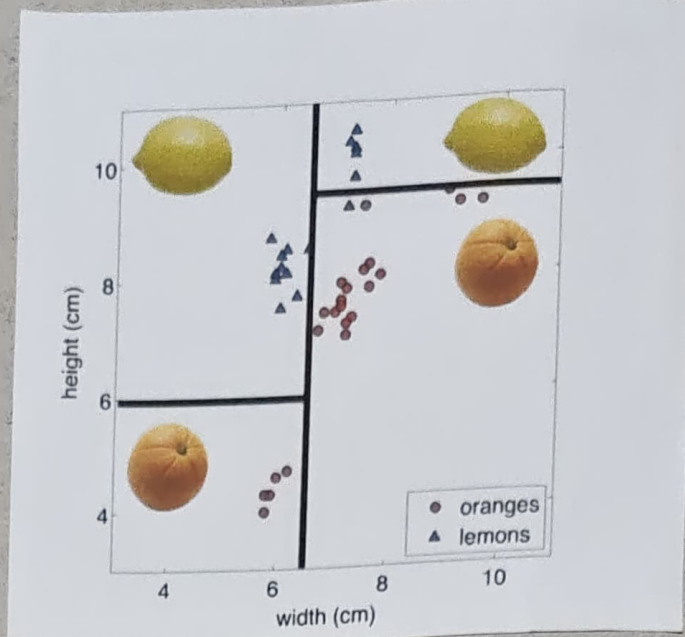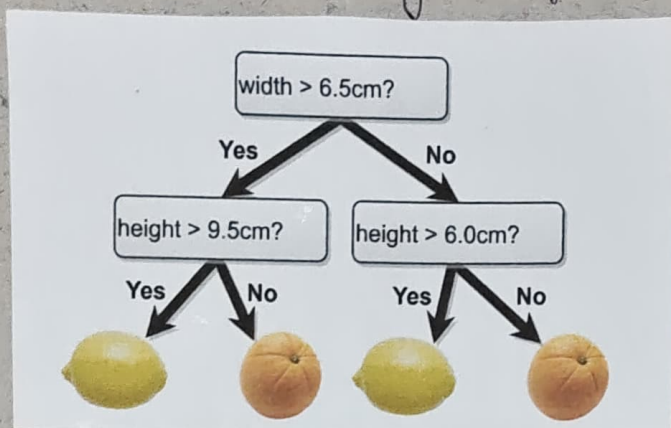# 9. Decision Tree Classifier (DTC)

Consider the simple dataset shown, comprises of __heights__ and __weights__ of lemons and Oranges.



Below DTC is built using the given dataset.



## Advantages
- interpretable, intuitive
- Popular in medical diagnosis application

How to choose the __best attribute__?

One can use principles of __information theory__ for this.

# ENTROPY.

Consider a source $Y$ with symbols $y_i$ and probabilities $p(y_i)$. The Entropy $H(Y)$ of this source is defined as

$$H(Y) = \sum_{i=1}^{n} p(y_i) \log_2 \frac{1}{p(y_i)}$$

Example:

$$p(0) = \frac{1}{4}, \quad p(1) = \frac{3}{4}, \quad H(Y) = ?$$

$$H(Y) = \frac{1}{4} \cdot \log_2 4 + \frac{3}{4} \log_2 \frac{4}{3}$$

$$= \frac{1}{2} + \frac{3}{4} \left( \frac{\ln\left(\frac{4}{3}\right)}{\ln(2)} \right)$$

