

Welcome back everyone

2 Questions

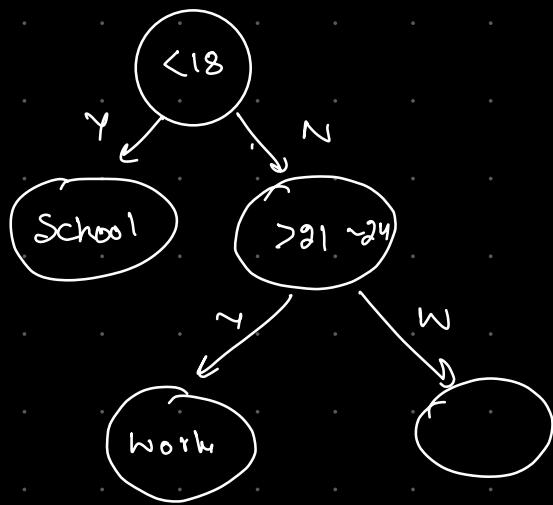
1. Classes in the morning from 10 - 1
2. If anyone has used a good
Condenser mic.

Decision Tree

1. Classification
2. Regressor

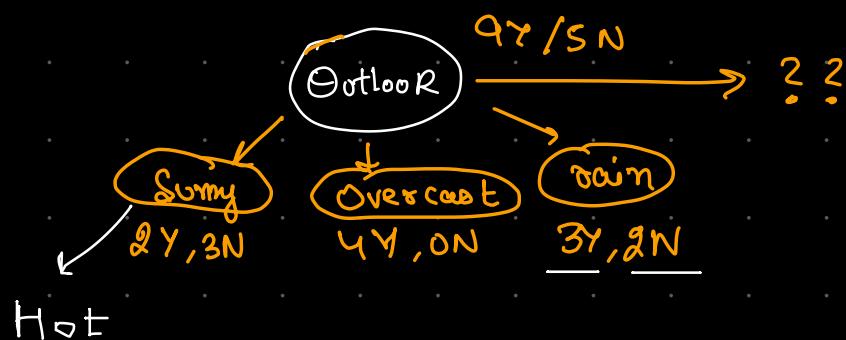
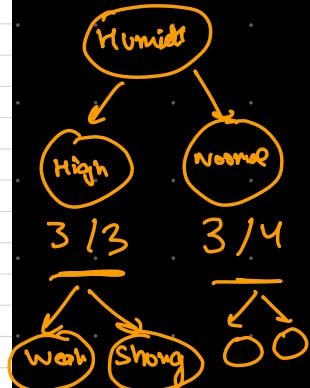
Nested if else clause

```
if (person < 18)
    _____ "College" School
elif person 18 - 22
    _____ 'College'
elif person 24 - 30
    _____ Job
else
    _____ Business.
```



Person \Rightarrow 14

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No



2 techniques in DT

① ID3

② CART *

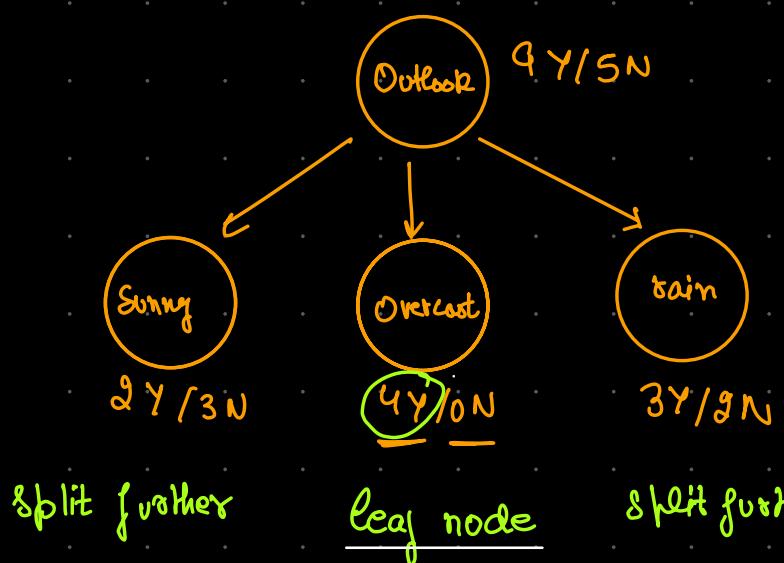
(binary)

$f_1 \ f_2 \ f_3 \ f_4 \text{ Play}$

ID3 and CART have same functionality

but CART \rightarrow binary classification

and ID3 \rightarrow not lead to binary classification



\Rightarrow We will continue to split until we get a leaf node

- 1 How will the model know that it is a leaf node?
- 2 How to select column on which we should divide first?

