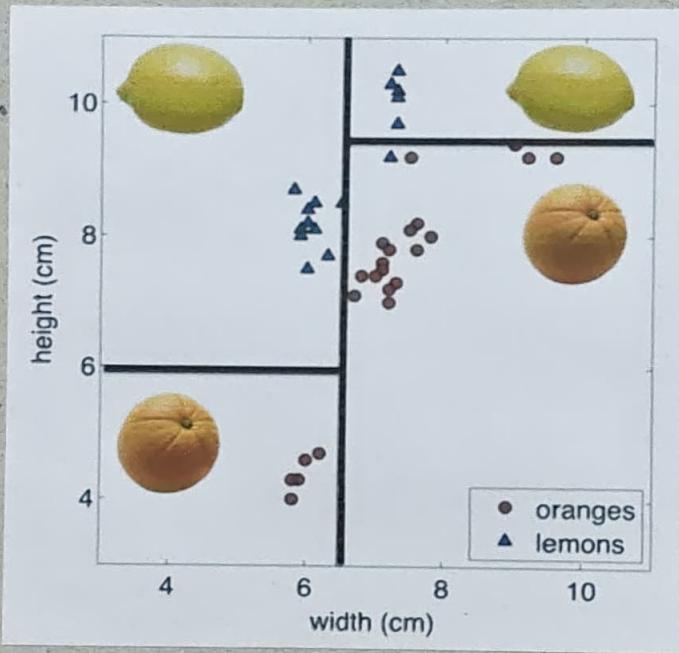


10. Decision Tree Classification (DTC)

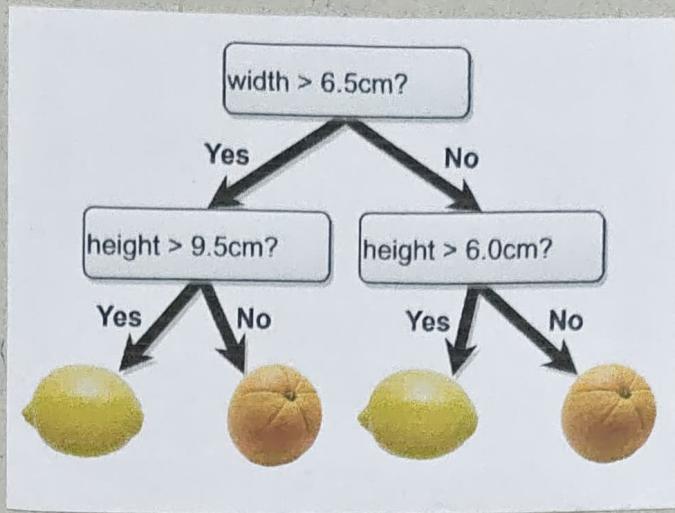
To understand DTC, let us consider the simple dataset, which comprises of Heights and Weights of Lemons and Oranges.



Now, Given a new object, based on their attributes/features, can we characterize / build Machine Learning algorithm which detects whether the object is Lemon or Orange?

Yes !

One can build Decision Tree Classifier (DTC) as below using the given dataset.



Advantages of DTC

- ① Interpretable, Intuitive (in contrast to Neural Nets)
- ② Popular in Medical diagnosis applications.

DTC's are well suited to model discrete outcomes. (ii)
Discrete decisions at every stage. To learn / build a DTC, one can use principles of Information theory to choose the best attribute.

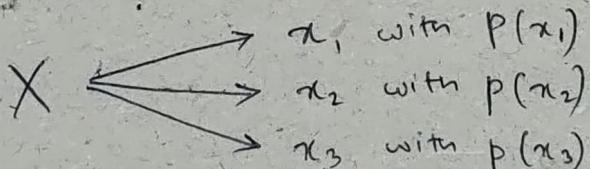
(ii) Which attribute to split?

Which attribute gives more information?

ENTROPY

Entropy is the intrinsic information present in a source.

Consider a source X with symbols x_i and probabilities $p(x_i)$.



The entropy $H(x)$ of this source is defined as

$$H(x) = p(x_1) \log_2 \left(\frac{1}{p(x_1)} \right) + p(x_2) \log_2 \left(\frac{1}{p(x_2)} \right) + \dots \\ = \sum_{i=1}^n p(x_i) \log_2 \left(\frac{1}{p(x_i)} \right)$$

$$\boxed{H(x) = - \sum_{i=1}^n p(x_i) \log_2 (p(x_i))}$$

This is the Entropy (or) Information content.

① Entropy is always positive.

