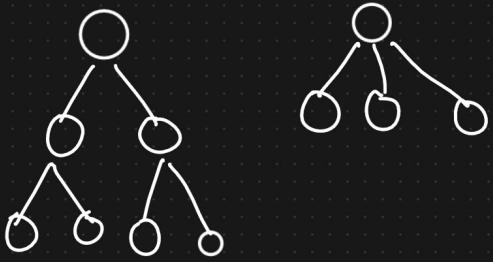


# Decision Tree Classifier

Decision Tree Classifier

→ ID3  
→ CART ✓



a) Entropy and Gini Index → Purity Split

b) Information Gain → features to select for

DT construction

$age = 14$

if ( $age \leq 15$ ):

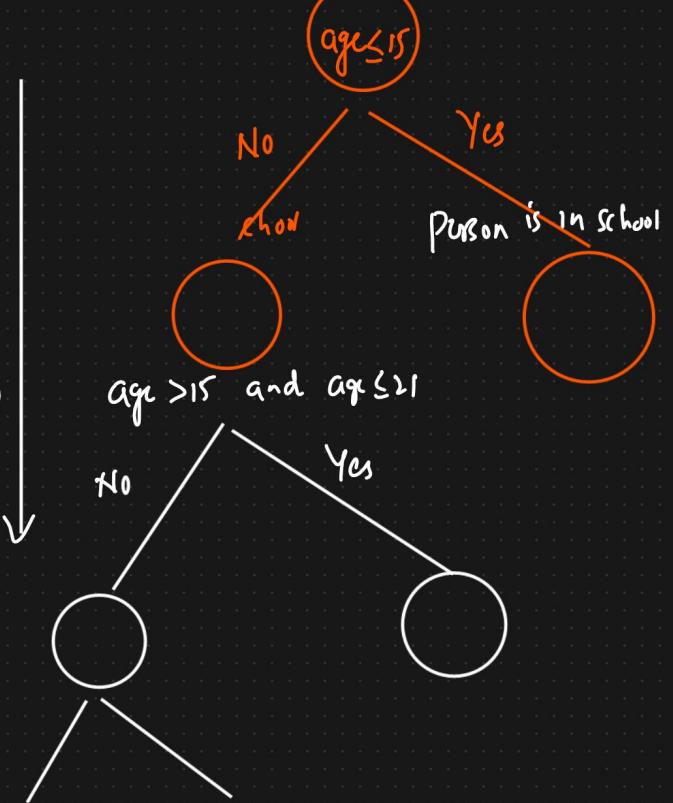
Print ("The person is in School")

elif ( $age > 15$  and  $age \leq 21$ ):

Print ("The person may be college")

else:

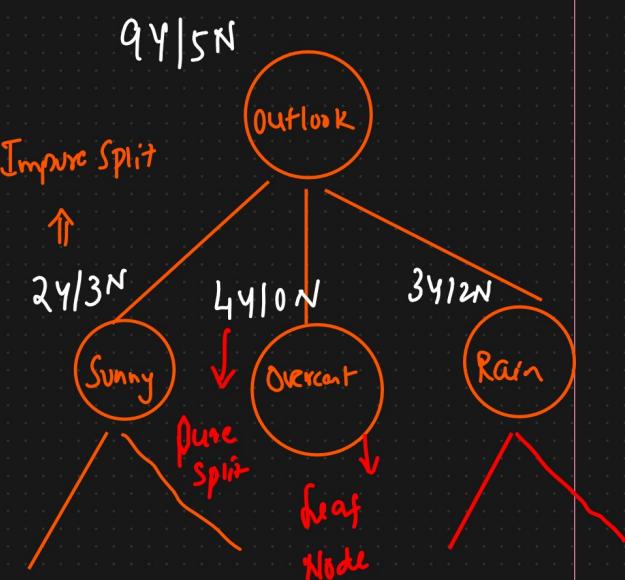
Print ("The person has passed")



## Data set

## Binary Classification

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny ✓	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast ✓	Hot	High	Weak	Yes
4	Rain ✓	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



① Punty → Pure or Impure Split

↳ Entropy  
↳ Gini Impurity }

② What feature you need Select for  
Splitting → Information Gain }

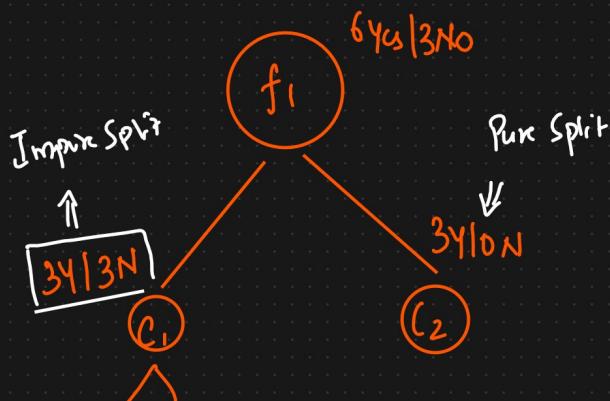
1

0

{Binary Classification}

1) Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

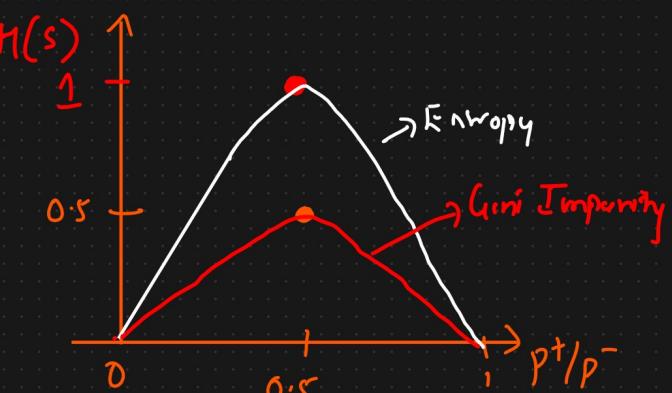


$$H(C_1) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

2) Gini Impurity

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



$\Rightarrow$  Impure Split

$$H(C_2) = -\frac{3}{3} \log_2 \frac{3}{3} - 0 \log_2 0$$

$\Rightarrow -1 \log_2 1 \Rightarrow 0 \Rightarrow$  Pure Split

(2) Min Impurity

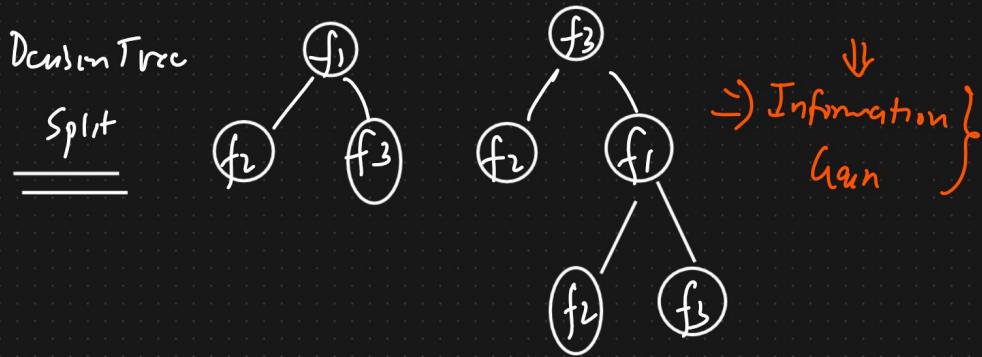
$$\begin{aligned} G \cdot I &= 1 - \sum_{i=1}^n (P_i)^2 \\ &= 1 - \left( (P_1)^2 + (P_2)^2 \right) \\ &= 1 - \left( \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right) \\ &= 0 \Rightarrow \text{Pure Split} \\ &\equiv 0.5 \Rightarrow \text{Impure Split} \end{aligned}$$

BY ION

$$\begin{aligned} &= 1 - \left( \left(\frac{3}{5}\right)^2 \right) \\ &= 1 - 1 \end{aligned}$$

$\equiv 0 \Rightarrow \text{Pure Split}$

$f_1 \quad f_2 \quad f_3$



Information Gain

$\text{Gain}(S, f_1) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$

Root Node

Entropy of the root node

$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$

$H(S) = 1 - \left( \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) = 0.971$

$H(S_1) = 1 - \left( \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = 1 - 1 = 0$

$H(S_2) = 1 - \left( \frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) = 0.918$

$H(S_3) = 1 - \left( \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = 1 - 1 = 0$

$\text{Gain}(S, f_1) = 0.971 - \frac{3}{5} \cdot 0 - \frac{2}{5} \cdot 0.918 = 0.971 - 0.367 = 0.604$

$f_1 \quad f_2 \quad f_3 \quad O/P$

$f_1 \rightarrow C_1$

$f_2 \rightarrow C_2$

$f_3 \rightarrow C_3$

$C_1 \rightarrow C_4$

$C_2 \rightarrow C_3$

$C_3 \rightarrow C_4$

$\leftrightarrow$

$0.604 = \text{Gain}$

$0.971 = H(S)$

$0.918 = H(S_2)$

$0.367 = H(S_1)$

$0 = H(S_3)$

$0 = H(S_4)$

$0.971 - 0.367 = 0.604$

$\text{Impure split}$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \quad H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

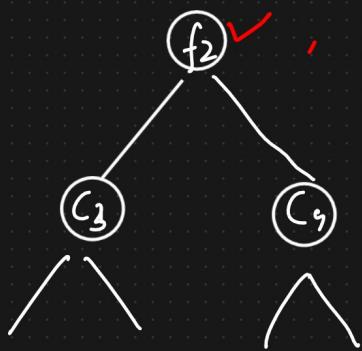
$\approx 0.94$

$$H(C_1) \approx 0.81$$

$$H(C_2) = 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\boxed{\text{Gain}(S, f_1) = 0.049}$$



$$\boxed{\text{Gain}(S, f_2) = 0.051} > \boxed{\text{Gain}(S, f_1) = 0.049}$$

Information  $\frac{\text{Gain}}{\text{Gain}}$  is Basically calculated.

Entropy  $\checkmark$  Vs Gini Impurity  $\checkmark$

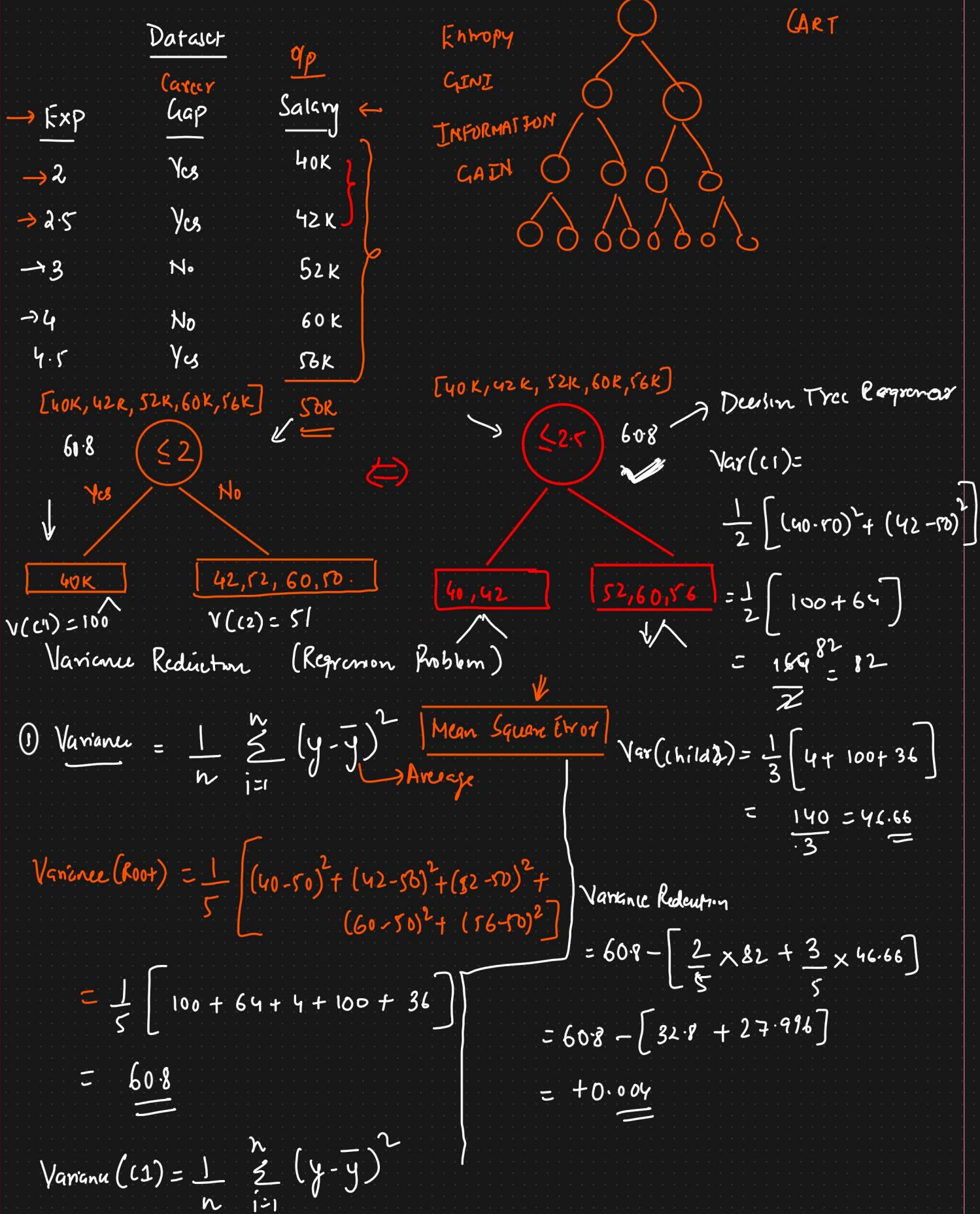
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_- \quad G.I. = 1 - \sum_{i=1}^n (P_i)^2 \Rightarrow$$

