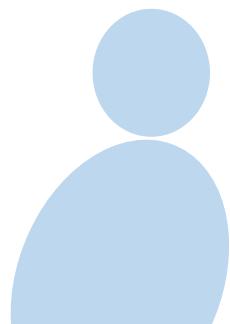


Elective Module:

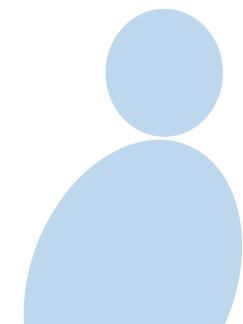
**Advanced ML
Techniques**



Chapter 10

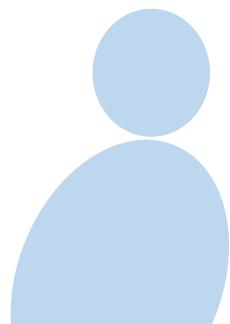
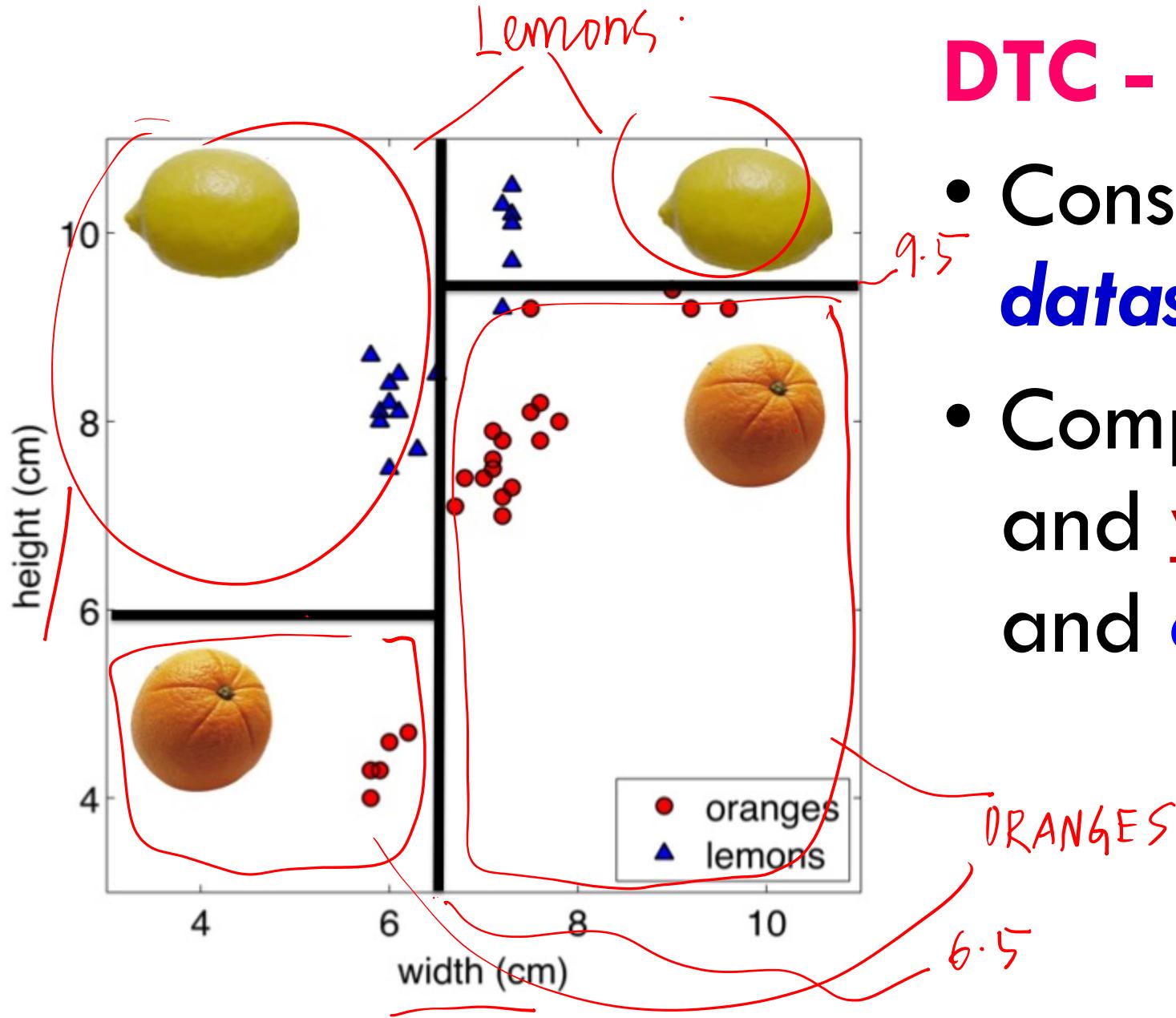
Decision Tree Classifiers

DTC



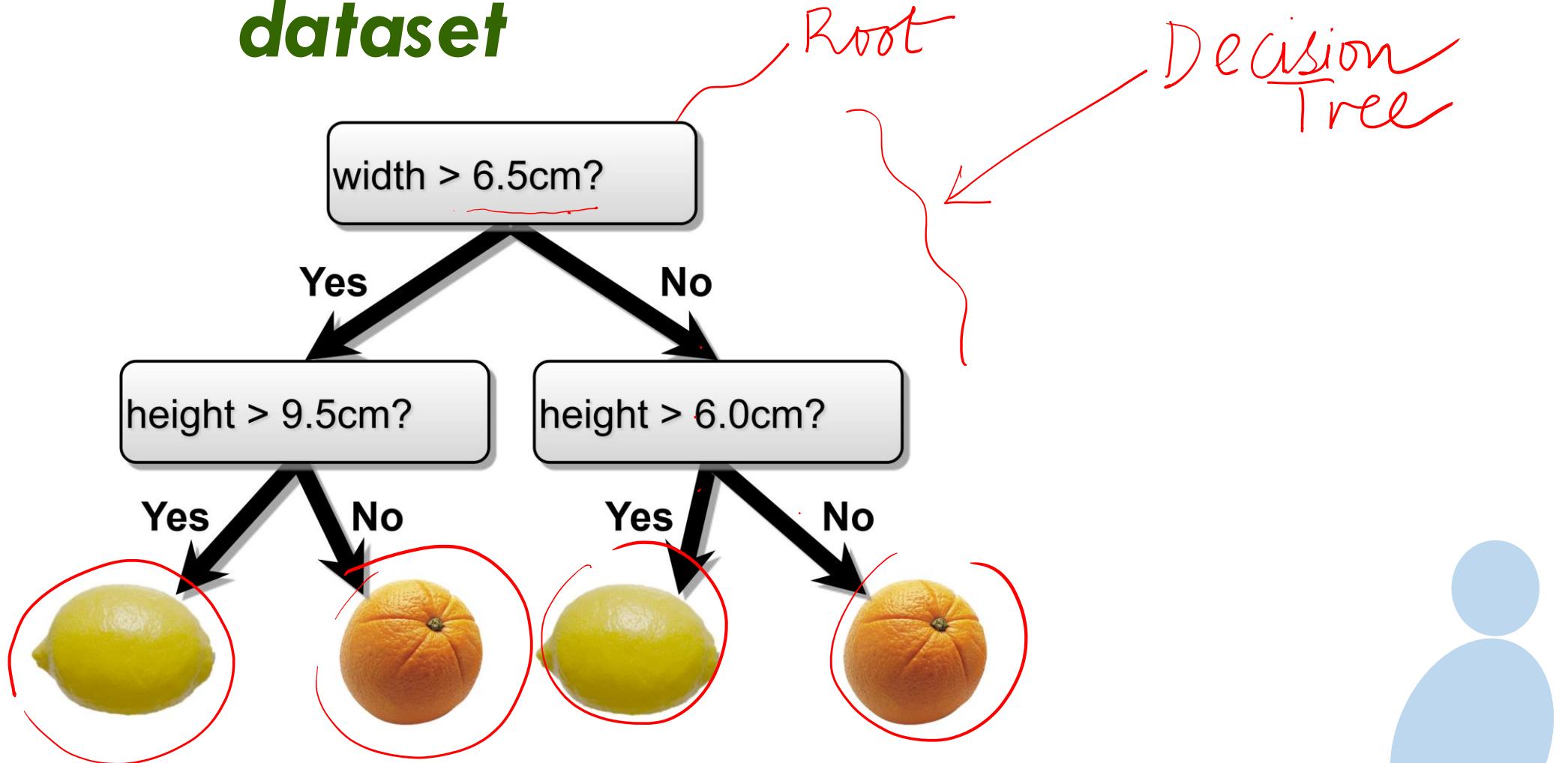
DTC - Example

- Consider the simple **dataset** shown
- Comprises of heights and weights of **lemons** and **oranges**



