

# Decision Tree Intuition

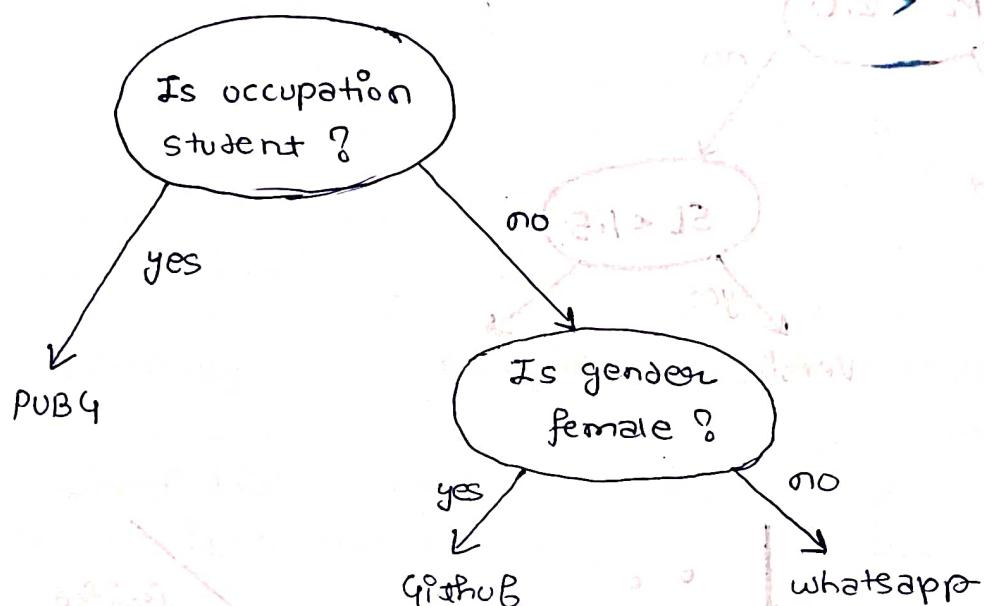
Example 1)

Gender	Occupation	Suggestion
F	Student	PUBG
F	Programmer	Github
M	Programmer	Whatsapp
F	Programmer	Github
M	Student	PUBG
M	Student	PUBG

```

if occupation == student
    print(PUBG)
else
    if gender == female
        print(Github)
    else
        print(Whatsapp)
    
```

Tree



Pseudo-code

- Begin with your training dataset, which should have some feature variables and classification or regression of.
- Determine the "best feature" in the dataset to split the data on (ie root node)

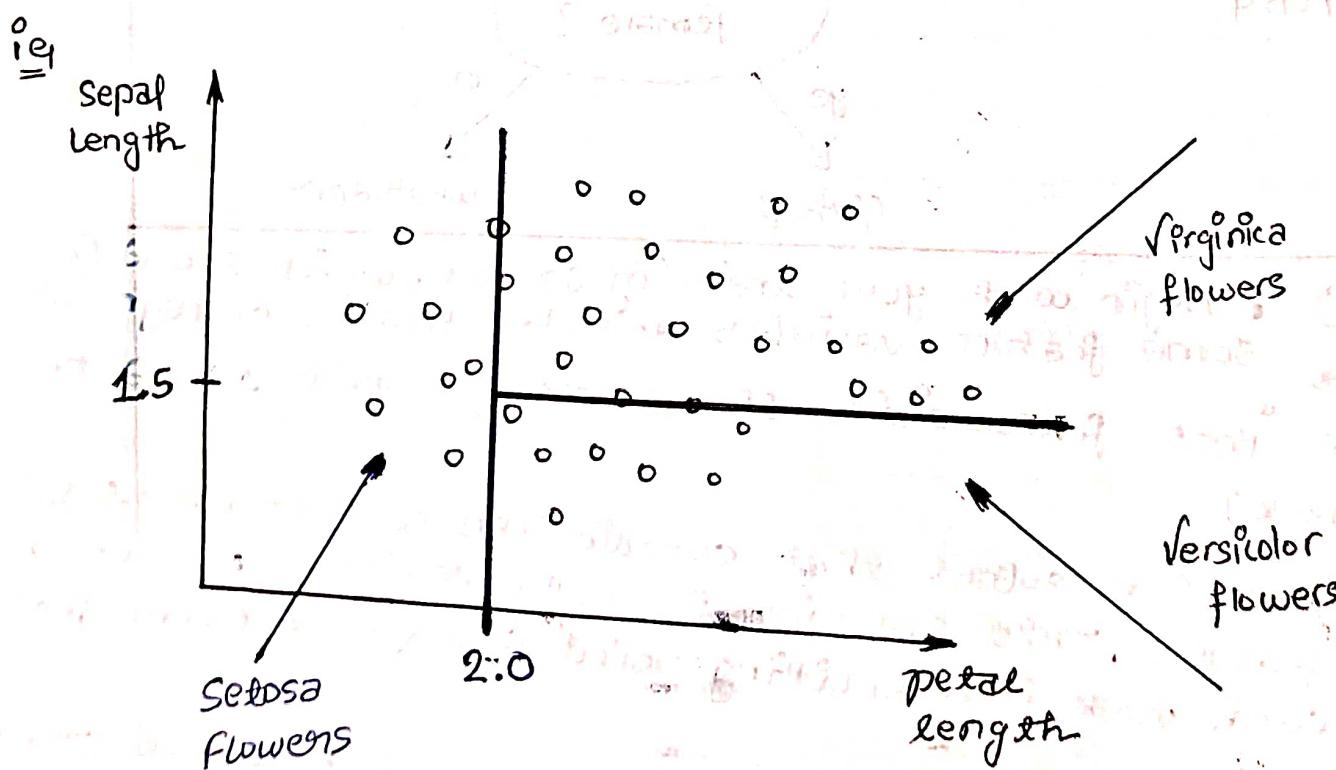
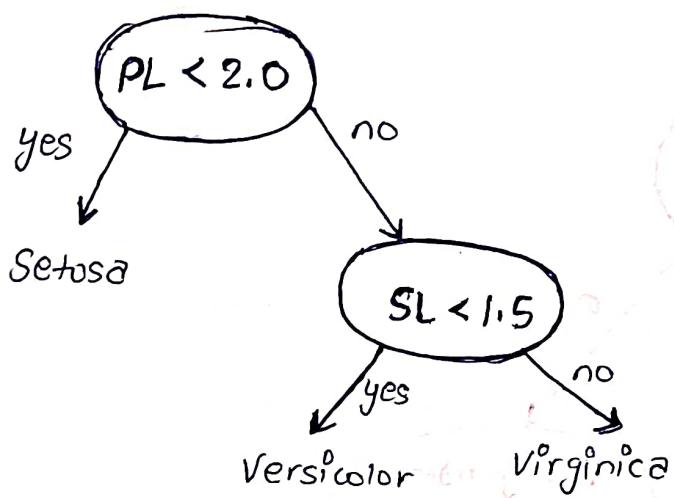
Split the data into subset that contain the correct values for the best feature. This splitting basically defines a node on the tree i.e each node is a splitting point based on a certain feature from our data

Recursively generate new tree nodes by using subset of data created from step 3

what if we have numerical data?