O/P = 3 categories

$$H(S) = -P_{C_1} \log_2 P_{C_1} - P_{C_2} \log_2 P_{C_2} - P_{C_3} \log_2 P_{C_3}$$

Whenever dataset is small  $\rightarrow$  Entropy  
large  $\rightarrow$  Gini Impurity

# Decision Tree Regression



$$= \frac{1}{n} (y_0 - \bar{y})^2$$

$$= 100$$

$$\text{Variance (C2)} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \frac{1}{4} \left[ (42 - 50)^2 + (52 - 50)^2 + (60 - 50)^2 + (56 - 50)^2 \right]$$

$$= \frac{1}{4} [64 + 4 + 100 + 36]$$

$$= 51$$

Variance Reduction  $\downarrow$

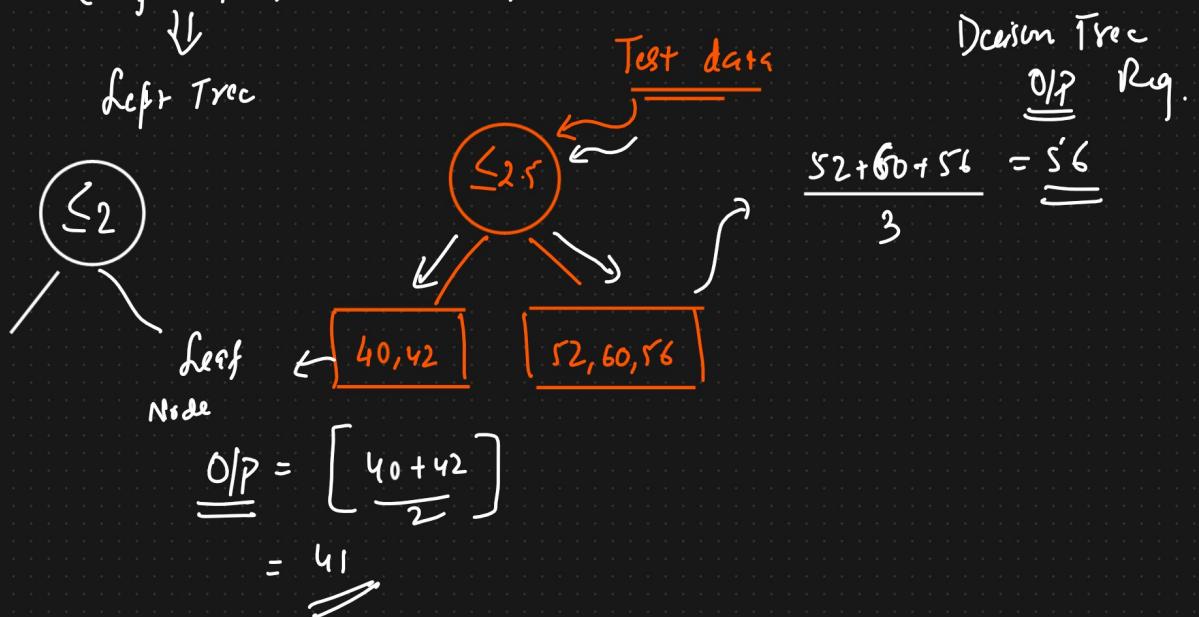
$$= \text{Var}(\text{Root}) - \sum w_i \text{Var}(\text{child})$$

$$= 60.8 - \left[ \frac{1}{8} \times 20 + \frac{4}{5} \times 51 \right]$$

$$= 60.8 - 20 - 40.8$$

Variance Reduction = 0

0 0.004  
 $\text{Variance Reduction (Left Split)} < \text{VR (Right Split)}$



# Decision Tree Split for Numerical Feature.

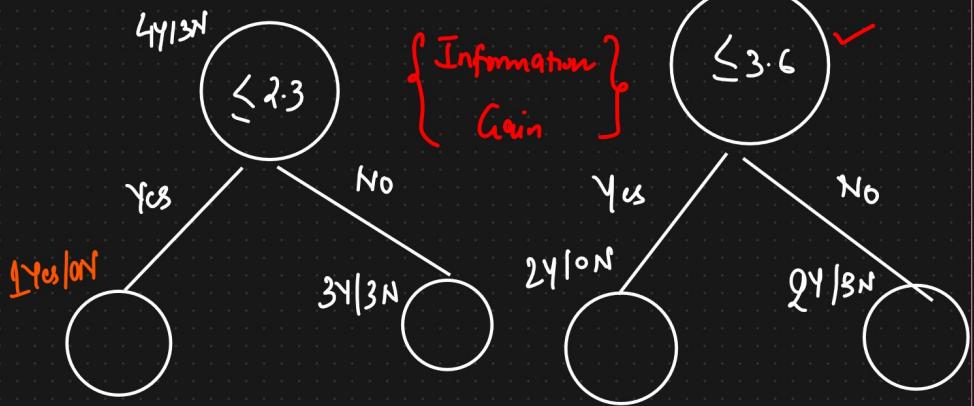
Day	Outlook	Temperature	Humidity	Wind	O/P	Play Tennis
1	Sunny	Hot	High	Weak	No	
2	Sunny	Hot	High	Strong	No	
3	Overcast	Hot	High	Weak	Yes	
4	Rain	Mild	High	Weak	Yes	
5	Rain	Cool	Normal	Weak	Yes	
6	Rain	Cool	Normal	Strong	No	
7	Overcast	Cool	Normal	Strong	Yes	
8	Sunny	Mild	High	Weak	No	
9	Sunny	Cool	Normal	Weak	Yes	
10	Rain	Mild	Normal	Weak	Yes	
11	Sunny	Mild	Normal	Strong	Yes	
12	Overcast	Mild	High	Strong	Yes	
13	Overcast	Hot	Normal	Weak	Yes	
14	Rain	Mild	High	Strong	No	

① Sort the feature value

f1	O/P
2.3	Yes
3.6	Yes
4	No
5.2	No
6.7	Yes
8.9	No
10.5	Yes

① Threshold = 2.3

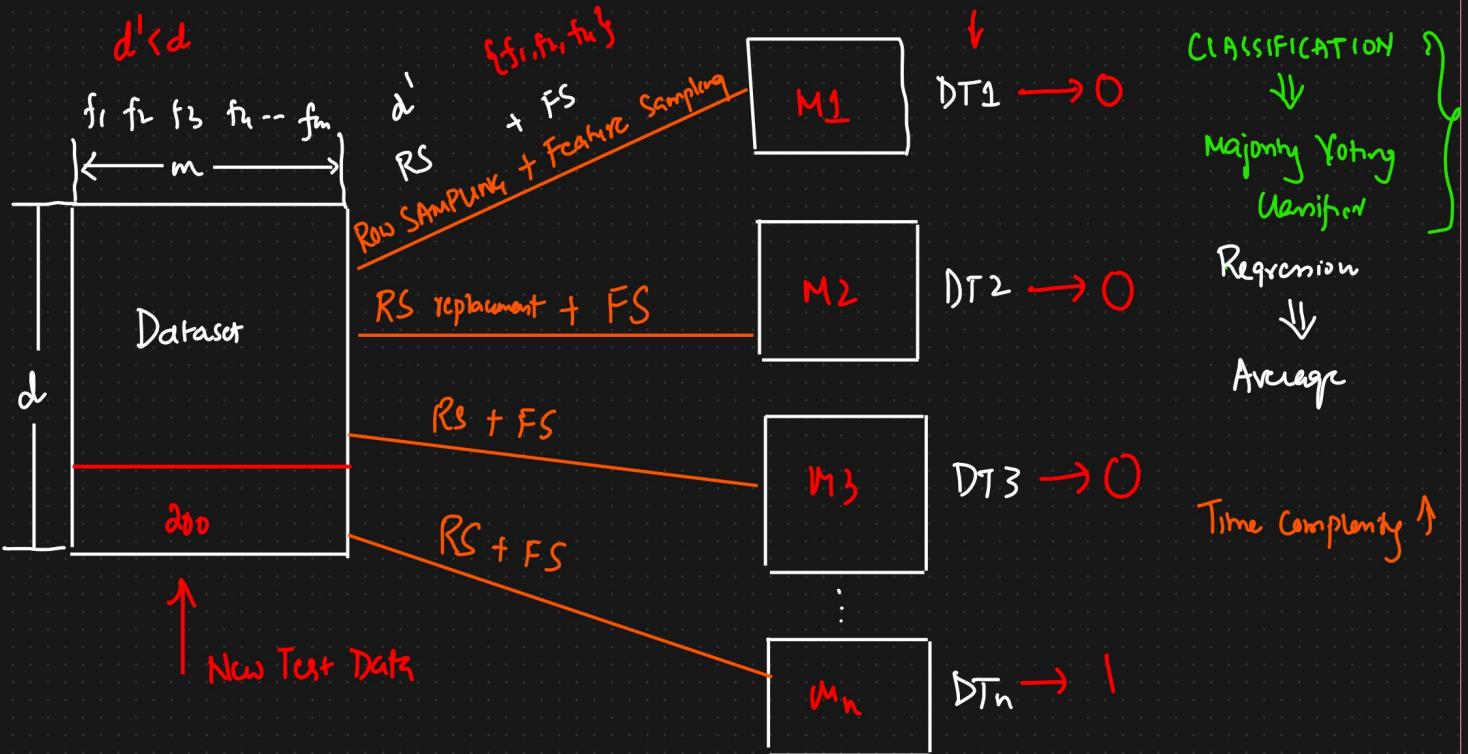
② Threshold = 3.6 ✓



Millions of records

(Time Complexity ↑↑)

# Random Forest CLASSIFICATION AND REGRESSIONS



CLASSIFICATION - MAJORITY VOTING CLASSIFIER

REGRESSION - Average O/P of the Models

Why should we use Random Forest instead of DT?

Decision Tree

Generalized Model → Low Bias

Overfitting

Random Forest

Train Acc  $\uparrow$  → Low Bias → Low Bias

Test Acc  $\downarrow$  → High Variance → Low Variance

# AdaBoost Machine Learning Algorithms

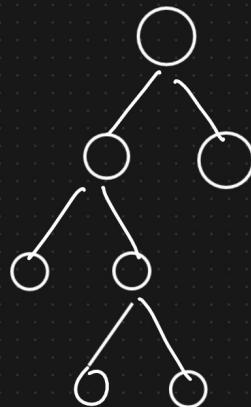
## Decision Trees

- {  
    } **Bagging**  
    {① Random Forest Classifier  
    ② Random Forest Regressor}

## Decision Tree :

Overfitting : Train Acc  $\uparrow$   
Test Acc  $\downarrow$

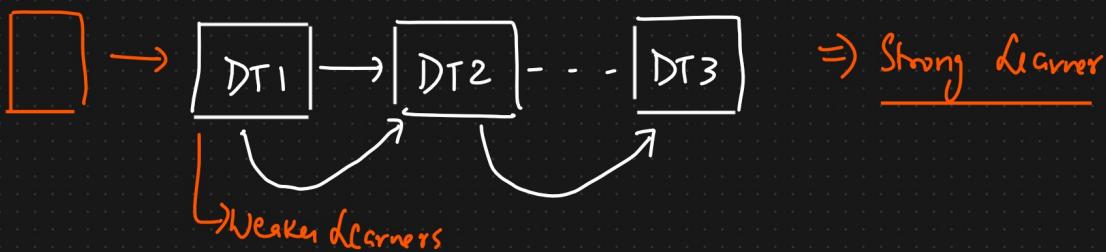
Low Bias  
High Variance



## Random Forest {     } **Bagging**

{  
    } Low Bias  
    { Low Variance }

## Boosting { Sequentially connected }



Weak Learner  $\rightarrow$  Haven't learnt much from the  
Training Dataset

Random Forest  $\rightarrow$  Majority Voting classifier  
Average of (0/p)

## Ensemble Techniques

{  
    } **weak learners**  
    **Boosting**

- ① AdaBoost
- ② Gradient Boosting
- ③ Xgboost

AdaBoost → Assignment weights to the weak learner

$M_1, \dots, M_n \rightarrow \underline{\text{Decision Tree Stumps}}$

$$f = \alpha_1(M_1) + \alpha_2(M_2) + \alpha_3(M_3) + \dots + \alpha_n(M_n)$$

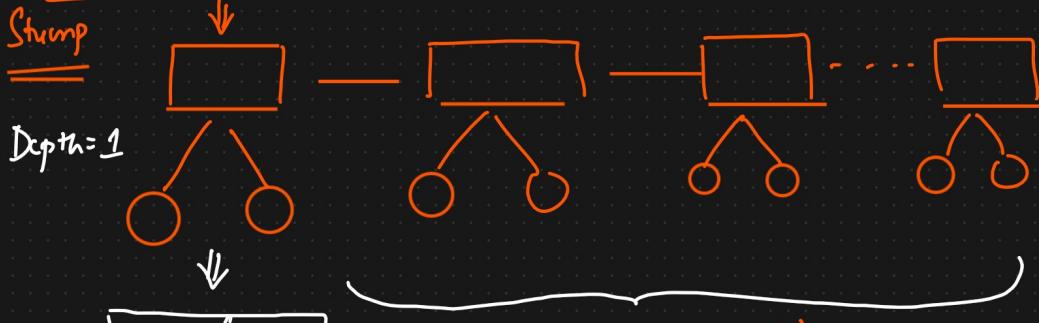
$\{\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_n\} \Rightarrow \text{weights}$