$$= 0.811 \text{ bits/symbol}$$

## DTC example

| Example | Input Attributes | | | | | | | | | | Goal |
|---------|-----|-----|-----|-----|-----|-------|------|-----|--------|-------|----------|
| | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | WillWait |
| $x_1$ | Yes | No | No | Yes | Some | \$\$\$ | No | Yes | French | 0–10 | $y_1$ = Yes |
| $x_2$ | Yes | No | No | Yes | Full | \$ | No | No | Thai | 30–60 | $y_2$ = No |
| $x_3$ | No | Yes | No | No | Some | \$ | No | No | Burger | 0–10 | $y_3$ = Yes |
| $x_4$ | Yes | No | Yes | Yes | Full | \$ | Yes | No | Thai | 10–30 | $y_4$ = Yes |
| $x_5$ | Yes | No | Yes | No | Full | \$\$\$ | No | Yes | French | >60 | $y_5$ = No |
| $x_6$ | No | Yes | No | Yes | Some | \$\$ | Yes | Yes | Italian | 0–10 | $y_6$ = Yes |
| $x_7$ | No | Yes | No | No | None | \$ | Yes | No | Burger | 0–10 | $y_7$ = No |
| $x_8$ | No | No | No | Yes | Some | \$\$ | Yes | Yes | Thai | 0–10 | $y_8$ = Yes |
| $x_9$ | No | Yes | Yes | No | Full | \$ | Yes | No | Burger | >60 | $y_9$ = No |
| $x_{10}$ | Yes | Yes | Yes | Yes | Full | \$\$\$ | No | Yes | Italian | 10–30 | $y_{10}$ = No |
| $x_{11}$ | No | No | No | No | None | \$ | No | No | Thai | 0–10 | $y_{11}$ = No |
| $x_{12}$ | Yes | Yes | Yes | Yes | Full | \$ | No | No | Burger | 30–60 | $y_{12}$ = Yes |

(Annotations pointing to column headers: Fri/sat → Fri, Hungry → Hun, Patrons → Pat, Reservation → Res, Estimated wait time → Est, Decision → WillWait)

Consider the table shown above.

customer decision to wait or not at restaurants

## Table columns:

| | |
|---|---|
| 1. | Alternate: whether there is a suitable alternative restaurant nearby. |
| 2. | Bar: whether the restaurant has a comfortable bar area to wait in. |
| 3. | Fri/Sat: true on Fridays and Saturdays. |
| 4. | Hungry: whether we are hungry. |
| 5. | Patrons: how many people are in the restaurant (values are None, Some, and Full). |
| 6. | Price: the restaurant's price range ($, $$, $$$). |
| 7. | Raining: whether it is raining outside. |
| 8. | Reservation: whether we made a reservation. |
| 9. | Type: the kind of restaurant (French, Italian, Thai or Burger). |
| 10. | WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60). |

## Problem 1 :



What is the entropy of the decision $Y \in \{Yes, No\}$ ?

$$P(Y = Yes) = \frac{1}{2} \quad , \quad P(Y = No) = \frac{1}{2}$$

$$H(Y) = H\left(\frac{1}{2}, \frac{1}{2}\right)$$

$$= \frac{1}{2} \log_2 2 + \frac{1}{2} \log_2 2$$

$$= 1 \text{ bits/symbol} .$$

## CONDITIONAL ENTROPY

Consider two sources : $Y$ with symbols $y_i$
$\qquad\qquad\qquad\qquad X$ with symbols $x_j$

The conditional Entropy $H(Y|X)$ is defined as

$$\sum_{j=1}^{m} p(x_j) \, H\left(Y \mid x = x_j\right)$$

Example: Cricket

WT — Winning Toss | WG — Winning Game
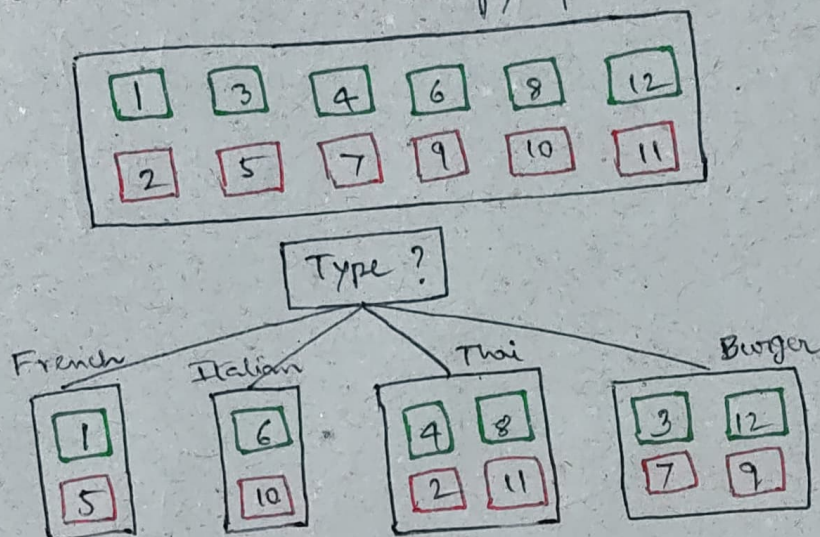LT — Losing Toss | LG — Losing Game.