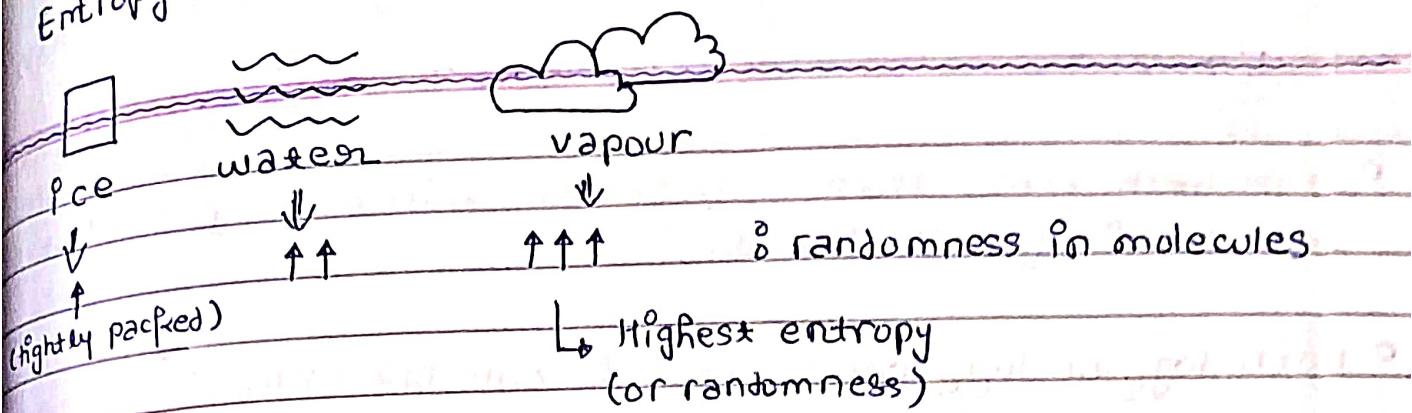
Petal Length	Sepal length	Size
1.34	0.34	Setosa
3.45	1.45	Versicolor
1.69	0.98	Setosa
2.56	1.79	Virginica
3.0	1.13	Versicolor
1.3	0.88	Setosa

- \* Working well with both either
- (i) Categorical data
- (ii) numerical data.

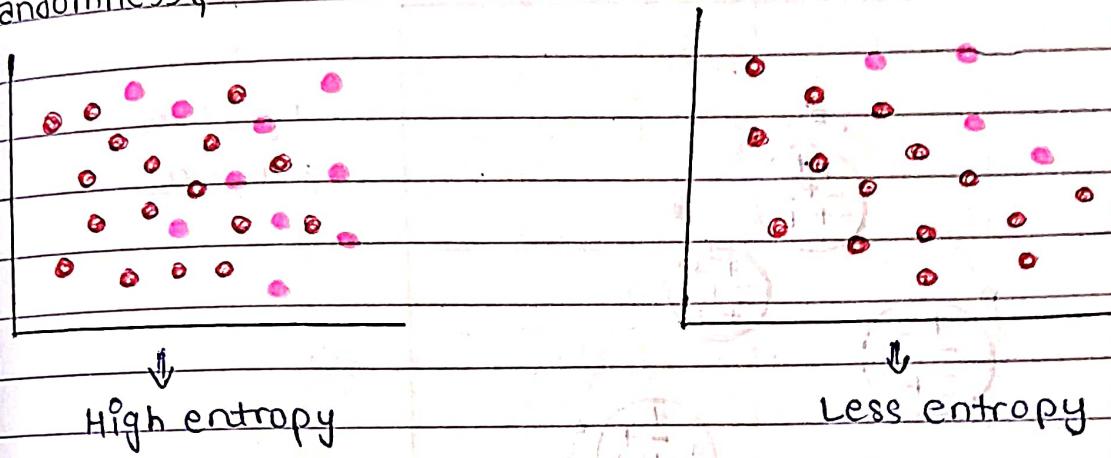


What is entropy?

Entropy is the measure of impurity / disorder



o Randomness,



More knowledge Less Entropy

Higher uncertainty Higher entropy

How to calculate Entropy?

$$E(S) = \sum_{i=1}^n -P_i \log_2 P_i$$

where ' $P_i$ ' is the frequentist probability of an element / class ' $i$ ' in our data.

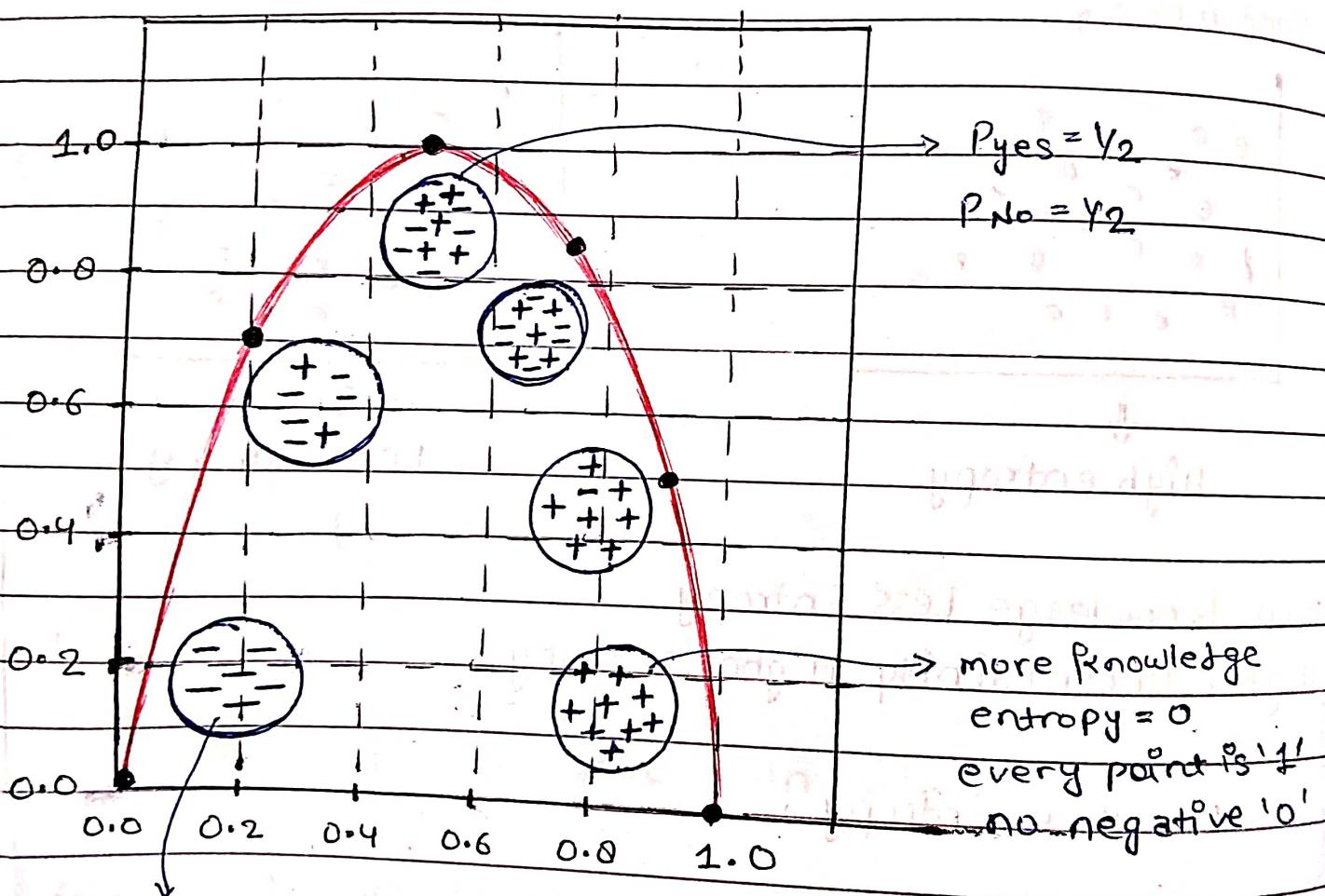
¶ Our data has two labels Yes and No.