→ CLASSIFICATION

→ REGRESSION

Decision

Stump



↳ Underfitting

↳ Train Acc ↓ 40%  
↳ Test Acc ↑ 45%

Decision Tree Stump

$$\left\{ \begin{array}{l} \text{High Bias} \\ \text{Low Variance} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \text{Low Bias} \\ \text{High Variance} \end{array} \right\}$$

AdaBoost Classifier Maths Indepth Intuition ① We create Decision Tree Stump

Salary	Credit	Approval
$\leq 50K$	B	No
$\leq 50K$	G	Yes
$\leq 80K$	G	Yes
$> 50K$	B	No
$> 50K$	G	Yes
$> 50K$	N	Yes
$\leq 50K$	N	No.

and we select the

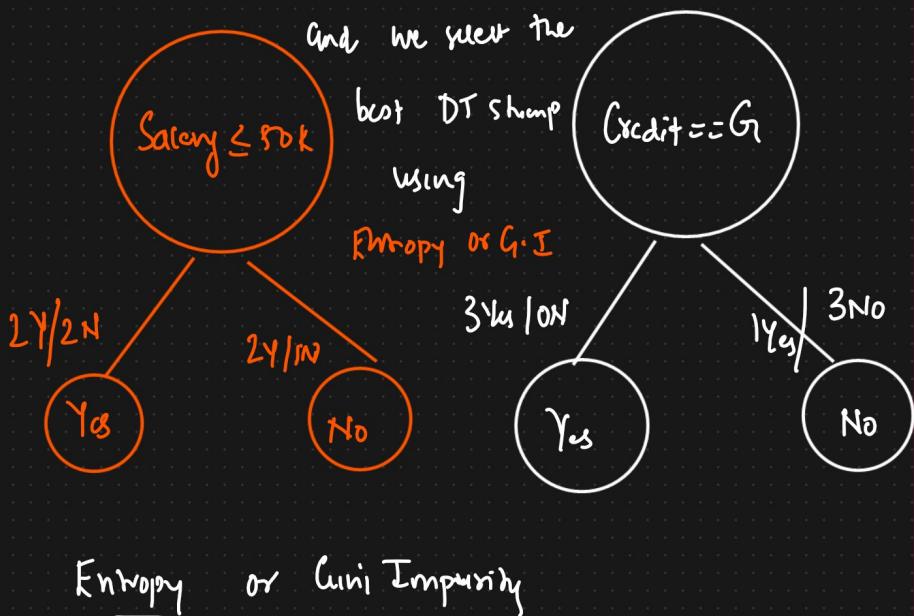
best DT stump

using

Entropy or G.I

$\text{Credit} == G_1$

3 Yes / 1 No  
1 Yes / 3 No



$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

(2) Sum of the Total Errors And Performance of Stump

Sum of all the  
① Total Error  
=  $\frac{1}{7}$

Salary	Credit	Approval	Sample Weights
$\leq 50K$	B	No.	$\frac{1}{7}$
$\leq 50K$	G	Yes	$\frac{1}{7}$
$\leq 50K$	G	Yes	$\frac{1}{7}$
$> 50K$	B	No	$\frac{1}{7}$
$> 50K$	G	Yes	$\frac{1}{7}$
$  > 50K$	N	Yes	$\frac{1}{7}$
$\leq 50K$	N	No.	$\frac{1}{7}$

② Performance of Stump =  $\frac{1}{2} \ln \left[ \frac{1 - TE}{TE} \right]$

$= \frac{1}{2} \ln \left[ \frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$

$= \frac{1}{2} \ln [6] \approx 0.896$

Performance of Stump  $\approx 0.896$

$$f = d_1(M_1) + d_2(M_2) + \dots + d_n(M_n)$$

$$d_1 = 0.896 \Rightarrow \text{Weight}$$

For correct classified points

- Performance of Stump

$$= \text{weight} * e^{-(0.896)}$$

$$= \frac{1}{7} * e^{-(0.896)}$$

$$= 0.058$$

For Incorrect classified

- Performance of Stump

$$= \text{weight} * e^{(0.896)}$$

$$= \frac{1}{7} * e^{(0.896)}$$

(3) Update the weights for correctly and Incorrectly classified points

Salary	Credit	Approval	Sample Weights	update wts	For correct classified points
$\leq 50K$	B	No.	$\frac{1}{7} \downarrow$	$0.058$	$= \text{weight} * e^{-0.058}$
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	$0.058$	$= \frac{1}{7} * e^{-0.058}$
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	$0.058$	$= 0.058$
$> 50K$	B	No	$\frac{1}{7} \downarrow$	$0.058$	
$> 50K$	G	Yes	$\frac{1}{7} \downarrow$	$0.058$	
$  > 50K$	N	Yes	$\frac{1}{7} \uparrow$	$0.349$	For Incorrect classified
$\leq 50K$	N	No.	$\frac{1}{7} \downarrow$	$0.058$	$= \text{weight} * e^{0.058}$

$$= 0.349$$

#### ④ Normalized Weights Computation And Assigning Bins

Salary	Credit	Approval	Update Wts	Normalized Weights	Bins Assignment
$\leq 50K$	B	No.	0.058	0.08	0 - 0.08
$\leq 50K$	G	Yes	0.058	0.08	0.08 - 0.16
$\leq 50K$	G	Yes	0.058	0.08	0.16 - 0.24
$> 50K$	B	No	0.058	0.08	0.24 - 0.32
$> 50K$	G	Yes	0.058	0.08	0.32 - 0.40
$  > 50K$	N	Yes	0.349	0.50	<span style="border: 1px solid orange; padding: 2px;">0.40 - 0.70</span>
$\leq 50K$	N	No.	0.058	0.08	0.50 - 0.58
				<u>0.697</u>	<u><math>\approx 1</math></u>



$\delta_1 = 0.896$  Prepare  
datapoints

#### ⑤ Select data points to send to Next Step

Salary	Credit	Approval	Bins Assignment
$\leq 50K$	B	No.	0 - 0.08
$\leq 50K$	G	Yes	0.08 - 0.16
$\leq 50K$	G	Yes	0.16 - 0.24
$> 50K$	B	No	0.24 - 0.32
$> 50K$	G	Yes	0.32 - 0.40
$  > 50K$	N	Yes	<span style="border: 1px solid orange; padding: 2px;">0.40 - 0.70</span>
$\leq 50K$	N	No.	0.50 - 0.58

① Iteration process selecting random value between 0 and 1

S	Credit	Approval	Random
$  > 50K$	N	Yes	0.50
$  \leq 50K$	G	Yes	0.10
$  > 50K$	N	Yes	0.60
$  > 50K$	N	Yes	0.75
$  \leq 50K$	G	Yes	0.24
$  > 50K$	B	No.	0.32
$  > 50K$	N	Yes	0.87

These records will be sent to Next DT Sharep.

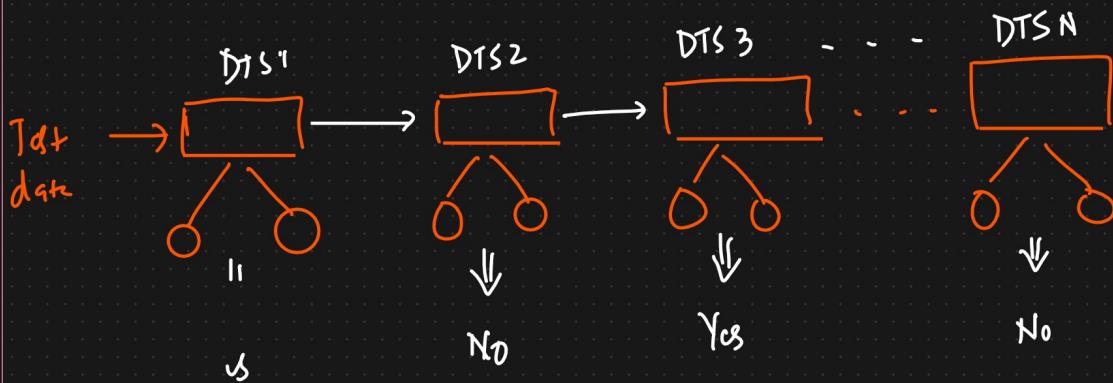
S	Credit	Approval	Sample weight	TE
>SDK	N	Yes	1/6	
<=SDK	G	Yes	1/6	Performance Shimp $\Rightarrow 0.65$
>SDK	N	Yes	1/6	
>SDK	N	Yes	1/6	
<=SDK	G	Yes	1/6	
>SDK	B	No	1/6	
>SDK	N	Yes	1/6	

$$f_1 = f_1(M_1) + f_2(M_2)$$

$$f_1 = 0.896 \quad f_2 = 0.65$$

## Final Prediction

Test data ( $\leq 50K$ , 6)



$$f_1 = 0.816 \quad f_2 = 0.650 \quad f_3 = 0.24 \quad f_n = -0.30$$

$$f = \alpha_1(M_1) + \alpha_2(M_2) + \alpha_3(m_3) + \dots + \alpha_4(M_4)$$

$$= 0.816(Y_{15}) + (0.650)(N_D) + 0.24(Y_{15}) - 0.30(N_D)$$

$$= \boxed{1.136} (\gamma_u) + 0.350 (N_0) \Rightarrow \underline{\underline{O/P}} : \underline{\underline{\gamma_u}}$$

Performance of say ( $\gamma_u$ ) = 1.136

Performance of say ( $N_0$ ) = 0.350

# GRADIENT BOOSTING ALGORITHM

## ① Regression

## ② Classification

### Regression Data

Exp	Degree	Salary	$\hat{y}$	$(y - \hat{y})$	$\hat{R}_1$	$R_2$	$\hat{y}$	$R_3$	$R_4$
→ 2	B.E.	50K	75K	-25K	-23	72.7	-22.7	-	
3	Masters	70K	75K	-5K	-3	74.7	-4.7	-	
5	Masters	80K	75K	5K	3	-	-	-	
6	Ph.D	100K	75K	25K	20	-	-	-	
$\hat{y} = 75K$									

### Steps

#### ① Create a Base Model

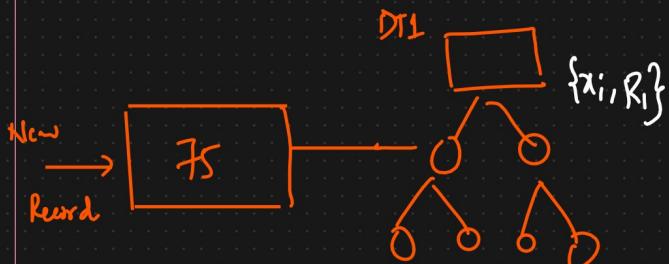
$\boxed{75}$

$$\text{Average} = \frac{50 + 70 + 80 + 100}{4} = 75$$

#### ② Compute Residuals, Error

#### ③ Construct a Decision Tree

Consider inputs  $x_i$  and O/p  $R_i$



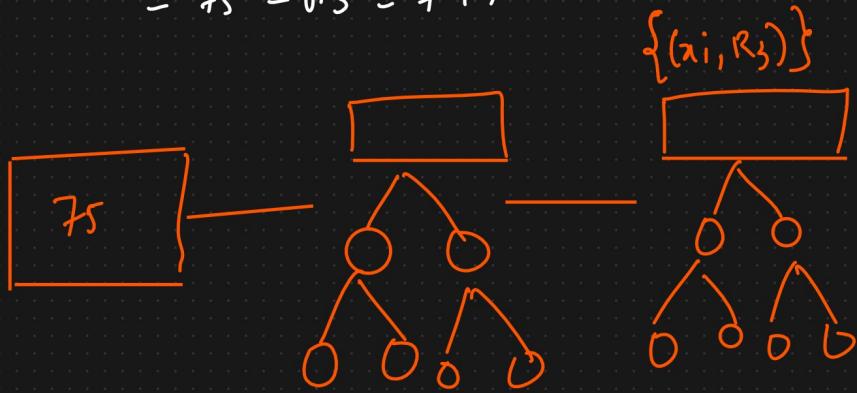
$$\text{Predicted O/p} \Rightarrow 75 + (-23) = 75 - 23 = \underline{\underline{52}} \quad \{ \text{Overfitting} \}$$

$$\text{Predicted O/P} \Rightarrow 75 + d_1(D_i) = 75 + (0.1)(-23)$$

$$\begin{aligned} d &= \text{Learning Rate} \\ d &= 0.1 \end{aligned}$$

$$\begin{aligned} &= 75 - 2.3 \\ &= 72.7 \end{aligned}$$

$$\begin{aligned} &\Rightarrow 75 + 0.1(-3) \\ &= 75 - 0.3 = 74.7 \end{aligned}$$



$$F(x) = L_0 h_0(x) + L_1(h_1(x)) + L_2(h_2(x)) + L_3(h_3(x)) + \dots + L_n(h_n(x))$$

Learning Rate  $\quad \delta = 0.1$

$$F(x) = \sum_{i=0}^n \delta_i h_i(x)$$

⇒ Final Function  
of Gradient  
Boosting

# Xgboost ML Algorithm (Classification)

$$\hat{y} = 0.5$$

Dataset

		$y$	$\downarrow$	$\wedge$	$\vee$
Salary ✓	Credit ✓	Approval	<u>R1</u>	<u>Y</u>	<u>R2</u>
$\rightarrow \{ \leq 50K \quad B \quad O \}$			-0.5	0.52	-0.48
$\{ \leq 50K \quad G \}$			1	0.5	0.58
$\{ \leq 50K \quad G \}$			1	0.5	-
$\{ > 50K \quad B \}$		O	-0.5	-	-
$\{ > 50K \quad G \}$		I	0.5	-	-
$\{ > 50K \quad N \}$		I	0.5	-	-
$\{ \leq 50K \quad N \}$		O	-0.5	-	-

Steps

① Construct a base Model

② Construct a Decision Tree  
with root.

