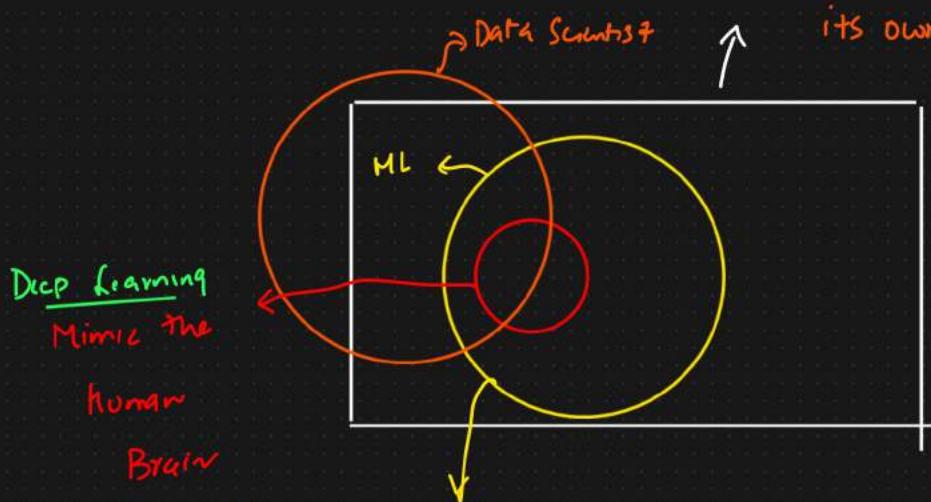


AI vs ML vs DL vs DS



AI → Smart Application that can perform its own task without any human intervention

Eg: Self Driving Car

# Robots

Alicks

- ① Statistics
  - ② Linear Algebra
  - ③ Calculus

{ Multi Layered NN }

## Eg: Object Detection

## Image recognition

CHATBOYS

## Recommendation

## Systems

It provides Stats tool to analyze, visualize, predictive models, forecasting

```

graph LR
    A[Amazon.in] --> B[Product]
    A --> C[Recommendation System]
    B --> D[Netflix]
    C --> D
    D --> E[Movies]
  
```

The diagram illustrates the architecture of a recommendation system. At the top, 'Amazon.in' is connected to two main components: 'Product' (indicated by a red bracket) and a 'Recommendation System' (also indicated by a red bracket). Below 'Amazon.in', there is a red bracket connecting 'Product' and the 'Recommendation System'. The 'Recommendation System' is then connected to 'Netflix' (indicated by a red bracket). Finally, 'Netflix' is connected to 'Movies' (indicated by a red bracket).

## Supervised, Unsupervised, Semi-Supervised, Reinforcement Learning

### Types of ML

- ① Supervised ML      → CLASSIFICATION
- REGRESSION
- ② Unsupervised ML
- ③ Semi-Supervised
- ④ Reinforcement Learning

- ① Supervised ML      → Classification
- Regression

(i) Dataset → O/p feature of the Data set

#### CLASSIFICATION

↓ O/p features → Dependent feature

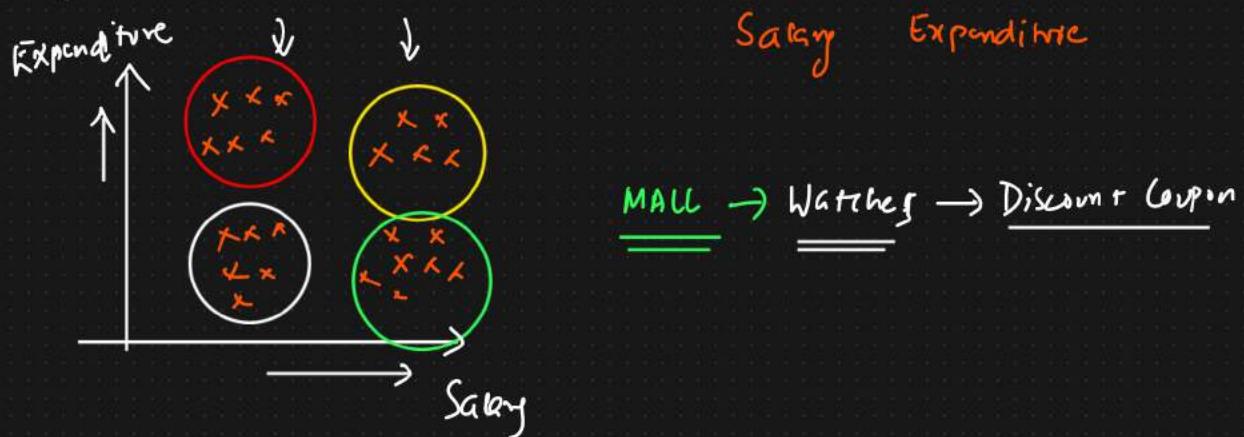
No. of hours played	No. of Study hours	Pars/Fail
8	2	Fail
7	3	Fail
6	4	Fail
5	5	Pass
4	6	Pass

Regression → O/p → continuous value

Size of house      No. of Room      Price of the House → continuous value

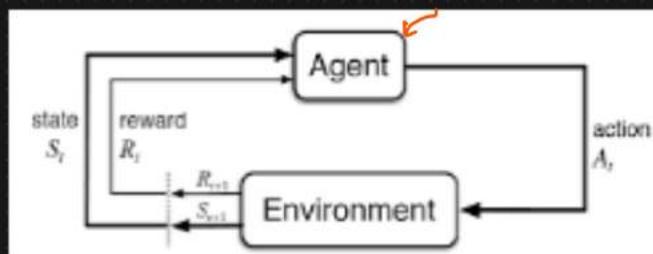
② Unsupervised ML  $\rightarrow$  No O/P  $\Rightarrow$  Clusters  $\rightarrow$  Group of Similar DATA

Eg: Customer Segmentation

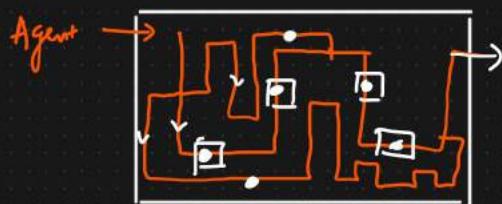


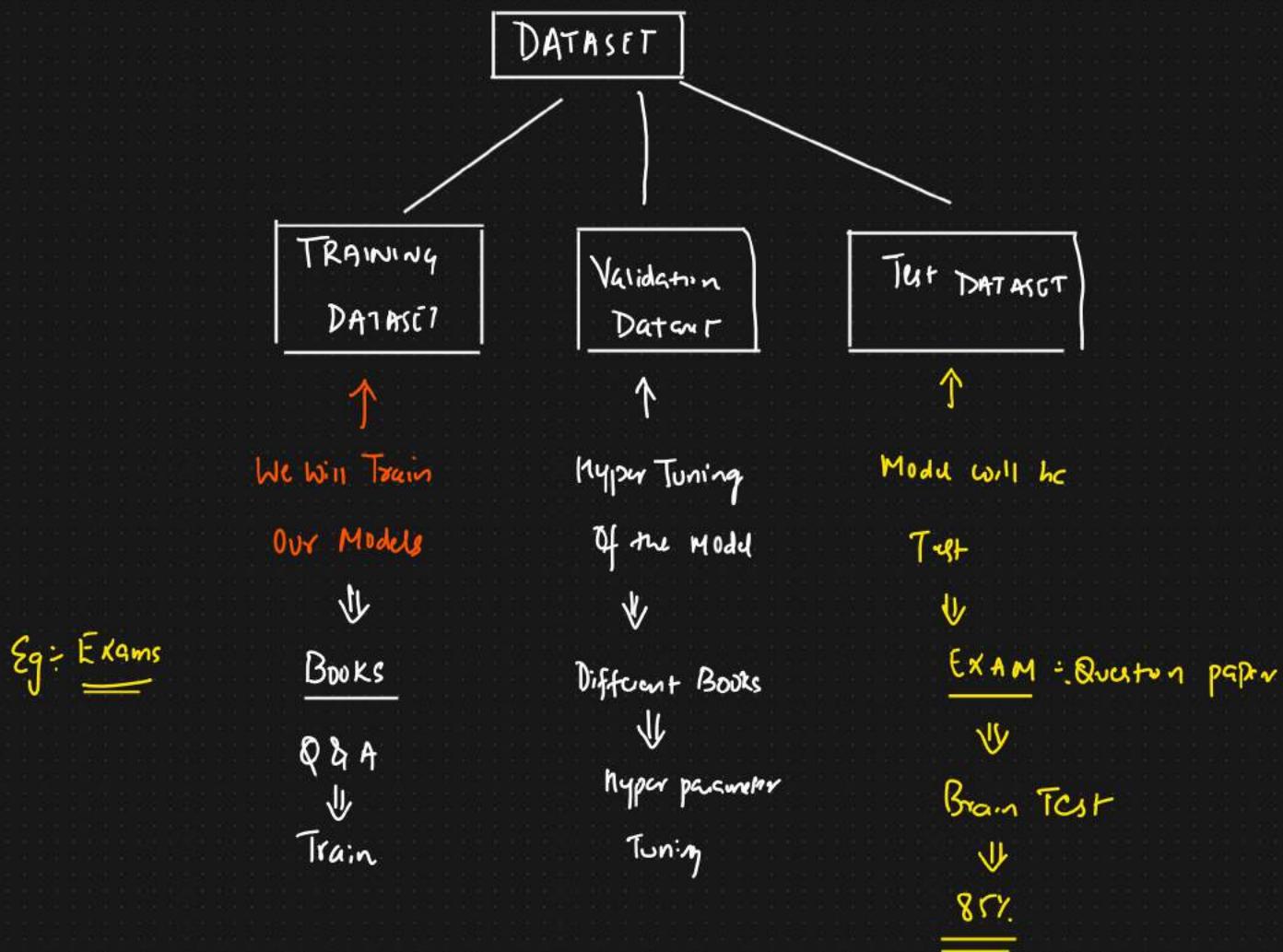
③ Semi Supervised : Supervised + Unsupervised

④ Reinforcement Learning :



Reinforcement learning is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. Reinforcement learning is one of three basic machine learning paradigms, alongside supervised learning and unsupervised learning.





## Machine Learning

- ① Model Performance - Accuracy ↑↑ → High
  - ② Overfitting, Underfitting
  - ③ Bias Vs Variance

## ① Overfitting Underfitting

## DATASET

- ① Books → Train → Model is Trained → Accuracy ↑↑ → 95% } Overfitting  
Exam → Test → Model is Tested → Accuracy ↓↓ → 60% } Low Bias  
High Variance

② Train → Accuracy ↓↓ → 55% } Underfitting  
Test → Accuracy ↓↓ → 50% } High Bias  
High Variance

## Generalized Model

↑

Train  $\rightarrow$  Acc  $\uparrow\uparrow$  }

Tcst → Acc ↑↑

1

{ low Bias  
low Variance }

## Feature Extraction

Feature Extraction is process of selecting and extracting the most important features from raw data.

ML Application → 1000 features



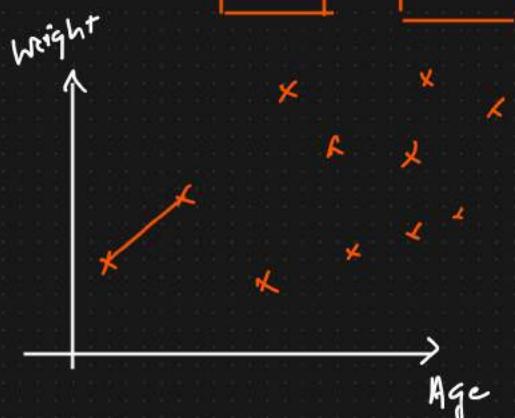
Most Important features



Machine Learning Algo.

### ① Feature Scaling

(y)	(Rg)	(cm)	BMI
Age	Weight	Height	
32	70	140cm	-
28	75	160cm	-
35	80	155cm	-



Normalize OR Standardize



$$Z\text{-score} = [\mu=0 \ \sigma=1]$$

$$Z\text{-score} = \frac{x_i - \bar{x}}{\sigma}$$



Standardization

Normalization [Min Max Scaler]

(0,1)

Unit Vector

② Feature Selection → We just pick the most Important feature

500 features → Top 10 features



ML Model Train

Correlation



Feature Selection

① Filter Method    ② Embedded Method

③ PCA {Principal Component Analysis}.

# Feature Scaling

① Standardization  $\rightarrow Z\text{-score}$

② Min Max Scaling {Normalization}

③ Unit vector.

① Standardization

Age

24

25

26

27

28

29

$$Z\text{-score}_i = \frac{x_i - \bar{x}}{\sigma}$$

Age'

$\mu=0$   $\sigma=1$

$\mu=0, \sigma=1$

② Normalization [Min Max Scaler]  $\rightarrow 0$  to  $1$

Age

24

25

26

27

28

29

30

TRANSFORMATION



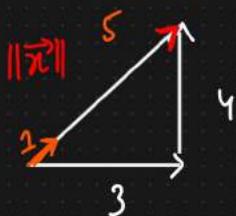
Age' [0-1]

$$x_{\text{Scaled}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

③ Unit Vector - Magnitude of 1

$$\vec{v} = (3, 4)$$

$$\|\vec{v}\| = \sqrt{(3)^2 + (4)^2} = \sqrt{25}$$
$$= 5$$



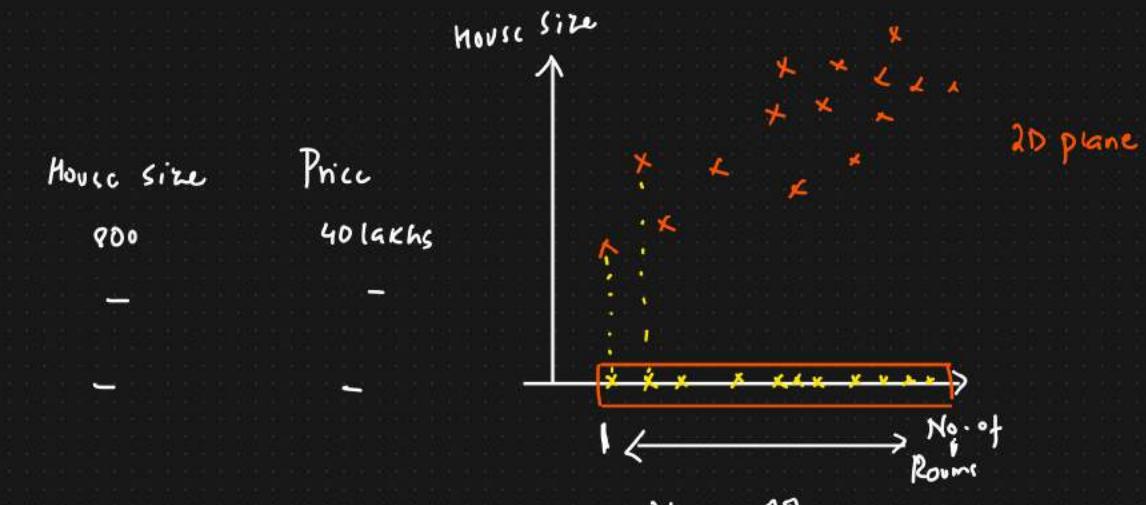
$$\hat{v} = \left( \frac{3}{\|\vec{v}\|}, \frac{4}{\|\vec{v}\|} \right) = \left( \frac{3}{5}, \frac{4}{5} \right)$$

$$\|\hat{v}\| = \sqrt{\left(\frac{3}{5}\right)^2 + \left(\frac{4}{5}\right)^2} = \sqrt{\frac{9+16}{25}}$$
$$= \sqrt{\frac{25}{25}} = 1$$

# PCA {Principal Component Analysis}

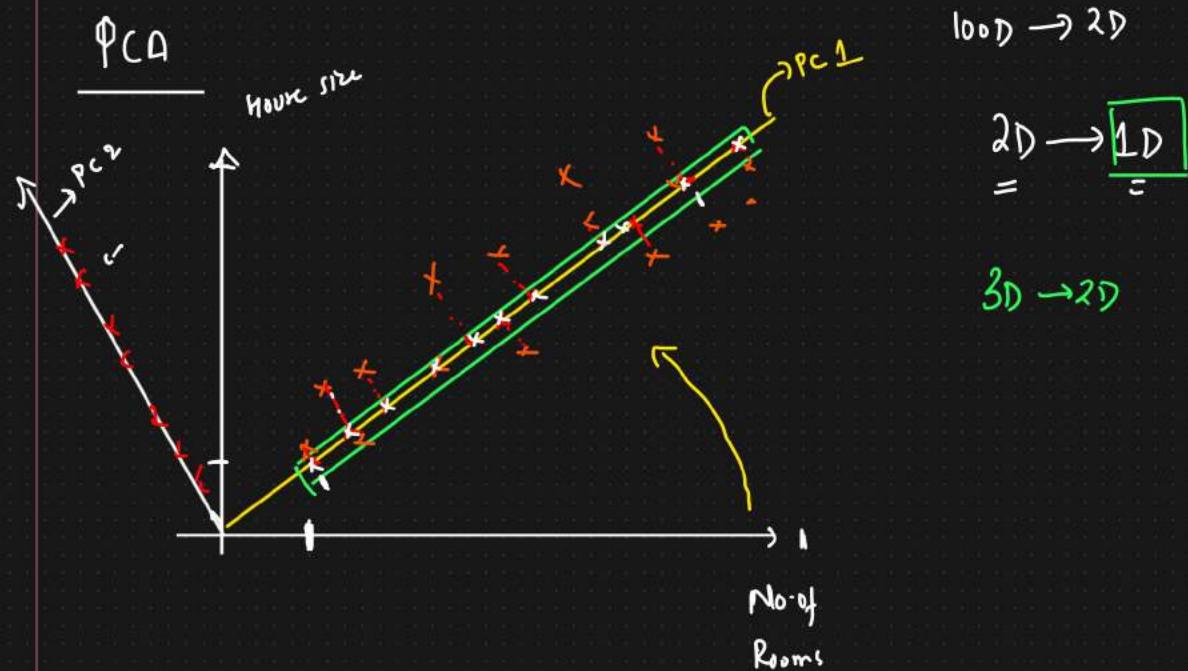
Dataset

No. of Rooms	House Size	Price
2	800	40 lakhs
1	-	-
-	-	-



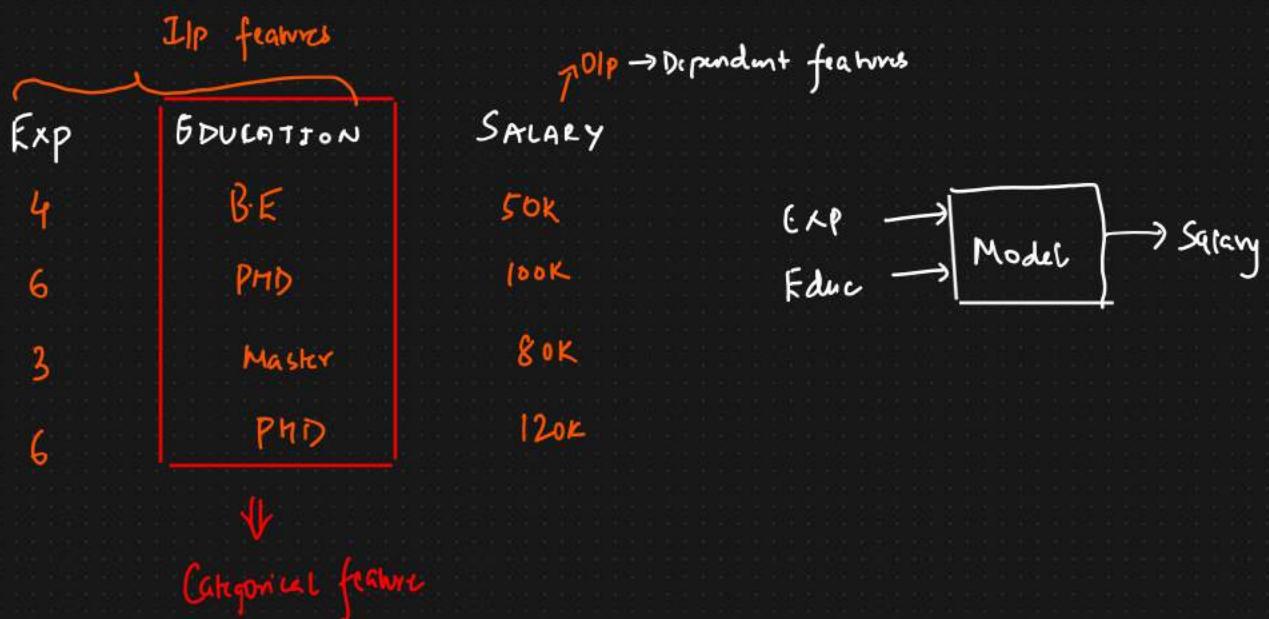
No. of Rooms      Price      Machine Learning Model

Problem : Data less



No. of Room    House size     $\boxed{[PC1 \quad Price]} \Rightarrow$  Model Train

# DATA ENCODING



Data Encoding → Categorical feature → Numerical  
↓  
ML Model  
↓  
TRAIN.