$$E(P) = -P_{\text{Yes}} \log_2 (P_{\text{Yes}}) - P_{\text{No}} \log_2 (P_{\text{No}})$$

for '3' labels (log<sub>3</sub>) Yes & No & maybe

$$E(D) = -P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}} - P_{\text{maybe}} \log_2 P_{\text{maybe}}$$

- Summary
- More the uncertainty  $\rightarrow$  Less Entropy.
  - For a 2 class problem the min Entropy is zero (0) & max Entropy is 1.0.
  - for ~~both~~ more than 2 classes the min entropy is 0 but may be greater than 1.
  - Both  $\log_2$  or  $\log_e$  can be used to calculate entropy.



More Knowledge

Entropy = 0

every point is '0'

no positive '1'

Entropy v/s Probability

## Entropy for continuous variables

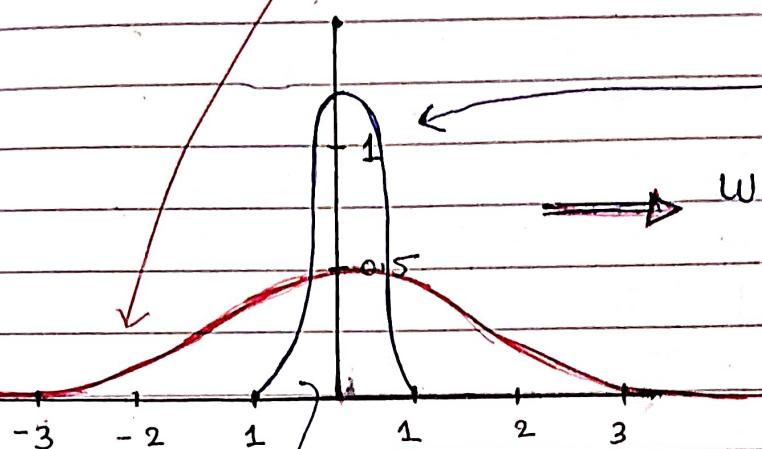
Q) Which of the below dataset have higher entropy?

Area	Built-in	Price
1200	1999	3.5
1800	2011	5.6
1400	2000	7.3
8000	0000	0000

Dataset 1

Area	Built-in	Price
2200	1989	4.6
800	2018	6.5
1100	2005	12.0
0000	0000	0000

Dataset 2



Whichever is less peaked  
↓  
High Entropy ...

Information Gain

- It's a metric used to train decision trees.

Specifically this metric measures the quality of a split.

The information gain is based on the decrease in entropy after a data-set is split on an attribute.

Constructing a decision tree is all about finding input that returns the highest information gain.

$$I = E\{ \text{parent} \} - \sum \text{Weighted Avg.} * E\{ \text{children} \}$$

## Decision Trees!

Day	Outlook	Temp	Humidity	Wind	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Strong	No
D6	Rain	Cool	Normal	Strong	Yes
D7	Overcast	Cool	Normal	Weak	No
D8	Sunny	Mild	High	Weak	Yes
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Weak	Yes
D13	Overcast	Hot	Normal	Strong	No
D14	Rain	Mild	High	Strong	No

Entropy measures homogeneity of examples,

• Entropy measures the impurity of a collection of examples  
 It depends from the distribution of the random variable  $p$ .

→  $S$  is a collection of training examples

→  $p_+$  the proportion of positive examples in  $S$

→  $p_-$  the proportion of negative examples in  $S$

- Entropy is the measure of disorder or impurity in a node

$$E = - \sum_{i=1}^n p_i \log_2 p_i$$

- node with a more variable composition considered to be more entropy than node with less variable composition
- Entropy of universal fact = 0;

High Entropy  $\rightarrow$  less homogenous / more impurity  
 less Entropy  $\rightarrow$  more homogenous / less impurity

### e.g. Attribute : Outlook

Values(Outlook) = Sunny, Overcast, Rain

$$S = [9+, 5-] \quad \text{entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \Rightarrow 0.94$$

$$S_{\text{Sunny}} \leftarrow [2+, 3-] \quad \text{entropy}(S_{\text{Sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \Rightarrow 0.971$$

$$S_{\text{Overcast}} \leftarrow [4+, 0-] \quad \text{Entropy}_{(\text{Overcast})} = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \Rightarrow 0$$

$$S_{\text{Rain}} \leftarrow [3+, 2-] \quad \text{Entropy}_{(\text{Rain})} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \Rightarrow 0.971$$

- Information gain, measure of change in Entropy

$$\text{Gain}(S, \text{outlook}) = \text{Entropy}(S) - \sum_{v \in \{\text{sunny, overcast, rain}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{outlook}) = \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{sunny}})$$

$$- \frac{4}{14} \text{Entropy}(S_{\text{overcast}}) - \frac{5}{14} \text{Entropy}(S_{\text{rain}})$$

$$\Rightarrow 0.94 - \frac{5}{14} (0.971) - \frac{4}{14} (0) - \frac{5}{14} (0.971) \Rightarrow 0.7464$$

Attribute : Temp

Values(Temp) : hot, mild, cold

$$S = [9+, 5-] \quad \text{Entropy}(S) \rightarrow -\frac{9}{14} \log_2^{\frac{9}{14}} - \frac{5}{14} \log_2^{\frac{5}{14}} \Rightarrow 0.94$$

$$\text{shot} \leftarrow [2+, 2-] \quad \text{Entropy}(\text{shot}) \rightarrow -\frac{2}{4} \log_2^{\frac{2}{4}} - \frac{2}{4} \log_2^{\frac{2}{4}} \Rightarrow 1.00$$

$$\text{mild} \leftarrow [4+, 2-] \quad \text{Entropy}(S_{\text{mild}}) \rightarrow -\frac{4}{6} \log_2^{\frac{4}{6}} - \frac{2}{6} \log_2^{\frac{2}{6}} \Rightarrow 0.9183$$

$$\text{cool} \leftarrow [3+, 1-] \quad \text{Entropy}(S_{\text{cool}}) \rightarrow -\frac{3}{4} \log_2^{\frac{3}{4}} - \frac{1}{4} \log_2^{\frac{1}{4}} \Rightarrow 0.8132$$

$$\text{Gain}(S, \text{Temp}) \Rightarrow \text{Entropy}(S) - \sum_{v \in \{\text{hot, mild, cold}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