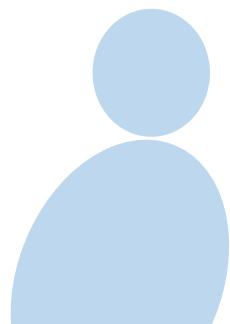
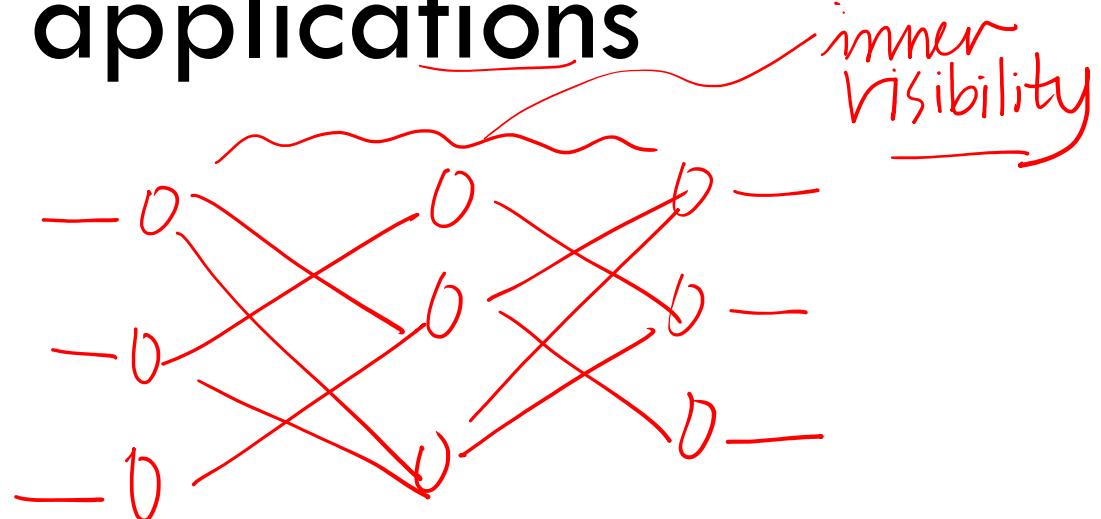
DTC - Example

- DTC below is built using the given **dataset**



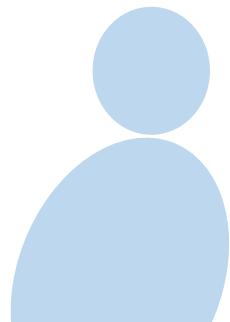
Decision Tree Classifiers

- Advantages: *interpretable, intuitive*
(in contrast to Neural Nets)
 - Popular in medical diagnosis
applications
- DTCS
very popular
in medicine*



Decision Tree Classifiers

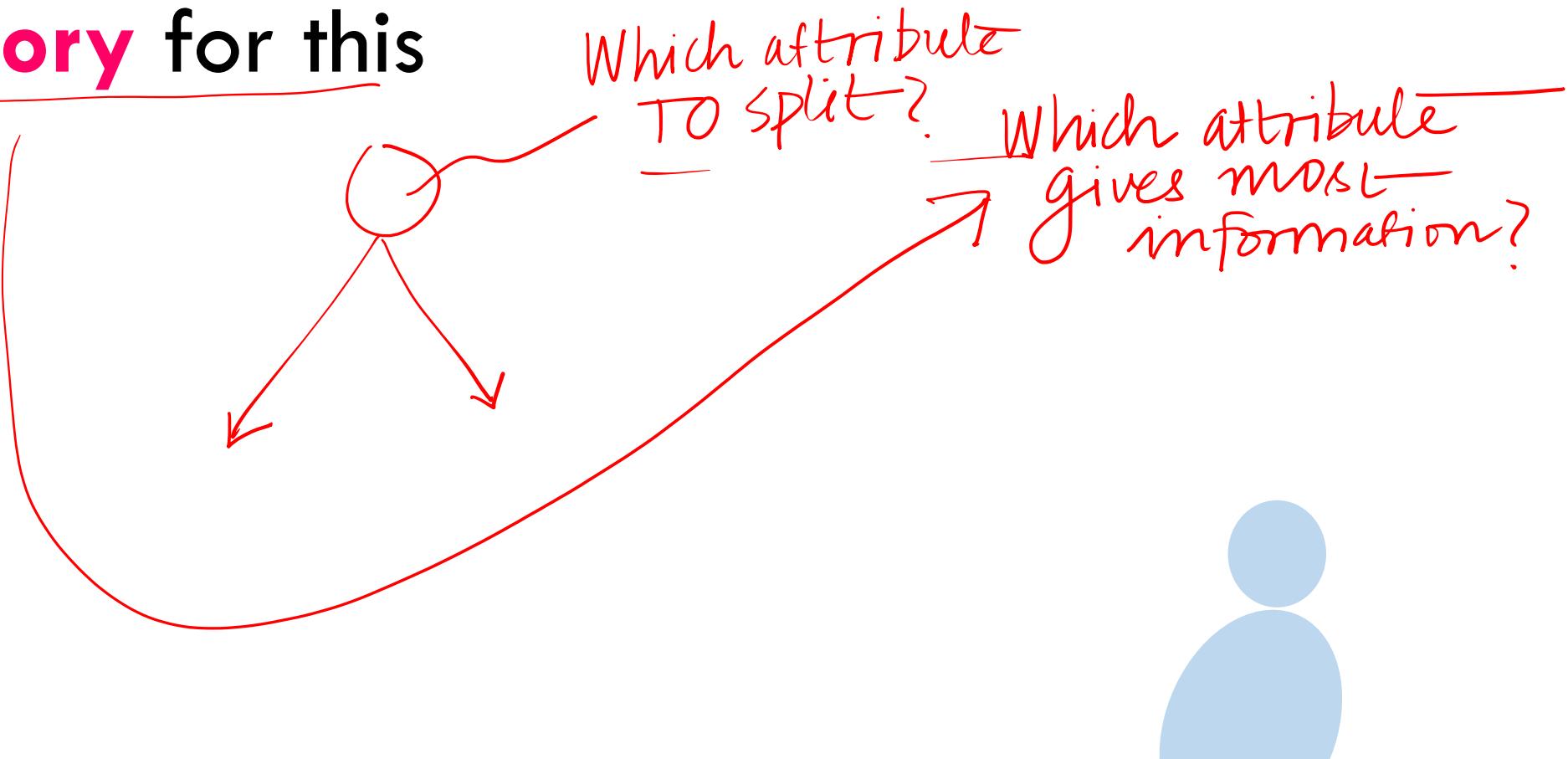
- DTCs are well-suited to model
discrete outcomes



Learning Decision Trees

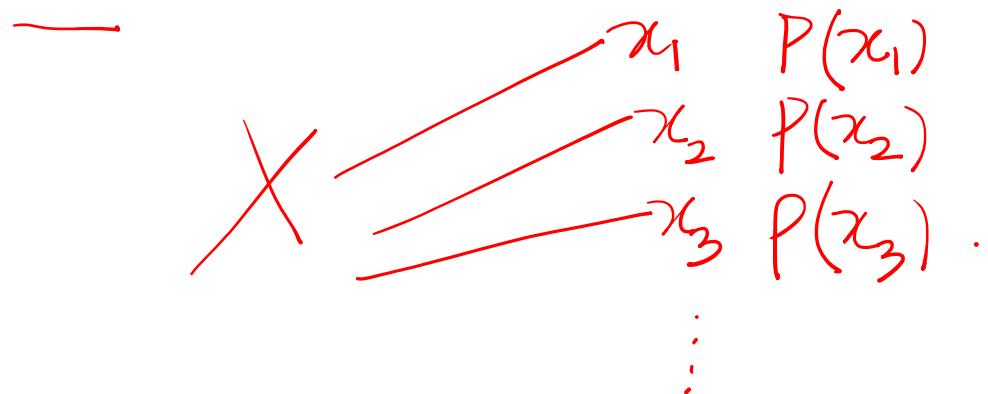
Build DTC

- How to choose the best attribute?
- One can use principles of information theory for this