③ Calculate Similarity Weight

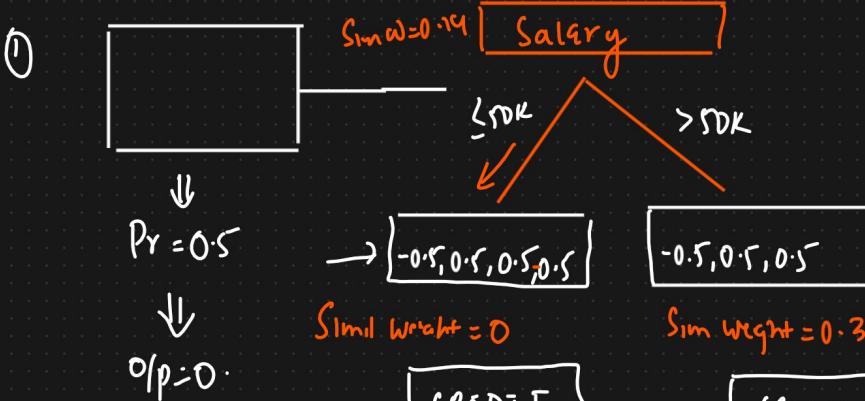
$$= \sum (\text{Residual})^2$$

$$\text{Conv value} \leftarrow \frac{\sum Pr(1-Pr)}{+ \lambda}$$

④ Calculate Gain

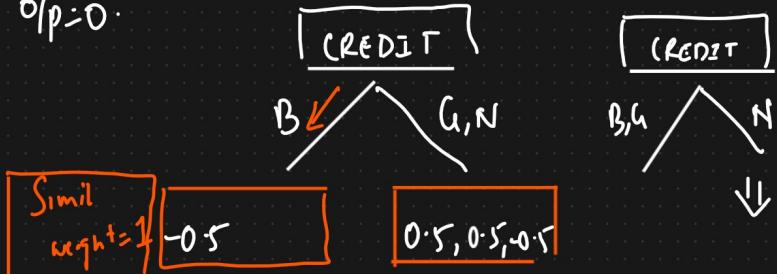
Hyperparameter

$$[-0.5, 0.5, 0.5, -0.5, -0.5, 0.5, 0.5]$$



$$\log(\text{odds}) = \log \left( \frac{P}{1-P} \right)$$

$$\log(\text{odds}) = \log \left( \frac{0.5}{0.5} \right) \\ = 0$$



$$\text{Similarity}(AC) = \frac{0.25}{0.25} = 1$$

$$\text{Similarity}(RC) = \frac{0.25}{0.75} = 0.33$$

$$\text{Gain} = 1 + 0.33 - 0 = \boxed{1.33}$$

|

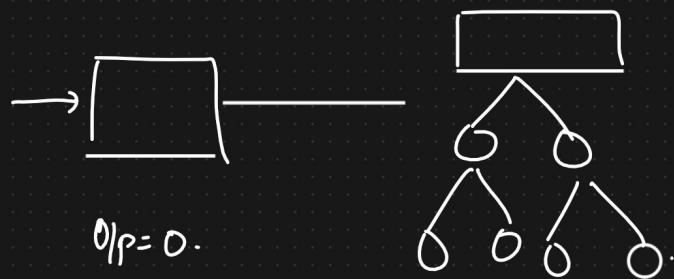
$$\text{Similarity weight} (d_C) = \frac{\sum (\text{Residual})^2}{\sum P_r (1-P_r)}$$

$$= \left[ (-0.5 + 0.5 + 1.5 - 0.5)^2 \right] = 0$$

$$1 \in [0.8(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)]$$

$$\text{Gain} = 0 + 0.33 - 0.14 = 0.21$$

Final O/P



$$\begin{matrix} \xrightarrow{\text{Now}} \\ D_{gtz} \end{matrix} = \sigma \left( 0 + d(1) \right) \quad d = 0.1$$

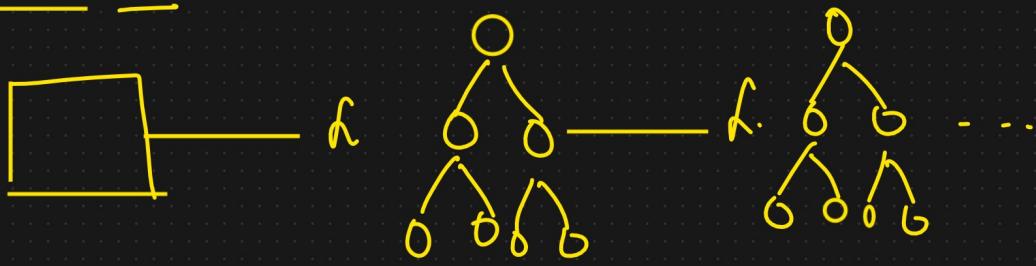
$$\begin{matrix} \xleftarrow{\text{logistic function}} \\ \xleftarrow{\text{Activation}} \end{matrix} \begin{matrix} \text{Sigmoid} \\ \text{Function} \end{matrix} = \sigma \left( 0 + 0.1(1) \right) = \frac{1}{1 + e^{-0.1}} = 0.52$$

$$\begin{matrix} \xleftarrow{\text{Second Record}} \\ \text{O/p} \end{matrix} = \sigma \left( 0 + d(0.33) \right) = \sigma \left( 0 + 0.1(0.33) \right)$$

$$\text{Similarity weight} (R_C) = \frac{(1.5 + 0.5 + 0.5)^2}{0.75} = \frac{4}{0.75} = 5.33$$

$$= \frac{1}{1+e^{-0.033}} = 0.508$$

Xg boost classifier



$$O/P = \sigma \left( \text{Basefimer} + f_1(DT_1) + f_2(DT_2) + f_3(DT_3) \right).$$

# Xgboost Regressor Mh Algorithm

Dataset {Regression} → 

<u>Exp</u>	<u>Gap</u>	Salary	<u>R<sub>1</sub></u>	<u><math>\hat{y}</math></u>	<u>R<sub>2</sub></u>
→ 2	Yes	40K	-11	49.9	-9.9
→ 2.5	Yes	42K	-9	49.9	-7.9
→ 3	No	52K	1	51.5	0.5
4	No	60K	9	51.5	8.5
4.5	Yes	62K	11	52.1	9.9
				≈ 51K	

$[51 + (0.1)(-10)] = 51 - 0.1 = 49.9$

$[51 + (0.1)(5^-)] = 51.5$

$[51 + 0.1(11)] = 51 + 1.1 = 52.1$

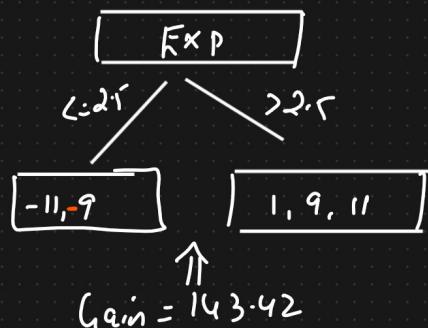
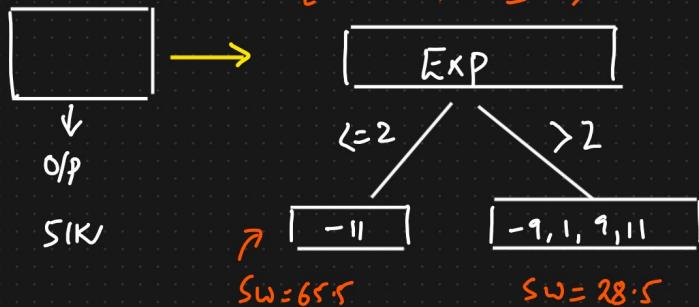
$$\text{Similarity weight} = \frac{\sum (\text{Residual})^2}{\sum p_x(1-p_x)}$$

Gain

## Steps

- ① Create a Base Model
- ② Residual Computation
- ③ Construct DT1 using  $\{x_i, R_i\}$

$$[-11, -9, 1, 9, 11] \Rightarrow SW = 0.16.$$



$$\text{Similarity weight} = \frac{\sum (\text{Residual})^2}{\text{No. of Residuals}}$$

$\lambda = 1$   $\rightarrow$  Hyperparameter

$$SW(\text{first child}) = \frac{121}{1+1} = 121/2 = 65.5 \text{ //}$$

$\boxed{\lambda \uparrow SW \downarrow}$

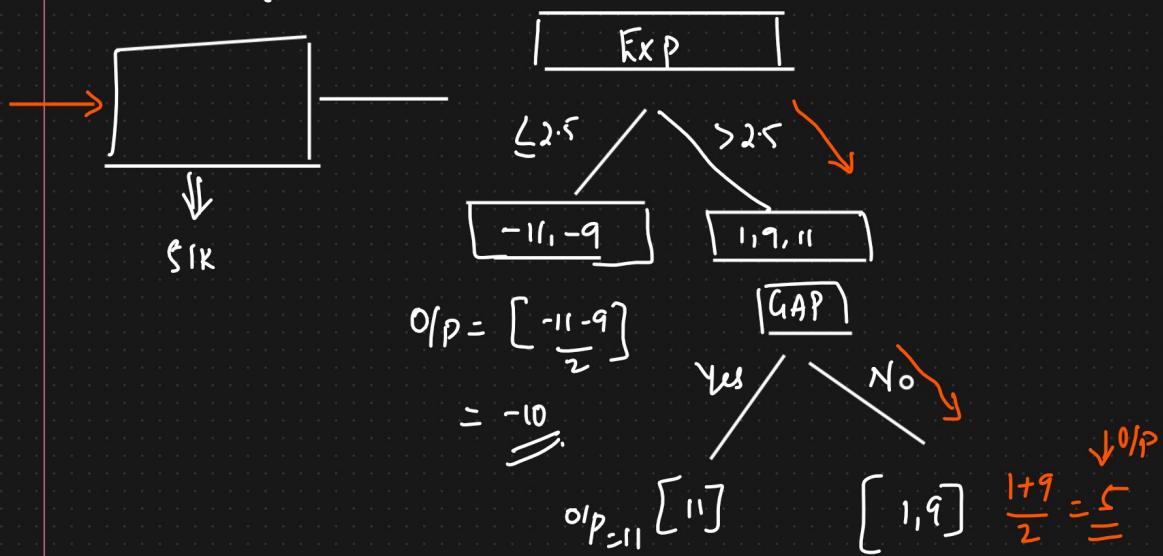
$$SW(\text{Right child}) = \frac{(-9+1+9+11)^2}{4+1}$$

$$= \frac{144}{5} = 28.5$$

### (5) Calculate Gain

$$\begin{aligned} \text{Gain} &= 65.5 + 28.5 - 0.16 \\ &= 98.34 \end{aligned}$$

Banc Farmer



$\alpha$ : Learning Rate     $\alpha = 0.1 \Rightarrow$  Hyperparameter

$$XGBoost \text{ Classifier} = \text{Banc Farmer} + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n)$$

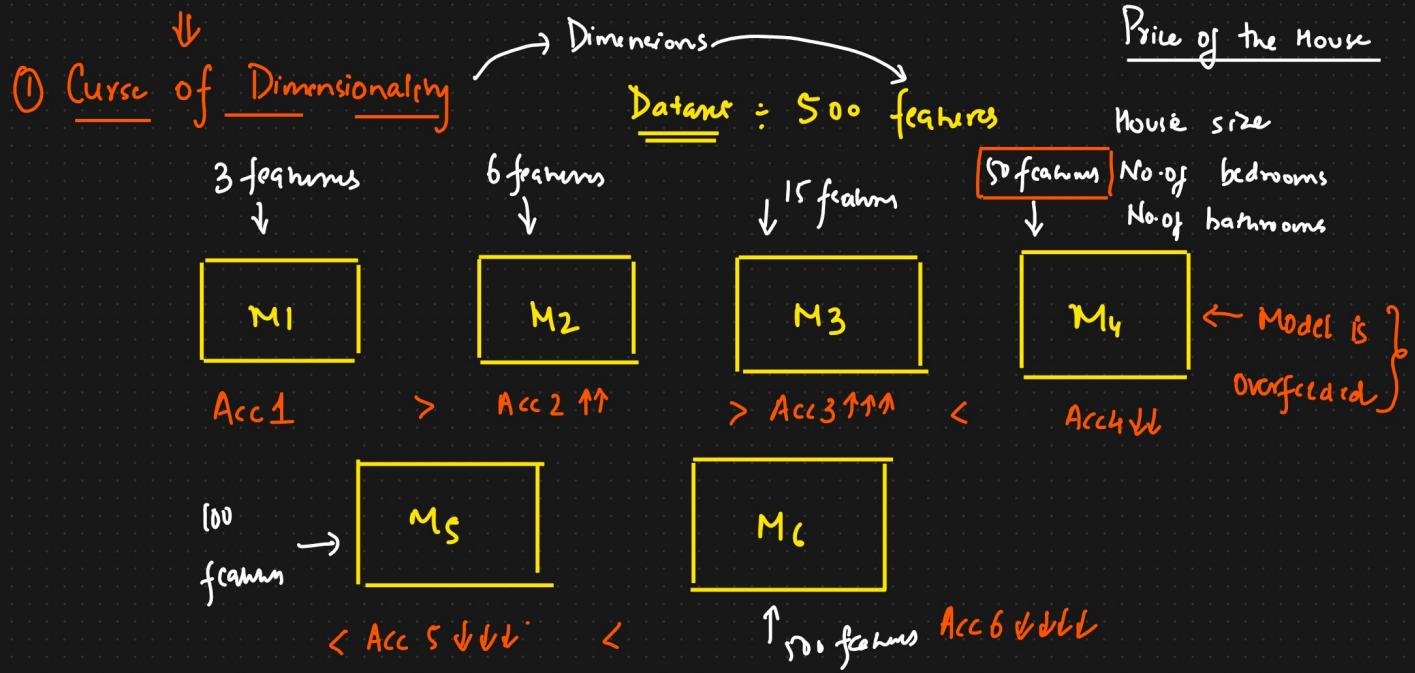
$$\begin{aligned} XGBoost \text{ Regressor} &= SIK + 0.1(5) \\ &= 51 + 0.5 \\ &= 51.5 \end{aligned}$$

$$\begin{aligned} \text{Smoothing weight} &= \frac{\sum (\text{Residual})^2}{\text{No. of Residuals} + \boxed{\lambda}} \end{aligned}$$