(2)

$$\text{Gain}(S, \text{Temp}) \Rightarrow \text{Entropy}(S) = \frac{4}{15} \text{Entropy}(S_{\text{hot}}) + \frac{6}{14} \text{Entropy}(S_{\text{mid}}) - \frac{4}{14} \text{Entropy}(S_{\text{cool}})$$

$$\Rightarrow 0.94 - \frac{4}{14} (1.0) - \frac{6}{14} (0.9183) - \frac{4}{14} (0.8113) \Rightarrow 0.0289$$

Attribute : Humidity

Values (Humidity) = High, Normal

$$S = [g+, 5-] \quad \text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \Rightarrow 0.94$$

$$S_{\text{High}} \leftarrow [3+, 4-] \quad \text{Entropy}_{(S_{\text{High}})} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9052$$

$$S_{\text{Normal}} \leftarrow [6+, 1-] \quad \text{Entropy}_{(S_{\text{Normal}})} = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \Rightarrow 0.5916$$

$$\text{Gain}(S, \text{Humidity}) = \text{Entropy}(S) - \sum_{v \in \{\text{High}, \text{Normal}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{Humidity}) \Rightarrow 0.94 - \frac{7}{14} (0.9052) - \frac{7}{14} (0.5916) = 0.1516$$

$\{D_1, D_2, \dots, D_{14}\}$

[9+, 5-]

outlook

Sunny

Overcast

Rain

$\{D_1, D_2, D_8, D_9, D_{11}\}$

$\{D_3, D_7, D_{12}, D_{13}\}$

$\{D_4, D_5, D_6, D_{10}, D_{14}\}$

[4+, 0-]

[2+, 3-]

[3+, 2-]

?

Yes

?

Day	Temp	Humidity	Wind	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Attribute : Temp

Values(Temp) = Hot, Mild, Cool

$$S_{\text{Sunny}} = [2+, 3-] \quad \text{Entropy}(S_{\text{Sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \Rightarrow 0.97$$

$$S_{\text{Hot}} = [0+, 2-]$$

$$S_{\text{Mild}} = [1+, 1-]$$

$$S_{\text{Cool}} = [1+, 0-]$$

$$\text{Entropy}(S_{\text{Hot}}) \downarrow$$

$$0.0$$

$$\text{Entropy}(S_{\text{Mild}}) \downarrow$$

$$1.0$$

$$\text{Entropy}(S_{\text{Cool}}) \downarrow$$

$$0.0$$

Attribute : Wind

Values (Wind) = Strong, Weak

$$S = [9+, 5-] \quad \text{Entropy}(S) \Rightarrow 0.94$$

$$\text{Strong} \leftarrow [3+, 3-] \quad \text{Entropy}(S_{\text{strong}}) = 1.0$$

$$\text{Weak} \leftarrow [6+, 2-] \quad \text{Entropy}(S_{\text{weak}}) \Rightarrow -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \Rightarrow 0.81$$

$$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \sum_{v \in \{\text{strong}, \text{weak}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \frac{6}{14} \text{Entropy}(S_{\text{strong}}) - \frac{8}{14} \text{Entropy}(S_{\text{weak}})$$

$$\text{Gain}(S, \text{wind}) = 0.94 - \frac{6}{14} (1.0) - \frac{8}{14} (0.81) \Rightarrow 0.0478$$

$$\text{Gain}(S, \text{outlook}) \rightarrow 0.2464$$

$$\text{Gain}(S, \text{temp}) \rightarrow 0.0289$$

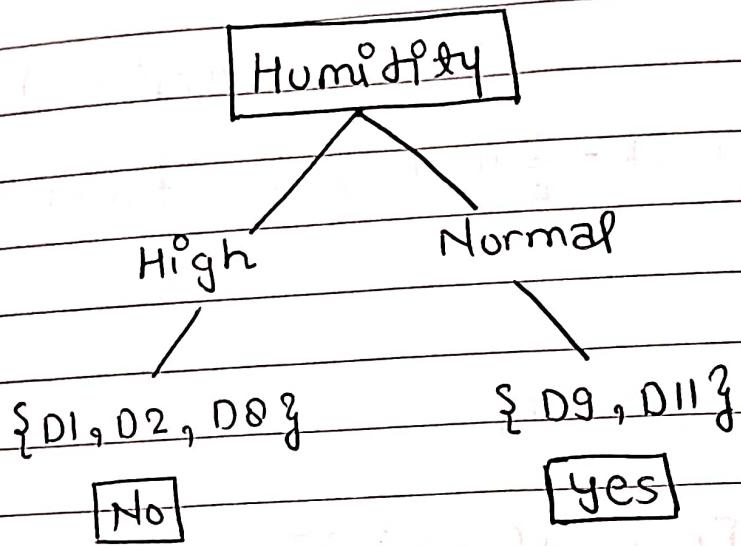
$$\text{Gain}(S, \text{humidity}) \rightarrow 0.1516$$

$$\text{Gain}(S, \text{wind}) \rightarrow 0.0478$$

$$Gain(S_{\text{sunny}}, \text{temp}) = 0.570$$

$$Gain(S_{\text{sunny}}, \text{Humidity}) = 0.97 \leftarrow \text{node}$$

$$Gain(S_{\text{sunny}}, \text{Wind}) = 0.0192$$



Day	Temp	Humidity	Wind	Play
D4	Mild	High	Weak	yes
D5	Cool	Normal	Weak	yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	yes
D14	Mild	High	Strong	No

Attribute: Humidity

Values(Humidity) = High, Normal

~~S<sub>sunny</sub> [3+, 2-]~~ Entropy  $\Rightarrow 0.97$   
~~(S<sub>sunny</sub>)~~

~~S<sub>high</sub> [1+, 1-]~~ Entropy  $(S_{\text{high}}) \Rightarrow 1.0$

~~S<sub>normal</sub> [2+, 1-]~~ Entropy  $(S_{\text{normal}}) \Rightarrow 0.9183$

$Gain(S_{\text{rain}}, \text{humidity}) \Rightarrow 0.0192$

$$Gain(S_{\text{sunny}}, \text{Temp}) = Entropy(S) - \sum_{v \in \{\text{hot}, \text{cool}, \text{mild}\}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$= Entropy(S) - \frac{2}{5} Entropy(S_{\text{hot}}) - \frac{2}{5} Entropy(S_{\text{mild}}) \\ - \frac{1}{5} Entropy(S_{\text{cool}})$$

$$Gain(S_{\text{sunny}}, \text{Temp}) \Rightarrow 0.97 - \frac{2}{5}(0.0) - \frac{2}{5}(1) - \frac{1}{5}(0.0) \Rightarrow 0.570$$