Our trees can be a lot diff. depending on how we start.

Purity of nodes \Rightarrow DT gonna check
whether the ^{node after} split is
pure or impure

① Entropy

② Gini Index.

For column to select first/prioritize, we
use something called Information
Gain.

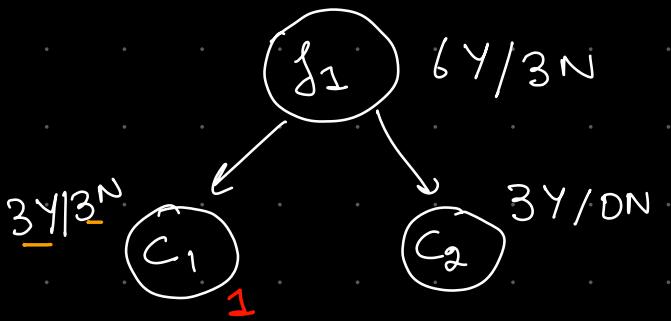
① Entropy.

$$H \text{ (of a node)} = - \sum p_i \log p_i$$

p_i = probability of getting the different categories of the node.

$$\begin{aligned} \text{Binary} &= -p_+ \log_2 p_+ - p_- \log_2 p_- \\ &\Leftrightarrow -p_0 \log_2 p_0 - p_1 \log_2 p_1 \end{aligned}$$

$$\begin{aligned} \text{Multiclass} &= -p_{c_1} \log_2 c_1 \\ &\quad - p_{c_2} \log_2 c_2 \\ &\quad - p_{c_3} \log_2 c_3 \end{aligned}$$



$$\begin{aligned}
 H(S)_{C_1} &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\
 &= -\frac{6}{6} \log_2 \frac{3}{6} \\
 &= 1
 \end{aligned}$$

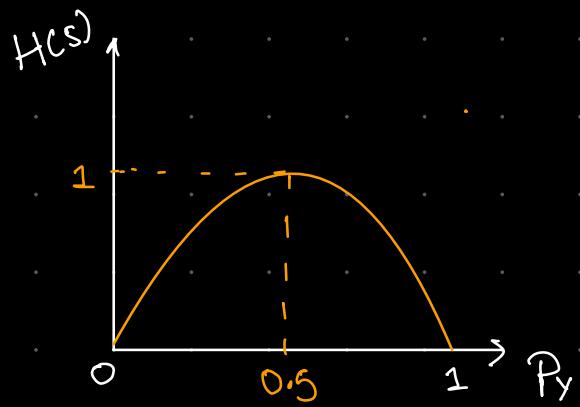
$$\begin{aligned}
 H(S)_{C_2} &= -\frac{3}{3} \log_2 \frac{3}{3} - 0 \\
 &= 0
 \end{aligned}$$



$$\begin{aligned}
 H(S)_{C_3} &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\
 &= 0.97
 \end{aligned}$$

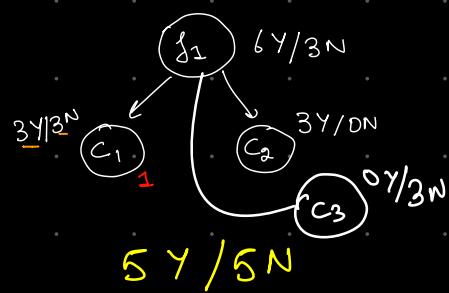
leaf node = Pure Split \Rightarrow 0 entropy