Entropy → intrinsic information

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



Entropy $\sum_i p(x_i) = 1$ Entropy = 0
 if $p(x_i) = 1$
 $p(x_j) = 0, j \neq i$

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

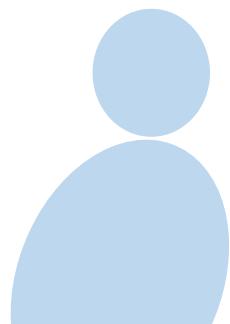
$$H(X) = p(x_1) \log_2 \left(\frac{1}{p(x_1)} \right) + p(x_2) \cdot \log_2 \left(\frac{1}{p(x_2)} \right) + \dots$$

$$= \sum_i p(x_i) \log_2 \frac{1}{p(x_i)}$$

$$= - \sum_i p(x_i) \log_2 (p(x_i))$$

Entropy maximum: $p(x_i) = \frac{1}{n}$ Equally probable symbols

Entropy or information content.
 $\underline{\text{Entropy} \geq 0}$



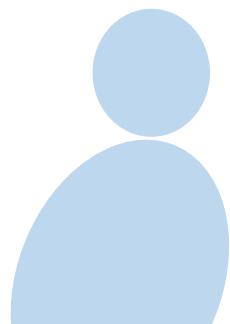
Entropy

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

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

Shannon:

Father of
information
Theory



Entropy-Example

- Example: Consider the binary event like or dislike ice cream (IC), (\bar{IC})

$$X = \{IC, \bar{IC}\}$$

$$P(IC) = \frac{3}{4}, P(\bar{IC}) = \frac{1}{4}$$

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

information
or entropy of Random variable

Like ice cream
Do NOT like ice cream

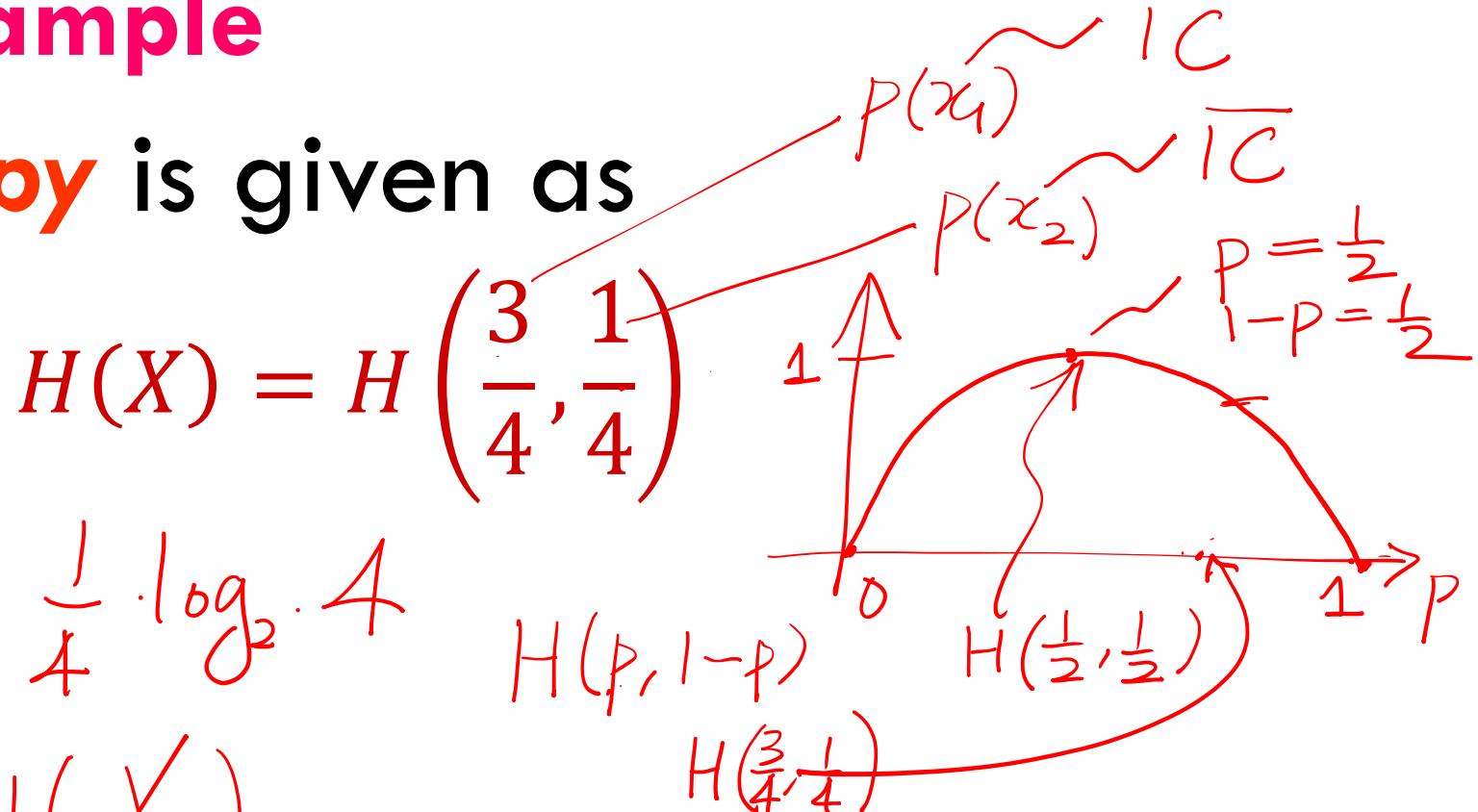
Entropy-Example

- The **entropy** is given as

$$= \frac{3}{4} \cdot \log_2 \frac{4}{3} + \frac{1}{4} \cdot \log_2 4$$

$$= 0.811 \approx H(X)$$

ice cream R ✓

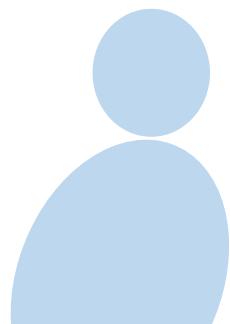


Entropy-Example

- The **entropy** is given as

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

$$= \frac{3}{4} \times \log_2 \frac{4}{3} + \frac{1}{4} \times \log_2 4 \approx 0.811$$



Conditional entropy

$$\underline{H(X|Y)}$$

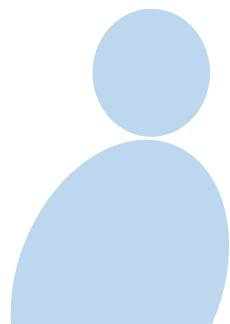
- Consider two sources: X with symbols x_i and Y with symbols y_j

$$\underline{H(X|Y)}$$

Conditional Entropy

$$\underline{X|Y}$$

Conditional
Entropy.



Conditional entropy

Weighing with probability.

- The **conditional entropy** $H(X|Y)$ is defined as

$$H(X|Y) = \sum_j P(y_j) H(X|Y=y_j)$$

Calculate entropy
of X for each value
of $Y = y_j$
conditional Entropy

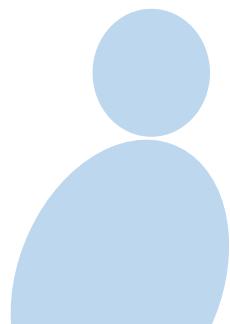
Weighted average

Conditional entropy

- The **conditional entropy** $H(X|Y)$ is defined as

$$Y \in \{y_1, y_2, \dots, y_m\}$$

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



Conditional entropy — Example:

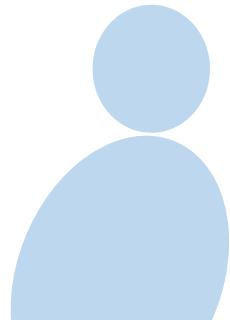
- Consider the table below showing **joint probabilities** of

$$X = \{\text{IC}, \overline{\text{IC}}\}, Y = \{\text{CHOC}, \overline{\text{CHOC}}\}$$

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

*Ice cream
Choc*

*like ice cream and chocolate
like chocolate but NOT
ice cream*



Conditional entropy

- What is $H(X|Y)$?