Attribute : Humidity

Values(Humidity) : High, Normal

$$S_{\text{sunny}} = [2+, 3-] \quad Entropy(S) = 0.97$$

$$S_{\text{high}} \leftarrow [0+, 3-] \quad Entropy(S_{\text{high}}) = 0.0$$

$$S_{\text{normal}} \leftarrow [2+, 0-] \quad Entropy(S_{\text{normal}}) = 0.0$$

$$Gain(S_{\text{sunny}}, \text{Humidity}) \Rightarrow 0.97$$

Attribute : Wind

Values(Wind) : Strong, Weak

$$S_{\text{sunny}} = [2+, 3-] \quad Entropy(S) = 0.97$$

$$S_{\text{strong}} \leftarrow [1+, 1-] \quad Entropy(S_{\text{strong}}) = 1.0$$

$$S_{\text{weak}} \leftarrow [1+, 2-] \quad Entropy(S_{\text{weak}}) = 0.9103$$

$$Gain(S_{\text{sunny}}, \text{Wind}) \Rightarrow 0.0192$$

Attribute : Wind

Values (wind) : strong, weak

$$S_{Rain} \leftarrow [3+, 2-] \quad \text{Entropy}(S) = 0.97$$

$$S_{Strong} \leftarrow [0+, 2-] \quad \text{Entropy}(S_{Strong}) = 0.0$$

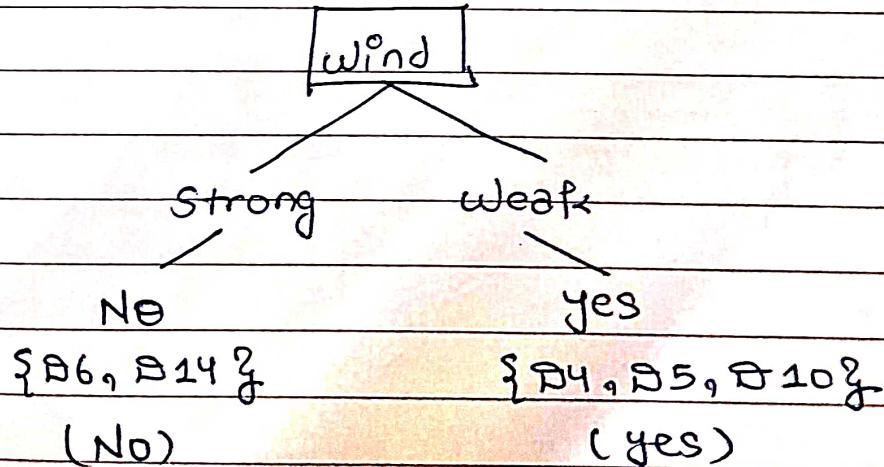
$$S_{Weak} \leftarrow [3+, 0-] \quad \text{Entropy}(S_{Weak}) = 0.0$$

$$\text{Gain}(S_{Rain}, \text{Wind}) = 0.97$$

$$\text{Gain}(S_{Rain}, \text{Temp}) = 0.0192$$

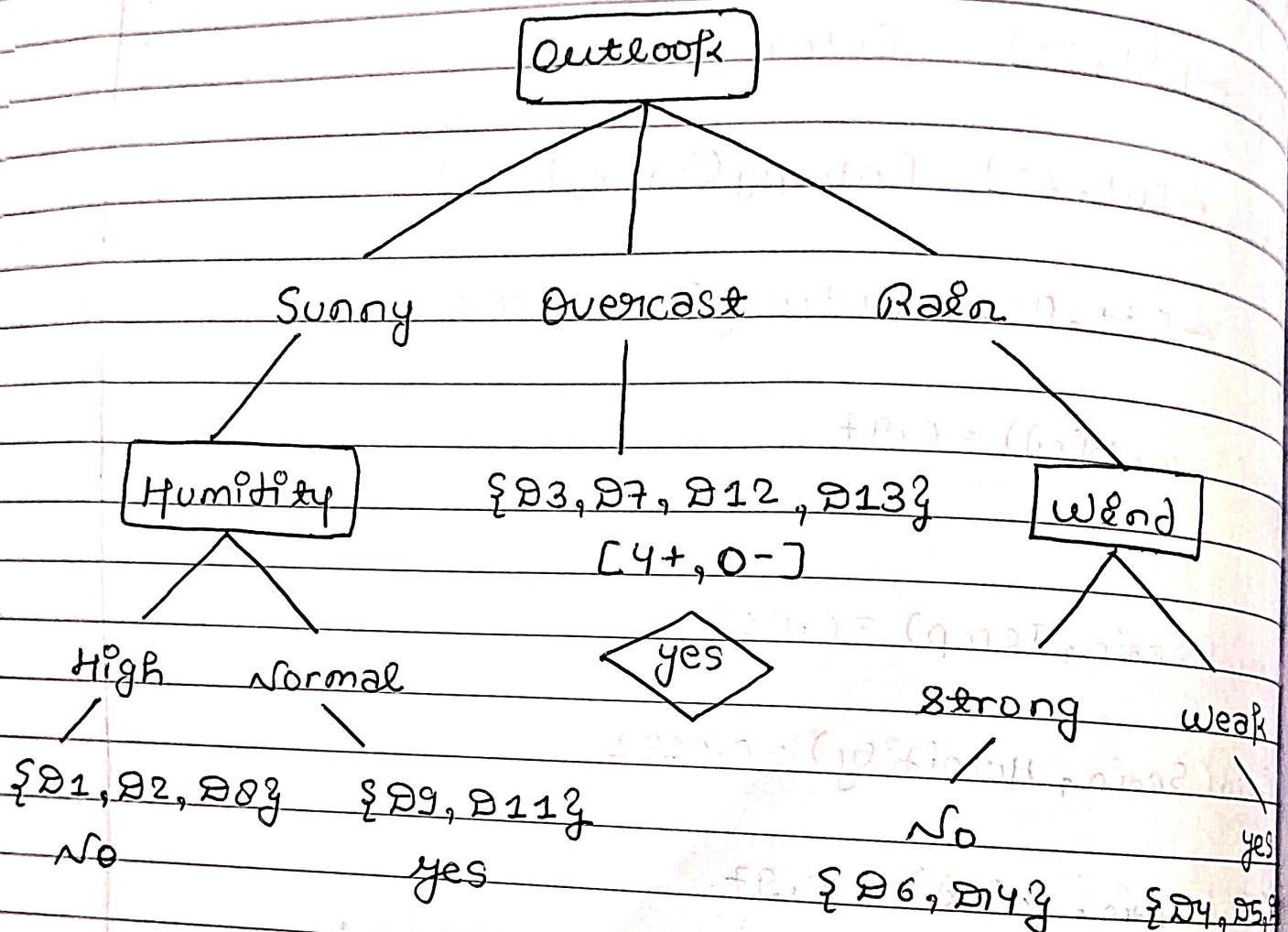
$$\text{Gain}(S_{Rain}, \text{Humidity}) = 0.0192$$