equal node = Highest entropy = 1



\Rightarrow more entropy

\Rightarrow equal class distribution



$3Y/0N$

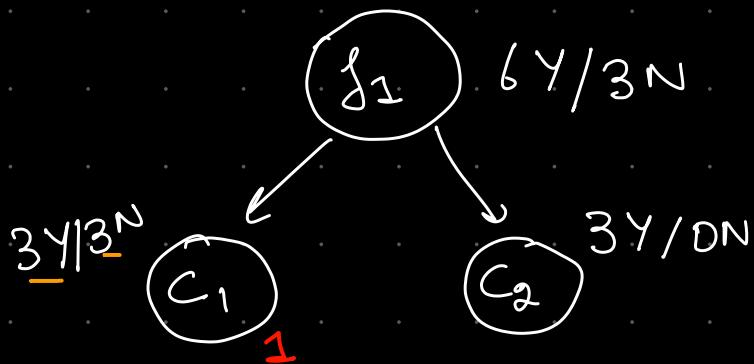
$\Rightarrow \delta_1$ impure, then we may split again.

Gini Index ($G.I.$)

Another purity check

$$G.I. = 1 - \sum_{i=1}^k (p_i)^2$$

$$= 1 - (p_+^2 + p_-^2)$$



$$G.I.C_1 = 1 - \left[\frac{3}{6}^2 + \frac{3}{6}^2 \right]$$

$$= \underline{0.5}$$

$$G.I.C_2 = 1 - \left[\frac{3}{3}^2 + \frac{0}{3}^2 \right]$$

$$= 0$$

C_4

$\frac{4Y}{8N}$ $G.I.C_4 = 4/9 = 0.44$

$8Y/2N$ $= 0.32$

$$\underline{87} / \underline{4N} = 0.44$$

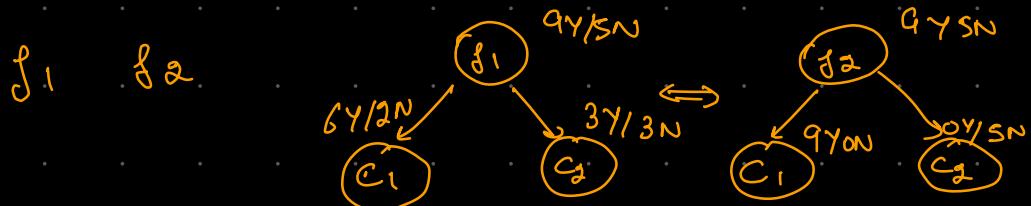


G.I. \Rightarrow longer dataset

Entropy \Rightarrow smaller dataset [log is expensive]

Information Gain

Why outlook and not wind, humidity etc?



I.G. helps in finalizing which node to split on.

It calculates the information which f will gain from f_1 & f_2 and divide on variable which give me more information.

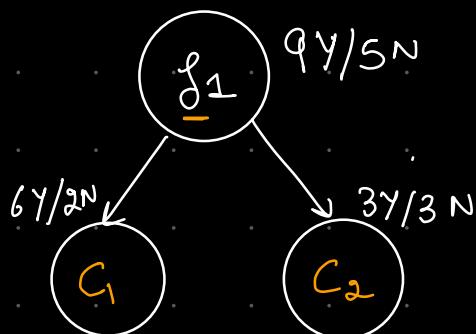
I. G₀ (feature):

$$H(\text{root}) = \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

\Downarrow

entropy of child nodes

Entropy but you can use G₀.I₀



$$H(f_1) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$

$$\frac{|S_v|}{|S|} = \frac{\text{total datapoint in Node}}{\text{total in Parent Node}}$$

$$\text{For nodes} = - \left[\frac{8}{14} \times \left[-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right] + \right.$$