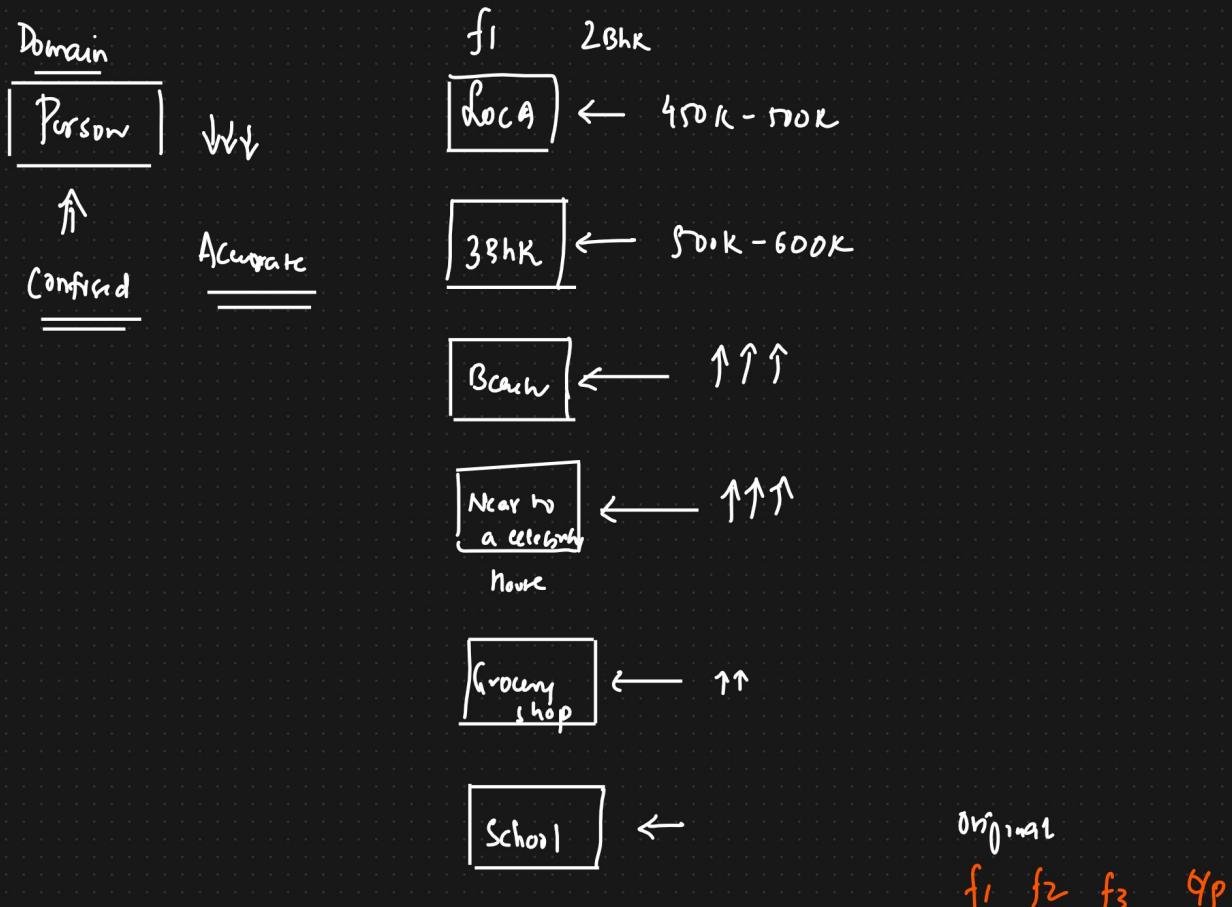
{ Regression }



## Principal Component Analysis (PCA) [Dimensionality Reduction]

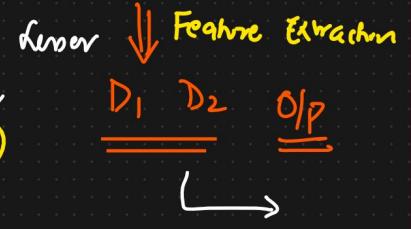


## ② Model performance Degrade



Two different ways to remove curse of Dimensionality

- ## ① Feature Selection      ② Dimensionality Reduction (PCA)



Imp features

## Y Feature Extraction

### Feature Selection Vs Feature Extraction

↳ Dimensionality Reduction

① Why Dimensionality Reduction?

- Ⓐ Prevent → Curse of Dimensionality
- Ⓑ Improve the performance of the model
- Ⓒ Visualize the data → understand the data

$\boxed{3d}$     $\boxed{2d}$

$\boxed{100d}$

$\downarrow$   
 $\boxed{3d}$  or  $\boxed{2d}$

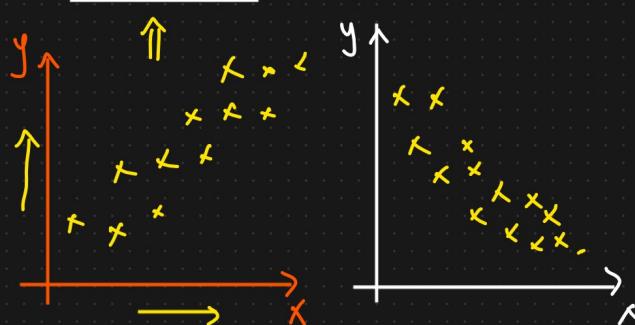
### Feature Selection

$$\begin{matrix} \text{JIP} & \text{OIP} \\ \boxed{X} \rightarrow y \\ - & - \end{matrix} \quad \text{tvc} \Rightarrow \begin{matrix} \downarrow \\ \boxed{X \uparrow \quad Y \uparrow \\ X \downarrow \quad Y \downarrow} \end{matrix}$$

$$\begin{matrix} \downarrow \\ \boxed{X \downarrow \quad Y \uparrow \\ X \uparrow \quad Y \downarrow} \end{matrix}$$

No relationship between X & Y

-ve y



$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1} = \begin{cases} +\text{ve} & \approx 0 \\ -\text{ve} & \approx 0 \end{cases}$$

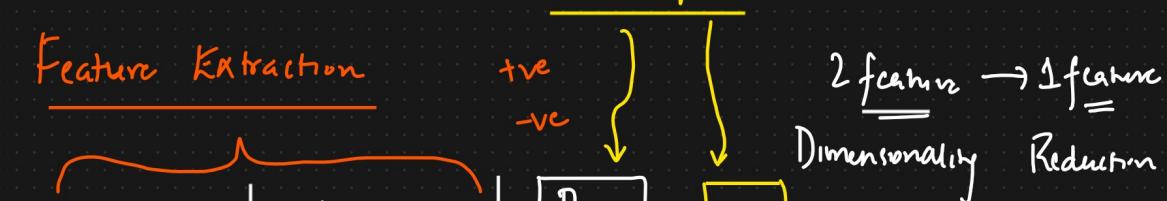
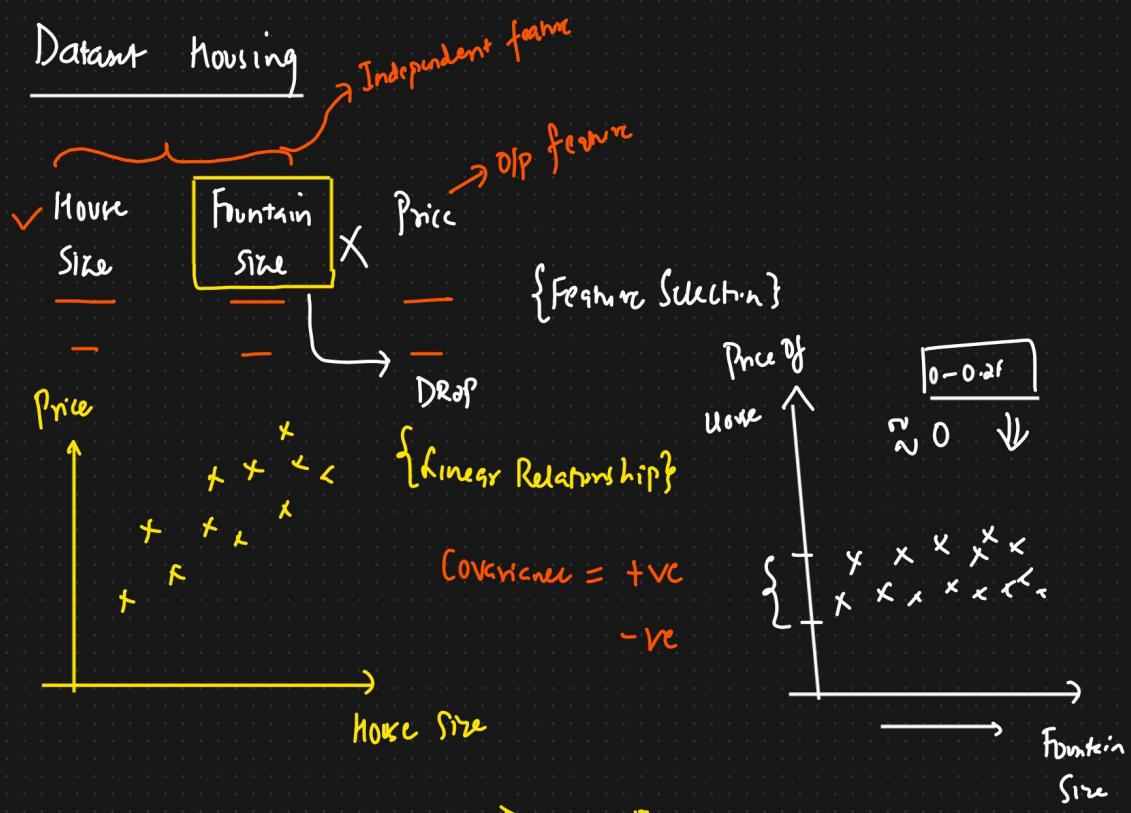
$\approx 0$  {No Relationship}

$\boxed{-\text{ve correlated}}$

$$\text{Pearson Correlation} = \frac{\text{Cov}(x, y)}{\sqrt{x^2 + y^2}} = \boxed{-1 \text{ to } 1}$$

$\left. \begin{array}{l} \text{The more towards the} \\ \text{Value of } +1 \text{ the} \end{array} \right\}$