$$\text{Gain}(S_{Rain}, \text{Wind}) = 0.97$$

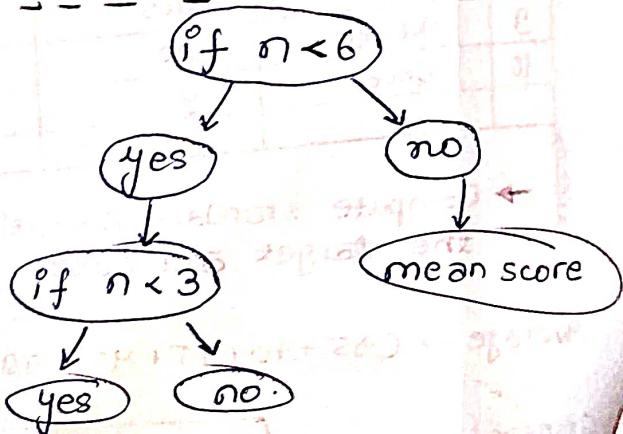
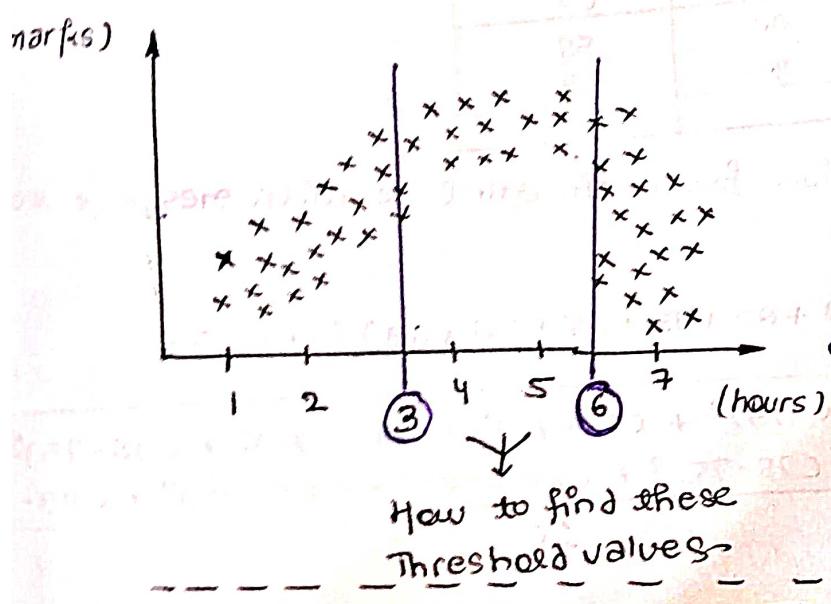
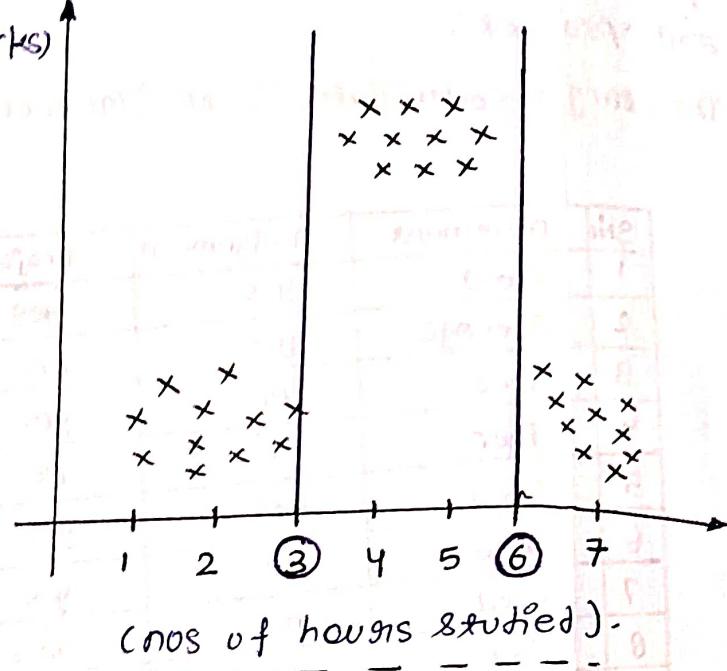
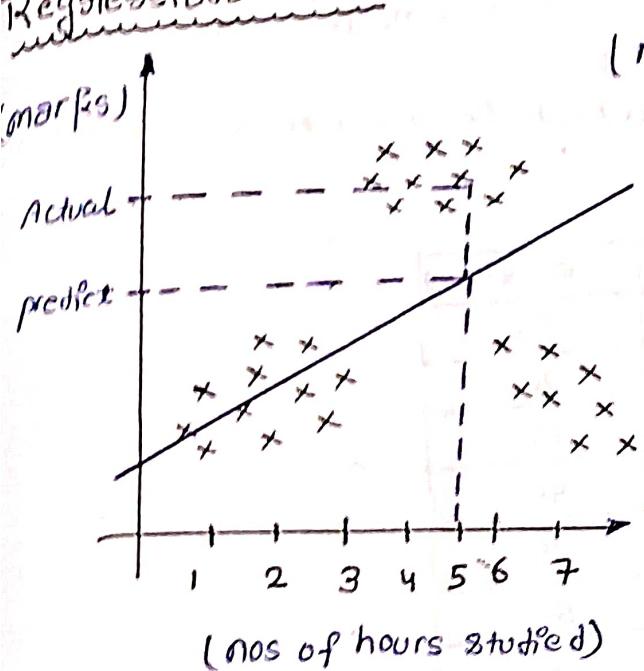


S 81, 82, 111, 814

-[9+5-]



# Regression Tree



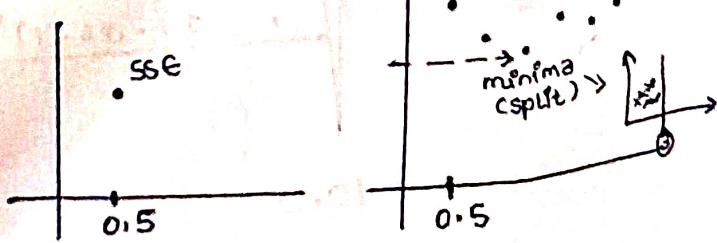
How to find these  
threshold values

- I) selecting a group, find mean
- II) perform splitting

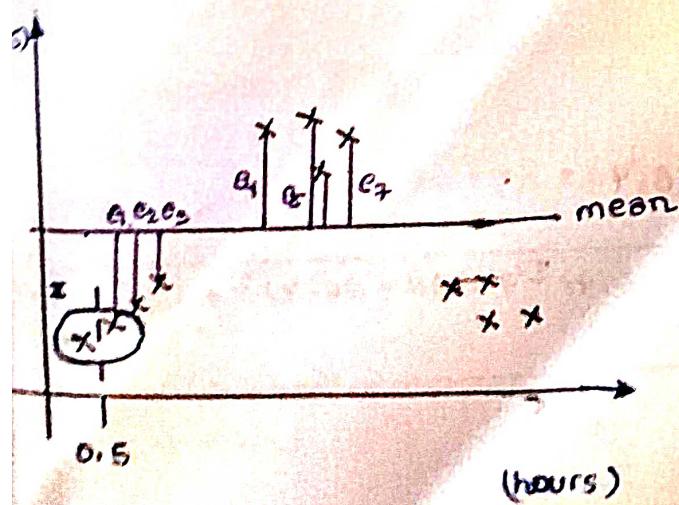
If  $R < 0.5$   
bc  $\rightarrow$  Marks      Remain mean

III)  $SSG \rightarrow e_1^2 + e_2^2 + \dots + e_n^2$

IV) plot



V) So pt by grouping others-



Construct a regression tree using the following data which consists of 10 data instances & 3 attributes 'Assessment', 'Assignment' and 'Project'.

The target attribute is the 'result' which is a continuous attribute.