## Types of Data Encoding

- ① Nominal / OHE (One Hot Encoding)
  - ② ORDINAL And LABEL ENCODING
  - ③ Target Guided Ordinal Encoding
- Categorical  
↓  
Numerical

## Nominal / One Hot Encoding Techniques

Nominal encoding is a technique used to transform categorical variables that have no intrinsic ordering into numerical values that can be used in machine learning models. One common method for nominal encoding is one-hot encoding, which creates a binary vector for each category in the variable.

House Price (DATASET)

Location → OHE

No. of Rooms	House Size	Location	Price	Bangalore	Delhi	Noida
		Bangalore		[ 1      0      0 ]		
		Noida		0      0      1		
		Delhi		0      1      0		
		Bangalore		1      0      0		
		Delhi		0      1      0		
		Noida				

Disadvantage [10 location]

① Sparse matrix [overfitting]

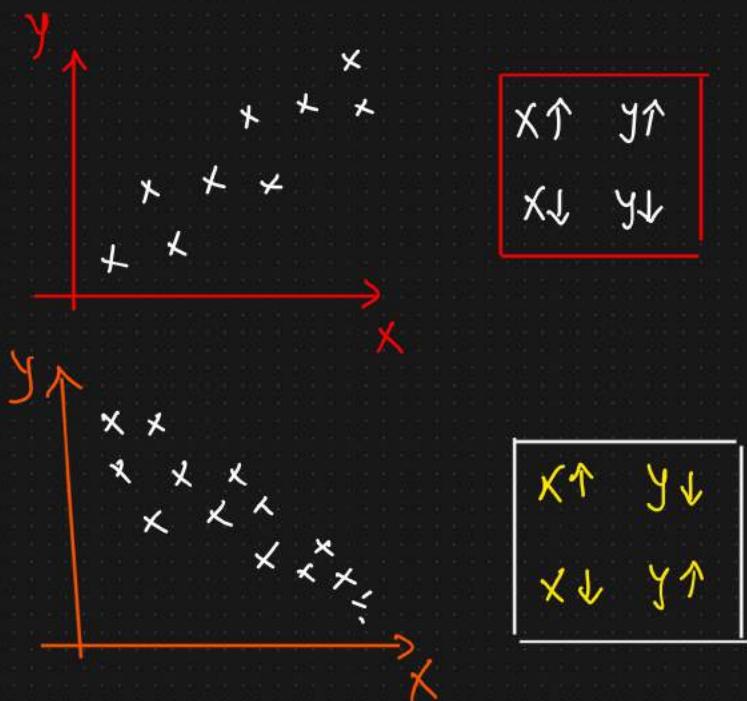
② 1000 features [1000 categories]  
↓

If will increase the number of the features

## Covariance And Correlation

[ Relationship between X and Y ]

X	Y	→	X↑ Y↑	↓↑ Size of House	↓ O/P Price of House ↑
2	3		X↓ Y↑		
4	5	→	X↓ Y↓		
6	7		X↑ Y↓		
8	9				



### Covariance

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

$$\begin{aligned} \text{Var}(x) &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{n-1} \end{aligned}$$

$x_i$  → Data points of x

$\bar{x}$  → Sample mean of x

$y_i$  → Data points of y

$\bar{y}$  → Sample mean of y

$\approx \text{Cov}(x, x) \Rightarrow \text{Spread}$

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

$X \uparrow$	$y \uparrow$	+ve Covariance
$X \downarrow$	$y \downarrow$	

$X \uparrow$	$y \downarrow$	-ve Covariance
$X \downarrow$	$y \uparrow$	

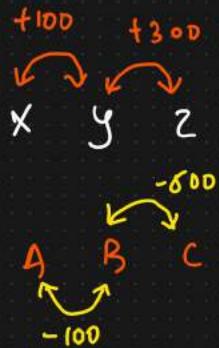
$X$	$y$
→ 2	3
→ 4	5
→ 6	7
$\bar{x} = 4$	$\bar{y} = 5$

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

$$= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{n-1}$$

$$= \frac{4+0+4}{2} = \frac{8}{2} = 4 \text{ tve value}$$

Positive  
Covarianu



$x$  &  $y$  are having a positive Covarianu

### Advantages

- ① Relationship between  $x$  and  $y$   
tve or -ve value

### Disadvantages

- ① Covariance does not have a specific limit value

## ② Pearson Correlation Coefficient [-1 to 1]

$$r_{x,y} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

- ① The more the value towards +1 the more +ve correlated it is  $(x, y)$
- ② The more the value towards -1 the more -ve correlated it is  $(x, y)$

## ③ Spearman Rank Correlation [-1 to 1]

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

$x$	$y$	$R(x)$	$R(y)$
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

## Feature Screen

+ve Size of House ↑	+ve No. of Rooms ↑	+ve Location ↑	≈ 0 No. of people staying	-ve Haunted	Price ↑
------------------------	-----------------------	-------------------	------------------------------	----------------	---------

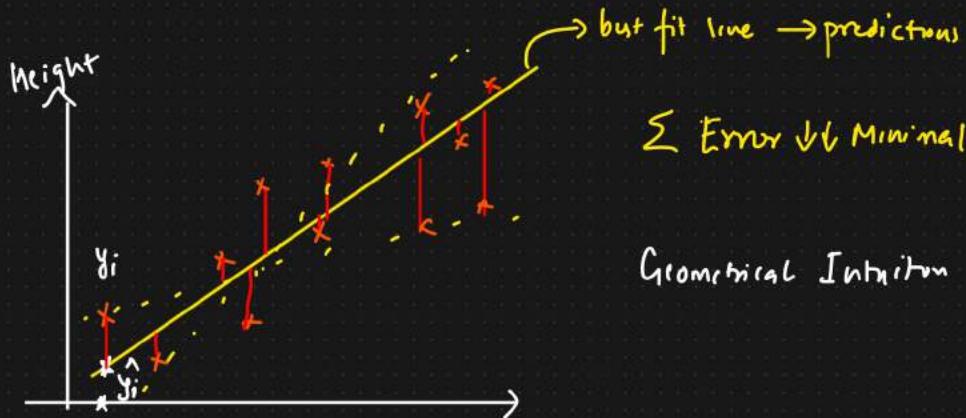
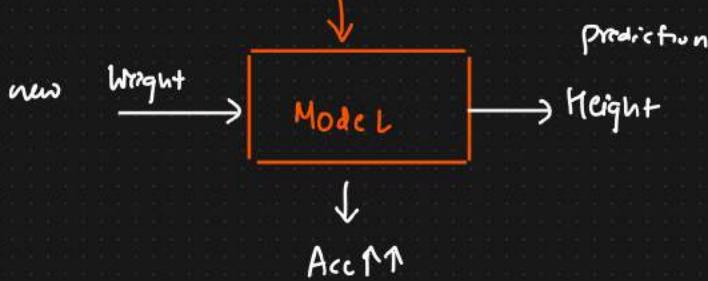
# Simple Linear Regression



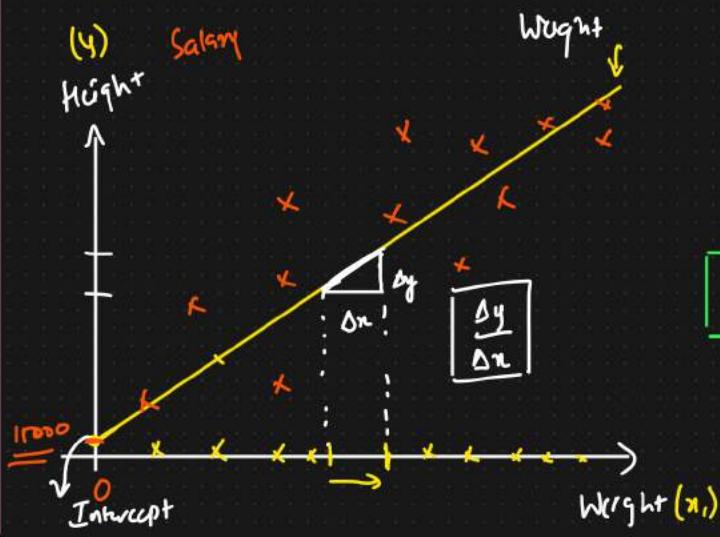
Data

	Independent feature	Dependent feature
Weight		Height $\hat{y}_p$
74		new weight $\rightarrow$ Height
80	170	
75	180	
-	175.5	
-	-	
-	-	
-	-	

TRAIN DATASET



Geometrical Intuition



$$\hat{y} = mx + c$$

$$\hat{y} = \beta_0 + \beta_1 x_i$$

$$h_{\theta}(x) = \underline{\theta_0} + \underline{\theta_1 x_i}$$

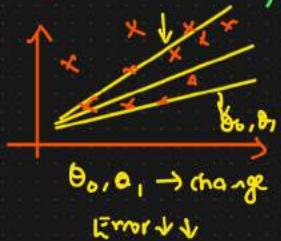
$\theta_0 \Rightarrow$  Intercept

$\theta_1 \Rightarrow$  Slope or Coefficient

Simple Linear Regression  
↑  
Model

↓

Hypothesis Testing



Experience  $\lambda_1 \Rightarrow$  Data point

Cost function [Error].

$$f(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n \left( \underset{\text{Actual}}{\downarrow} y_i - \underset{\text{predicted}}{\downarrow} \hat{y}_i \right)^2 \quad \left[ \text{Mean Squared Error} \right]$$

$n$ : no. of datapoints

$y_i$ : Actual value

$\hat{y}_i$ : predicted value

Final Aim [In order to get best fit line]

$$\text{Minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



$$\hat{y}_i = \theta_0 + \theta_1 x_i$$

Let's consider  $\theta_0 = 0$

$$\hat{y}_i = \theta_1 x_i$$

Let  $\theta_1 = 1$

Let  $\theta_1 = 0.5$

Dataset		
x	y	$\hat{y}_i$
1	1	1
2	2	2
3	3	3

Predicted

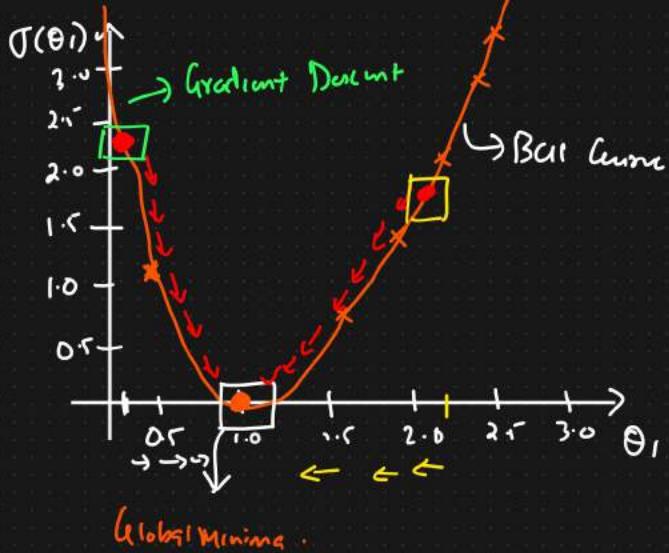
$$\left\{ \begin{array}{lll} x=1 & \hat{y}_i = 1(1) & \hat{y}_i = 0.5 \\ & = 1 & = 0.5 \\ x=2 & \hat{y}_i = 1(2) & \hat{y}_i = 1.0 \\ & = 2 & = 1.0 \\ x=3 & \hat{y}_i = 1(3) & \hat{y}_i = 1.5 \\ & = 3 & = 1.5 \end{array} \right.$$

Cost fn       $\theta_1 = 1$

$$J(\theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x)_i)^2$$

$$= \frac{1}{3} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(1) = 0$$



Cost fn       $\theta_1 = 0.5$

$$J(\theta_1) = \frac{1}{3} [(1-0.5)^2 + (2-1)^2 + (3-1.5)^2]$$

$$\underline{J(\theta_1)} = \underline{\underline{1.16}}$$

Cost fn       $\theta_1 = 0$

$$J(\theta_1) = \frac{1}{3} [(1-0)^2 + (2-0)^2 + (3-0)^2]$$

$$J(\theta_1) = 4.66$$

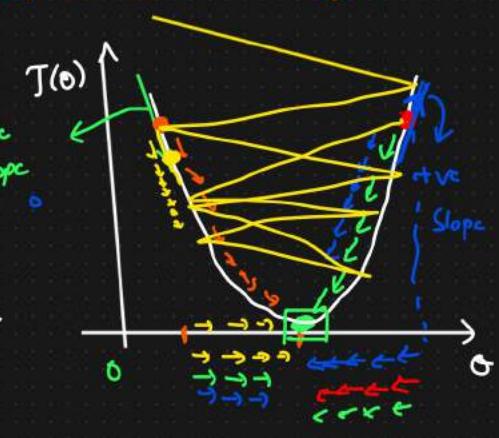
Convergence Algorithm      {Optimize the change of  $\theta_0, \theta_1$  to Global Minima}.

Repeat until Convergence

$$\left\{ \begin{array}{l} j=0,1 \\ \downarrow \downarrow \end{array} \right.$$

$$\theta_j : \theta_j - \alpha \left[ \frac{\partial J(\theta_j)}{\partial \theta_j} \right] \Rightarrow \text{Determine Slope of}$$

$$\alpha = 0.01$$



Learning  
Rate :

a point

$$\textcircled{1} \quad \theta_1 = \theta_1 - \alpha (-\text{ve})$$

$$\theta_{1\text{new}} = \theta_{1\text{old}} + (\text{value})$$

$$\theta_{1\text{new}} > \theta_{1\text{old}}$$

Learning Rate  $\Rightarrow$  Speed of Convergence

$$\textcircled{2} \quad \theta_{1\text{new}} = \theta_{1\text{old}} - \alpha (+\text{ve})$$

$$\theta_{1\text{new}} = \theta_{1\text{old}} - (\text{value})$$

$$\theta_{1\text{new}} < \theta_{1\text{old}}$$

### Conclusion

Repeat until Convergence

{

$$\theta_j : \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

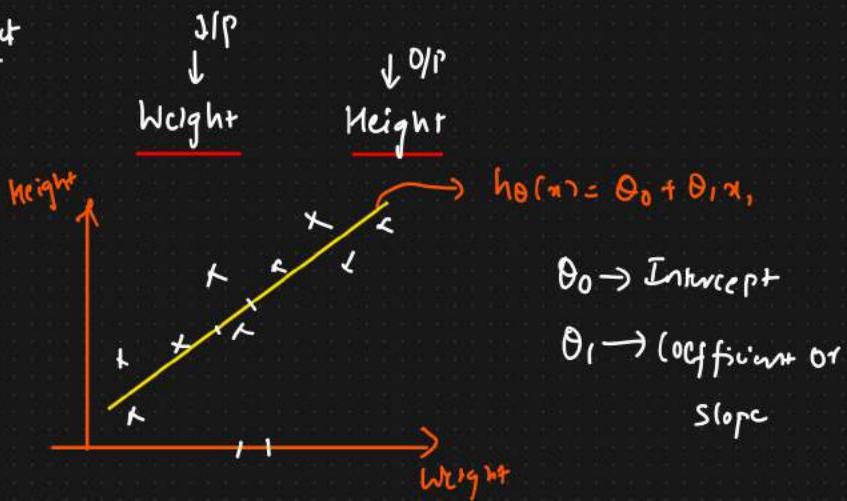
$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2$$



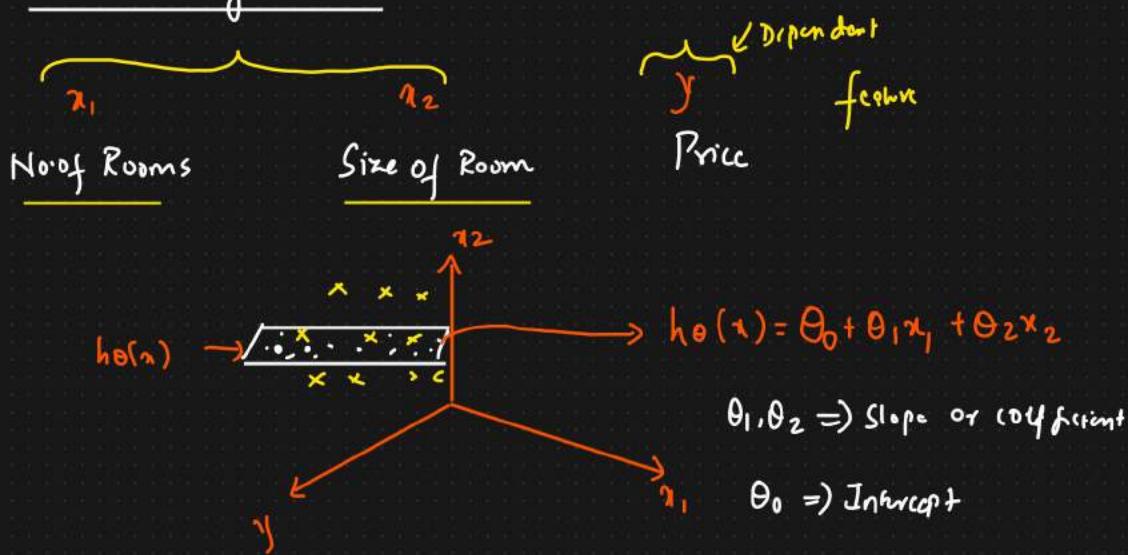
Mean Square Error

## Multiple Linear Regression

Dataset



House Pricing Dataset

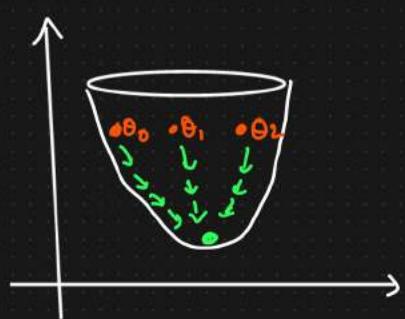
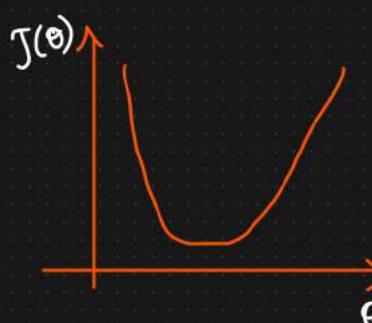


Simple Linear Regression

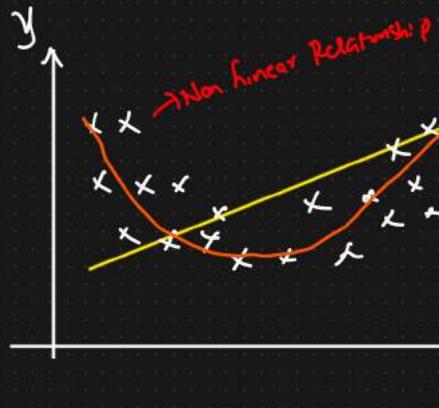
Generic Equation Multiple Linear Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$



# Polynomial Regression



## Simple Linear Regression

$$h_0(x) = \theta_0 + \theta_1 x,$$

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

↓  
Multiple Linear Regression

## Polynomial Regression

→ Simple Polynomial Reg.