(i) $H(x) \geq 0$.

② Entropy is maximum when all symbols are equiprobable.

$$(i) p(x_i) = \frac{1}{n}$$

③ Entropy is zero, only if one of the symbols have probability 1 and rest of the symbols have probability 0.

(i) Entropy = 0 if $p(x_i) = 1$

$$\text{and } p(x_j) = 0, j \neq i$$

Remember, Sum of all ten probabilities has to be 1.

$$(i) \sum_{i=1}^n p(x_i) = 1.$$

"Shannon", who is the Father of Information Theory, first came up with the concept of Entropy.

Example :

Consider the binary event of people who like or dislike Ice Cream, I_C, \bar{I}_C .

$$\text{Random variable, } X = \{ I_C, \bar{I}_C \}$$

Probability that people like Ice cream, $P(I_C) = \frac{3}{4}$

Probability that people don't like Ice cream, $P(\bar{I}_C) = 1 - P(I_C)$

$$= 1 - \frac{3}{4}$$

$$= \frac{1}{4}$$

Now, the Information / Entropy of the Random Variable X is given as

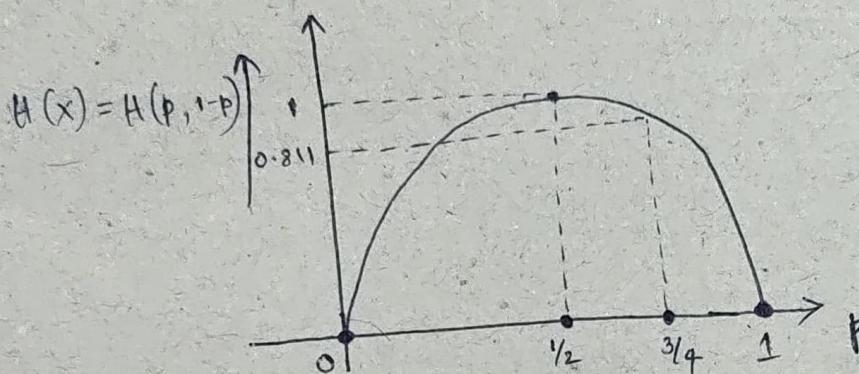
$$H(X) = H(P(I_c), P(\bar{I}_c))$$

$$= H\left(\frac{3}{4}, \frac{1}{4}\right)$$

$$= \frac{3}{4} \log_2\left(\frac{4}{3}\right) + \frac{1}{4} \log_2(4)$$

$$H(x) = 0.811$$

Plot the Entropy as a function of the probability of the source.



$$\begin{aligned} p=0 &\Rightarrow H(0,1)=0 \\ p=1 &\Rightarrow H(1,0)=0 \\ p=\frac{1}{2} &\Rightarrow H\left(\frac{1}{2}, \frac{1}{2}\right)=1 \\ p=\frac{3}{4} &\Rightarrow H\left(\frac{3}{4}, \frac{1}{4}\right)=0.811 \end{aligned}$$

The Entropy $H(x)$ reaches the peak when all the symbols have equal probability. (ii) when $p = \frac{1}{2}$, $H(x) = H\left(\frac{1}{2}, \frac{1}{2}\right) = 1$.

CONDITIONAL ENTROPY

Consider two sources :

- X with symbols x_i
- Y with symbols y_j

The conditional Entropy, $H(X|Y)$

(i) The information in one Random Variable, given the other Random Variable is defined as

$$H(X|Y) = \sum_{j=1}^m p(y_j) H(X|Y=y_j).$$

(ii) First calculate entropy of X for each value of $Y=y_i$, then take the weighted average of all these entropies, by weighing with respective probabilities $p(y_j)$ and take the sum.

For every possible value that the variable Y can take $y \in \{y_1, y_2, \dots, y_m\}$, we calculate the entropy of X and then we take a weighted average of all these entropies by weighing with the respective probabilities of each value of the Random Variable Y .

Week 6 : Session 4

Example :

Consider the table below showing Joint Probabilities of $X = \{\text{IC}, \overline{\text{IC}}\}$, $Y = \{\text{CHOC}, \overline{\text{CHOC}}\}$.

	IC	$\overline{\text{IC}}$
CHOC	$\frac{1}{2}$	$\frac{1}{8}$
$\overline{\text{CHOC}}$	$\frac{1}{4}$	$\frac{1}{8}$

Eg. Probability that people like Icecream and don't like Chocolate.