more +ve correlated X & Y is



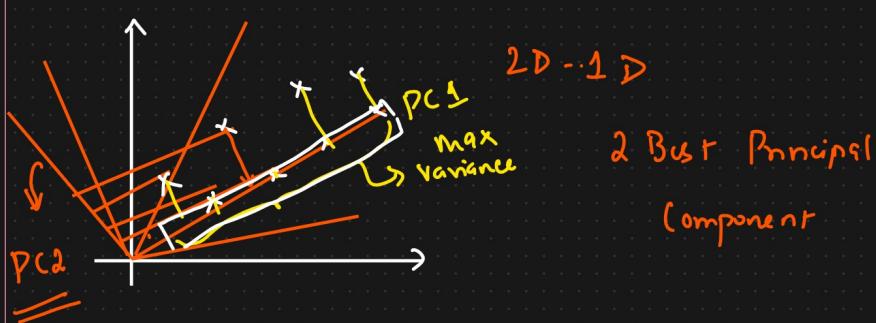
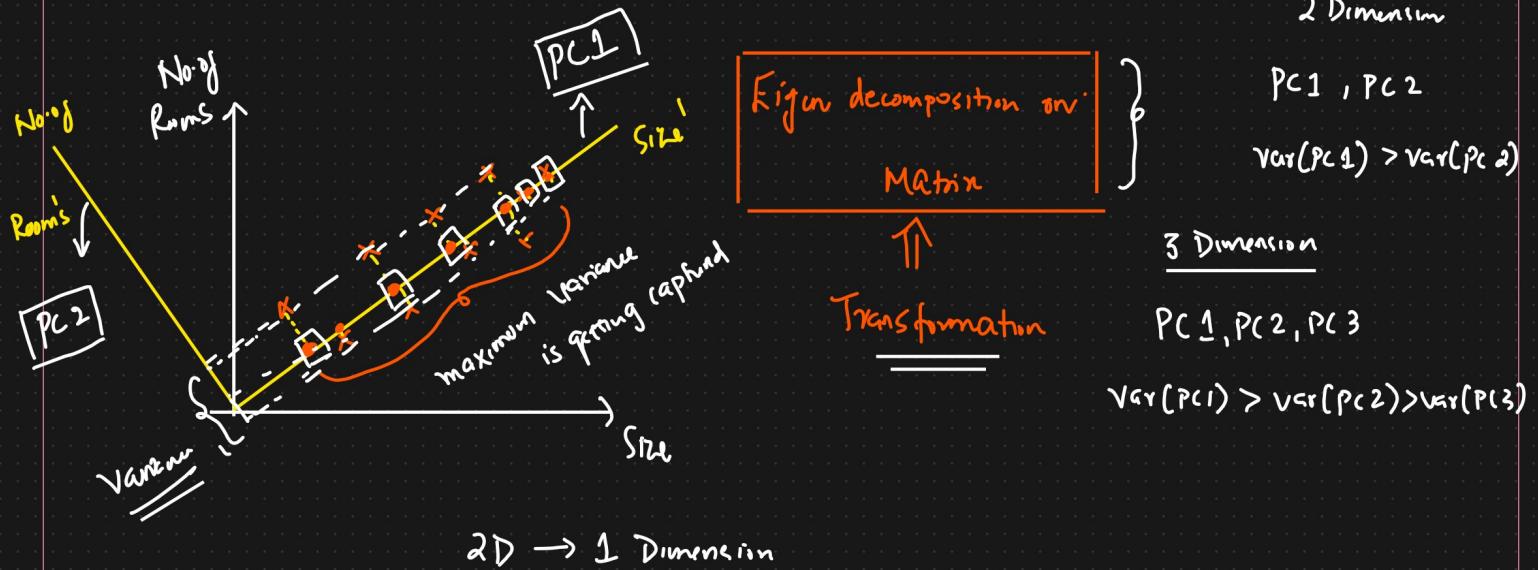
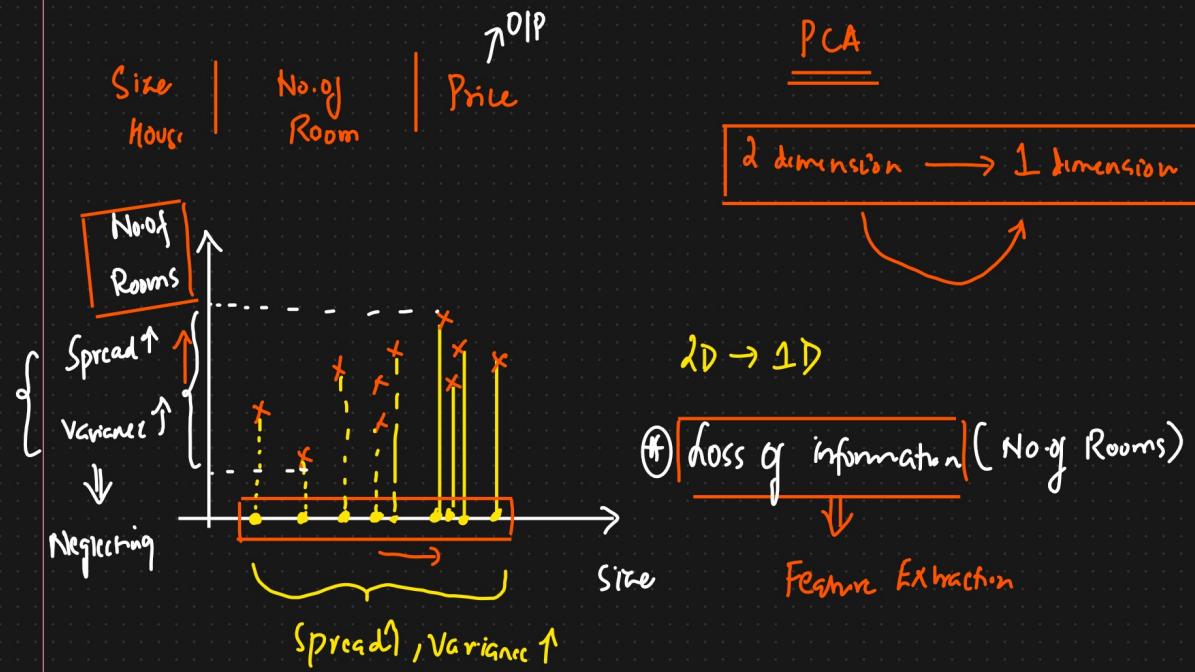
↓ ↓ Transformation To extract New feature



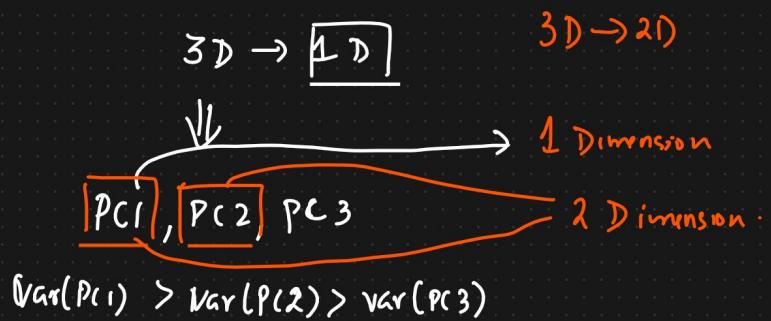
# PCA Geometric Intuition

{ Dimensionality Reduction }

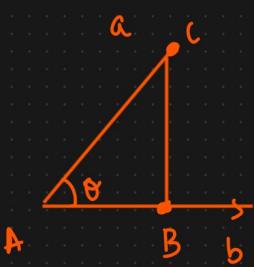
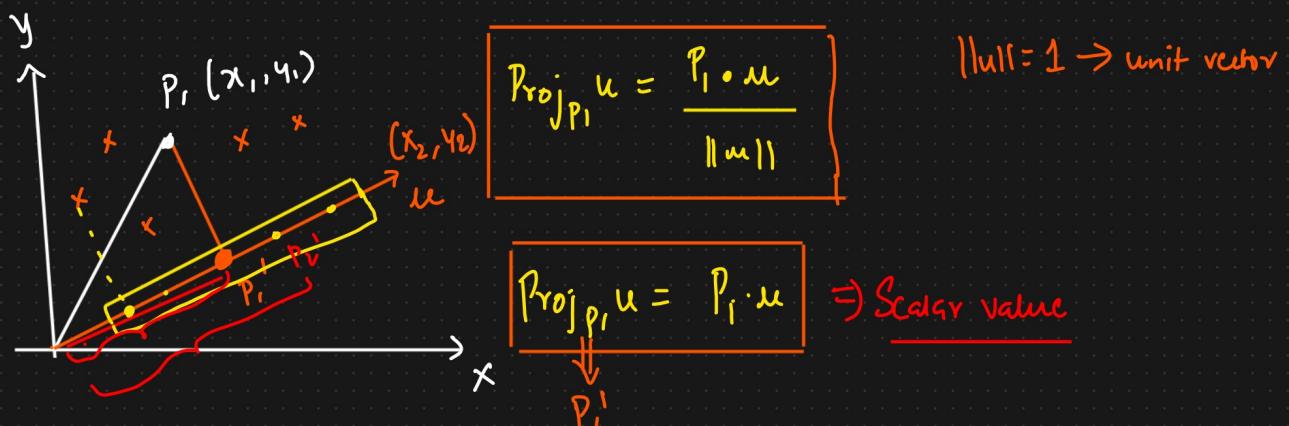
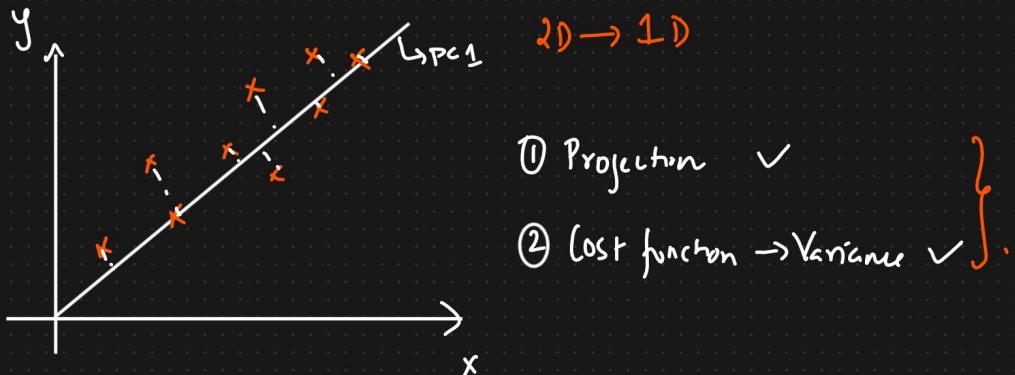
## Moving Dataset



To get the best Principal Component which captures maximum variance



### Maths Intuition behind PCA Algorithm



$$[P_0^T, P_1^T, P_2^T, P_3^T, P_4^T, \dots, P_n^T]$$

$\Downarrow$   
Scalar values  
 $\Downarrow$   
Variance

$$\boxed{P_0^1, P_1^1, P_2^1, P_3^1, P_4^1, \dots, P_n^1}$$

$$\downarrow$$

$$x_0^1, x_1^1, x_2^1, x_3^1, x_4^1, \dots, x_n^1$$

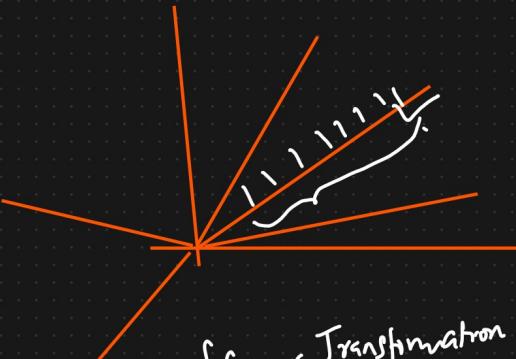
$$\text{Max Variance} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

{ goal: Find the best unit vector which captures maximum variance }.

$\downarrow$

Cost function

Eigen vectors And Eigen values.



① Covariance Matrix between features

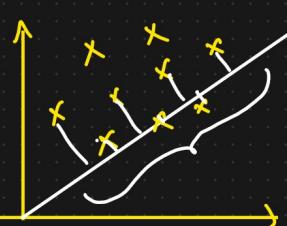
{Linear Transformation of matrix}

② Eigen vectors and Eigen values will be found out from this covariance matrix

$$A v = \lambda v$$

③ Eigen vector  $\rightarrow$  Eigen value  $\rightarrow$  magnitude of the Eigen vector  $\rightarrow$  Capture the maximum variance

Eigen vectors And Eigen values [Linear Transformation]



[Eigen decomposition of covariance Matrix]



Eigen vector & Eigen values

$$[ ] * [v] = \lambda * v$$

↓  
Eigen  
Value

$$\boxed{A * v = \lambda * v}$$

↑ v ↓



Eigen vector → Maximum magnitude



Eigen vector → Max Magnitude



Principal Component



Max Eigen Vector



Max Var

Best Principal Component → PC1

Steps to calculate Eigen value and vectors

---

① Covariance of features

$$\boxed{(X, Y)} \quad Z$$

$\downarrow$

$$x_i$$

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$2 \times 2$      $x \quad y$

$$A = \begin{matrix} X & \boxed{\begin{array}{|c|c|} \hline \text{Var}(x) & \text{Cov}(x,y) \\ \hline \text{Cov}(y,x) & \text{Var}(y) \\ \hline \end{array}} \\ Y \end{matrix}$$

$\text{Cov}(x, x) = \text{Var}(x)$   
 $\text{Cov}(y, y) = \text{Var}(y)$

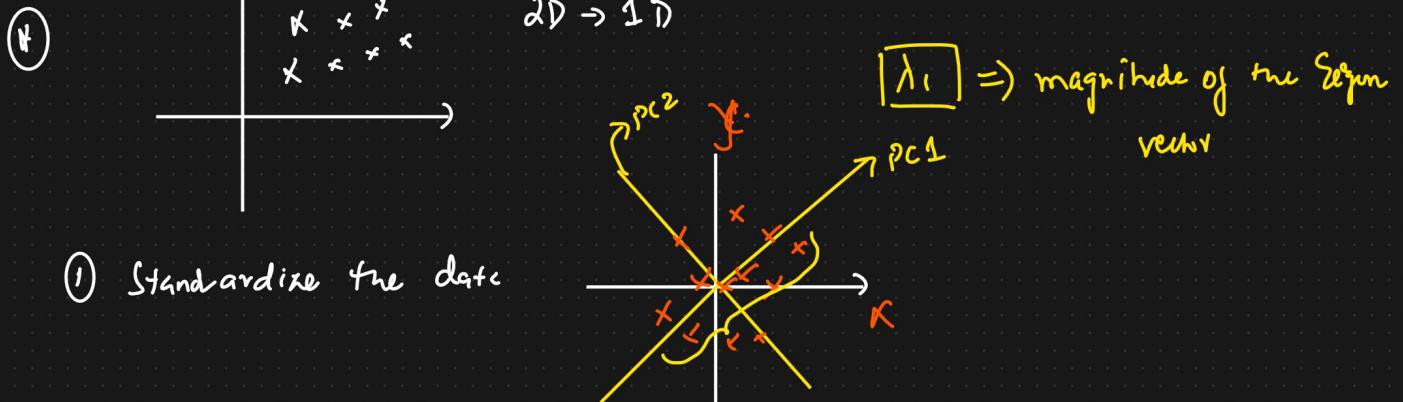
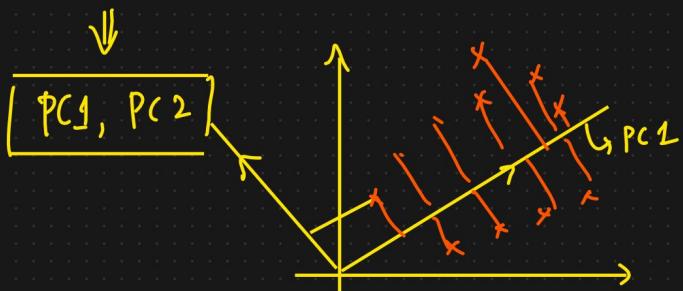
$x \quad y \quad z$