→ Multiple Polynomial Reg.

↓  
Polynomial Degrees

### Simple Polynomial Regression

polynomial degree = 0

$$h_0(x) = \theta_0 x^0 = \theta_0 \overset{\circ}{x_i} \Rightarrow \text{Constant}$$

polynomial degree = 1

$$h_0(x) = \theta_0 \overset{\circ}{x_i} + \theta_1 x_i^1 \Rightarrow \text{Simple Linear Reg.}$$

polynomial degree = 2

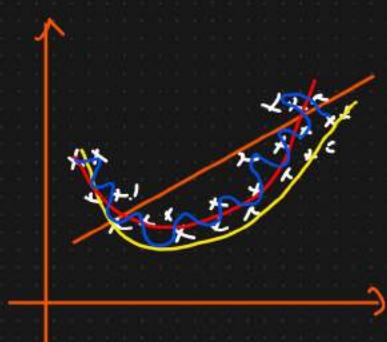
$$h_0(x) = \theta_0 \overset{\circ}{x_i} + \theta_1 x_i^1 + \theta_2 x_i^2$$

$$h_0(x) = \theta_0 \overset{\circ}{x_i} + \theta_1 x_i^1 + \theta_2 x_i^2 + \theta_3 x_i^3$$

⋮

polynomial degree = n

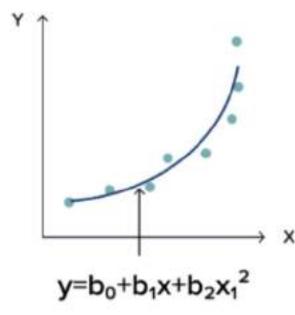
$$h_0(x) = \theta_0 \overset{\circ}{x_i} + \theta_1 x_i^1 + \theta_2 x_i^2 + \dots + \theta_n x_i^n$$



Simple linear model



Polynomial model



## Multiple Polynomial Regression

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \{ \text{Multiple Linear Regression} \}$$

Polynomial Degree = 2

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_3^2$$

# Performance Metrics Used In Regression

① R Squared

② Adjusted R Squared

① R Squared

$$R^{\text{Squared}} = 1 - \frac{SS_{\text{Res}}}{SS_{\text{Total}}} \quad \begin{array}{l} \rightarrow \text{Error} \\ \{ \text{Average of } y_i \text{ vs } \hat{y}_i \} \rightarrow \text{Error} \end{array}$$

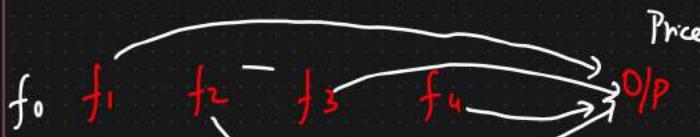
$SS_{\text{Res}}$  = Sum of square Residual {Error}

$SS_{\text{Total}}$  = Sum of Square Total

$$R^{\text{Squared}} = 1 - \frac{\sum_{i=1}^n (y_i - h_{\theta}(x))^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \Rightarrow \text{Small} \quad \begin{array}{l} \nearrow SS_{\text{Res}} \\ \searrow \text{Big} \end{array}$$

$$R^{\text{Squared}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

R Squared ranges between 0 to 1

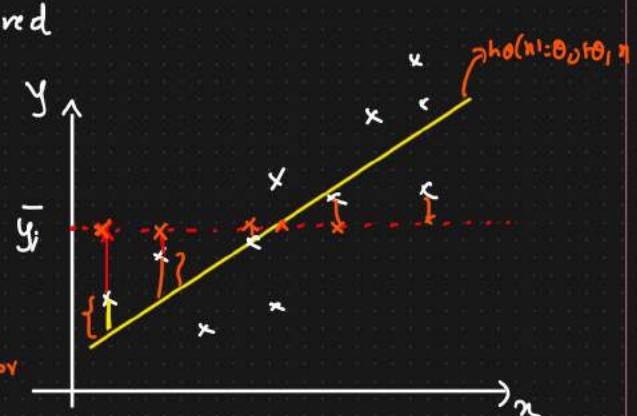


Price

$R^{\text{Squared}} = 75\% = 0.75$

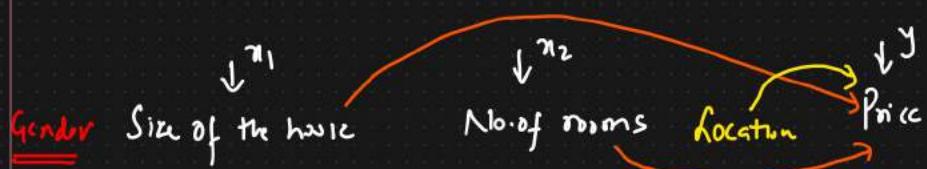
$R^{\text{Squared}} = 82\%$

$R^{\text{Squared}} = 90\% \uparrow \uparrow \uparrow$



$$R^2_{\text{Squard}} = 92\%, \text{PA}$$

## ② Adjusted R squared



$$R_{-} \text{ squared} = 85\% = 0.85$$

$$R^2 - \text{Squared} = 90\% = 0.90$$

$$\text{Adjusted } R \text{ square} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

$N = \text{no. of dots/parts}$

$$R^2 = R \text{ Squared}$$

P = No. of Independent features.

$$R^L = 0.8 \quad N=11 \quad p=2$$

P=3 → feature

$$\text{Adjusted } R^2 = 1 - \left[ \frac{(0.2)(10)}{11-2-1} \right] = 0.75$$

$$R^2 \gg \text{Adj} \nu_{\text{fid}} R^2$$

$$R^2 = 80\% \quad \text{Adjusted } R^2 = 75\%$$

$$P=3 \quad R^2 = 85\% \quad \text{Adjusted } R^2 = 78\%$$

$P=4$   $\uparrow R^2 = 87\%$  Adjusted  $R^2 = 76\%$   $\downarrow \downarrow$

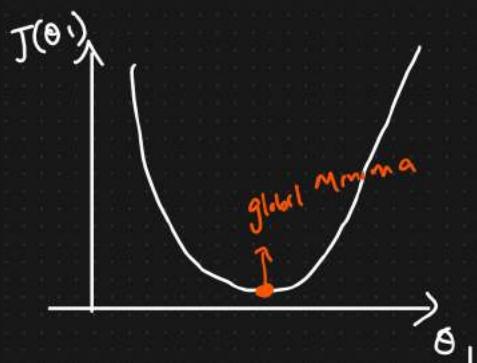
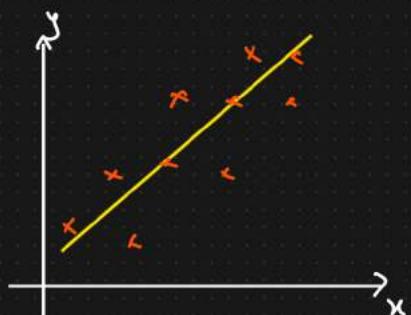
Independent factor is not that important

MSE, MAE, RMSE [Cost function]

① Mean Squared Error (MSE) ✓

② Mean Absolute Error (MAE)

③ Root Mean Squared Error (RMSE)



$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2$$



Mean Squared Error

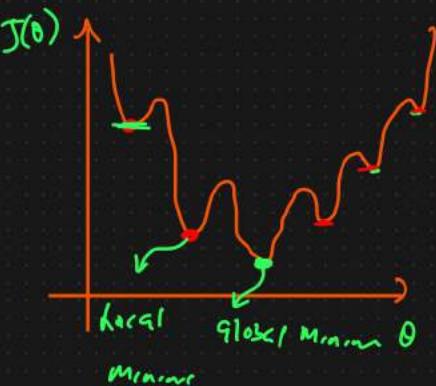
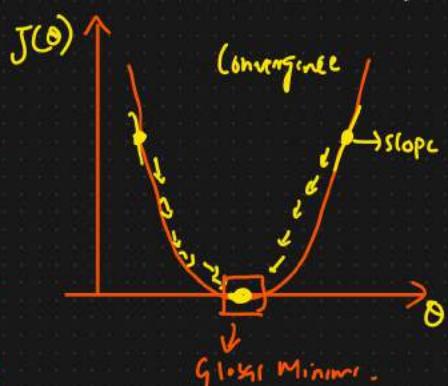
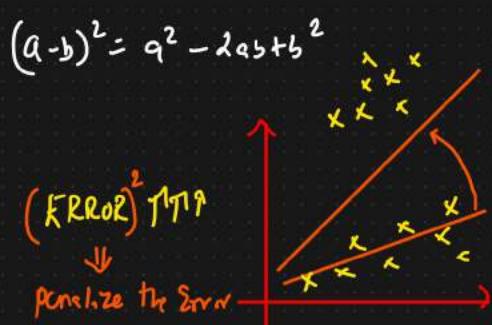
① Mean Squared Error (MSE)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2$$

$$\boxed{MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \rightarrow \text{Quadratic Equation}$$

↙ ↘ Convex function

Non Convex function



Advantage

Disadvantage

① Equation is differentiable

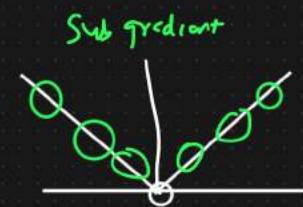
① Not Robust to outliers

② It has only one local  
or global minima.

② It is not in the  
same unit

## ② Mean Absolute Error (MAE)

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



### Advantage

① Robust to outliers

### Disadvantage

② It will be in the same unit

① Convergence usually takes more time

## ③ RMSE (Root Mean Squared Error)

$$\begin{aligned} RMSE &= \sqrt{MSE} \\ &= \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2} \end{aligned}$$

### Advantages

① Same Unit

### Disadvantage

① Not Robust to outliers.

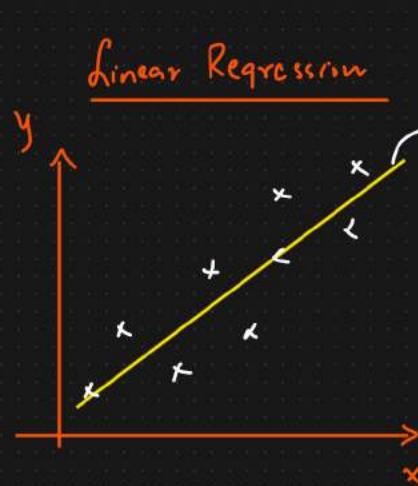
② Differentiable

## Note: Linear Regression

Performance Metrics =  $R^2$  and Adjusted  $R^2 \Rightarrow$  Acc of model

Cost function  $\rightarrow$  Error  $\rightarrow$  MSE, MAE, RMSE

# Ridge, Lasso And Elasticnet Regression



Independent



$$J(\theta)$$

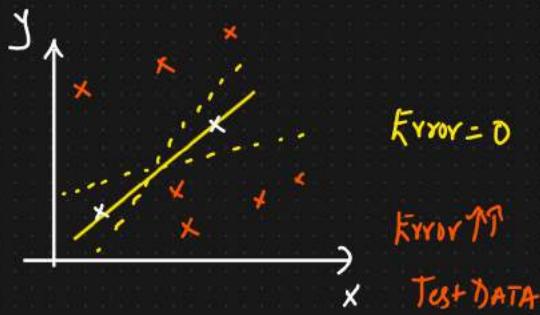
Gradient Descent

Optimal Minima



$$\text{Cost fn: } \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad [\text{Mean Squared Error}]$$

## ① Ridge Regression (L2 Regularization) → Reducing Overfitting



Overfitting

TRAIN → Acc ↑↑ → R<sup>2</sup>

TEST → Acc ↓↓ → R<sup>2</sup>

$$h_\theta(x) = \theta_0 + \theta_1 x$$

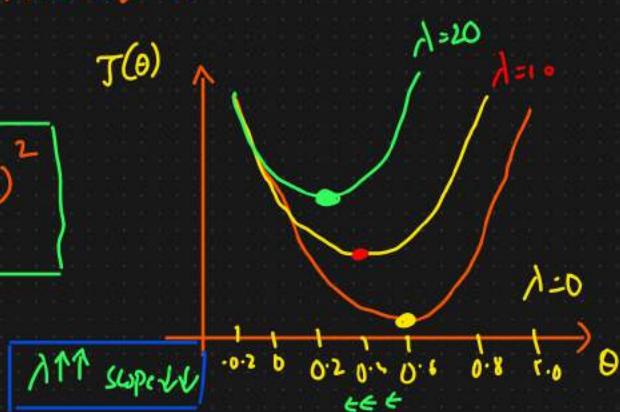
$$\lambda = 1$$

$$J(\theta)$$

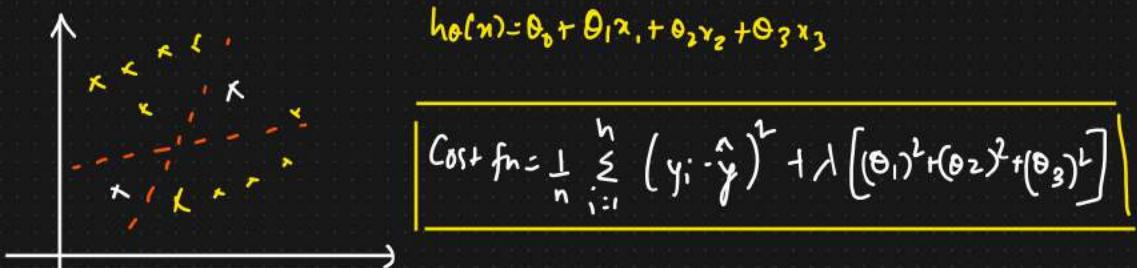
$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n (\text{slope}_i)^2}$$

↓  
0

Hypersparameter



$$= 0 + 1 [(\theta_1)^2] \leftarrow \text{Penalize the cost function}$$



② Lasso Regression (d, Regularization) → Feature Selection

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \text{feature remove}$$

$$= 0.52 + 0.65 \underline{\underline{x_1}} + 1.5 \underline{\underline{x_2}} + \boxed{0.2} \cancel{x_3}$$

↓

feature

Selection

$\underline{\underline{x_1}} \rightarrow 1 \text{ unit}$

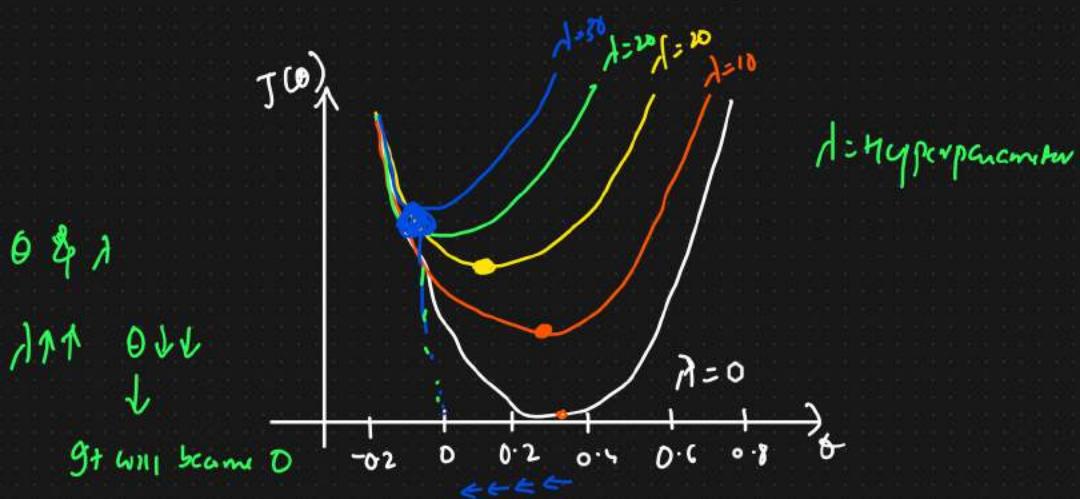
$\underline{\underline{x_2}} \rightarrow 1.5 \text{ unit}$

$x_3 \rightarrow 4 \text{ unit}$

$y \rightarrow 0.2 x_3$

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n |\text{slope}|}$$

$$= \text{Error} + 1 [|\theta_0| + |\theta_1| + |\theta_2|]$$



### ③ ElasticNet Regression

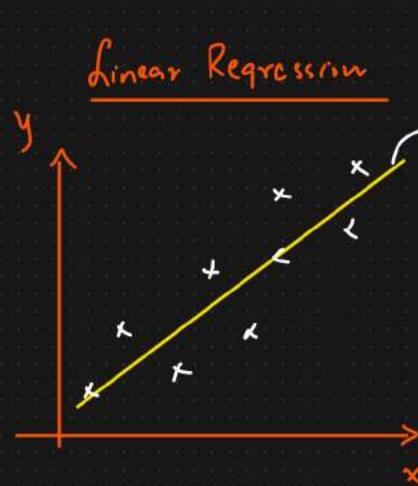
→ Reducing Overfitting → Ridge

→ Feature Selection → Lasso

$$\text{Cost fn} = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\downarrow \text{MSE}} + \lambda_1 \underbrace{\sum_{i=1}^n (\text{slope})^2}_{\downarrow \text{Reducing overfitting}} + \lambda_2 \underbrace{\sum_{i=1}^n |\text{slope}|}_{\downarrow \text{Feature selection}},$$

$\lambda_1, \lambda_2$  {Hyperparameter Tuning}.