What is the Conditional Entropy $H(X|Y)$?

$$\begin{aligned}
 H(X|Y) &= \sum_{j=1}^m P(y_j) H(X|Y=y_j) \\
 &= P(\text{CHOC}) \times H(X|Y=\text{CHOC}) \\
 &\quad + P(\overline{\text{CHOC}}) \times H(X|Y=\overline{\text{CHOC}})
 \end{aligned}$$

Entropy in X , given person likes chocolate

Entropy in X , given person doesn't like chocolate

- From the table, we get

$$P(\text{CHOC}) = \frac{1}{2} + \frac{1}{8} = \frac{5}{8}$$

$$P(\overline{\text{CHOC}}) = \frac{1}{4} + \frac{1}{8} = \frac{3}{8}$$

- Now, $H(X|Y=\text{CHOC}) = H(P(\text{IC}|\text{CHOC}), P(\overline{\text{IC}}|\text{CHOC}))$

$$P(\text{IC}|\text{CHOC}) = \frac{P(\text{IC} \cap \text{CHOC})}{P(\text{CHOC})} = \frac{\frac{1}{2}}{\frac{1}{2} + \frac{1}{8}} = \frac{4}{5}$$

$$P(\overline{\text{IC}}|\text{CHOC}) = \frac{P(\overline{\text{IC}} \cap \text{CHOC})}{P(\text{CHOC})} = \frac{\frac{1}{8}}{\frac{1}{2} + \frac{1}{8}} = \frac{1}{5}$$

$$\text{Therefore, } H(X|Y = \text{CHOC}) = H\left(\frac{4}{5}, \frac{1}{5}\right)$$

$$= \frac{4}{5} \log_2\left(\frac{5}{4}\right) + \frac{1}{5} \log_2(5) \\ = 0.722.$$

$$\textcircled{O} \text{ Now, } H(X|Y = \overline{\text{CHOC}}) = H\left(P(\text{IC}|\overline{\text{CHOC}}), P(\overline{\text{IC}}|\overline{\text{CHOC}})\right)$$

$$P(\text{IC}|\overline{\text{CHOC}}) = \frac{P(\text{IC} \cap \overline{\text{CHOC}})}{P(\overline{\text{CHOC}})} = \frac{\frac{1}{4}}{\frac{1}{4} + \frac{1}{8}} = \frac{2}{3}$$

$$P(\overline{\text{IC}}|\overline{\text{CHOC}}) = \frac{P(\overline{\text{IC}} \cap \overline{\text{CHOC}})}{P(\overline{\text{CHOC}})} = \frac{\frac{1}{8}}{\frac{1}{4} + \frac{1}{8}} = \frac{1}{3}$$

$$\text{Therefore, } H(X|Y = \overline{\text{CHOC}}) = H\left(\frac{2}{3}, \frac{1}{3}\right)$$

$$= \frac{2}{3} \log_2\left(\frac{3}{2}\right) + \frac{1}{3} \log_2(3) \\ = 0.918.$$

Finally, the conditional Entropy is calculated as

$$H(X|Y) = \frac{5}{8} \times 0.722 + \frac{3}{8} \times 0.918$$

$$H(X|Y) = 0.7955$$

Information Gain (IG)

The Information Gain (IG) is used to choose the best attribute. The Information Gain is defined as

$$IG(X|Y) = H(X) - H(X|Y)$$

$$\Rightarrow IG(X|Y) = 0.811 - 0.7955 = 0.0155.$$

\nwarrow Information gained by observing Y

This is also known as the Mutual Information (MI).

Larger the Mutual Information,

More the information Conveyed ...

Week 7 Session 1

DTC Feature Selection

- Choose the feature / attribute , that maximizes the Information Gain !

$$IG(x|y) = H(x) - H(x|y)$$

- For each attribute Y , we compute Information Gain , and choose the attribute that maximizes the Information Gain.

Example

Consider the table shown below , which says the decisions of each and every customer , to WAIT / NOT WAIT at the restaurants .

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

The customer's decisions are based on the attributes as defined below .

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).

Out of these 10 attributes/features, which one to choose to construct the DTC?