$$WT = \frac{1}{4} \left\{ WG = \frac{5}{6}, \ LG = \frac{1}{6} \right\}$$

$$LT = \frac{3}{4} \left\{ WG = \frac{1}{5}, \ LG = \frac{4}{5} \right\}$$

$$\Rightarrow \frac{1}{4} \times \left\{ \frac{5}{6} \log_2 \frac{6}{5} + \frac{1}{6} \log_2 6 \right\}$$

$$+ \frac{3}{4} \times \left\{ \frac{1}{5} \log_2 5 + \frac{4}{5} \log_2 \frac{5}{4} \right\}$$

## Problem 2 :

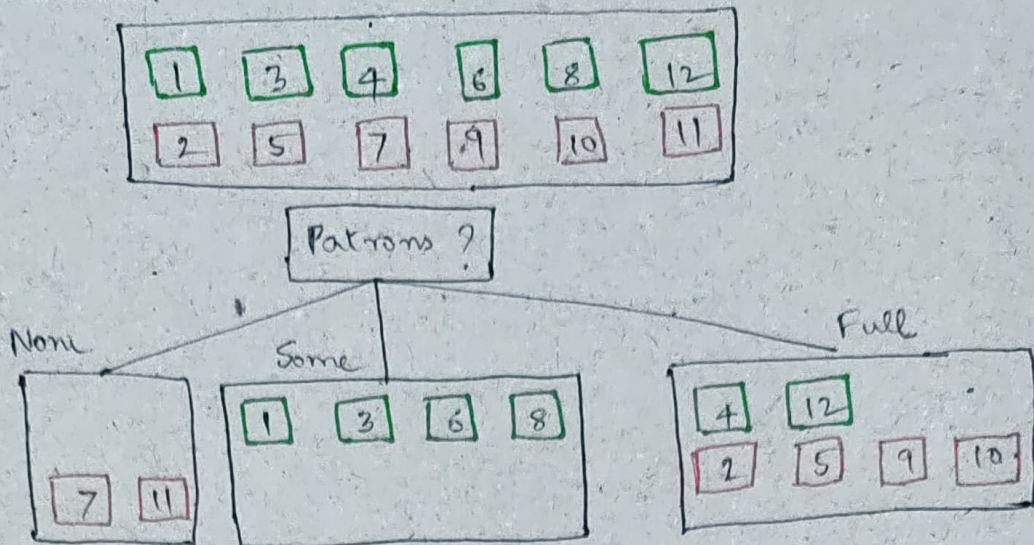What is the conditional entropy of the TYPE feature?



$$P(Fr) \times H(Y|Fr) + P(It) \times H(Y|It)$$
$$+ P(Th) \times H(Y|Th)$$
$$+ P(Bu) \times H(Y|Bu)$$

$$= \left( \frac{2}{12} \times 1 \right) + \left( \frac{2}{12} \times 1 \right) + \left( \frac{4}{12} \times 1 \right) + \left( \frac{4}{12} \times 1 \right)$$

$$= 1.$$

# Problem 3 :



Patrons ?

None — Some — Full

What is the conditional entropy of the PATRONS feature?

$$P(None) * H(Y \mid None) + P(Some) * H(Y \mid Some)$$
$$+ P(Full) * H(Y \mid Full)$$

$$= \left(\frac{2}{12} \times 0\right) + \left(\frac{4}{12} \times 0\right) + \left(\frac{1}{2} \times H\left(\frac{1}{3}, \frac{2}{3}\right)\right)$$

$$= 0.459.$$

## INFORMATION GAIN

The Information Gain (IG) is defined as

$$IG(x) = H(Y) - H(Y \mid x).$$

Choose the feature that maximizes the IG!

· Example :
- IG for the TYPE feature is given as

$$IG(TYPE) = H(Y) - H(Y \mid x = TYPE)$$
$$= 1 - 1 = 0$$

- IG for the PATRONS feature is given as

$$IG(PATRONS) = H(Y) - H(Y \mid x = PATRONS)$$
$$= 1 - 0.459 = 0.541.$$

# DTC example.

Since $IG(Patron) = 0.541 > IG(Type) = 0$

Therefore, we choose PATRONS as the feature to split.

## Final DTC

```
                    Patrons ?
           None      Some        Full
          [NO]      [Yes]     Hungry ?
                          NO           Yes
                        [NO]         Type ?
                  French    Italian    Thai      Burger
                  [Yes]     [NO]    Fri/Sat ??    [Yes]
                                    NO      Yes
                                  [NO]     [Yes]
```