# Ridge, Lasso And Elasticnet Regression



Independent



$$J(\theta)$$

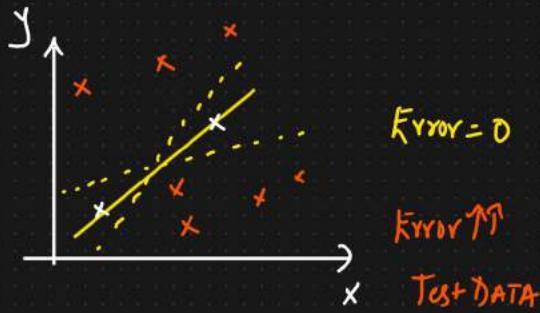
Gradient Descent

Optimal Minima



$$\text{Cost fn: } \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad [\text{Mean Squared Error}]$$

## ① Ridge Regression (L2 Regularization) → Reducing Overfitting



Overfitting

TRAIN → Acc ↑↑ → R<sup>2</sup>

TEST → Acc ↓↓ → R<sup>2</sup>

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$\lambda = 1$$

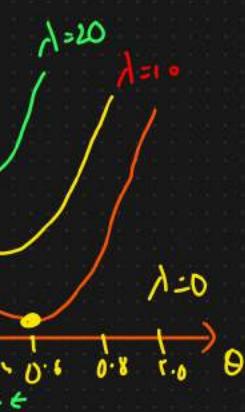
$$J(\theta)$$

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n (\text{slope}_i)^2$$

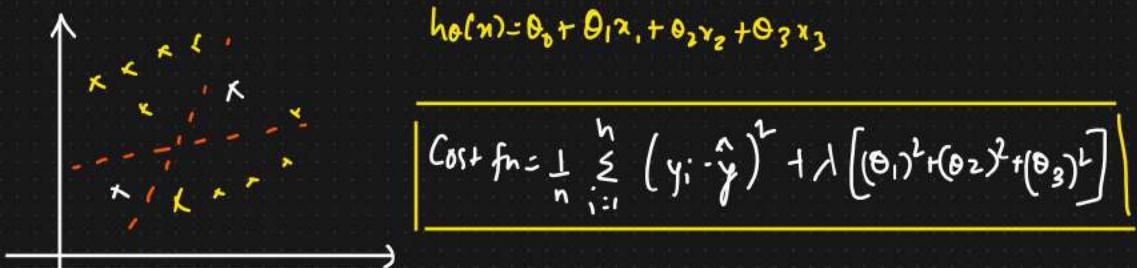
↓  
0

Hypersparameter

$\lambda \uparrow \uparrow \text{ slope} \downarrow \downarrow$



$$= 0 + 1 [(\theta_1)^2] \leftarrow \text{Penalize the cost function}$$



② Lasso Regression (d, Regularization) → Feature Selection

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \text{feature remove}$$

$$= 0.52 + 0.65 \underline{\underline{x_1}} + 1.5 \underline{\underline{x_2}} + \boxed{0.2} \cancel{x_3}$$

↓

feature

Selection

$\underline{\underline{x_1}} \rightarrow 1 \text{ unit}$

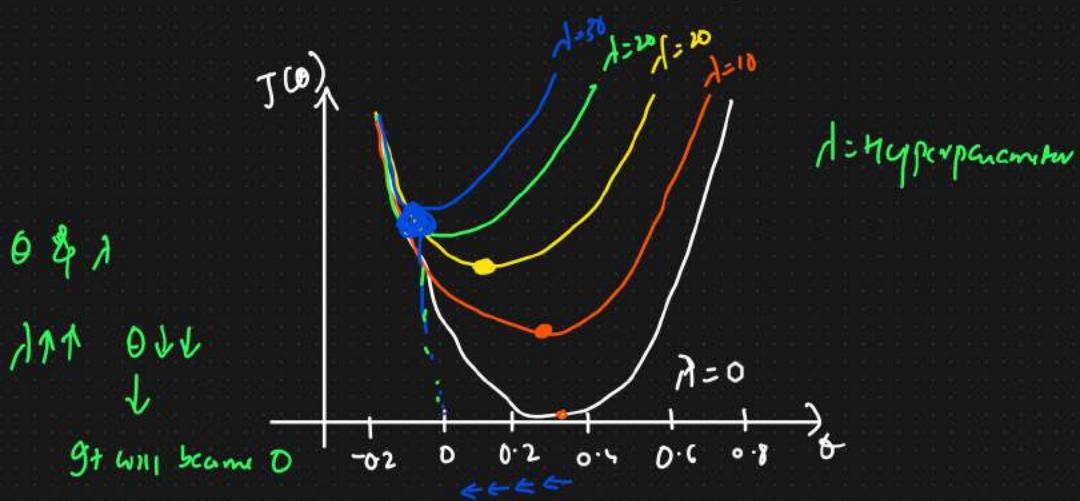
$\underline{\underline{x_2}} \rightarrow 1.5 \text{ unit}$

$x_3 \rightarrow 4 \text{ unit}$

$y \rightarrow 0.2 x_3$

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n |\text{slope}|}$$

$$= \text{Error} + 1 [|\theta_0| + |\theta_1| + |\theta_2|]$$



### ③ ElasticNet Regression

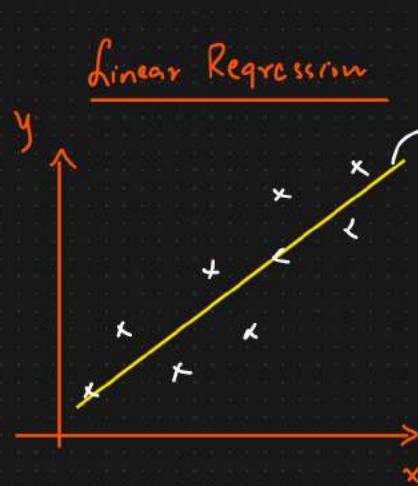
→ Reducing Overfitting → Ridge

→ Feature Selection → Lasso

$$\text{Cost fn} = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\downarrow \text{MSE}} + \lambda_1 \underbrace{\sum_{i=1}^n (\text{slope})^2}_{\downarrow \text{Reducing overfitting}} + \lambda_2 \underbrace{\sum_{i=1}^n |\text{slope}|}_{\downarrow \text{Feature selection}},$$

$\lambda_1, \lambda_2$  {Hyperparameter Tuning}.

# Ridge, Lasso And Elasticnet Regression



Independent



$$J(\theta)$$

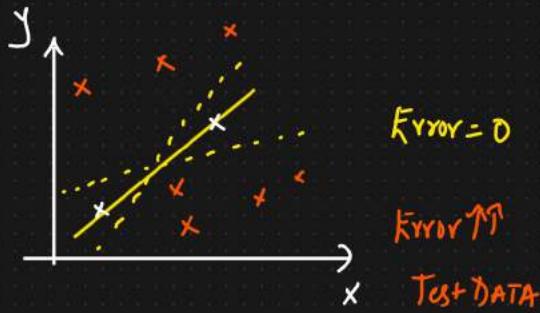
Gradient Descent

Optimal Minima



$$\text{Cost fn: } \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad [\text{Mean Squared Error}]$$

## ① Ridge Regression (L2 Regularization) → Reducing Overfitting



Overfitting

TRAIN → Acc ↑↑ → R<sup>2</sup>

TEST → Acc ↓↓ → R<sup>2</sup>

$$h_\theta(x) = \theta_0 + \theta_1 x$$

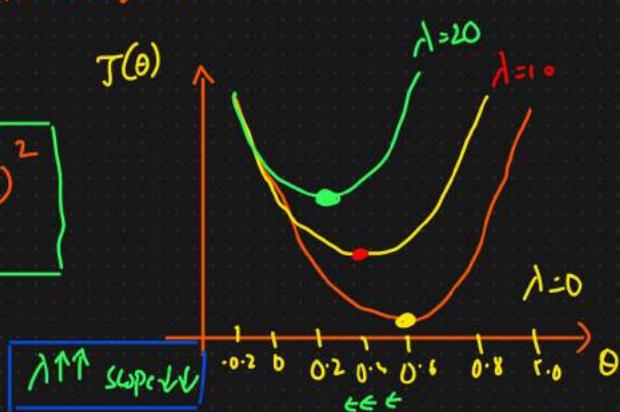
$$\lambda = 1$$

$$J(\theta)$$

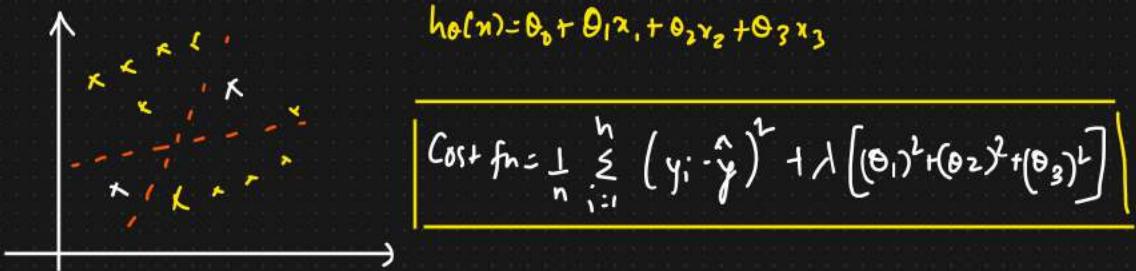
$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n (\text{slope}_i)^2}$$

↓  
0

Hypersparameter



$$= 0 + 1 [(\theta_1)^2] \leftarrow \text{Penalize the cost function}$$



② Lasso Regression ( $\ell_1$  Regularization)  $\rightarrow$  Feature Selection

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n |\text{Slope}|}$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \text{feature remove}$$

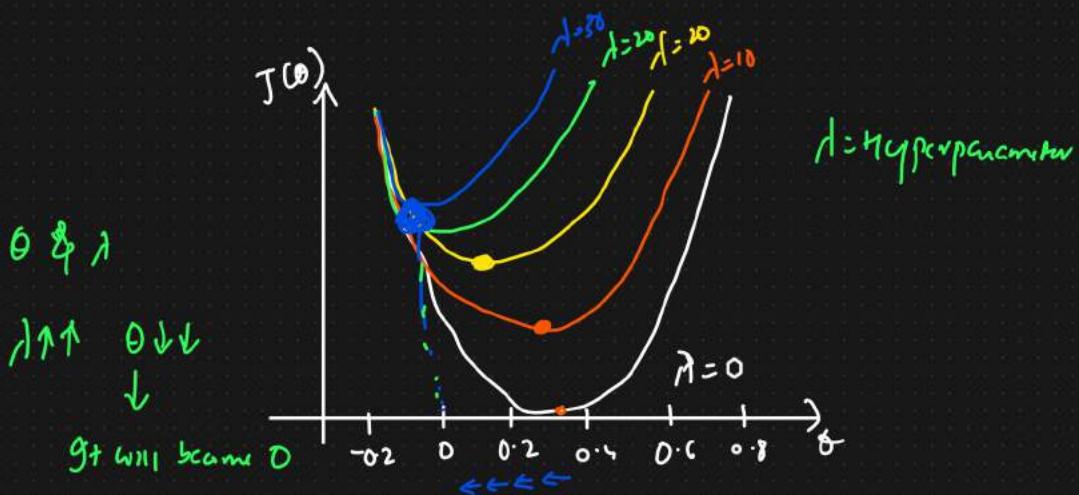
$$= 0.52 + 0.65 x_1 + 1.5 x_2 + \cancel{0.2 x_3}$$

↓  
 feature selection  
 prediction

$x_2 \rightarrow 1 \text{ unit}$   
 $y \rightarrow 1.5 x_2$   
 $x_3 \rightarrow 4 \text{ unit}$   
 $y \rightarrow 0.2 x_3$

$$\text{Cost fn} : \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \boxed{\lambda \sum_{i=1}^n |\text{Slope}|} \quad \lambda = 1$$

$$= \text{Error} + 1 [\lvert \theta_0 \rvert + \lvert \theta_1 \rvert + \lvert \theta_2 \rvert + \lvert \theta_3 \rvert]$$



### ③ ElasticNet Regression

→ Reducing Overfitting → Ridge

→ Feature Selection → Lasso

$$\text{Cost fn} = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\downarrow \text{MSE}} + \lambda_1 \underbrace{\sum_{i=1}^n (\text{slope})^2}_{\downarrow \text{Reducing overfitting}} + \lambda_2 \underbrace{\sum_{i=1}^n |\text{slope}|}_{\downarrow \text{Feature selection}},$$

$\lambda_1, \lambda_2$  {Hyperparameter Tuning}.

# Logistic Regression

To solve classification problem

- Binary classification
- Multiclass classification

## Dataset

Independent feature

No. of play hours

9

8

7

6

5

4

3

2

Dependent or op feature

Pass/Fail of y<sub>j</sub>.  $\hat{y} \rightarrow$  predicted

Fail 0

Fail 0

Fail 0

Fail 0

Pass 1

Pass 1

Pass 1

Fail 0

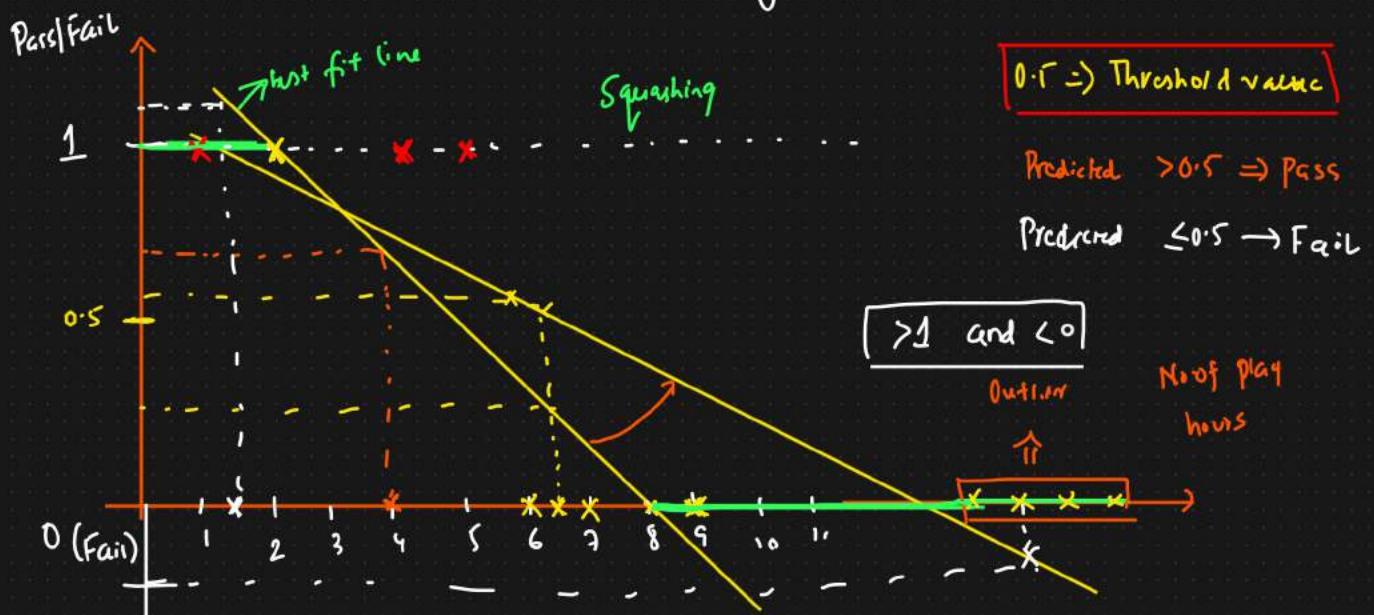
No. of play hours

TRAIN

Model

Accuracy M

Can we solve this classification problem using Regression?



- ① Best fit line changes because of outliers  $\rightarrow$  prediction goes wrong

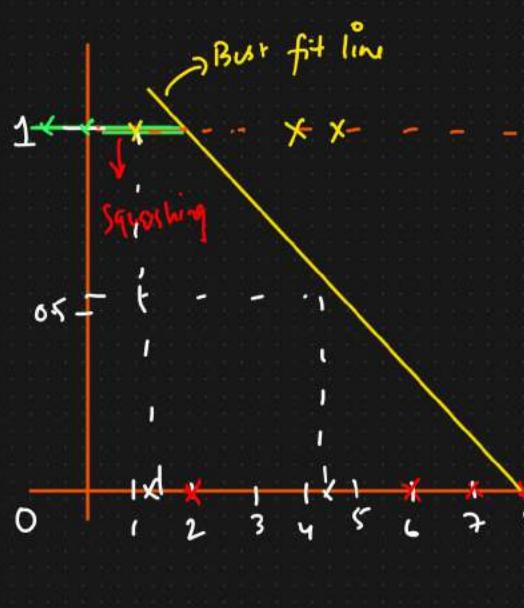
② The outcome comes  $>1$  and  $<0$  also

To solve this problem we use Logistic Regression

$$\boxed{0 \text{ to } 1} \Rightarrow \text{Squashing}$$

Technique

## How Logistic Regression Solves Classification Problem



$$L = \boxed{h_{\theta}(x) = \theta_0 + \theta_1 x_1} \rightarrow \text{Best fit line}$$

$\Downarrow$   
[Sigmoid Activation function]

$0 \text{ to } 1$

$$\boxed{f(x) = \frac{1}{1+e^{-x}}} \Rightarrow \boxed{0 \text{ to } 1}$$

$$h_{\theta}(x) = f(\theta_0 + \theta_1 x_1)$$

$$\boxed{h_{\theta}(x) = \frac{1}{1+e^{-x}}} \Rightarrow L = \theta_0 + \theta_1 x_1 \Rightarrow$$

$$\boxed{h_{\theta}(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x_1)}}}$$

$\hookrightarrow$  Logistic Regression

## Linear Regression Cost function

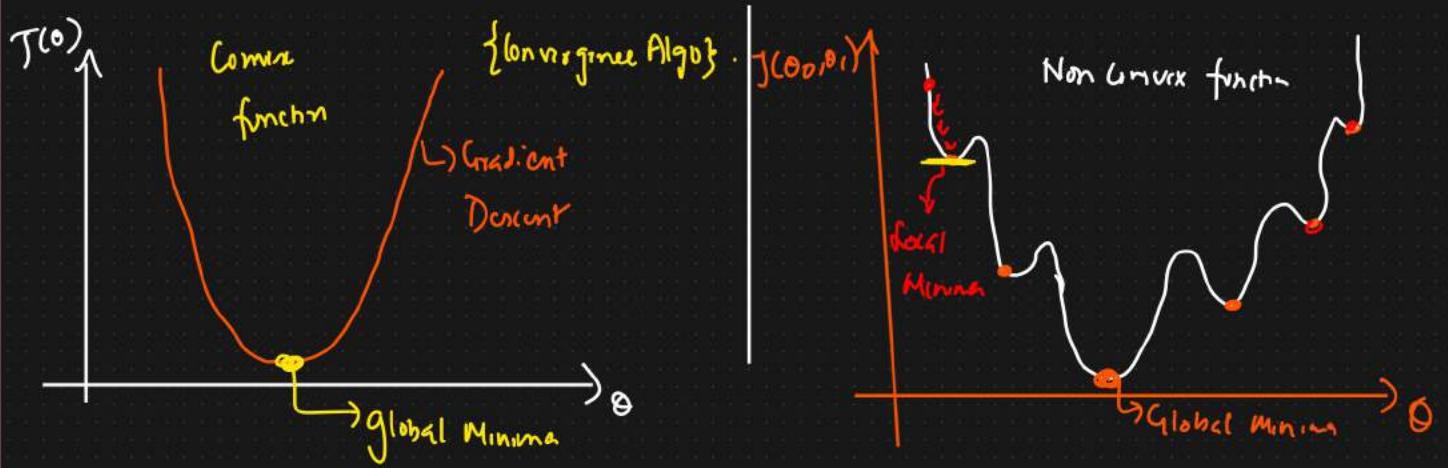
$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x,$$