SNo	Assessment	Assignment	Project	Result(%)
1	Good	yes	yes	95
2	Average	yes	no	70
3	Good	no	yes	75
4	Poor	no	no	45
5	Good	yes	yes	90
6	Average	no	yes	80
7	Good	no	no	75
8	Poor	yes	yes	65
9	Average	no	no	50
10	Good	yes	yes	89

→ Compute standard deviation for each attribute with respect to the target attribute.

$$\text{Average} \rightarrow (95 + 70 + 75 + 45 + 90 + 80 + 75 + 65 + 50 + 89) / 10 \rightarrow 75$$

Standard Deviation →

$$\sqrt{\frac{(95-75)^2 + (70-75)^2 + (75-75)^2 + (45-75)^2 + (90-75)^2 + (80-75)^2 + (75-75)^2 + (65-75)^2 + (50-75)^2 + (89-75)^2}{10}}$$

$$\Rightarrow 16.55$$

\* Let, Attribute,

Assessment = Good

$$\text{Average} \rightarrow (95 + 75 + 90 + 75 + 89) \rightarrow 86.4$$

Standard Deviation →

$$\sqrt{\frac{(95-86.4)^2 + (75-86.4)^2 + (90-86.4)^2 + (75-86.4)^2 + (89-86.4)^2}{5}}$$

$$\Rightarrow 10.9$$

Let Attribute,  
Assessment = Average

$$\text{Average} \rightarrow (70 + 80 + 50) \rightarrow 69.3$$

$$\text{std} \rightarrow \sqrt{\frac{(70 - 69.3)^2 + (80 - 69.3)^2 + (50 - 69.3)^2}{3}} \\ \rightarrow 11.01$$

Let Attribute,  
Assessment = Poor

$$\text{Average} \rightarrow (45 + 65) / 2 \rightarrow 55$$

$$\text{std} \rightarrow \sqrt{\frac{(45 - 55)^2 + (65 - 55)^2}{2}}$$

$$\rightarrow 14.14$$

Assessment	Std. Dev	Data Instances
Good	10.9	5
Average	11.01	3
Poor	14.14	2

Weighted Standard Deviation for Assessment

$$\rightarrow (\frac{5}{10}) \times 10.9 + (\frac{3}{10}) \times 11.01 + (\frac{2}{10}) \times 14.14 \rightarrow 11.58$$

Standard deviation reduction for Assessment

$$\rightarrow 16.55 - 11.58 \rightarrow 4.97$$

Let Attribute,  
Assignment = Yes

Assignment = Yes

$$\text{Average} \rightarrow (95 + 70 + 90 + 65 + 89) / 5 \rightarrow 83.4$$

$$\text{Std. Dev} \rightarrow 14.98$$

Let Attribute,

Assignment = No

$$\text{Average} \rightarrow (75 + 75 + 80 + 75 + 50) / 5 \rightarrow 66.6$$

$$\text{Std. Dev} \rightarrow 24.7$$

construct a regression

Standard Deviation for Assignment.

	Standard Deviation	Data instance
Assignment	std.dev.	5
Yes	14.98	5
No	14.7	

15

Good  
Poor



$$\begin{aligned} \text{Weighted standard deviation for assignment} \\ = \left(\frac{5}{10}\right) \times 14.98 + \left(\frac{5}{10}\right) \times 14.07 \rightarrow 14.84 \\ = \left(\frac{5}{10}\right) \times 14.98 + \left(\frac{5}{10}\right) \times 14.07 \rightarrow 14.84 \\ \text{Std. dev. prediction for assignment} \rightarrow 10.71. \end{aligned}$$

Std. dev. prediction for assignment

	Standard Deviation	Data instance
Project	std.dev.	6
Project	12.6	4
Project	13.39	4
Project	No	

$$\begin{aligned} \text{Weighted std. dev. for Project} \\ = \left(\frac{6}{10}\right) \times 2.6 + \left(\frac{4}{10}\right) \times 13.39 \\ \rightarrow 12.92 \end{aligned}$$

Std. dev. prediction for assessment

$$16.55 - 12.92 \rightarrow 3.63$$

Attribute	Standard Dev. Prediction
Assessment	10.71
Assignment	10.71
Project	3.63

make request  
as note

SN	Assessment	Assignment	Project	Result(%)
4	Poor	No	No	45
8	Poor	Yes	Yes	65
2	Average	Yes	No	70
6	Average	No	Yes	80
9	Average	No	No	50

Average

SN	Assessment	Assignment	Project	Result(%)
2	Average	Yes	No	70
6	Average	No	Yes	80
9	Average	No	No	50

SN	Assessment	Assignment	Project	Result(%)
4	Poor	No	No	45
8	Poor	Yes	Yes	65
2	Average	Yes	No	70

## ★ function to calculate gini-impurity

# example labels

labels = [0, 1, 1, 0, 1, 0].

def gini-impurity(labels):

total\_samples = len(labels)

label\_counts = np.unique(labels, → [3, 3])

return counts = True) [1].

→ [0.5 0.5]

probabilities = ~~label\_counts~~

label\_counts / total\_samples

gini = 1 - np.sum(probabilities \*\* 2)

Feature	Root	Branch A	Branch B	Branch C
1	0.5	0.5	0.5	0.5
2	0.5	0.5	0.5	0.5
3	0.5	0.5	0.5	0.5

# Classification & Regression Trees!

CGPA	Interactive	Practical Knowledge	Common skills	Job offers
$\geq 9$	yes	Very Good	Good	yes
$\geq 8$	no	Good	Moderate	yes
$\geq 9$	no	Average	Poor	no
$< 8$	no	Average	Good	no
$\geq 8$	yes	Good	Moderate	yes
$\geq 9$	yes	Good	Moderate	yes
$< 8$	yes	Good	Poor	no
$\geq 9$	no	Very good	Good	yes
$\geq 8$	yes	Good	Good	yes
$\geq 8$	yes	Average	Good	yes

Step 1: Calculate the Gini Index for the dataset

The target attribute "job offer" has 7 instances as 'yes' & 3 instances as 'no'.

$$\text{Gini-Index}(T) = 1 - \sum_{i=1}^m p_i^2$$

$$\begin{aligned}\text{Gini-Index}(T) &= 1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 \\ &\Rightarrow 1 - 0.49 - 0.09 \\ &\Rightarrow 1 - 0.58\end{aligned}$$

$$\text{Gini-Index}(T) \Rightarrow 0.42$$

Step 2: Compute Gini-Index for each of the attribute and each of the subset in the attribute

Categories of CGPA

CGPA	Job Offer = "Yes"	Job Offer = "No"
$\geq 9$	3	1
$\geq 8$	4	0
$< 8$	0	2

possible subsets

$2^3$

↓

$\{ \}, \{ \geq 9 \}, \{ \geq 8 \}, \{ \geq 9, \geq 8 \}, \{ \geq 9, > 8 \}, \{ \geq 8, > 8 \}$   
 $\{ \geq 9, > 8 \} \text{ and } \{ \geq 8, > 8 \}$

$$\circ \text{Gini-Index}(T) = 1 - \sum_{i=1}^m p_i^2$$

$$\circ \text{Gini-Index}(T, A) = \frac{|S_1|}{|T|} \text{Gini}(S_1) + \frac{|S_2|}{|T|} \text{Gini}(S_2)$$

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ \geq 9, \geq 8 \}) \Rightarrow 1 - \left(\frac{7}{8}\right)^2 - \left(\frac{1}{8}\right)^2 \\ \Rightarrow 1 - 0.7006 \Rightarrow 0.2194$$