$x$			

$$Y = \begin{bmatrix} f_1 & f_2 \end{bmatrix}$$

$$A \cdot v = \lambda \cdot v$$

$$\lambda_1, \lambda_2 \rightarrow \text{Eigen values}$$



(2) Covariance Matrix of  $X \& Y$

$$A = \begin{matrix} X & Y \\ X & Y \end{matrix} \quad 2 \times 2$$

(3) Find out Eigen vectors And value

$$A v = \lambda v$$

$$\lambda_1, \lambda_2 \rightarrow \text{Eigen values.}$$

$\downarrow \quad \downarrow$

PC1 PC2



## Unsupervised Machine Learning

# Supervised ML



Input Data Dependent or  
↓  
O/P features  
 I/O  
 $f_1 \quad f_2 \quad f_3 \quad f_4$       O/P

Unsupervised ML → Clustering {group your data into similar clusters}

## Data set

Age      Experience      Salary      No      O/P

## Cluster 9



$\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$      $\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$      $\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$      $\Rightarrow$  Unsupervised ML  
 $\Downarrow$   
 $\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$      $\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$      $\left\{ \begin{array}{c} - \\ - \\ - \end{array} \right.$   
 Group or cluster  
 this data

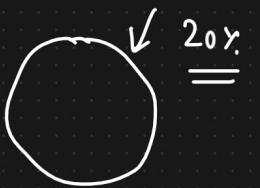
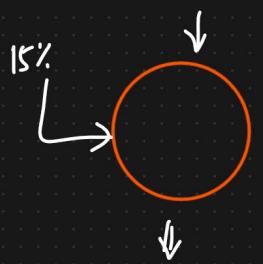
A diagram showing two circles, one yellow and one green, positioned side-by-side on a grid background. The yellow circle is on the left and the green circle is on the right.

{ - - - }

eg: Customer Segmentation

Apple      Product       $\downarrow$       Data      Salary      Spending-Score

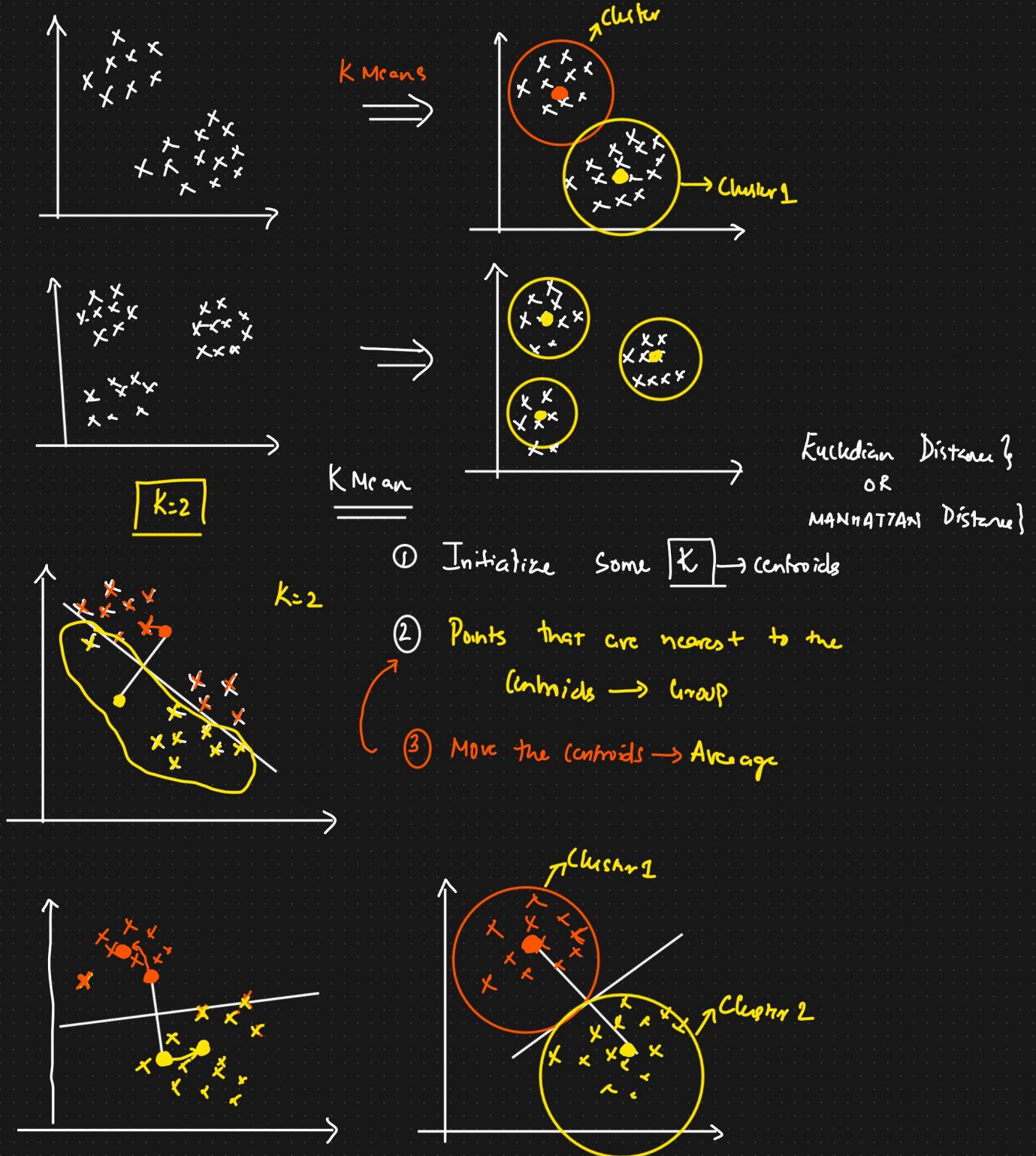
↳ New Product  $\leftarrow$  Discount



## Unsupervised ML

- ① K Means Algorithm
  - ② Hierarchical Clustering
  - ③ DBScan Clustering
  - ④ Silhouette Scoring
- } validate
- =

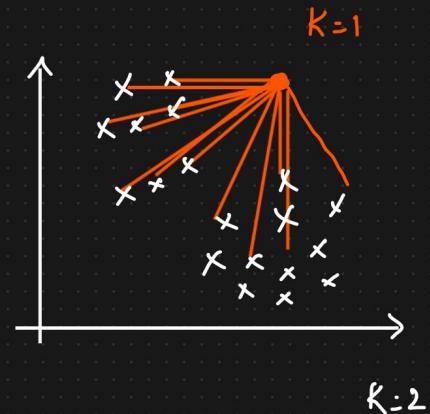
# K Means Clustering Algorithm



How do we select the K value?

WCSS = Within Cluster Sum of Squares

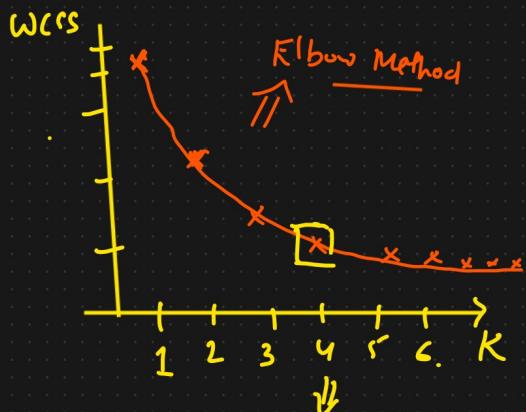
Initialize K=1 to 20



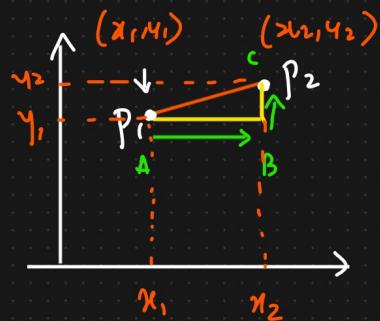
$$WCSS = \sum_{i=1}^n \left( \text{distance between Points to nearest Centroid} \right)^2$$



$WCSS \downarrow$



### # Euclidean Distance



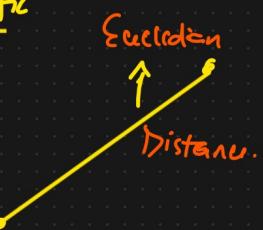
$$\text{Euclidean dist} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Manhattan dist} = |x_2 - x_1| + |y_2 - y_1|$$

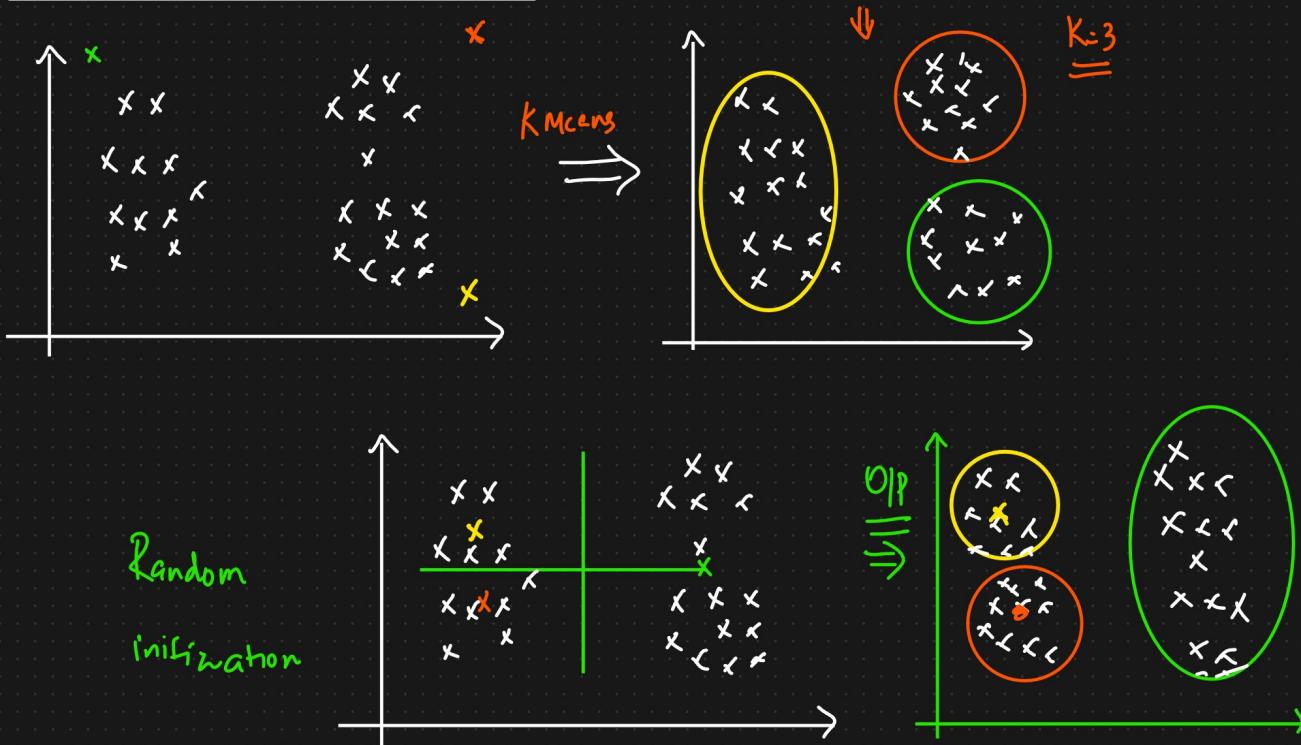
IRON MAN  $\rightarrow$  U.S



Air Traffic

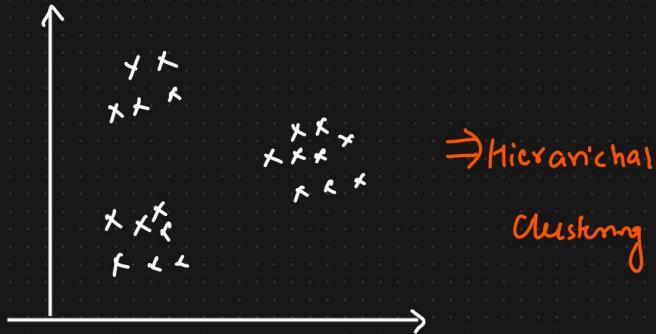


# Random Initialization TRAP (Kmeans++)



## Kmeans++ Initialization Technique

# Hierarchical Clustering

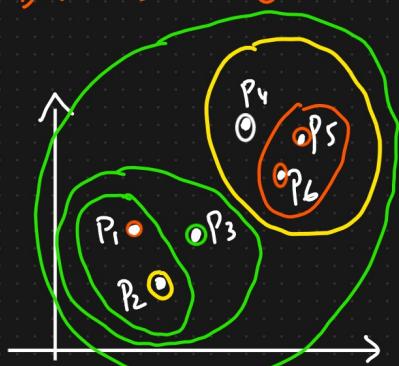


K=3      No centroids

## HC

- ① Agglomerative
- ② Divisive

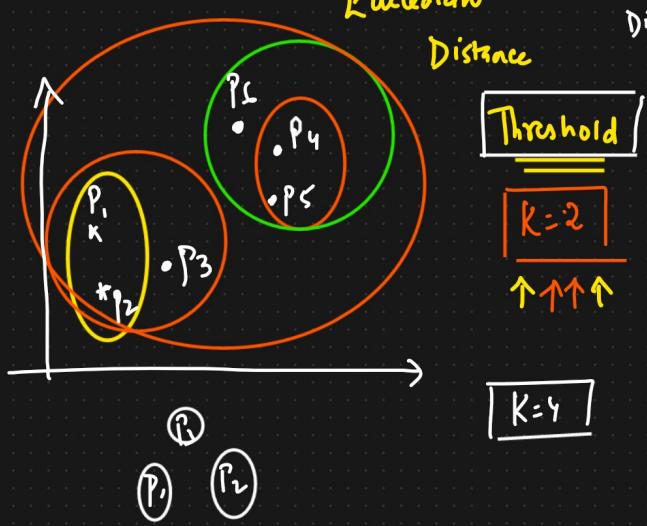
⇒ Geometric Intuition.