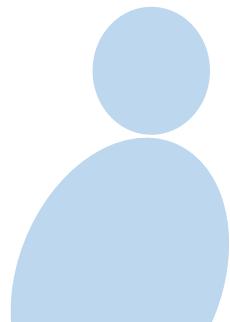
$$H(X|Y) = P(\text{CHOC}) \times H(X|Y = \text{CHOC})$$

$$+ P(\text{CH}\bar{\text{O}}\text{C}) \times H(X|Y = \text{CH}\bar{\text{O}}\text{C})$$

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

$$P(\text{CH}\bar{\text{O}}\text{C}) = \frac{1}{4} + \frac{1}{8} = \frac{3}{8}$$

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



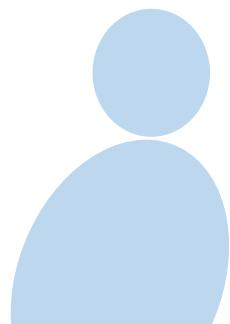
Conditional entropy

- What is $H(X|Y)$?

$$H(X|Y)$$

$$= P(\text{CHOC}) \times H(X|\text{CHOC}) + P(\overline{\text{CHOC}}) H(X|\overline{\text{CHOC}})$$

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



Conditional entropy

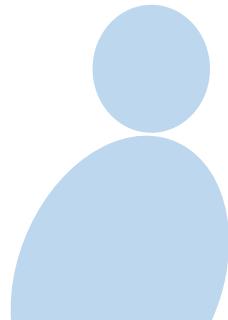
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\bullet H(X|Y = \text{CHOC}) = ?$$

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

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

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

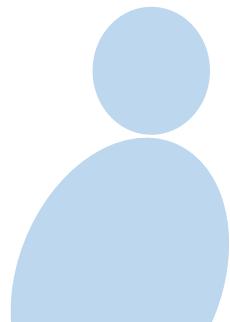


Conditional entropy

- $H(X|Y = \text{CHOC}) = ?$

$$P(\text{IC}|\text{CHOC}) = \frac{4}{5}, P(\overline{\text{IC}}|\text{CHOC}) = \frac{1}{5}$$

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



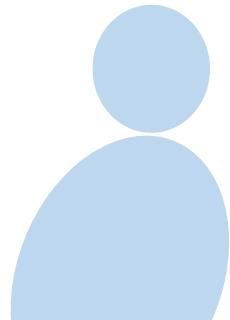
Conditional entropy

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

$$= \frac{4}{5} \log_2 \frac{5}{4} + \frac{1}{5} \cdot \log_2 5 = 0.722$$

Entropy in X
given person likes chocolate

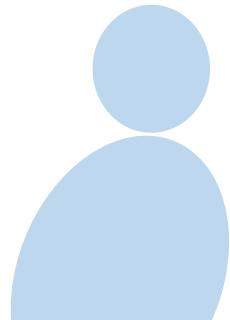
	IC	$\overline{\text{IC}}$
CHOC	$\frac{1}{2}$	$\frac{1}{8}$
$\overline{\text{CHOC}}$	$\frac{1}{4}$	$\frac{1}{8}$



Conditional entropy

$$\begin{aligned} H(X|\text{CHOC}) &= H\left(\frac{4}{5}, \frac{1}{5}\right) \\ &= \frac{4}{5} \times \log_2\left(\frac{5}{4}\right) + \frac{1}{5} \times \log_2(5) = 0.722 \end{aligned}$$

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



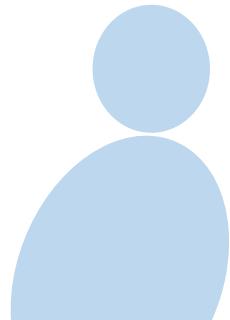
Conditional entropy

$$\bullet H(X|Y = \overline{\text{CHOC}}) = ?$$

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

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

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

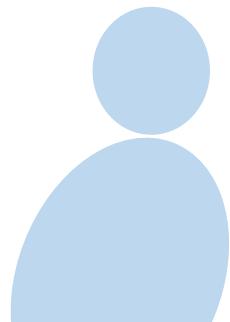


Conditional entropy

- $H(X|Y = \overline{\text{CHOC}}) = ?$

$$P(\text{IC}|\overline{\text{CHOC}}) = \frac{2}{3}, P(\overline{\text{IC}}|\overline{\text{CHOC}}) = \frac{1}{3}$$

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

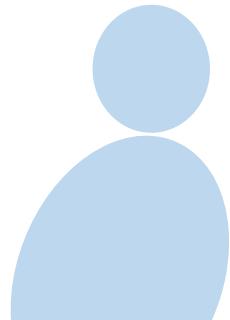


Conditional entropy

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

$$= \frac{2}{3} \cdot \log_2 \frac{3}{2} + \frac{1}{3} \cdot \log_2 3 = 0.918$$

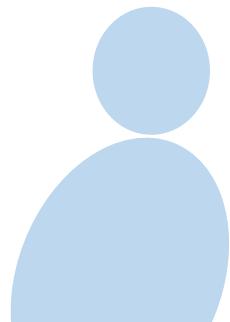
	IC	$\overline{\text{IC}}$
CHOC	$\frac{1}{2}$	$\frac{1}{8}$
$\overline{\text{CHOC}}$	$\frac{1}{4}$	$\frac{1}{8}$



Conditional entropy

$$\begin{aligned} H(X|\overline{\text{CHOC}}) &= H\left(\frac{2}{3}, \frac{1}{3}\right) \\ &= \frac{2}{3} \times \log_2\left(\frac{3}{2}\right) + \frac{1}{3} \times \log_2(3) \\ &= 0.918 \end{aligned}$$

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

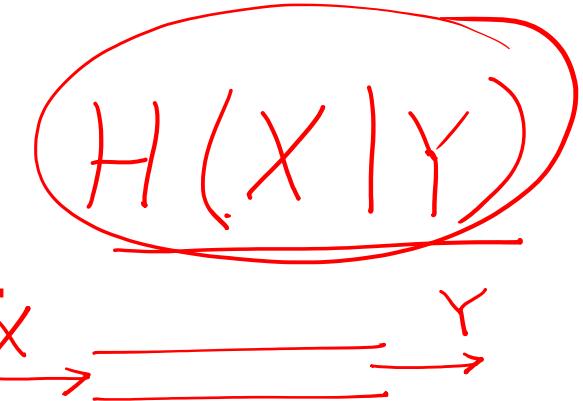


Conditional entropy

- Finally, the conditional entropy is given as,

$$\underline{H(X|Y)}$$

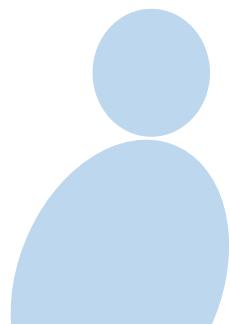
$$\begin{aligned} &= P(\text{CHOC}) \times H(X|\text{CHOC}) + p(\overline{\text{CHOC}}) \\ &\quad \times H(X|\overline{\text{CHOC}}) \\ &= \frac{5}{8} \times 0.722 + \frac{3}{8} \times 0.918 \\ &\approx 0.7955 \end{aligned}$$



Conditional entropy

- Finally, the conditional entropy is given as,

$$\begin{aligned}H(X|Y) &= P(\text{CHOC}) \times H(X|\text{CHOC}) + p(\overline{\text{CHOC}}) \\&\quad \times H(X|\overline{\text{CHOC}}) \\&= \frac{5}{8} \times 0.722 + \frac{3}{8} \times 0.918 = 0.7955\end{aligned}$$



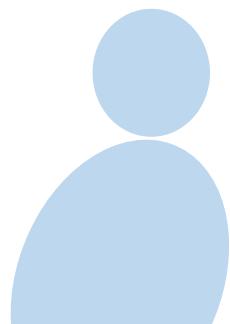
Information gain

- The information gain (IG) is defined as

$$IG(X|Y) = H(X) - H(X|Y) \geq 0$$

$$\begin{aligned} IG(X|Y) &= H(X) - H(X|Y) \\ &= 0.811 - 0.7955 \\ &= 0.0155 \end{aligned}$$

information gained by observing Y.

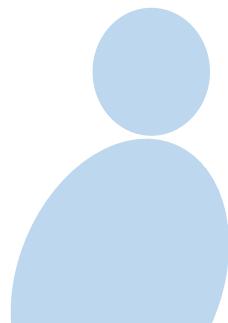


Information gain

- The **information gain (IG)** is defined as

$$\begin{aligned} \text{IG}(X|Y) &= H(X) - H(X|Y) \\ &= 0.811 - 0.7955 = 0.0155 \end{aligned}$$

.



Mutual information

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

Mutual information

Larger mutual info
⇒ more information conveyed.

DTC Feature Selection

- Choose the feature...that maximizes the **information gain!**

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

Choose attribute
that maximizes information
gain.

