times \frac{5}{14}$$

<div align="center">Strong</div>

$$= {}^+0.57 \times \{-0.75 \times (-0.41) - 0.25 \times (-2)\} +$$
$$\quad 0.42 \times \{-0.5 \times (-1) - 0.5 \times (-1)\}$$

$$= 0.57 \times 0.8075 + 0.42 \times 1$$

$$= 0.460 + 0.42$$

$$= 0.88$$

$$Gain(Wind) = Info(D) - Info_{Wind}(D)$$
$$= 0.92 - 0.88$$
$$= 0.04$$

So, Gain of
Outlook = 0.2438
Temparature = 0.028
Humidity = 0.135
Wind = 0.04

→ Outlook এর gain বেশি, তাই সেটা হবে decision tree এর root.

Ans-