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ < 8 \}) \Rightarrow 1 - (0/2)^2 - (2/2)^2 \\ \Rightarrow 1 - 1 \Rightarrow 0$$

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ C \geq 9, > 8 \}, < 8 \}) \quad (\text{CT}) \text{ robust}$$

$$\rightarrow \left(\frac{8}{10}\right) 0.2194 + \left(\frac{2}{10}\right) 0 \Rightarrow 0.17552$$

↓ similarly

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ \geq 9, < 8 \}) \Rightarrow 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2$$

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ \geq 8 \}) \Rightarrow 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 \\ \Rightarrow 1 - 1 \Rightarrow 0$$

$$* \text{Gini-Index}(T, \text{CGPA} \in \{ \geq 9, < 8 \}, g \geq 8 \})$$

$$\Rightarrow \left(\frac{6}{10}\right) \times 0.5 + \left(\frac{4}{10}\right) \times 0 \Rightarrow 0.3 -$$

$$Gini\text{-Index}(T, CGPA \in \{>=8, <8\}) \Rightarrow 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 \quad (1)$$

$$\Rightarrow 1 - 0.55 \Rightarrow 0.445$$

$$Gini\text{-Index}(T, CGPA \in \{>=9, <9\}) \Rightarrow 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2$$

$$\Rightarrow 1 - 0.625 \Rightarrow 0.375$$

\*  $Gini\text{-Index}(T, CGPA \in \{>=8, <8\}, >=9) \Rightarrow$

$$\Rightarrow \left(\frac{6}{10}\right)(0.445) + \left(\frac{4}{10}\right)(0.375) \Rightarrow 0.417$$

Step III) Choose the best splitting subset which has ~~minimum~~ maximum  $Gini\text{-Index}$  for an attribute

The subset

$$\{>=9, >=8, <8\}$$

has lowest ~~highest~~  $Gini\text{-Index}$  value

as 0.1755 is

chosen as Best splitting subset.

$Gini\text{-Index}$  of CGPA

Subsets	$Gini\text{-Index}$	
$>=9, >=8$	$<8$	0.1755
$>=9, <8$	$>=8$	0.3
$>=8, <8$	$>=9$	0.417

Step IV) Compute  $\Delta Gini$  or Best splitting subset of that attribute

$$\Delta Gini(CGPA) = Gini(T) - Gini(T, CGPA)$$

$$\rightarrow 0.42 - 0.1755 \rightarrow 0.2445$$

Repeat the same process for the remaining attributes in the dataset such as for interaction, practical know & communication skills

## Categories for practical knowledge

Practical knowledge	Job offer = Yes	Job offer = No
Very good	2	0
good	4	1
Average	1	2

- Gini-Index CT, Prac. know.  $\in \{\text{Very good, good}\}$

$$\Rightarrow 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \Rightarrow 0.7544$$

- Gini-Index CT, Prac. know.  $\in \{\text{Average}\}$   $= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2$

$$\Rightarrow 1 - 0.55 \Rightarrow 0.445$$

- Gini-Index CT, Prac. know.  $\in \{\text{Very good, good, Avg}\}$

$$\rightarrow \left(\frac{7}{10}\right)^2 \times 0.2456 + \left(\frac{3}{10}\right) \times 0.445 \Rightarrow 0.3054$$

- Gini-Index CT, Prac. know.  $\in \{\text{Very good, Average}\}$

$$\Rightarrow 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \Rightarrow 1 - 0.52 \Rightarrow 0.48$$

- Gini-Index CT, Prac. know.  $\in \{\text{Good}\}$

$$\Rightarrow 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2 \Rightarrow 1 - 0.68 \Rightarrow 0.32$$

- Gini-Index CT, Prac. know.  $\in \{\text{Very good, Avg, Good}\}$

$$\rightarrow \left(\frac{5}{10}\right) \times 0.48 + \left(\frac{5}{10}\right) \times 0.32 \Rightarrow 0.40$$

- Gini-Index CT, Prac. know.  $\in \{\text{Good, Avg}\}$

$$\Rightarrow 1 - \left(\frac{5}{10}\right)^2 - \left(\frac{3}{10}\right)^2 \Rightarrow 1 - 0.5312 \Rightarrow 0.4688$$

- Gini-Index CT, Prac. know.  $\in \{\text{Good}\}$

$$\Rightarrow 1 - \left(\frac{2}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \Rightarrow 1 - 1 \Rightarrow 0$$

- Gini-Index CT, Prac. know.  $\in \{\text{Good, Average, Very good}\}$

$$\rightarrow \left(\frac{8}{10}\right) \times 0.4688 + \left(\frac{2}{10}\right) \times 0 \Rightarrow 0.3750$$

## Gini-Index for practical knowledge

Subsets		Gini-Index
(Very good, Good)	Average	0.3054
(Very good, Average)	Good	0.40
(Good, Average)	Very Good	0.3750

$$\Delta \text{Gini}(\text{practical knowl}) = \text{Gini}(T) - \text{Gini}(T, \text{practical knowl}) \\ \Rightarrow 0.42 - 0.3054 \Rightarrow 0.1146.$$

Categories for communication skills

Communication skills	Total	No
Good	4	1
Moderate	3	0
Poor	0	2

Gini-Index for subset commnno skills

Subsets		Gini-Index
{Good, Moderate}	Poor	0.1755
{Good, Poor}	Moderate	0.3429
{Moderate, Poor}	Good	0.40

$$\Delta \text{Gini}(\text{Commnno skills}) = \text{Gini}(T) - \text{Gini}(T, \text{Commnno skills}) \\ \Rightarrow 0.42 - 0.1755 \\ \Rightarrow 0.2445$$

Gini Index &  $\Delta \text{Gini}$  for all attributes

Attributes	Gini-index	$\Delta \text{Gini}$
CGPA	0.1755	0.2445
Interactiveness	0.368	0.1052
Practical knowl	0.3054	0.1146
Communication skills	0.1755	0.2445

Subset CGPA = {C >= 9, > 0} , {< 0}

is the best splitting subset

from table,

CQPA

< 8

Job Offer = No

sno	CQPA	Interactivity	Prac know	commo skills	Job
1	$\geq 9$	yes	Very good	Good	yes
2	$\geq 8$	no	good	Moder	yes
3	$\geq 9$	no	Average	Poor	no
5	$\geq 8$	yes	good	Moder	yes
6	$\geq 9$	yes	good	Moder	yes
8	$\geq 9$	no	Very	Good	yes
9	$\geq 8$	yes	good	Good	yes
10	$\geq 8$	yes	Average	Good	yes

Age	Gender	Education	Experience
20-25	Male	High School	1-2 years
26-30	Female	College	3-5 years
31-35	Male	Post Grad	6-10 years
36-40	Female	Post Grad	11+ years

related variables are not significant

Age Group	Gender	Education
20-25	Male	High School
26-30	Female	College
31-35	Male	Post Grad
36-40	Female	Post Grad

1. Data analysis CQPA - CT PA = (left side minus right side)  
Diff = 20.5 - 20.5 = 0  
Diff = 0

Age Group	Gender	Education	Experience
20-25	Male	High School	1-2 years
26-30	Female	College	3-5 years
31-35	Male	Post Grad	6-10 years
36-40	Female	Post Grad	11+ years