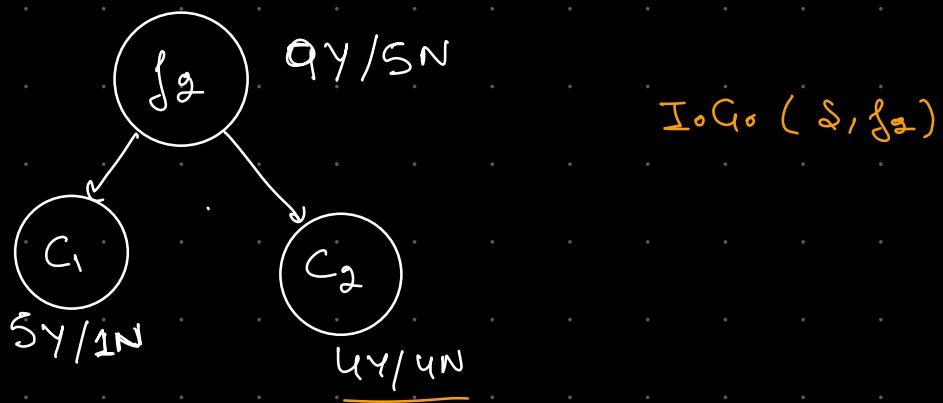
\Downarrow
 C_1

$$+ \left. \frac{6}{14} (1) \right]$$

\Downarrow
 C_2

$$I.G_0(S, f_1) = 0.045$$

* Above is for feature f_1



$$H(\text{root} | f_2) = 0.94$$

$$\Rightarrow \begin{aligned} f_2 &= 0.94 - \left[\frac{6}{14} \left[\frac{5}{6} \log \frac{5}{6} + \frac{1}{6} \log \frac{1}{6} \right] \right. \\ &\quad \left. + \frac{8}{14} [1] \right] \\ &\quad \text{C1} \\ &\quad \text{C2} \end{aligned}$$

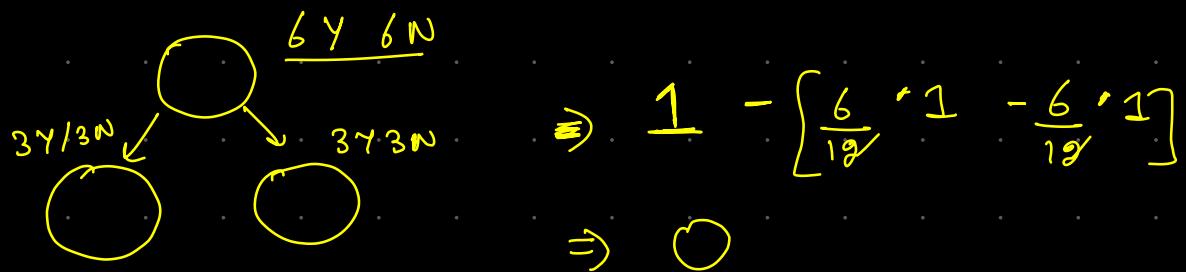
$$\Rightarrow 0.59$$

$$I.G. (S, f_1) = 0.45$$

$$I.G. (S, f_2) = 0.59$$

$\Rightarrow f_2$ split is better than f_1 split

so we are gonna first split
using f_2

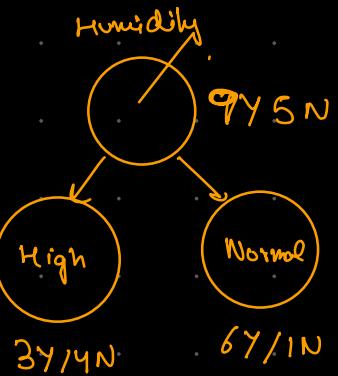


$f_1 \ f_2 \ f_3 \ f_4$

$f_1 \quad f_2$
 $f_2 \quad f_1$
 $f_3 \quad f_3$
 $f_4 \quad f_4$

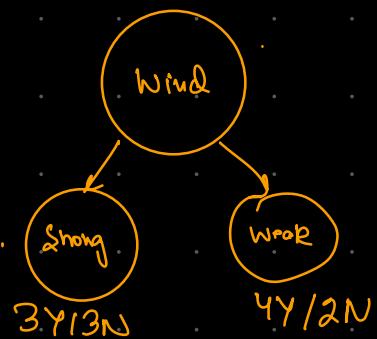
(4!)