## Logistic Regression Cost function

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$h_{\theta}(x) = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$$



### Log Loss

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

$h_\theta(x) = \frac{1}{1+e^{-(\theta_0+\theta_1 x)}}$

$\Downarrow$

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)). \quad \Rightarrow \text{Convex function}$$

if  $y=1$

$$= -\log(h_\theta(x)) \Rightarrow y=1$$

$$\text{if } y=0 \Rightarrow -\log(1-h_\theta(x))$$

Final Aim :

Minimize Cost function  $J(\theta_0, \theta_1)$  by changing  $\theta_0$  &  $\theta_1$

Convergence Algorithm

Repeat

{

$$\theta_j : \theta_j - \alpha \frac{\partial J(\theta, \theta_1)}{\partial \theta}$$

}

j=0,1

## Logistic Regression With Regularization Parameters

### Cost function

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)). & \text{if } y=0 \end{cases}$$

Reduce Overfitting  
↓

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) + \lambda_2 \text{Regularization}$$

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) + \lambda_1 \text{Regularization}$$

$$\hat{J}(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) + \lambda_2 \text{Reg} + \lambda_1 \text{Reg.}$$

- $\lambda_2$  Regularization  $\Rightarrow$  Reduce Overfitting

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) + \lambda \sum_{i=1}^n (\text{Slope})^2$$

- $\lambda_1$  Regularization  $\Rightarrow$  Feature Selection

$\lambda \Rightarrow$  Hyperparameter

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) + \lambda \sum_{i=1}^n |slope|$$

## ElasticNet

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) + \lambda_1 \sum_{i=1}^n (slope)^2 + \lambda_2 \sum_{i=1}^n |slope|$$

C & λ Relationship

$$\boxed{C = \frac{1}{\lambda}} \quad \Rightarrow \quad \boxed{\lambda = \frac{1}{C}}$$

# Performance Metrics, Accuracy, Precision, Recall And F-Beta

## Topics to be covered

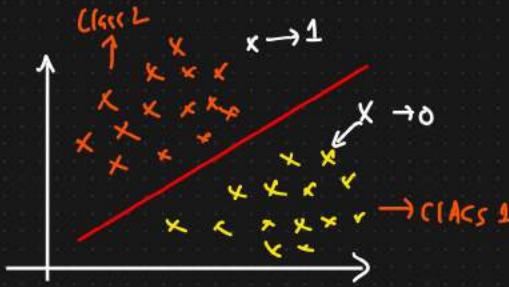
① Confusion Matrix ✓

② Accuracy ✓

③ Precision ✓

④ Recall ✓

⑤ F-Beta Score



Dataset	Actual		Predicted	
	$x_1$	$x_2$	$y$	$\hat{y}$ ← Model Prediction
	—	—	0	1   → Wrong Prediction

→ — — 1 1 → Correct Prediction

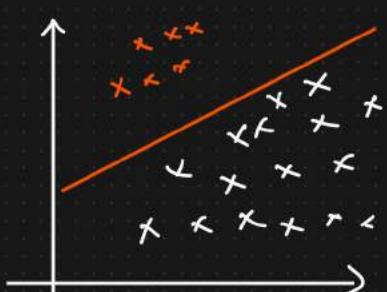
① Confusion Matrix

		Actual Values			
		1	0	—	— 0 → " "
1	3	2	—	— 1	1 → " "
0	1	1	—	— 0	1 → Wrong Prediction
			—	— 1	0

		Actual	
		1	0
Predicted	1	TP	FP
	0	FN	TN

$$\text{Model Acc.} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$= \frac{3 + 1}{3 + 2 + 1 + 1} = \frac{4}{7} = 0.571 \approx 57.1\%$$



④ DATASET → Imbalanced dataset

↓  
1000 datapoints  $\begin{cases} 900 \rightarrow 1 \\ 100 \rightarrow 0 \end{cases}$  → Imbalanced dataset

Dumb Model → 0/p → 1 ⇒

$$\boxed{\text{Accuracy} = 90\%}$$

Imbalanced dataset

In this scenario we cannot use Accuracy performance

② Precision =  $\frac{TP}{TP+FP}$  } Out of all the actual values how many are correctly predicted

	1	0
1	TP	FP
0	FN	TN

FP → Important  
FP ↓

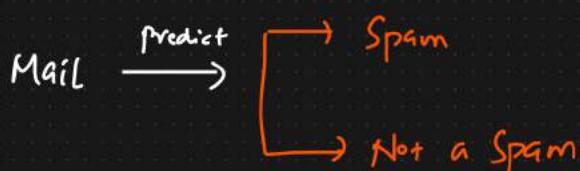
③ Recall =  $\frac{TP}{TP+FN}$  ⇒ FN ↓

⇒ Out of all the predicted values how many are correctly predicted

	1	0
1	TP	FP
0	FN	TN

With actual values

Use Case 1 ⇒ Spam classification



	1	0
1	TP	FP
0	FN	TN

Mail → Spam } Good Scenario  
Model → Spam }

FP is Important

FP ↓↓↓  
0 ∈ Mail → Not a Spam } Blunder  
1 ∈ Model → Spam }

1 ∈ Mail → Spam }  
0 ∈ Model → Not a Spam } ⇒ FN

## PRECISION PERFORMANCE METRICS.

Use Case 2  $\Rightarrow$  FN is Important

To predict whether a person has diabetes or not

$\downarrow$

① Actual  $\rightarrow$  Diabetes } Good  
Model  $\rightarrow$  Diabetes }  $\rightarrow$

Diabetic	No Diabetic	Actual
TP	FP	Diabetes
FN $\downarrow\downarrow$	TN	
No Diabetes		No Diabetic

$\downarrow$

② Actual  $\rightarrow$  Diabetes }  $\Rightarrow$  FN  $\downarrow\downarrow \Rightarrow$  Important  
Model  $\rightarrow$  No. Diabetes } Blunder

③ Actual  $\rightarrow$  No Diabetes } FP  $\Rightarrow$  Wrong prediction  
Model  $\rightarrow$  Diabetes }

④ Actual  $\rightarrow$  No Diabetes }  
Model  $\rightarrow$  No Diabetes } Correct

$\downarrow$   
RECALL

Assignment : ① Tomorrow the stock will crash or not

Reduce FP  $\downarrow\downarrow$  or FN  $\downarrow\downarrow$

$$\textcircled{4} \quad \underline{F\text{-Beta Score}} = (1 + \beta^2) \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

① If  $FP$  &  $FN$  are both important

$$\beta = 1$$

$\overbrace{\text{Mean}}$   
 $\overbrace{\text{Harmonic}}$

$$F1 \text{ Score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

② If  $FP$  is more important than  $FN$

$$\beta = 0.5$$

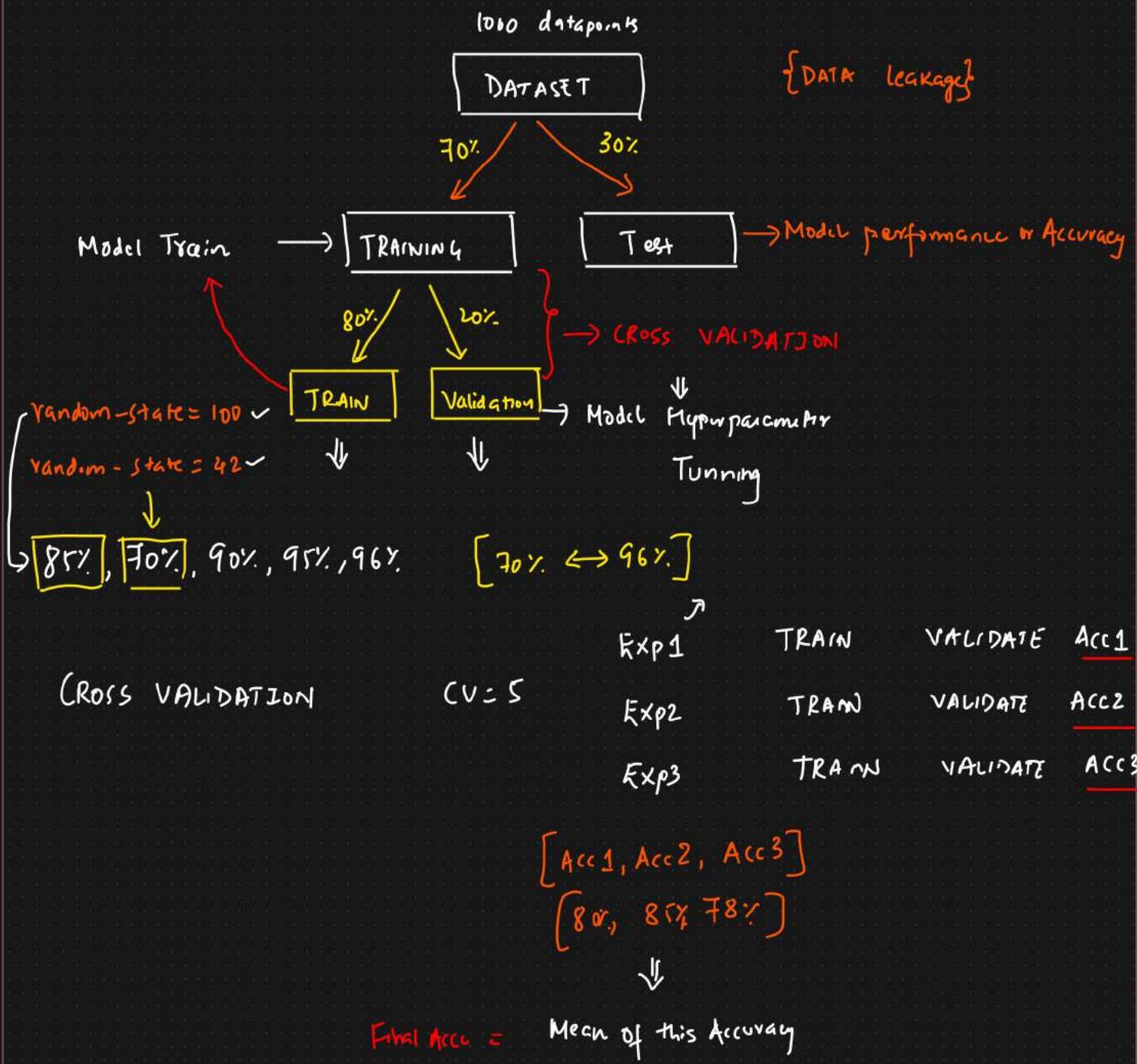
$$F0.5 \text{ Score} = (1 + 0.25) \frac{P * R}{P + R}$$

③ If  $FN \gg FP$

$$\beta = 2$$

$$F2 \text{ Score} = (1 + 4) \frac{P * R}{P + R}$$

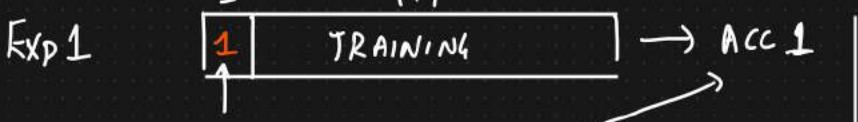
# CROSS VALIDATION AND ITS TYPES



## Types of Cross Validation

- ① leave One Out Cross Validation (100 cv)

TRAINING DATA → 500 Records



TRAINING → Model TRAIN

VALIDATION → Model Predict



⇒ Mean of All these Accuracies



### Disadvantage

- ① Time Comsuming is huge for training Big datasets
- ② Model Overfit → TRAINING Acc  $\uparrow\uparrow$   
New data → Validation Acc  $\downarrow\downarrow$

### Leave p out CV

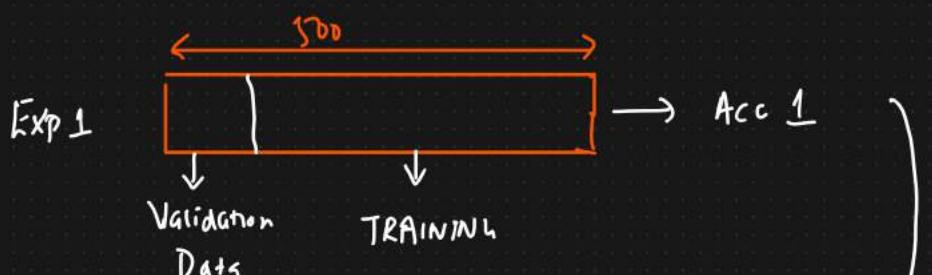
$p = 10, 20, 30, \dots, n_0$  → Hyperparameter

### K Fold Cross Validation

$K = 5$

$n = 500$

$$\text{Validation size} \uparrow \frac{500}{5} = 100$$



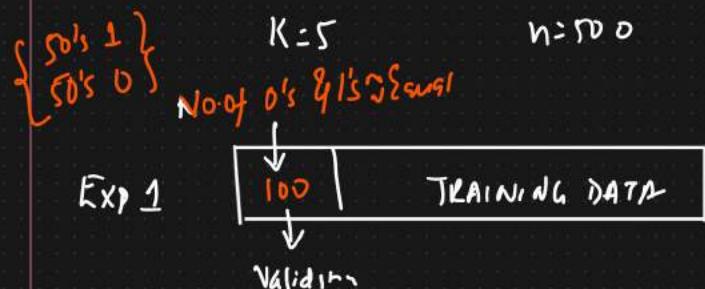
Average of all the Accuracy.



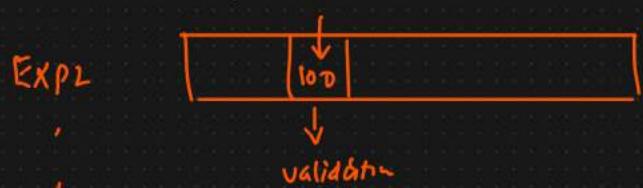
Exp 5  $\left[ \begin{array}{c|c} \text{TRAIN} & \text{Validation} \end{array} \right] \rightarrow \text{Acc} \int$

$f_1$	$f_2$	$y$
4		0 }
		0 }
		0 }
		1 }
		,
		0 }

#### ④ Stratified K Fold CV {Imbalanced DATASET}



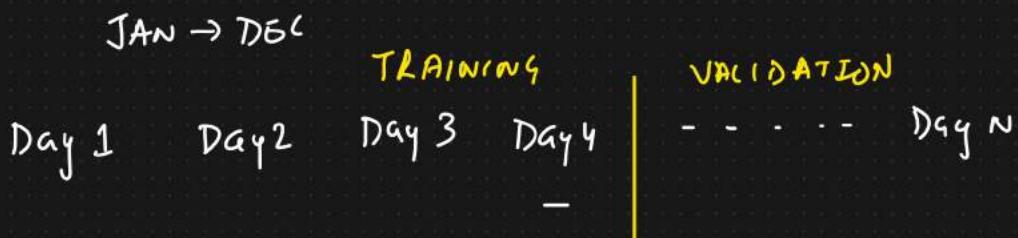
$n=500$   $\left\{ \begin{array}{l} 350 \rightarrow 1 \\ 150 \rightarrow 0 \end{array} \right\}$   $\left\{ \begin{array}{l} 50's \ 0 \\ 50's \ 1 \end{array} \right\}$  !



#### ⑤ Time Series CV

Amazon review Sentiment Analysis

Product A



Time Series Application

# Hypoparameter Tuning

↳ Finding the best parameters while training the model

① GridSearchCV

② RandomizedSearchCV

① GridSearchCV [ Grid Search + CROSS VALIDATION ] [ CV=5 ]

logistic Regression

penalty['l1', 'l2', 'elasticnet', None],

solver['lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky', 'sag', 'saga']

CV  $\Rightarrow$  K folds CV

k=5

TRAIN AND VALIDATION = 5

$\Downarrow$

Average Accuracy  $\Rightarrow$

model = LogisticRegression( elasticnet, newton-cg )

$\Downarrow$

model.predict(nclo-data)



Increase the Model Performance or Accuracy

Disadvantage

① Time Complexity increases for Training the Model

② RandomizedSearchCV  $\leftarrow$

$n\_iter = 10$  +  $CV = 5$

10 different combination +  $CV = 5$

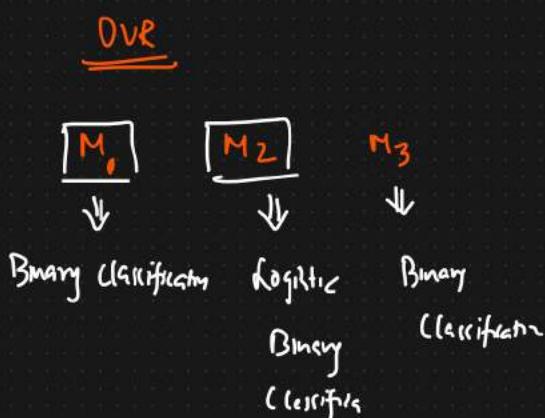
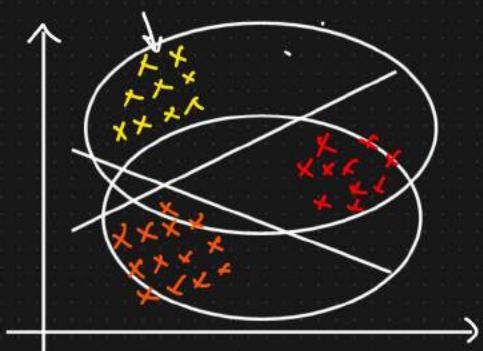
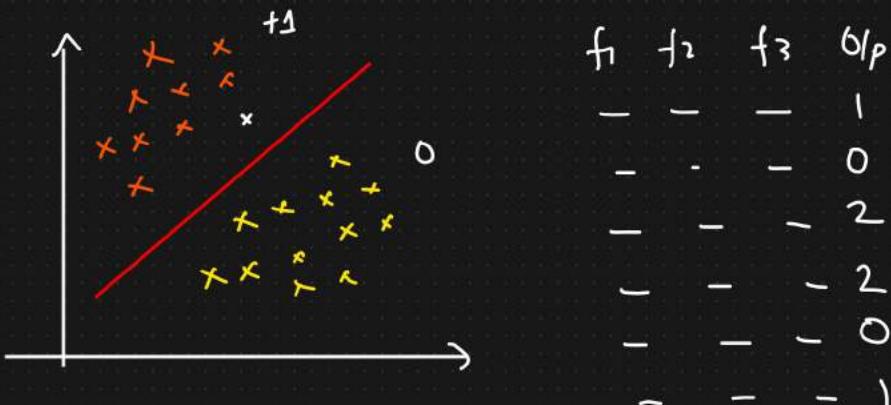
## Advantage

- ① Time Complexity Decrease

# Logistic Regression For Multiclass Classification

① One Versus Rest

② Multinomial



f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>	b/p	0 <sub>1</sub>	0 <sub>2</sub>	0 <sub>3</sub>
-	-	-	1	0	1	0
-	-	-	0	1	0	0
-	-	-	2	0	0	1
-	-	-	2	0	0	1
-	-	-	0	1	0	0
-	-	-	1	0	1	0

I/P feature

O/P

M<sub>1</sub>  $\Rightarrow$  f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>

0<sub>1</sub>  $\rightarrow$  Binary classification

M<sub>2</sub>  $\Rightarrow$  f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>

0<sub>2</sub>  $\rightarrow$  Binary classification

M<sub>3</sub>  $\Rightarrow$  f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>

0<sub>3</sub>  $\rightarrow$  " "

New Data point =  $\left[ \begin{array}{ccc} M_1 & M_2 & M_3 \\ 0.20, & 0.35, & 0.45 \end{array} \right] \Rightarrow 1$

03

↳ ap

Multinomial Probabilities.