Attribute

Conditional Entropy

DTC Example

Attributes
Features.

- Consider the table shown below
- Customer decisions** to wait or not at restaurants

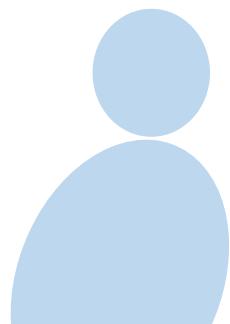
1.	Alternate: whether there is a suitable alternative restaurant nearby.
2.	Bar: whether the restaurant has a comfortable bar area to wait in.
3.	Fri/Sat: true on Fridays and Saturdays.
4.	Hungry: whether we are hungry.
5.	Patrons: how many people are in the restaurant (values are None, Some, and Full).
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).

Weekend

Hungry

Price Range?

Cuisine
TYPE



Data Set

Yes/No

DTC Example

- Table columns

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	Yes	Thai	0–10	$y_8 = \text{Yes}$	
x_9	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$	
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	$y_{12} = \text{Yes}$	

Alternative

NOT Friday

Saturday

Yes hungry

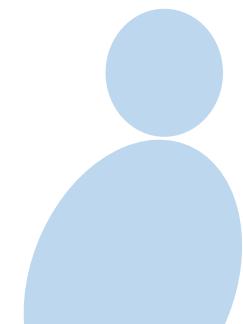
No Rain

Cuisine

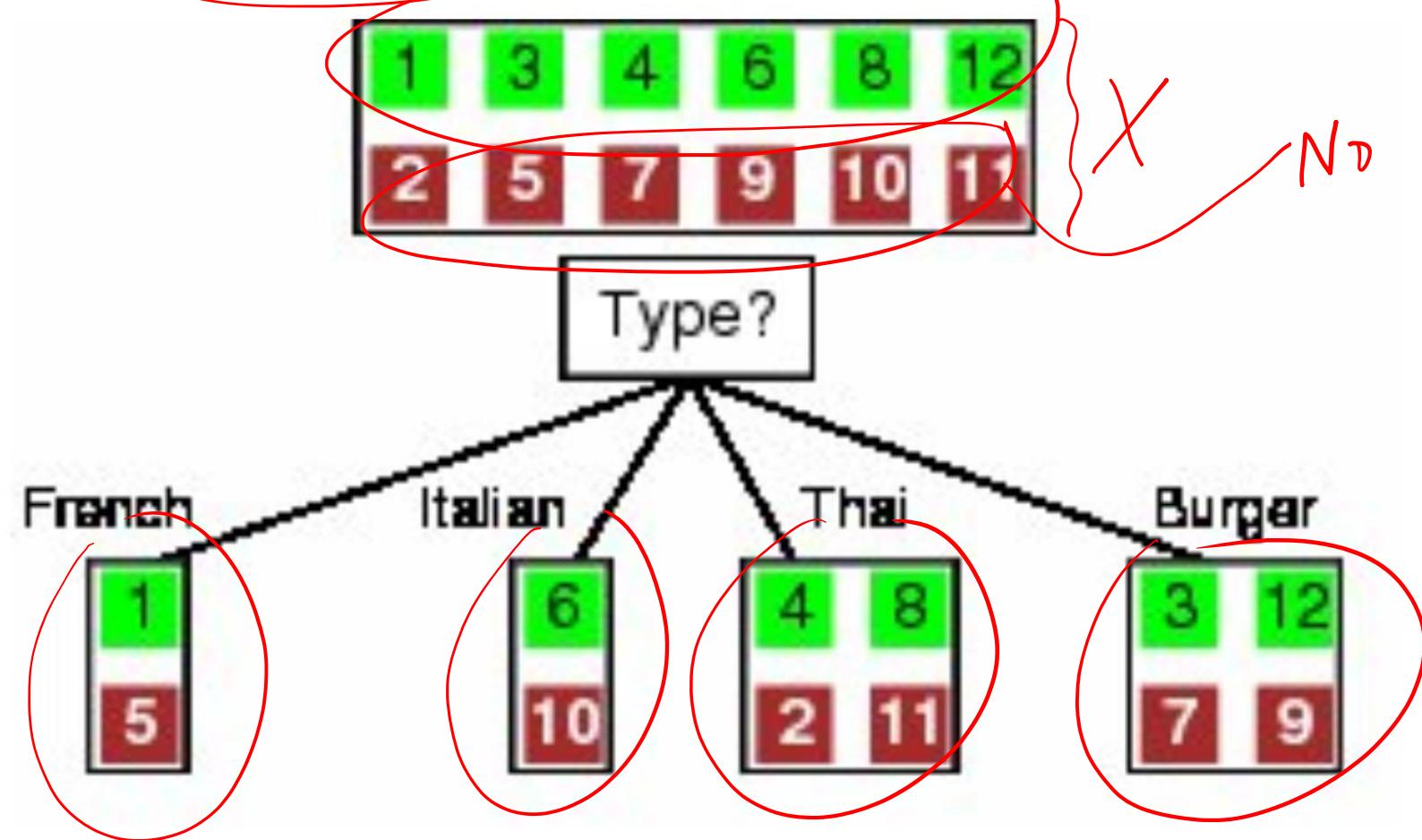
Wait : Yes !

10 Attributes

TYPE or PATRONS?
Which is better for DTC?



IG for Type



Yes: $\Pr(Y_{18}) = \frac{1}{2}$
 $\Pr(ND) = \frac{1}{2}$

IG for Type

Information
Gain

- IG for the **TYPE feature** is given as

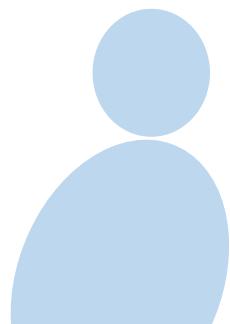
$$H(X) - H(X|TYPE)$$



IG for Type

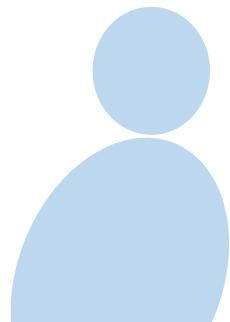
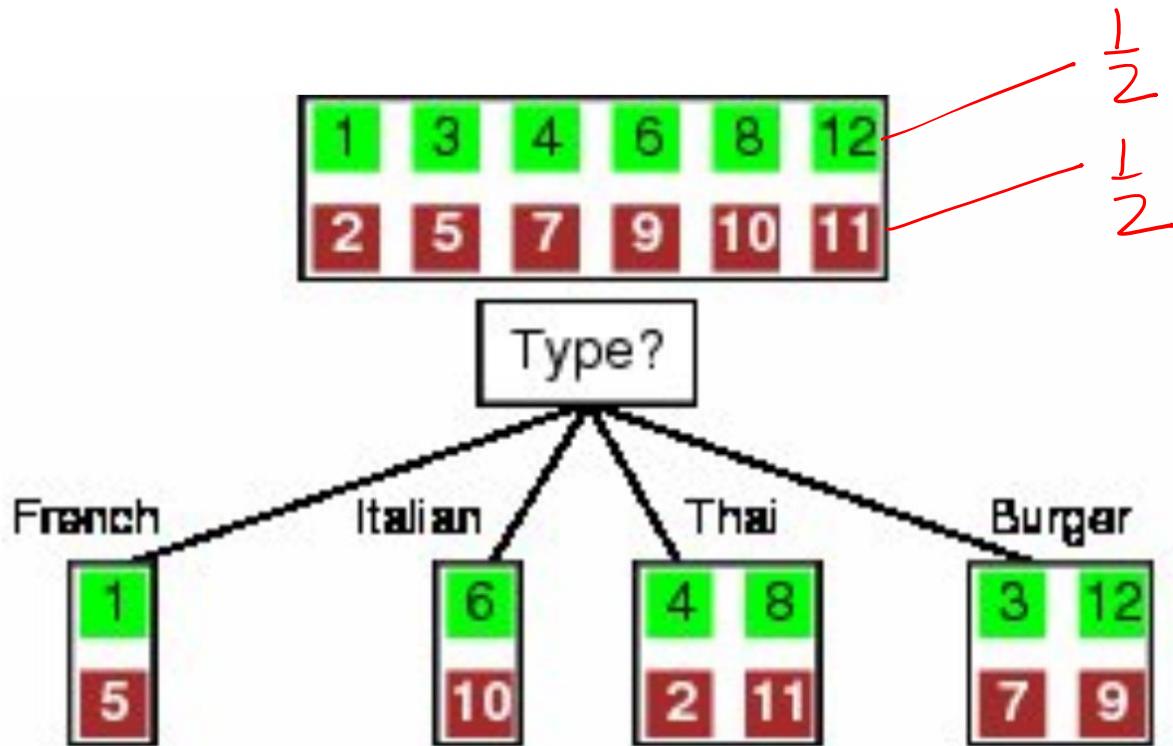
- IG for the **TYPE feature** is given as

$$\text{IG}(TYPE) = H(X) - \underline{H(X|TYPE)}$$



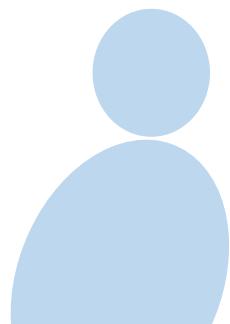
IG for Type

$$H(X) = H\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{1}{2} \log_2 2 + \frac{1}{2} \log_2 2$$
$$= \frac{1}{2} \times 1 + \frac{1}{2} \times 1 = 1$$