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No



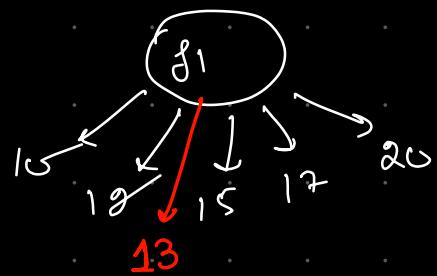
I.G. (Humidity)

2.G. (Wind)



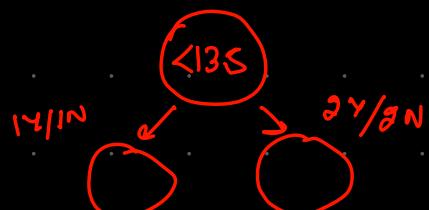
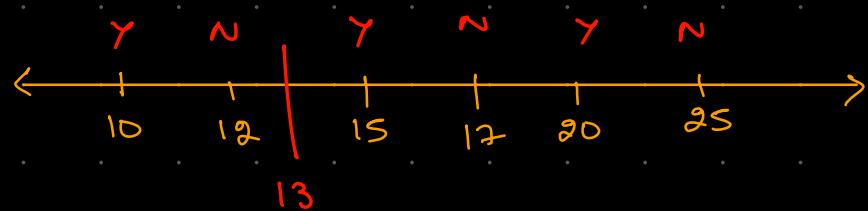
How to deal with numeric features?

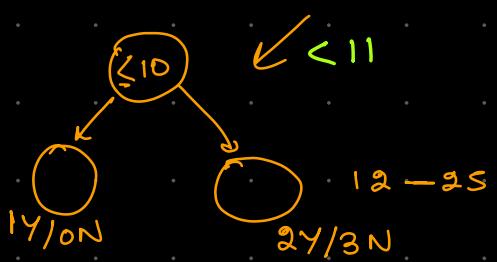
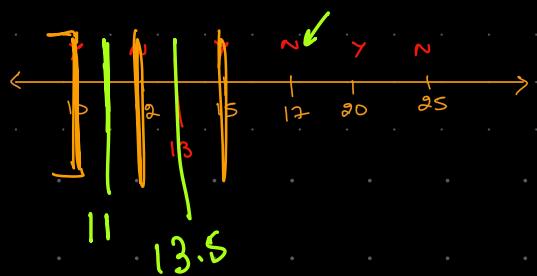
f_1	O/P
10	1
12	0
15	1
17	0
20	0
25	1



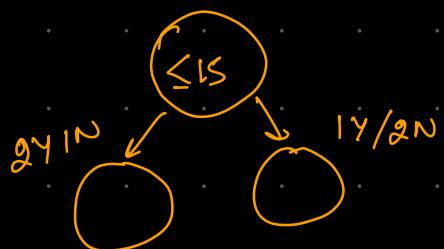
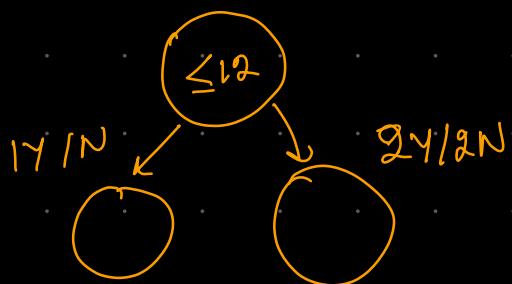
test
data 13

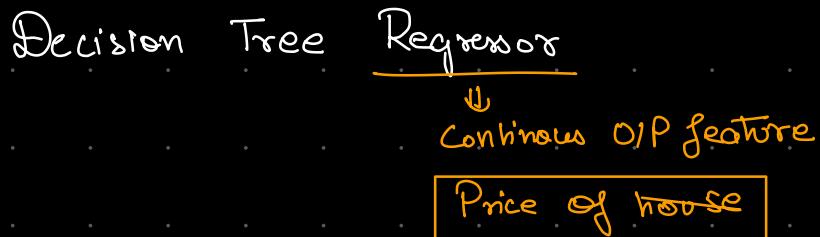
13.5





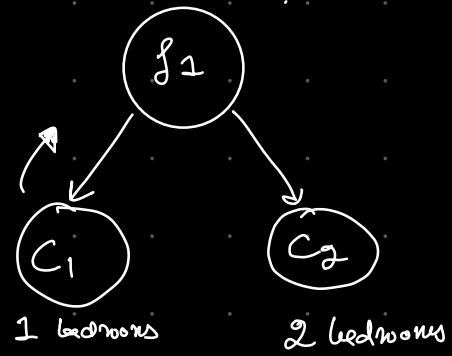
I.G.





(Bedrooms)

	f_1	f_2	<u>Price</u>
1			20
1			24 }
1			28 }
2			14
2			18
2			20



MSE / MAE