# Decision Tree

① Decision Tree Classifier [classification]

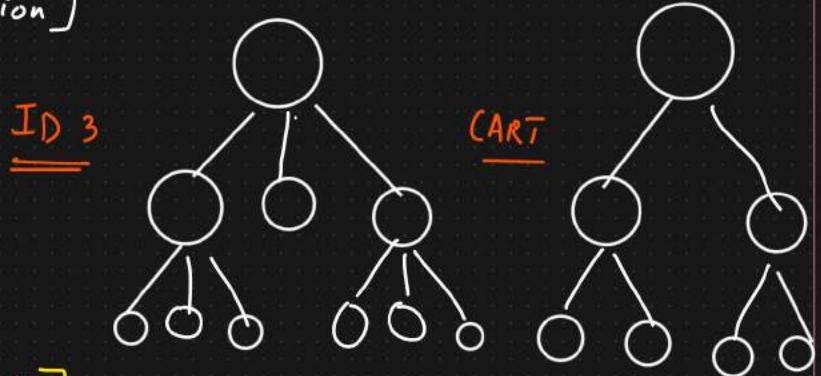
② Decision Tree Regressor [Regression]

Decision Tree Classifier =

Two types

① ID3 [Iterative Dichotomiser 3]

② CART ✓ [Classification And Regression Tree]



$age = 14$

if ( $age \leq 15$ ) :

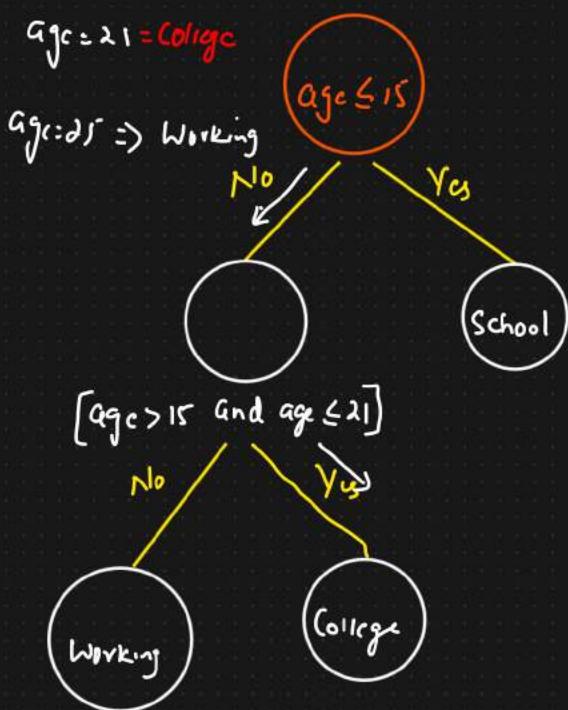
Print ("School").

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

Print ("College")

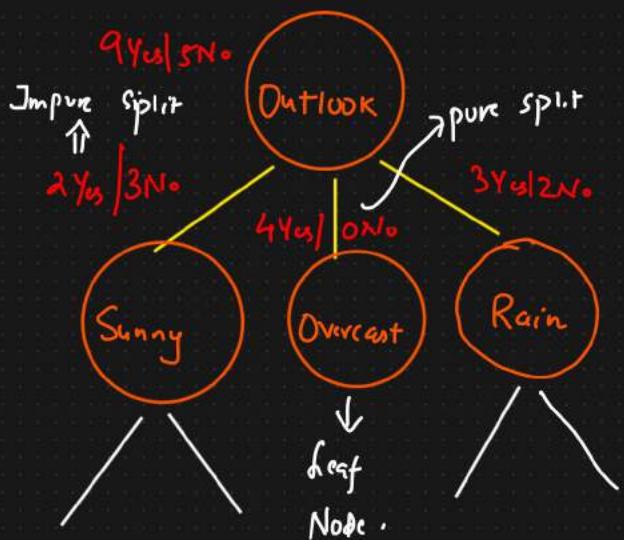
else:

Print ("WORKing")



Dataset → Predict Play Tennis OR Not

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 Purity check : Pure Split or Impure Split

Entropy      } Measure of ✓  
Gini Impurity      } purity.

② What feature you need to select to start the split? → Information Gain

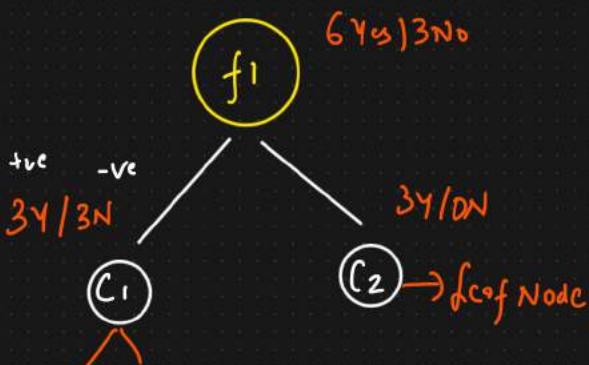
Binary classification

① Entropy

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

$P_+$  = probability of positive category

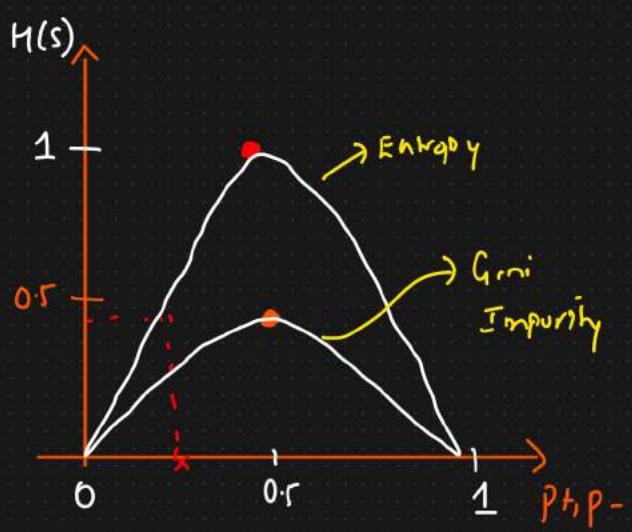
$P_-$  = probability of negative category



$$\begin{aligned} H(C_1) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\ &= -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right) \end{aligned}$$

② Gini Impurity

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



$$= -\frac{1}{2} \log_2(1/2) - (1/2) \log_2(1/2).$$

$\therefore 1 \Rightarrow$  Impure Split

$$H(c_2) = -\frac{3}{3} \log_2(3/3) - 0/3 \log_2(0/3)$$

$\therefore 0 \Rightarrow$  Pure Split

$c_1 \ c_2 \ c_3$

Yes | No | May Be

Multiclass

$$H(s) = -P_{c_1} \log_2 P_{c_1} - P_{c_2} \log_2 P_{c_2} - P_{c_3} \log_2 P_{c_3}$$

② Gini Impurity

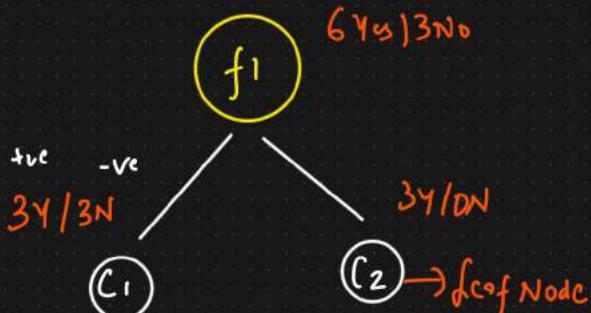
$$GI = 1 - \sum_{i=1}^n (P_i)^2$$

$$= 1 - \left[ (P_+)^2 + (P_-)^2 \right]$$

$$= 1 - \left[ \left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right]$$

$$= 1 - \left[ \frac{1}{4} + \frac{1}{4} \right]$$

$$= \frac{1}{2} = 0.5 \Rightarrow \text{Impure Split}$$



$$= 1 - \left[ \left(\frac{3}{3}\right)^2 + \left(0/3\right)^2 \right]$$

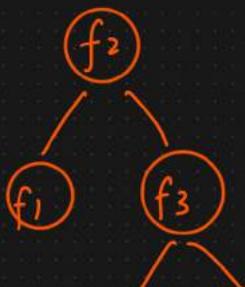
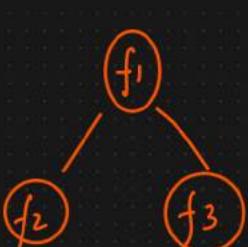
$$= 1 - 1$$

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

② What feature you need to select to

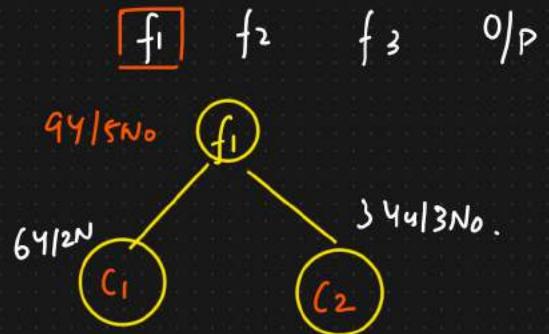
start the split?  $\rightarrow$  Information Gain

$f_1 \ f_2 \ f_3 \ \text{Op}$



Information Gain Entropy of the root node

$$\text{Gain}(S, f_1) = \boxed{H(S)} - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$



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

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

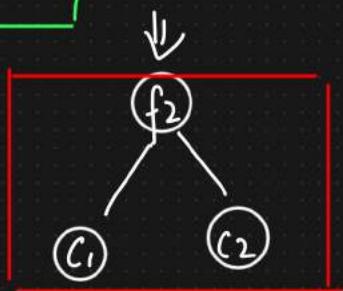
$$\approx 0.94 //$$

$$H(C_1) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \frac{2}{8} \approx \boxed{0.81}$$

$$H(C_2) = 1 //$$

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

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



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

Information gain is more when we split using  $f_2$ .

## Entropy Vs Gini Impurity

Whenver dataset is small  $\rightarrow$  Entropy }  
dataset is large  $\rightarrow$  GiniImpurity }.

# Decision Tree For Numerical Split

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

$f_2$        $f_1$       O/P

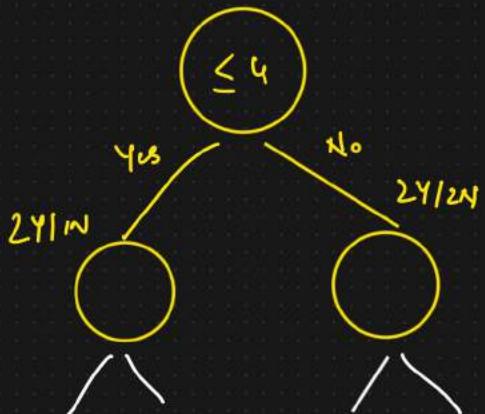
- {  $\rightarrow 2.3$  Yes  
   {  $\rightarrow 3.6$  Yes  
   {  $\rightarrow 4$  No  
   {  $\rightarrow 5.2$  No  
   {  $\rightarrow 6.7$  Yes  
   {  $\rightarrow 7.8$  No  
   {  $\rightarrow 9.0$  Yes

① Threshold = 2.3

② Threshold = 3.6

① Sorting the feature

③ Threshold = 4



Millions of Records

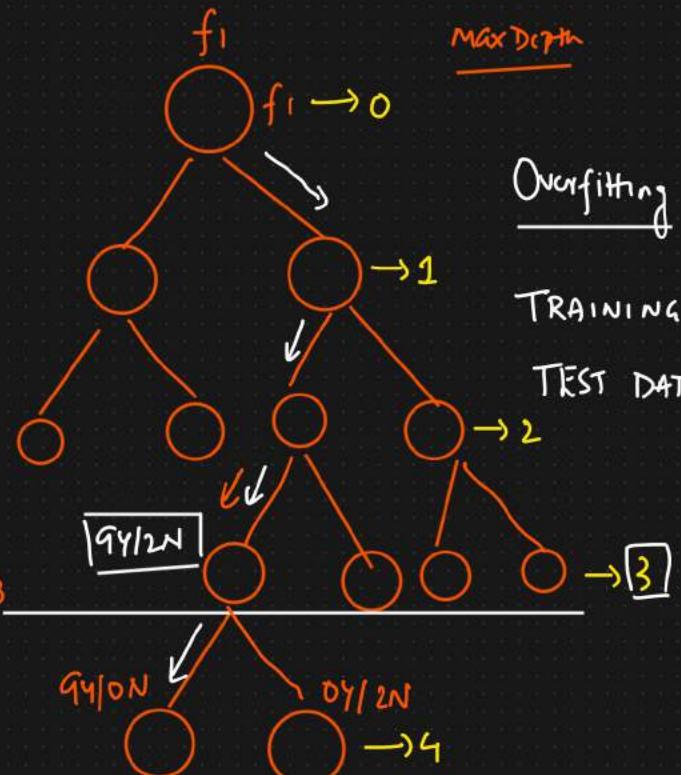
| Time Complexity ↑↑ |

# Decision Tree Post Pruning And Pre Pruning [Reduce Overfitting]

Training Data

- ① Post Pruning
- ② Pre Pruning

To Reduce Overfitting.



Max Depth

Overfitting

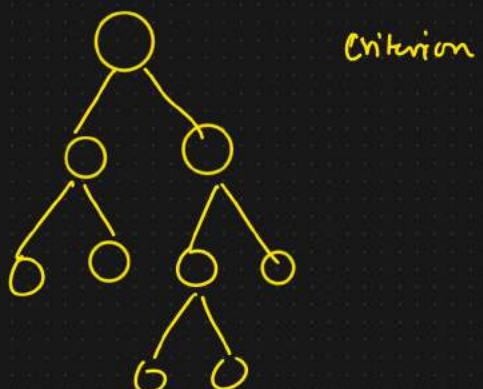
TRAINING ACC  $\uparrow\uparrow$

TEST DATA ACC  $\downarrow\downarrow$

- ① Post Pruning
- ② Decision Tree Construct
- ③ Prune the Tree

② Pre Pruning = Hypoparameter Tuning

$\Downarrow$   
Large Dataset



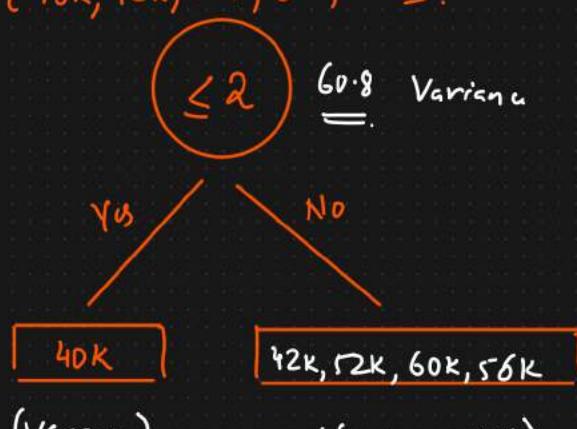
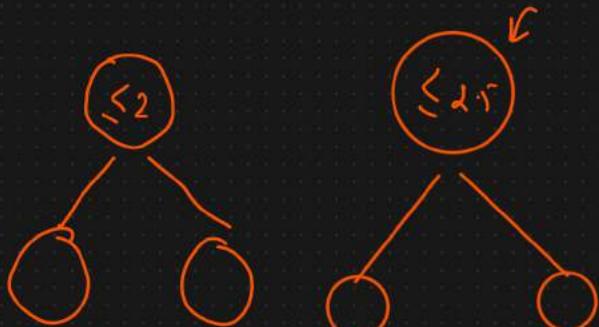
# Decision Tree Regressor

Dataset

Exp	Career Gap	Salary
→ 2	Yes	40K
→ 2.5	Yes	42K
3	No	52K
4	No	60K
4.5	Yes	56K

$$\frac{50K}{5} = \bar{y}$$

[40K, 42K, 52K, 60K, 56K].