## Steps

- ① For each point initially will consider it as a separate cluster
- ② Find the nearest point and create a new cluster
- ③ Keep on doing the same process until we get a single cluster

## Dendrogram

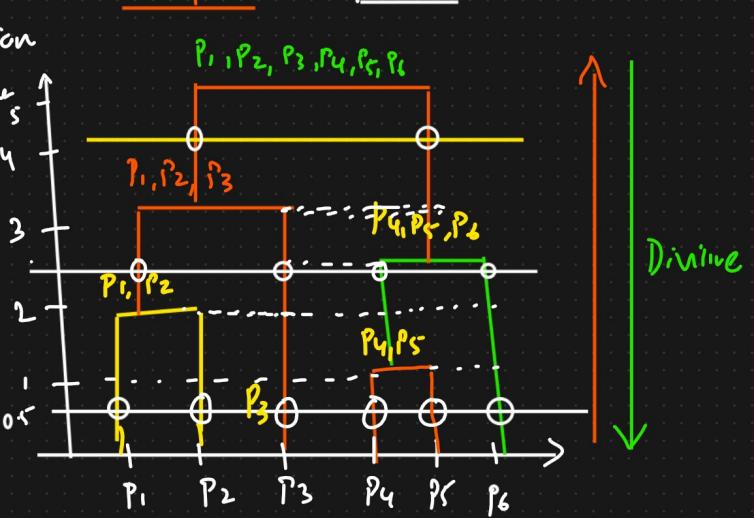


## Dendrogram

K=2

K=2

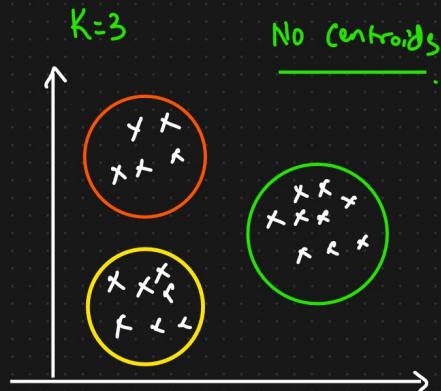
Agglomerative



(K) Select the longest vertical line such that }  
no horizontal line passes through it }

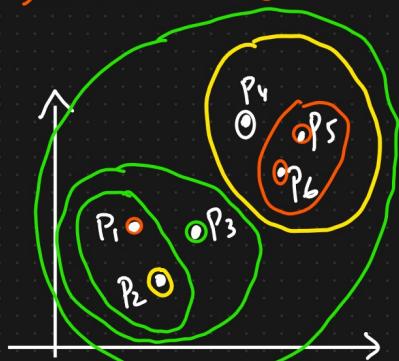
Threshold { Euclidean Distance } .

# Hierarchical Clustering



## HC

- ① Agglomerative
  - ② Divisive
- ]  $\Rightarrow$  Geometric Intuition.

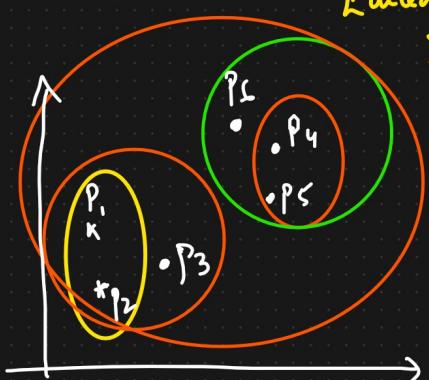


## Steps

- ① For each point initially will consider it as a separate cluster
- ② Find the nearest point and create a new cluster
- ③ Keep on doing the same process until we get a single cluster

## Cosine Similarity

## Dendrogram



No. of clusters

Euclidean Distance

Threshold

$K=2$   
↑↑↑↑

$K=4$

Fusion

Distance

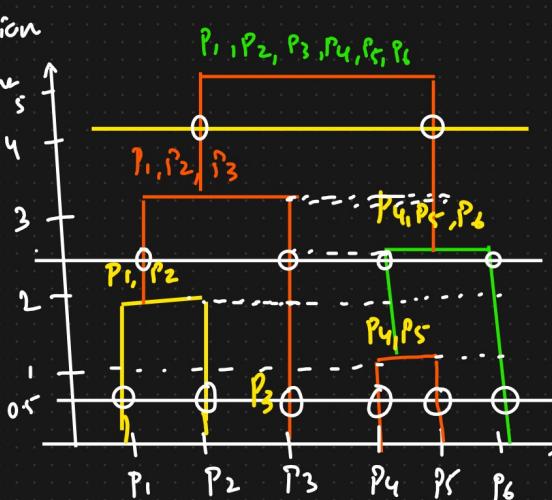
## Dendrogram

$K=2$

$K=2$

Agglomerative

Divisive



④ Select the longest vertical line such that }  
no horizontal line passes through it }

Threshold { Euclidean Distance } .

## K Means Vs Hierarchical Clustering

Scalability ✓ And Flexibility ✓

- ① Dataset size → Huge → K Means  
Small → Hierarchical Clustering



- ② Kmean → Numerical data  
Hierarchical clustering → Variety of data.

$$\Rightarrow \text{Cosine Similarity}$$

- ③ Centroids → Elbow method → No. of centroids  
→ No. of clusters

# Silhouette (clustering)

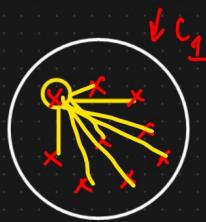
①

For data point  $i \in C_I$  (data point  $i$  in the cluster  $C_I$ ), let

$$a(i) = \frac{1}{|C_I| - 1} \sum_{j \in C_I, i \neq j} d(i, j) \quad \checkmark$$

be the mean distance between  $i$  and all other data points in the same cluster, where  $|C_I|$  is the number of points belonging to cluster  $i$ , and  $d(i, j)$  is the distance between data points  $i$  and  $j$  in the cluster  $C_I$  (we divide by  $|C_I| - 1$  because we do not include the distance  $d(i, i)$  in the sum). We can interpret  $a(i)$  as a measure of how well  $i$  is assigned to its cluster (the smaller the value, the better the assignment).

$a(i)$



②

We then define the mean dissimilarity of point  $i$  to some cluster  $C_J$  as the mean of the distance from  $i$  to all points in  $C_J$  (where  $C_J \neq C_I$ ).

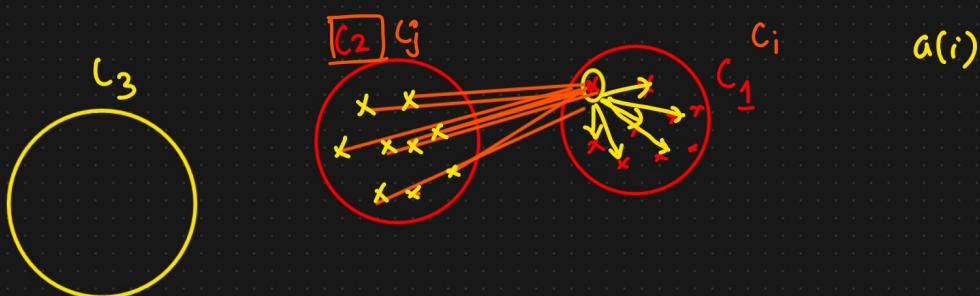
For each data point  $i \in C_I$ , we now define

$$b(i) = \min_{J \neq I} \frac{1}{|C_J|} \sum_{j \in C_J} d(i, j) \Rightarrow b(i)$$

$a(i) \ll b(i)$

to be the *smallest* (hence the `min` operator in the formula) mean distance of  $i$  to all points in any other cluster, of which  $i$  is not a member. The cluster with this smallest mean dissimilarity is said to be the "neighboring cluster" of  $i$  because it is the next best fit cluster for point  $i$ .

$c_i$        $a(i)$



## ② Silhouette Score

We now define a *silhouette* (value) of one data point  $i$

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |C_I| > 1$$

$$\boxed{-1 \text{ to } 1}$$

and

$$s(i) = 0, \text{ if } |C_I| = 1$$

Which can be also written as:

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

More near to 1 better

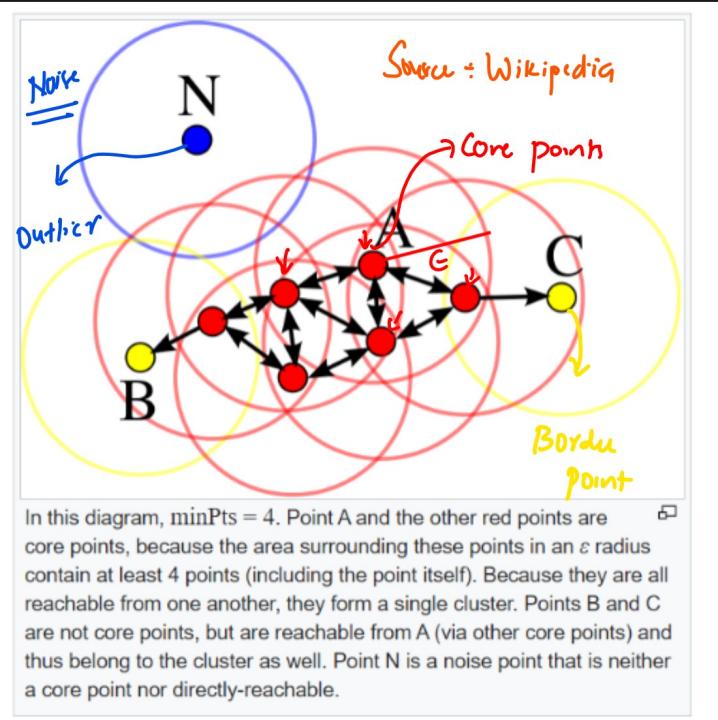
Clustering model we

have created

From the above definition it is clear that

$$\boxed{-1 \leq s(i) \leq 1}$$

# DBSCAN Clustering.



● → Core point  
● → border point  
● → Outlier

$\left. \begin{array}{l} \text{Non linear} \\ \text{Clustering} \end{array} \right\}$

$$\text{minpts} = 4 \quad \epsilon = \text{radius}$$

### Core point

- ① No. of points within the  $\epsilon$  Should be greater  $\geq 4$



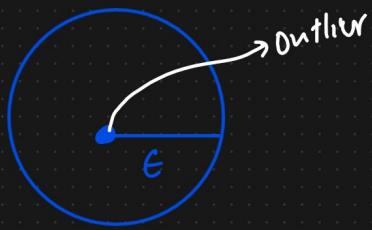
$$\text{minpts} \geq 4$$

### Border point

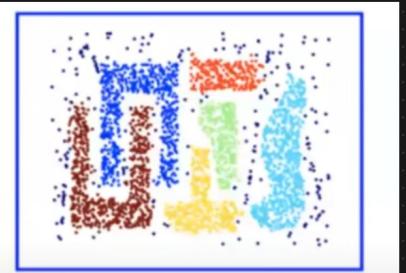
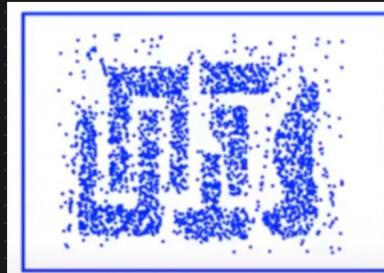
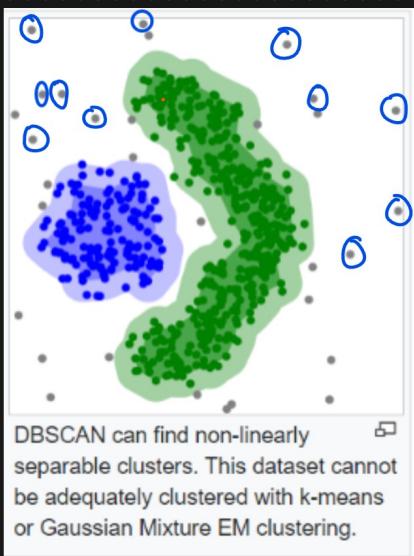
no. of data points within this radius will be less than  $\text{minpts}$



### Outlier (Noise)



## Some Examples after we apply DBScan Clustering



The left image depicts a more traditional clustering method that does not account for multi-dimensionality. Whereas the right image shows how DBSCAN can contour the data into different shapes and dimensions in order to find similar clusters.