IG for Type

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



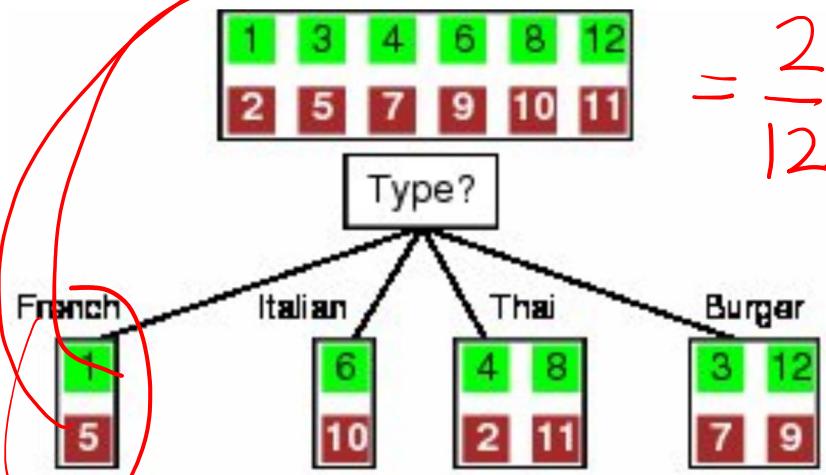
IG for Type

Conditional Entropy / type

$$H(X|TYPE) = P(\text{Fr}) \times H(X|\text{Fr}) + P(\text{It}) \times H(X|\text{It}) + P(\text{Th}) \times H(X|\text{Th}) + P(\text{Bu}) \times H(X|\text{Bu})$$

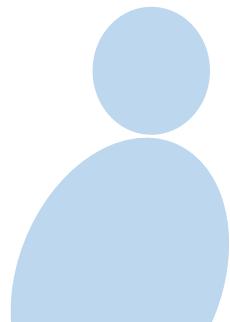
$$= \frac{2}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right)$$

No
Yes



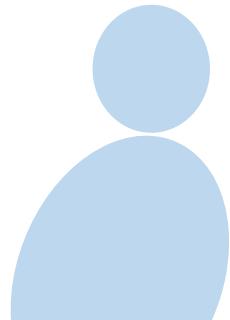
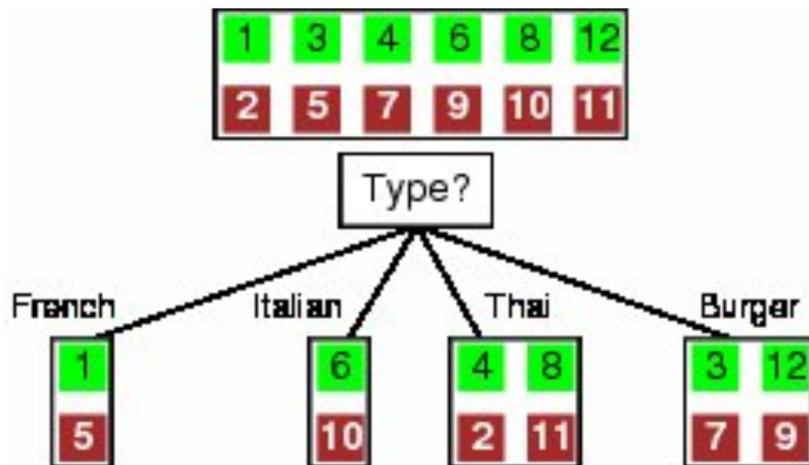
$$= \frac{2}{12} + \frac{2}{12} + \frac{4}{12} + \frac{4}{12} = 1$$

$$H(X|TYPE)$$



IG for Type

$$\begin{aligned} & \underline{H(X|TYPE)} \\ &= P(\text{Fr}) \times H(X|\text{Fr}) + P(\text{It}) \times H(X|\text{It}) + P(\text{Th}) \\ &\quad \times H(X|\text{Th}) + P(\text{Bu}) \times H(X|\text{Bu}) \\ &= \frac{2}{12} \times 1 + \frac{2}{12} \times 1 + \frac{4}{12} \times 1 + \frac{4}{12} \times 1 \\ &= \underline{\underline{1}} \end{aligned}$$



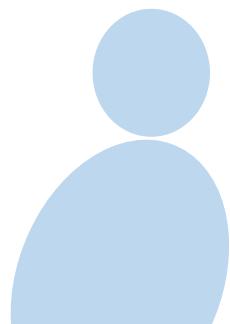
IG for Type

- IG for the **TYPE feature** is given as

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

$$= | - | = 0$$

NOT suitable for DTC



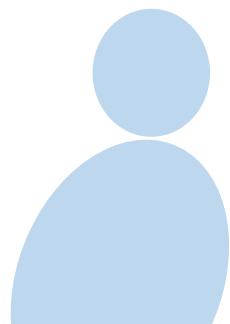
IG for Type

- IG for the **TYPE feature** is given as

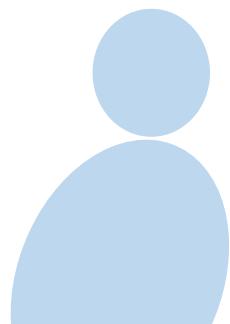
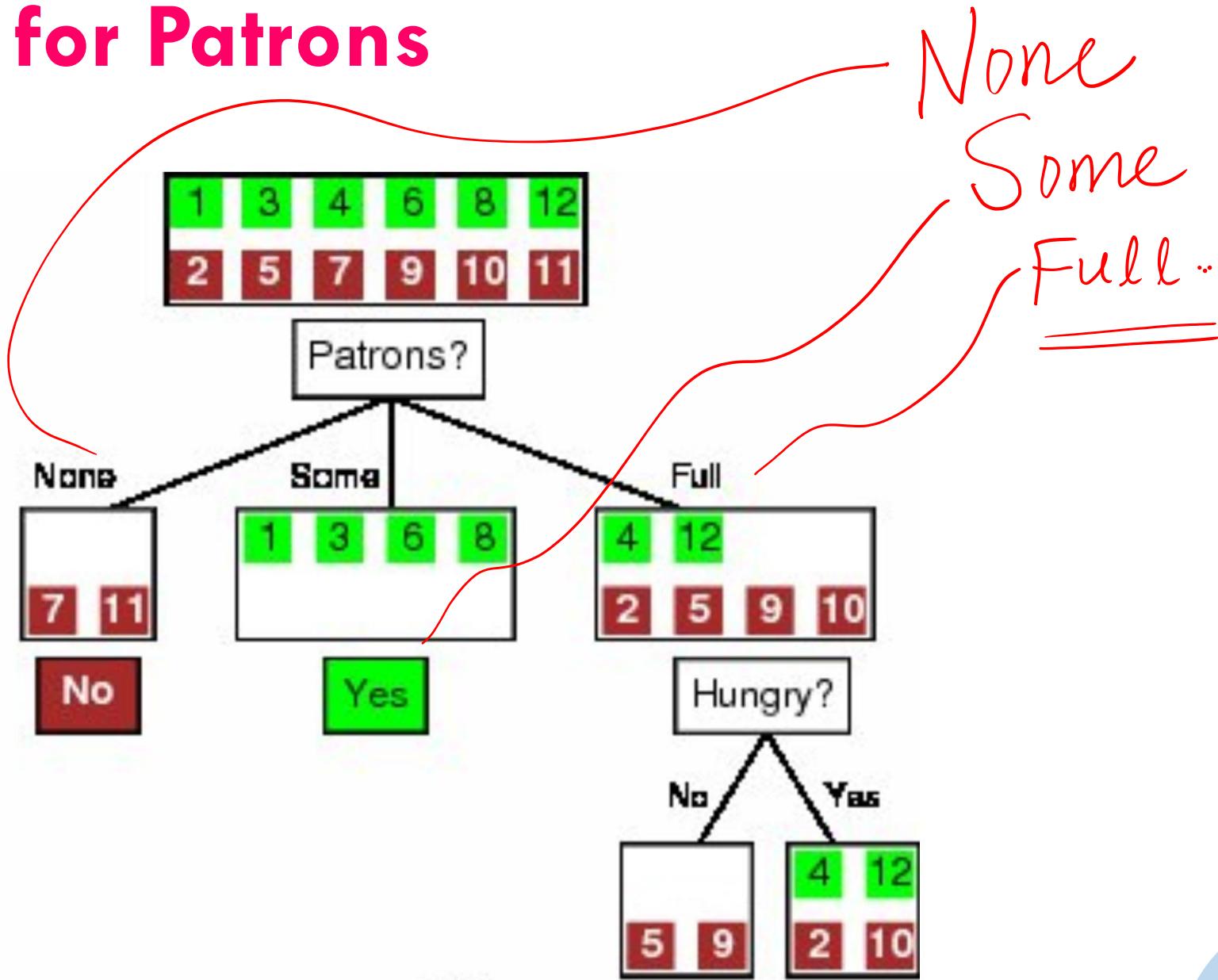
$$\begin{aligned} \text{IG}(TYPE) &= H(X) - H(X|TYPE) \\ &= 1 - 1 = 0 \end{aligned}$$