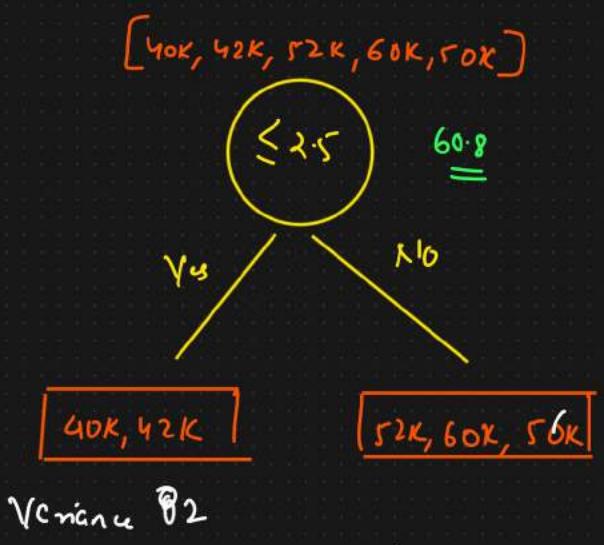


100 (Variance)  
Variance Reduction

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \begin{matrix} \text{Mean Squared} \\ \text{Error} \end{matrix}$$

$$\text{Variance (Root)} = \frac{1}{5} \left[ (40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right]$$

$$= \frac{1}{5} [100 + 64 + 4 + 100 + 36]$$



Variance Right

$$\begin{aligned} \text{Var(Lift)} &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \frac{1}{2} \left[ (40-50)^2 + (42-50)^2 \right] \\ &= \frac{1}{2} [100 + 64] \\ &= \frac{164}{2} = 82\% \end{aligned}$$

Var(Right)

$$= \frac{1}{3} [4 + 100 + 36]$$

$$= 60.8 \quad \underline{\underline{=}} \quad = \frac{140}{3} = 46.66$$

$$\text{Variance (Left)} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ = \frac{1}{4} [(40 - 50)^2]$$

$$\text{Variance Reduction} = 60.8 - \left[ \frac{2}{5} * 100 + \frac{3}{5} * 46.66 \right] \\ = 0.004 \quad \underline{\underline{=}}$$

$$\text{Variance (Left)} = 100$$

$$\text{Variance (Right)} = \frac{1}{4} \left[ (42 - 50)^2 + (52 - 50)^2 + (60 - 50)^2 + (56 - 50)^2 \right] \\ = \frac{1}{4} [64 + 4 + 100 + 36]$$

$$= 51 \quad \underline{\underline{=}} \quad w_i \in \mathcal{L} = \frac{1}{5} \quad w_i(r) = 4/5$$

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

$$= 60.8 - \left[ \frac{1}{5} * 100 + \frac{4}{5} * 51 \right]$$

$$= 60.8 - 20 - 40.8$$

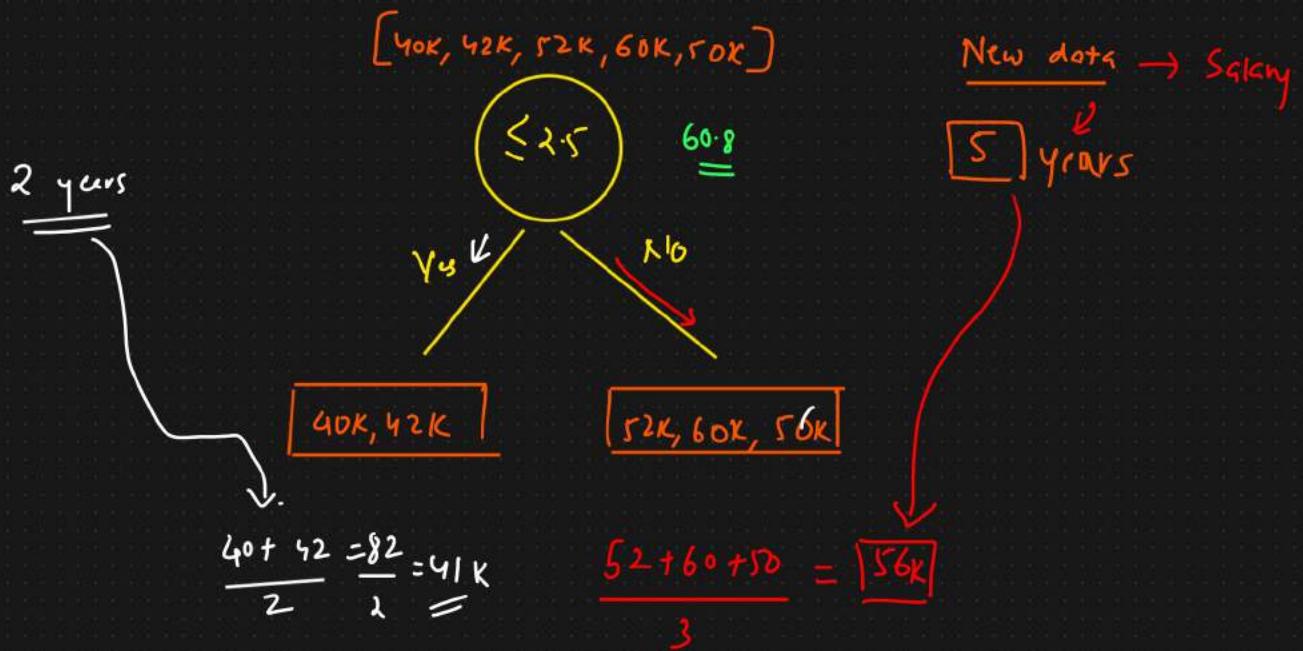
$$\boxed{\begin{array}{l} \text{Variance} = 0 \\ \text{Reduction} \end{array}}$$

○

0.004

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

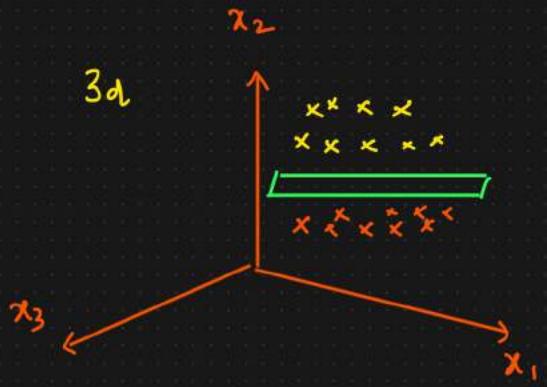
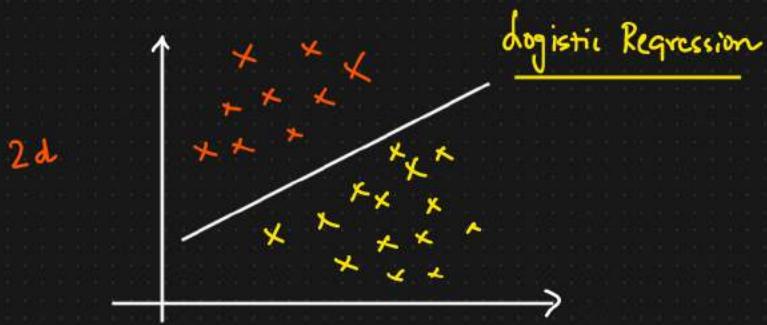




# Support Vector Machines ML Algorithm

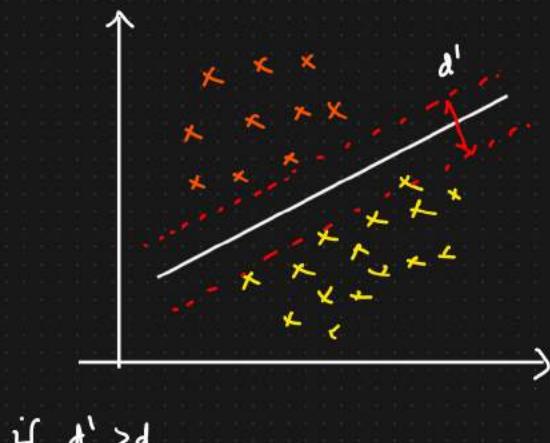
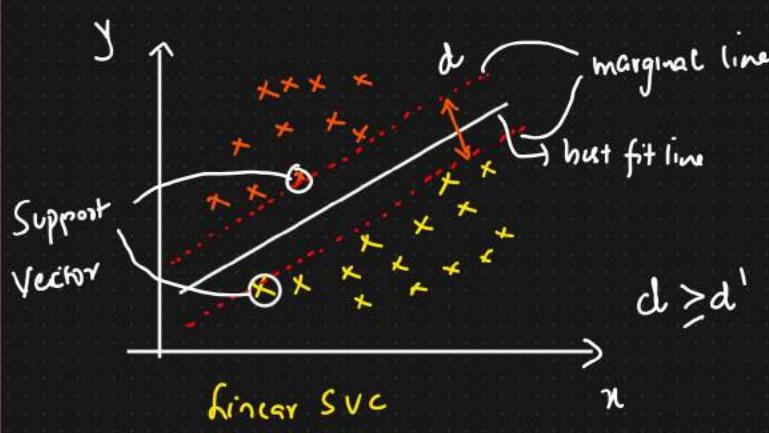
① SVC (Support Vector Classifier) → classification

② SVR (Support Vector Regressor) → Regression

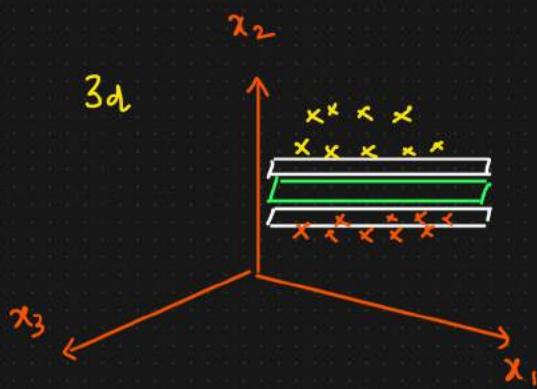


① Support Vector Classifier (SVC)

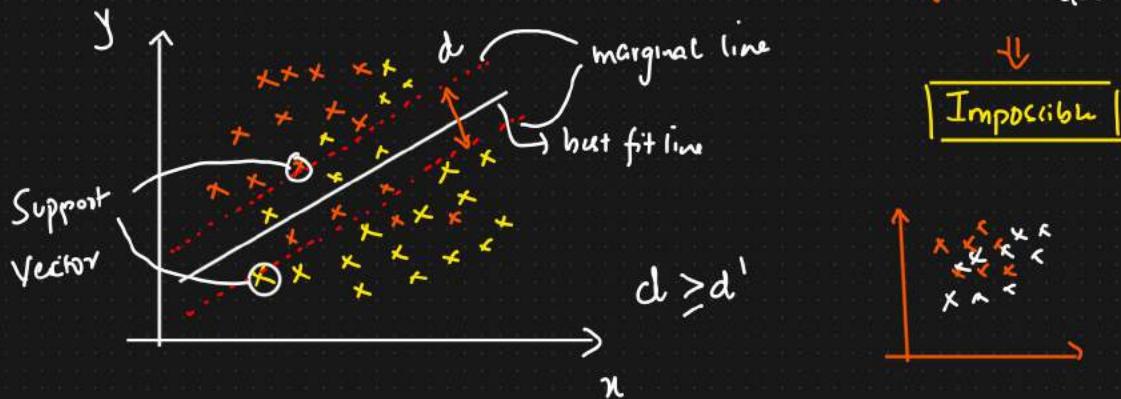
$d$  = marginal plane distance



distance is maximum

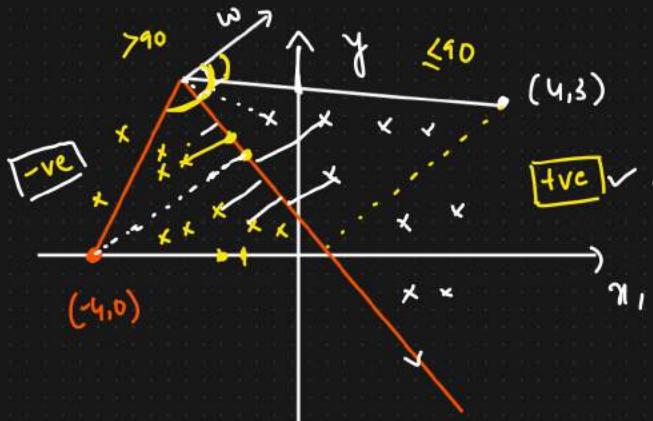


# Soft Margin And Hard Margin In SVC



② Soft Margin : Some data point are misclassified [Error]

③ Support Vector Machines (SVC) Maths Intuition



Equation of a straight line

$$\hat{y} = mx + c \Leftrightarrow ax + by + c = 0$$

$$h_0(x) = \theta_0 + \theta_1 x_1, \quad by = -ax - c$$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$y = \frac{-a}{b} x - \frac{c}{b}$$

$$w^T x + b = 0$$

$$y = b + [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$w^T = [w_1 \ w_2 \ w_3] \cdot x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

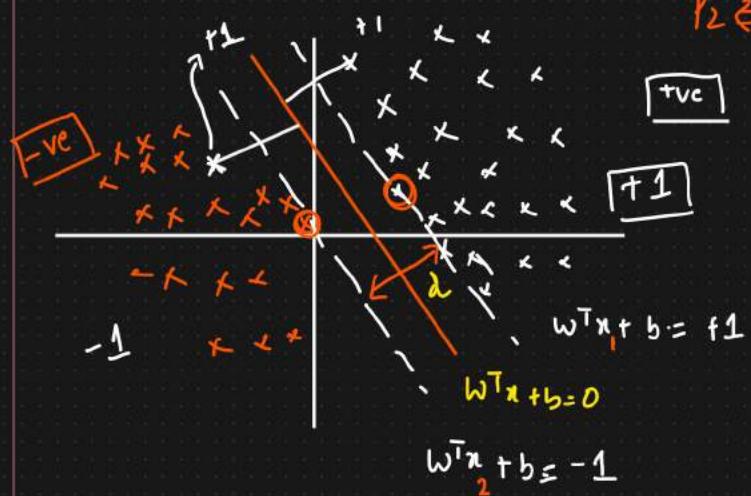
$$w^T x = [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$y = w^T x + b \Rightarrow y = mx + c$$

$$ax + by + c = 0$$

$$\boxed{w^T x + b = 0}$$

Marginal plane & SVC



$$P_1 \in x_1 \Rightarrow (x_1, x_2)$$

$$P_2 \in x_2 \Rightarrow (x_1, x_2)$$

$$w^T x_1 + b = +1$$

$$w^T x_2 + b = -1$$

$$\frac{(-)}{(-)} \quad \frac{(+)}{(+)} \quad \frac{(+)}{(-)}$$

$$\vec{w} \leftarrow \frac{w^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|} \Rightarrow \begin{array}{l} \text{distance} \\ \text{between} \\ \text{Marginal} \\ \text{plane.} \end{array}$$

Maximize

Maximize  
 $w, b$

$$\frac{2}{\|w\|} \Rightarrow \text{Distance between Marginal plane}$$

Constraint such that

$$y_i \begin{cases} +1 & \text{if } w^T x + b \geq 1 \\ -1 & \text{if } w^T x + b \leq -1 \end{cases}$$

$\Downarrow$

For all correctly classified data points

$$\boxed{y_i * w^T x + b \geq 1}$$

## Modified Cost function of SVC

$$\underset{w,b}{\text{Maximize}} \quad \frac{2}{\|w\|} \Rightarrow \underset{w,b}{\text{Minimize}} \quad \frac{\|w\|}{2}$$

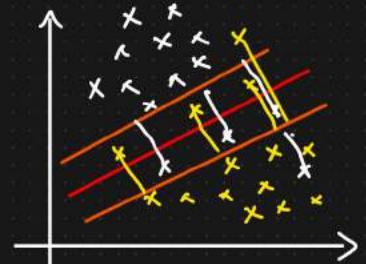
Constraint such that

$$y_i \begin{cases} +1 & \text{if } w^T x + b \geq 1 \\ -1 & \text{if } w^T x + b \leq -1 \end{cases}$$

## Cost function of Soft Margin SVC

$$\text{Cost fn} = \underset{w,b}{\text{Min}} \quad \frac{\|w\|}{2} + \underbrace{\left[ C_i \sum_{i=1}^n \xi_i \right]}_{\text{hyperparameter}} \Rightarrow \text{hinge loss}$$

$$C_i = 5$$



Summation of the distance  
of incorrect data points

{How many points to the marginal  
we can consider plane}

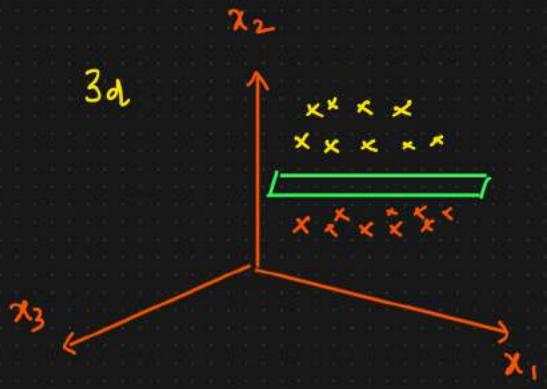
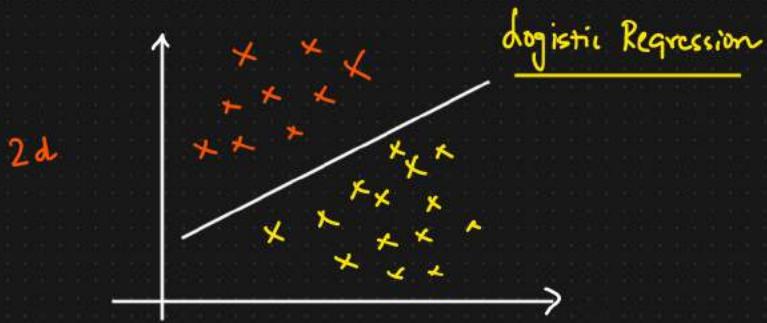
for misclassification}

$$\boxed{C = \frac{1}{\lambda}}$$

# Support Vector Machines ML Algorithm

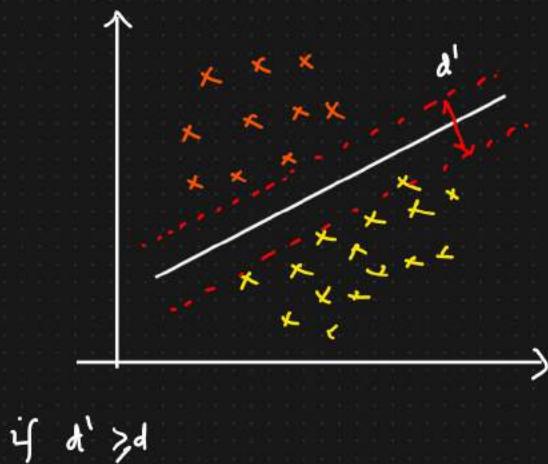
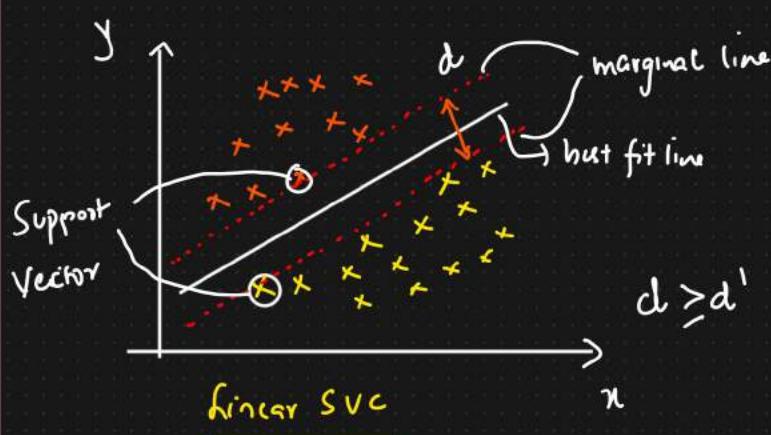
① SVC (Support Vector Classifier) → classification

② SVR (Support Vector Regressor) → Regression

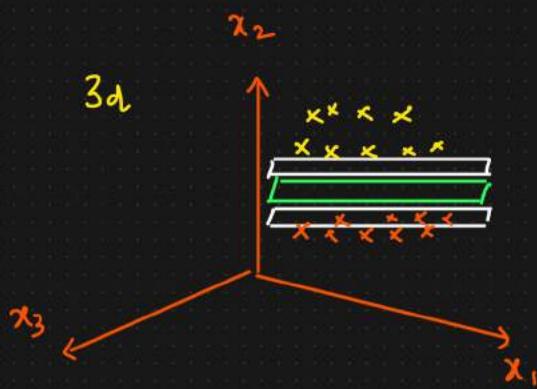


① Support Vector Classifier (SVC)