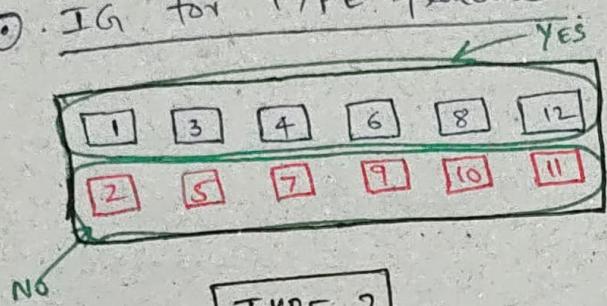
For this example, we consider the two attributes

- TYPE
- PATRONS

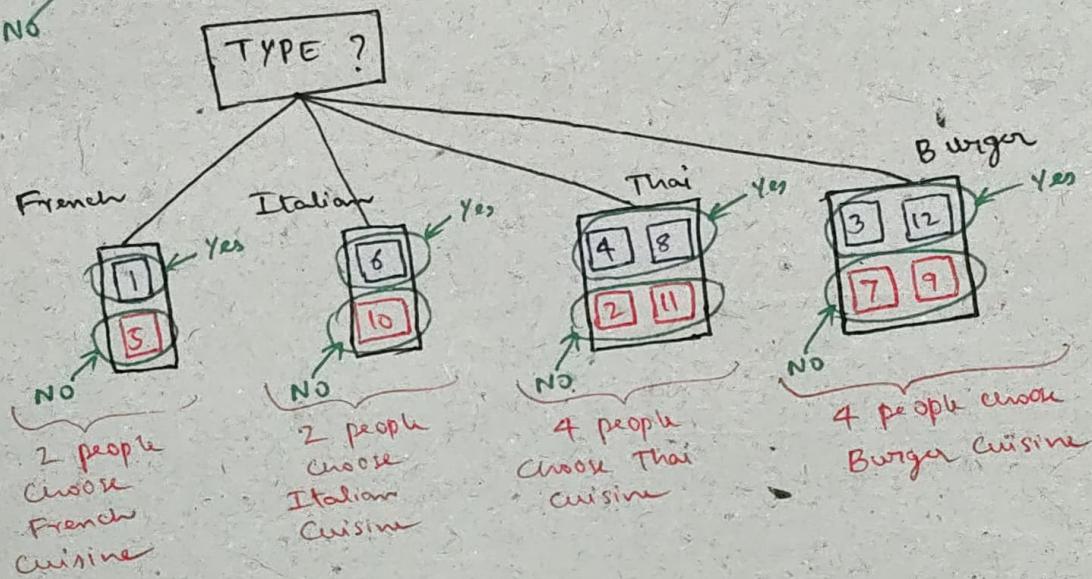
Out of these two, which is better for DTC?

Naturally, we choose the one which has the maximum IG.

① IG for TYPE feature



Out of 12, half people choose to WAIT, half people choose NOT to WAIT.



In any of the cuisines, half people choose to wait, and half people choose NOT to wait.

The IG for TYPE feature is given as

$$IG(\text{TYPE}) = H(x) - H(x|\text{TYPE}).$$

$$\Rightarrow H(x) = H\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{1}{2} \log_2(2) + \frac{1}{2} \log_2(2) = 1.$$

$$\Rightarrow H(x|\text{TYPE}) = Pr(\text{Fr}) H(x|\text{Fr}) + Pr(\text{It}) H(x|\text{It}) \\ + Pr(\text{Th}) H(x|\text{Th}) + Pr(\text{Bu}) H(x|\text{Bu})$$

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

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

$$= \frac{2}{12} + \frac{2}{12} + \frac{4}{12} + \frac{4}{12} = 1.$$

Therefore, the IG for the TYPE feature

$$IG(\text{TYPE}) = 1 - 1 = 0$$

Inference :