~~1~~

Information Gain = 0 .



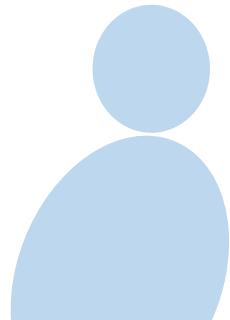
IG for Patrons



IG for Patrons

- IG for PATRONS feature is given as follows

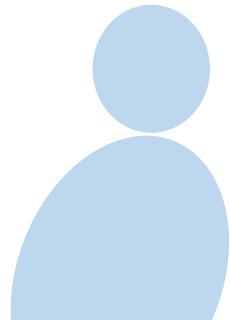
$$H(X) - \underline{H(X|PATRONS)}.$$



IG for Patrons

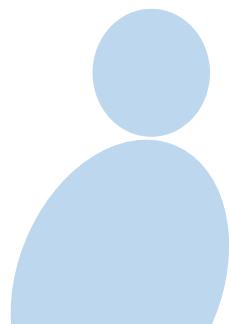
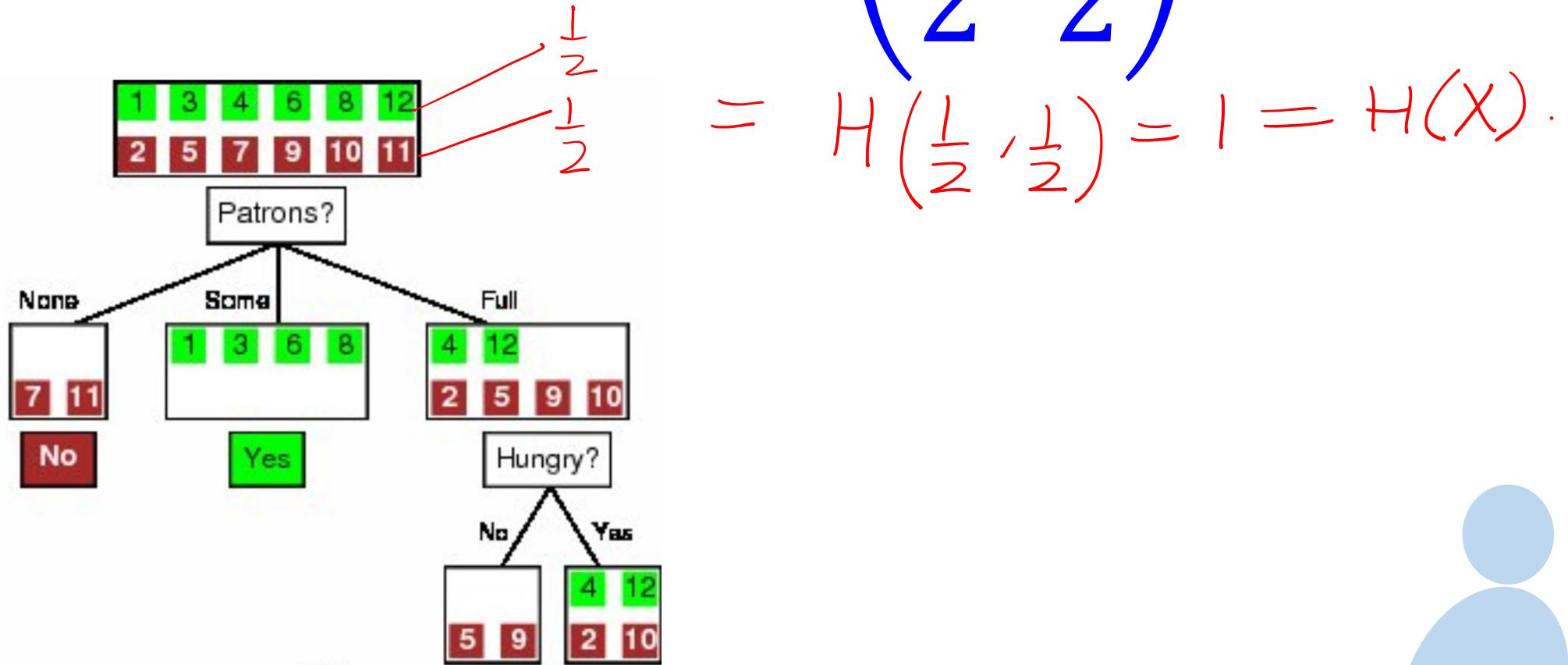
- IG for PATRONS feature is given as follows

$$H(X) - H(X|PATRONS)$$



IG for Patrons

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



IG for PATRONS

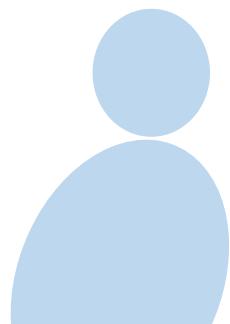
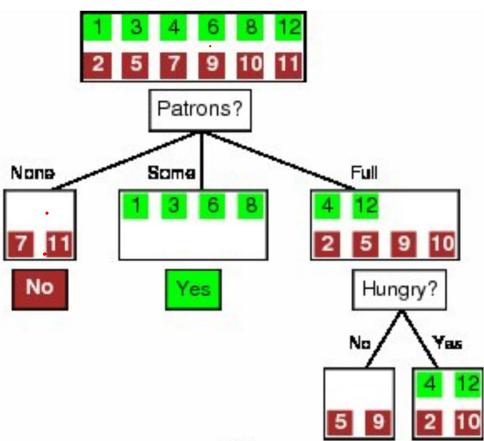
$$H(X|PATRONS)$$

$$= P(\text{None}) \times \cancel{H(X|\text{None})} + P(\text{Some}) \\ \times H(X|\text{Some}) + P(\text{Full}) \times H(X|\text{Full})$$

$$= \frac{2}{12} \times \cancel{H(0,1)}_0 + \frac{4}{12} \times \cancel{H(1,0)}_0 + \frac{6}{12} \times H\left(\frac{1}{3}, \frac{2}{3}\right)$$

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

$$\approx 0.46 \sim H(X|PATRONS)$$



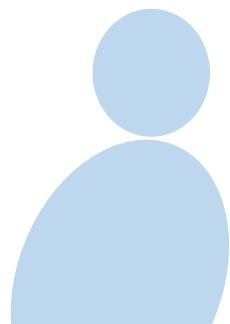
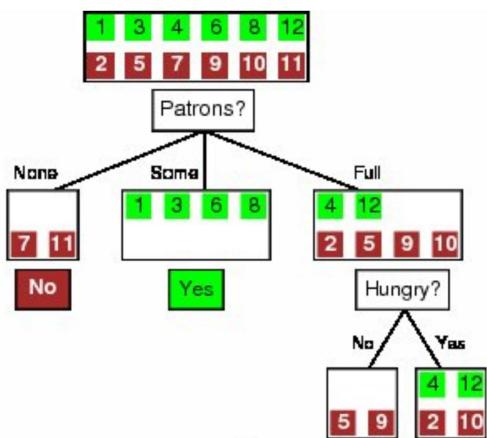
~~IG for PATRONS~~

$$H(X|Y = \text{PATRONS})$$

$$= P(\text{None}) \times H(X|\text{None}) + P(\text{Some}) \\ \times H(X|\text{Some}) + P(\text{Full}) \times H(X|\text{Full})$$