What is the information that we gained, by knowing the Cuisine ? **NOTHING.**

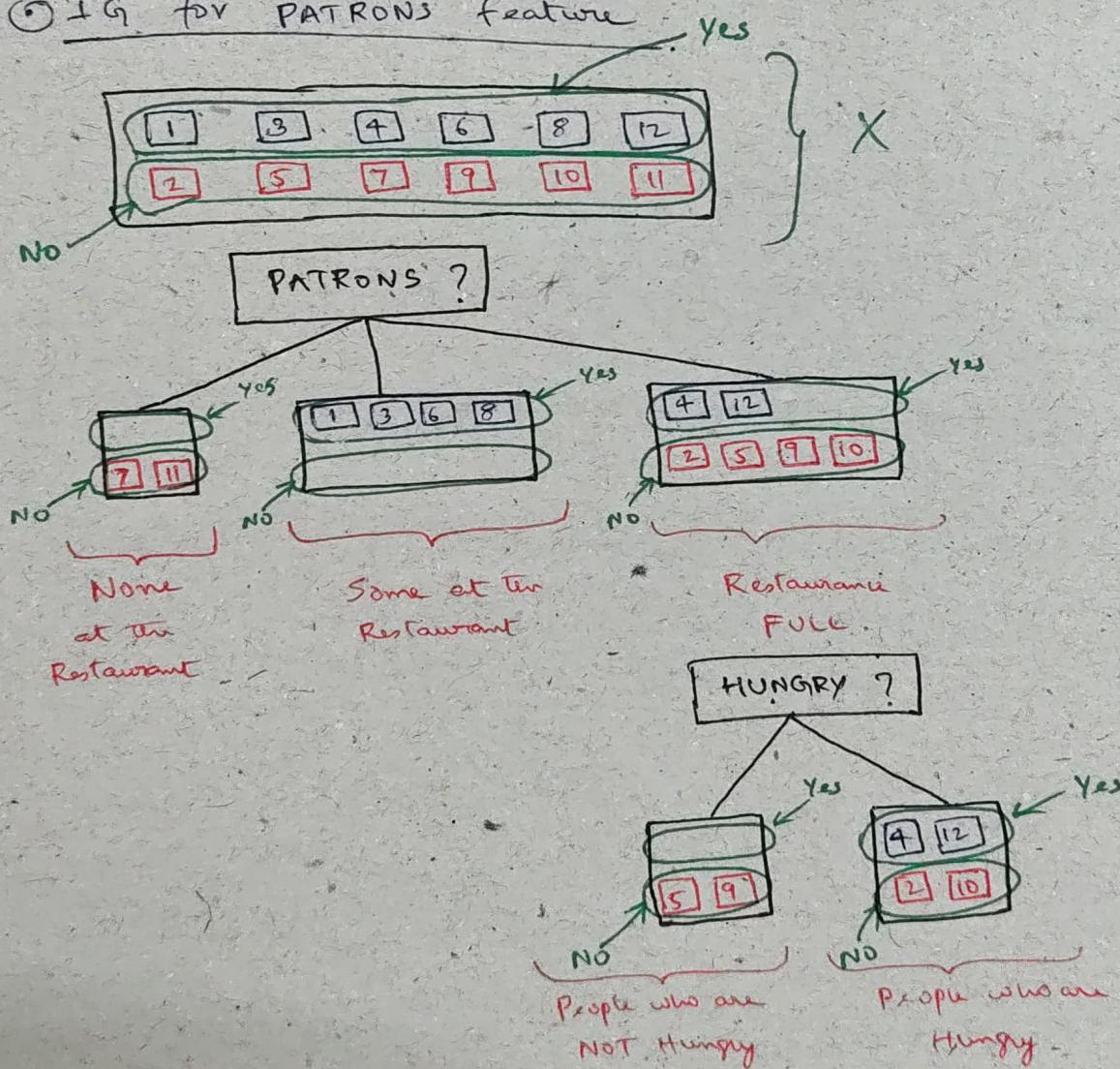
So, ① Probability of persons Waiting = $\frac{1}{2}$

② After knowing the Cuisine, }
Probability of persons waiting } = $\frac{1}{2}$

So, knowing the Cuisine is not helping us. (ie) There is no Information Gain even if we get to know what is the Cuisine (French / Italian / Thai / Burger).

There is no additional information that the Cuisine is giving us about the decision of the customer to wait / Not. Therefore, the TYPE feature is not suitable for the DTC.

③ IG for PATRONS feature



The IG for PATRONS feature is given as

$$IG(\text{PATRONS}) = H(x) - H(x \mid \text{PATRONS})$$

WKT, $H(x) = 1$.

$$\begin{aligned} \text{Now, } H(x \mid \text{PATRONS}) &= \Pr(\text{None}) H(x \mid \text{None}) \\ &\quad + \Pr(\text{Some}) H(x \mid \text{Some}) \\ &\quad + \Pr(\text{Full}) H(x \mid \text{Full}) \\ &= \frac{2}{12} H(0, 1) + \frac{4}{12} H(1, 0) \\ &\quad + \frac{6}{12} H\left(\frac{1}{3}, \frac{2}{3}\right) \\ &= \frac{1}{2} \left(\frac{1}{3} \log_2(3) + \frac{2}{3} \log_2\left(\frac{3}{2}\right) \right) \\ &= 0.46. \end{aligned}$$

Therefore, the IG for the PATRONS feature

$$IG(\text{FEATURE}) = 1 - 0.46 = 0.54$$

- Since $IG(\text{PATRONS}) > IG(\text{TYPE})$, we choose PATRONS as the feature to split.

- Final DTC.