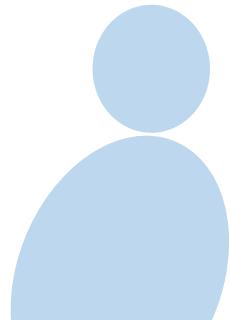
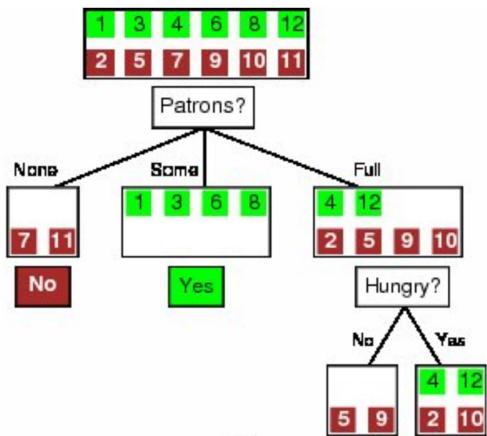
$$= \frac{2}{12} \times 0 + \frac{4}{12} \times 0 + \frac{1}{2} \times H\left(\frac{1}{3}, \frac{2}{3}\right)$$

$$= 0.46$$



IG for Type

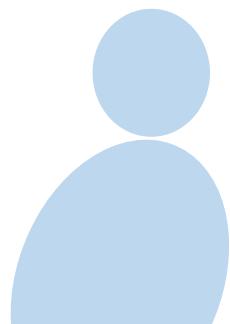
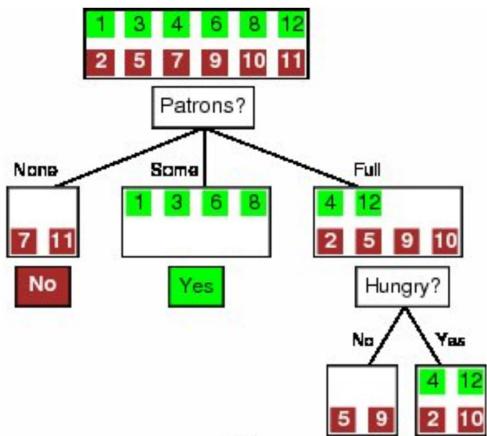
$$H(X) - H(X|PATRONS)$$
$$= | - 0.46 \neq 0.54$$



IG for Type

$$\begin{aligned} & H(X) - H(X|PATRONS) \\ &= 1 - \left(\frac{2}{12} \times 0 + \frac{4}{12} \times 0 + \frac{1}{2} \times H\left(\frac{1}{3}, \frac{2}{3}\right) \right) \\ &= 1 - 0.46 = 0.54 \end{aligned}$$

斜线

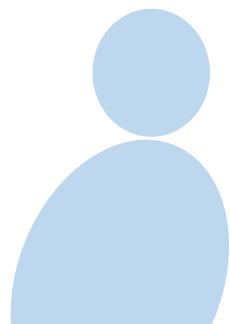


IG for Type

- Since

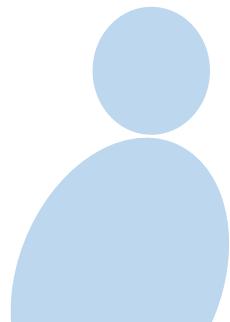
$$\text{IG(PATRONS)} = 0.54 > 0 = \text{IG(TYPE)}$$

Choose PATRONS
as attribute for DTC.

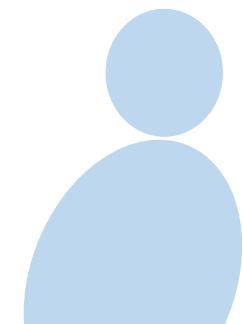
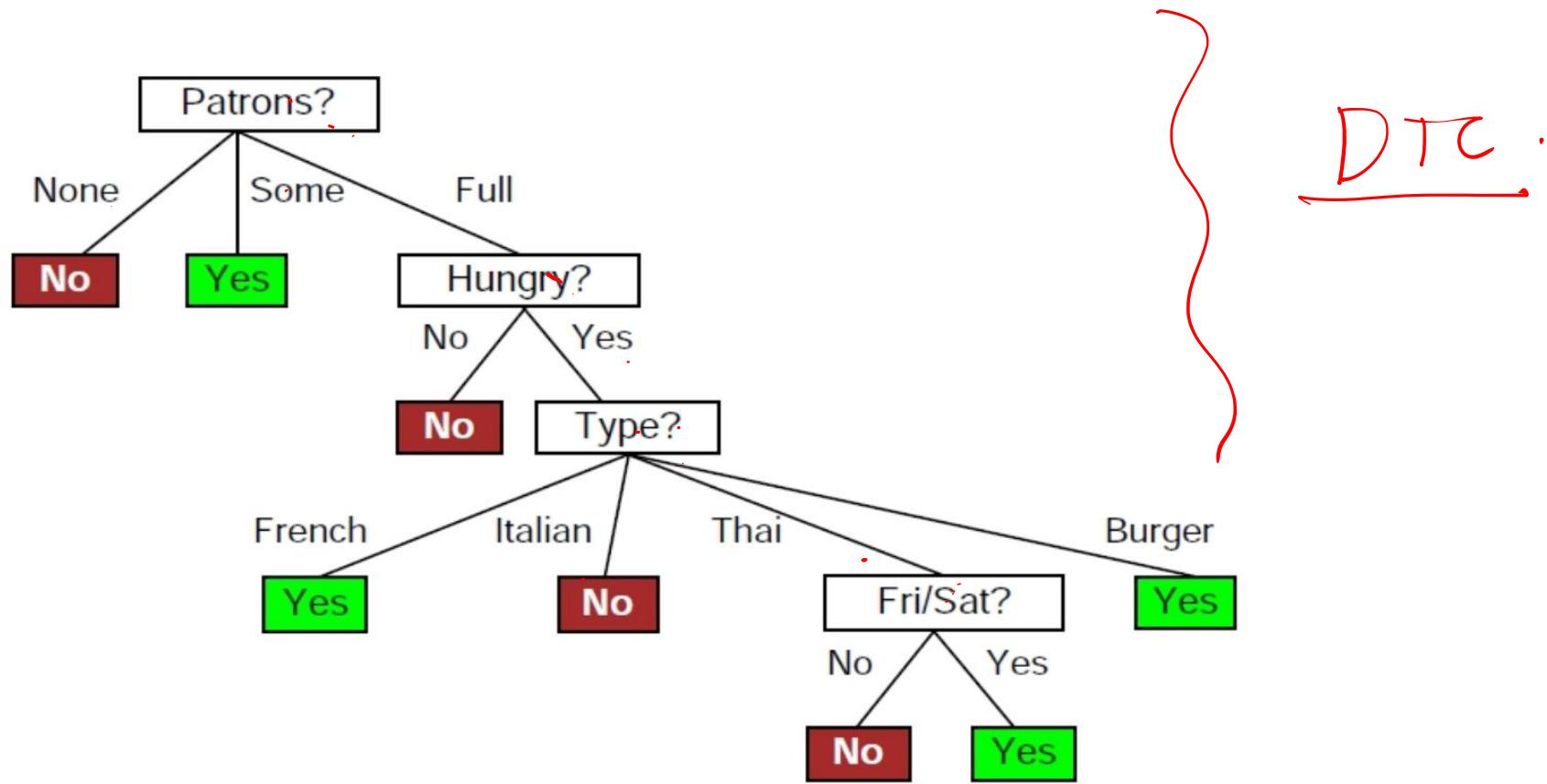


IG for Type

- We choose **PATRONS** as
the feature to split



Final DTC



Instructors may use this white area (14.5 cm / 25.4 cm) for the text.
Three options provided below for the font size.

Font: Avenir (Book), Size: 32, Colour: Dark Grey

Font: Avenir (Book), Size: 28, Colour: Dark Grey

Font: Avenir (Book), Size: 24, Colour: Dark Grey

Do not use the space below.