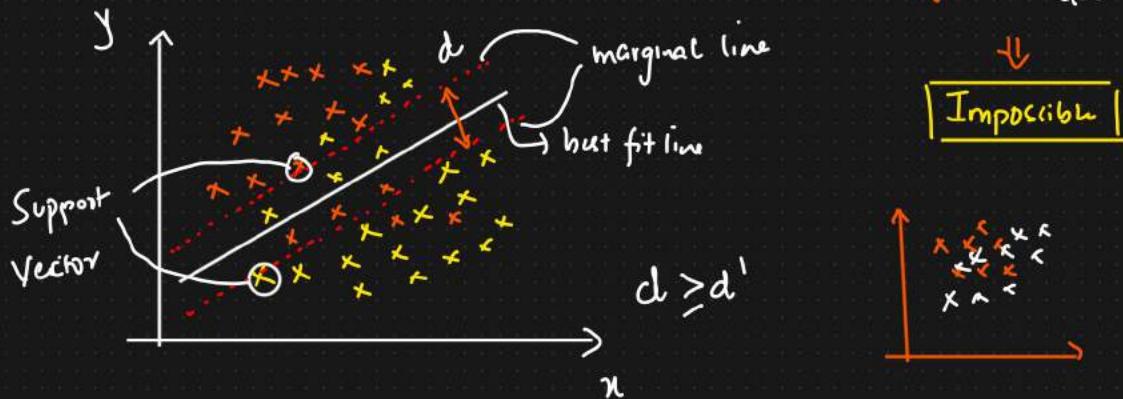
$d$  = marginal plane distance



distance is maximum

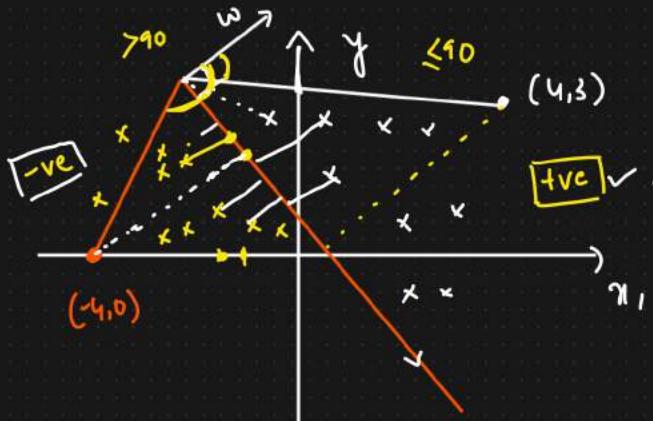


# Soft Margin And Hard Margin In SVC



② Soft Margin : Some data point are misclassified [Error]

③ Support Vector Machines (SVC) Maths Intuition



Equation of a straight line

$$\hat{y} = mx + c \Leftrightarrow ax + by + c = 0$$

$$h_0(x) = \theta_0 + \theta_1 x_1, \quad by = -ax - c$$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$y = \frac{-a}{b} x - \frac{c}{b}$$

$$w^T x + b = 0$$

$$y = b + [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$w^T = [w_1 \ w_2 \ w_3] \cdot x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

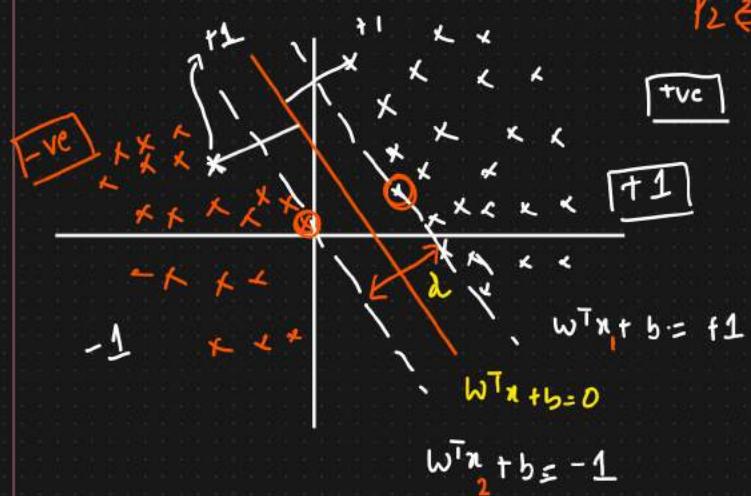
$$w^T x = [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$y = w^T x + b \Rightarrow y = mx + c$$

$$ax + by + c = 0$$

$$\boxed{w^T x + b = 0}$$

Marginal plane & SVC



Cost function

Maximize  
 $w, b$

$$\frac{2}{\|w\|} \Rightarrow \text{Distance between Marginal plane}$$

Constraint such that

$$y_i \begin{cases} +1 & \text{if } w^T x + b \geq 1 \\ -1 & \text{if } w^T x + b \leq -1 \end{cases}$$

$\Downarrow$

For all correctly classified data points

$$y_i * w^T x + b \geq 1$$

## Modified Cost function of SVC

$$\underset{w,b}{\text{Maximize}} \quad \frac{2}{\|w\|} \Rightarrow$$

$$\underset{w,b}{\text{Minimize}} \quad \frac{\|w\|}{2}$$

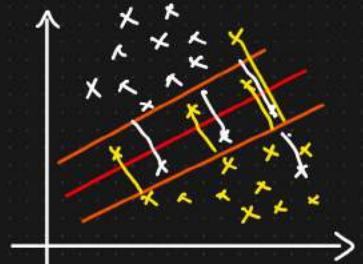
Constraint such that

$$y_i \begin{cases} +1 & \text{if } w^T x + b \geq 1 \\ -1 & \text{if } w^T x + b \leq -1 \end{cases}$$

## Cost function of Soft Margin SVC

$$\text{Cost fn} = \underset{w,b}{\text{Min}} \quad \frac{\|w\|}{2} + \left[ C_i \sum_{i=1}^n \xi_i \right] \Rightarrow \text{Hinge Losses}$$

$$C_i = 5$$

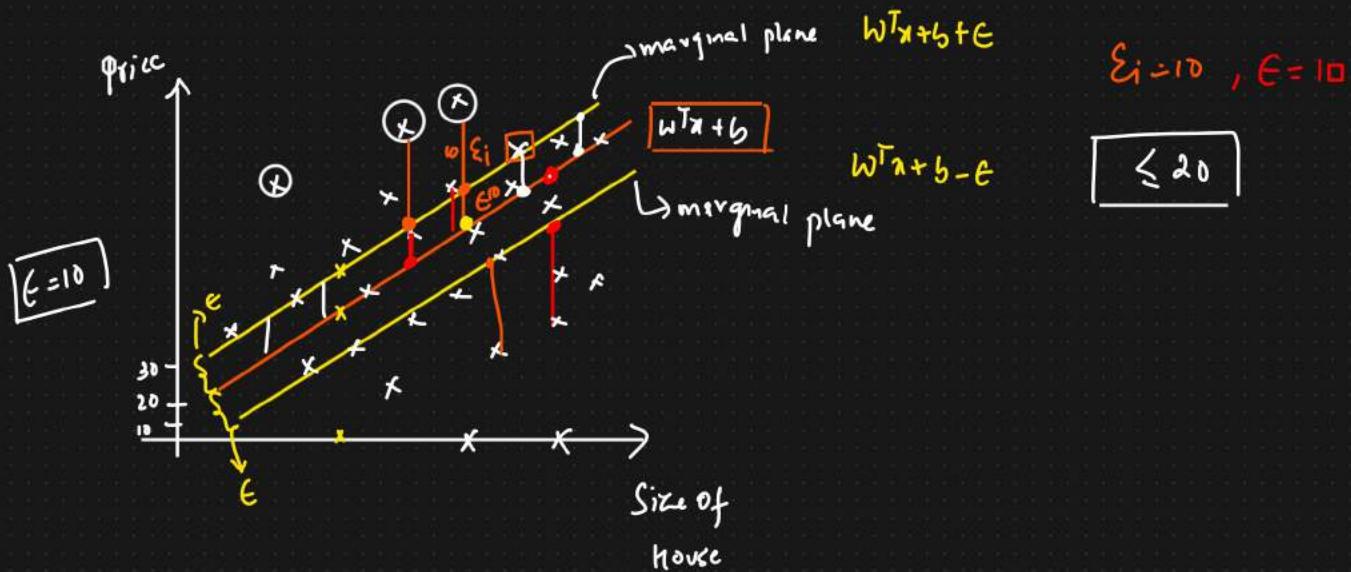


hyperparameter  
 $C$  → How many points to the marginal plane for misclassification

$$\boxed{C = \frac{1}{\lambda}}$$

# Support Vector Regressor (SVR)

$\epsilon$  = Marginal Error



## Cost fn

$$\underset{w, b}{\text{Min}} \quad \frac{\|w\|}{2} + \left[ C \sum_{i=1}^n \xi_i \right] \Rightarrow \text{Hinge loss}$$

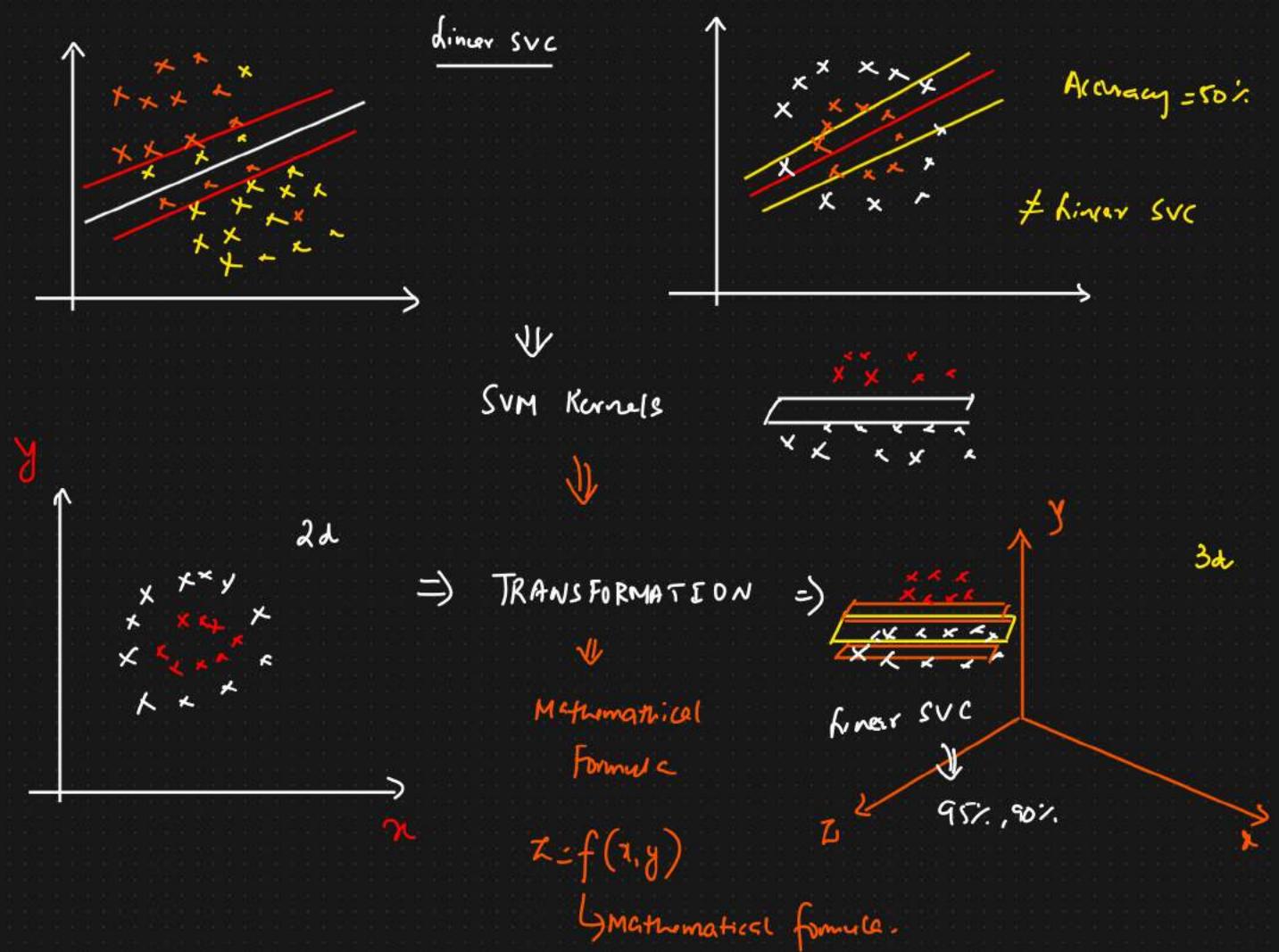
## Constraint

$$|y_i - w^T x_i| \leq \epsilon + \xi_i$$

$\epsilon$  Marginal Error

$\xi_i$  = Error above the Margin

## SVM Kernels



## Data set

X	Y	$Z = x^2$
2	Yes	4
3	No	9
4	Yes	16
-	Yes	-
-	-	-

$\Rightarrow$  1 dimension

$\Rightarrow$

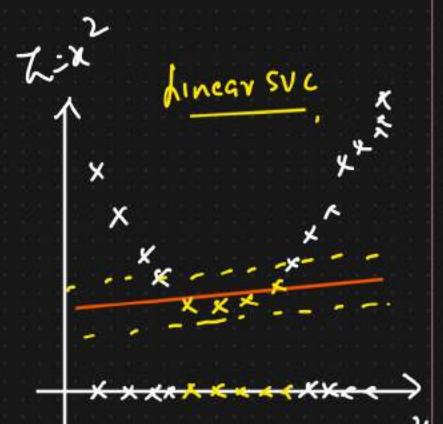
$x \rightarrow x^2$

[1d - 2d]

SVM Kernel

↓  
Data Transformation

$$Z = x^2$$



- ① Polynomial Kernel
- ② RBF Kernel
- ③ Sigmoid Kernel

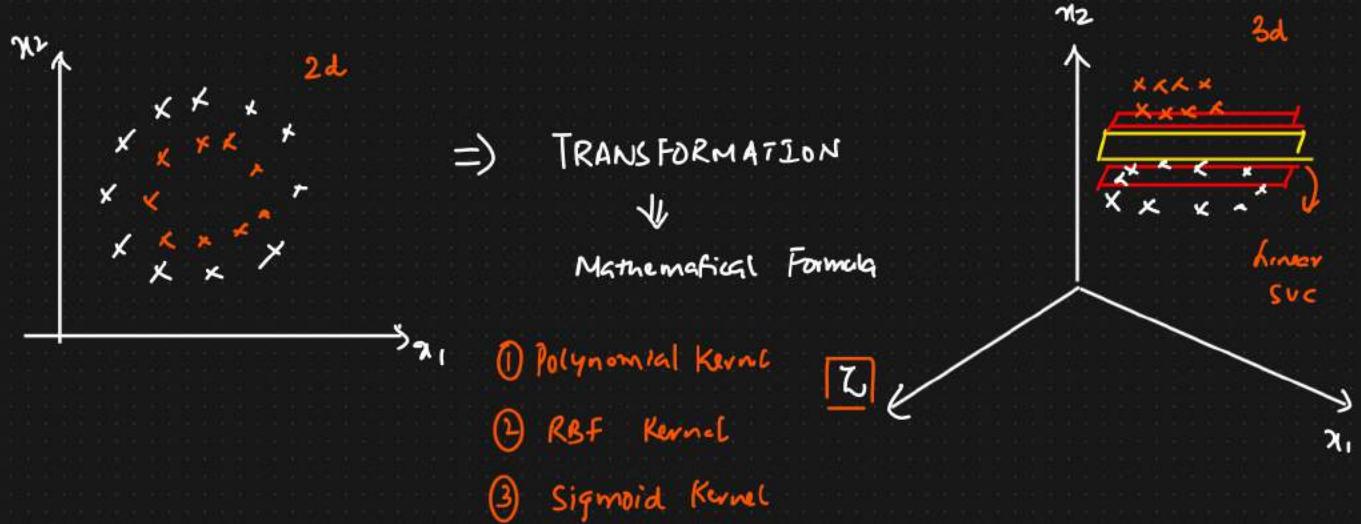
$\Rightarrow$  Transformation  $\Rightarrow$  Mathematical formula

# SVM Kernels

① Polynomial Kernel

② RBF Kernel

③ Sigmoid Kernel



① Polynomial Kernel

$$x = [x_1, x_2]$$

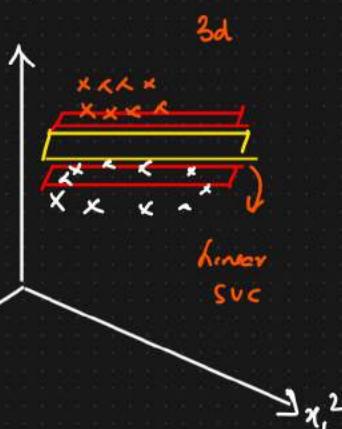
$$y = [x_1, x_2]$$

$$f(x, y) = \left( \underline{x^T y + c} \right)^d \quad |c=1|$$

$$\boxed{2d} \quad \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} x_1, x_2 \end{bmatrix}$$

$$x_1 \quad x_2 \quad y = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix} \quad \begin{matrix} \downarrow \\ \text{Transformation} \end{matrix} \quad x_1, x_2$$

$$\begin{matrix} \boxed{3d} \\ x_1^2 \quad x_1 x_2 \quad x_2^2 \end{matrix} \quad y \quad \Rightarrow \text{Polynomial Kernel}$$



② RBF Kernel

↓

Formula

③ Sigmoid Kernel

↓

Formula

## Naive Bayes Algorithm (Classification)

① Probability [Independent And Dependent Events]

② Baye's Theorem

③ Naive Baye's Math Intuition.

① Probability

Independent Events

Rolling a Dice  $\{1, 2, 3, 4, 5, 6\}$

$$Pr(1) = \frac{1}{6} \quad Pr(2) = \frac{1}{6} \quad Pr(3) = \frac{1}{6}$$

$= \quad =$

Dependent Events

① What is the probability of first removing  
a orange marble and then a yellow marble?

0	0	0
0	0	

①  $\rightarrow P(O) = \frac{3}{5} \rightarrow 1^{\text{st}} \text{ Event}$

0	0
0	0

②  $\rightarrow P(Y) = \frac{2}{4} \rightarrow 2^{\text{nd}} \text{ Event}$

$P(Y|O) = \frac{2}{4}$   $\Rightarrow$  Conditional Probability

$$Pr(O \text{ and } Y) = P(O) * \underline{P(Y|O)} \Rightarrow \text{Conditional Probability}$$

$$= \frac{3}{5} * \frac{2}{4} = \boxed{\frac{3}{10}}$$

$$\Pr(A \text{ and } B) = \Pr(A) * \Pr(B/A)$$

## Bayes Theorem

$$\Pr(A \text{ and } B) = \Pr(B \text{ and } A)$$

$$\Pr(A) * \Pr(B/A) = \Pr(B) * \Pr(A/B)$$

$$\Pr(A/B) = \frac{\Pr(A) * \Pr(B/A)}{\Pr(B)}$$

Bayes Theorem

$\Pr(A|B)$  = Probability of Event A given B has occurred

$\Pr(A)$  = Probability of Event A

$\Pr(B)$  = Probability of Event B

$\Pr(B/A)$  = Probability of Event B given A has occurred.

### DATASET

$x_1$	$x_2$	$x_3$	$\downarrow \text{Predict}$
-	-	-	Yes
-	-	-	No
-	-	-	Yes
-	-	-	No
-	-	-	Yes

$$\Pr(Y/(x_1, x_2, x_3)) = \frac{\Pr(Y) * \Pr(x_1, x_2, x_3/Y)}{\Pr(x_1, x_2, x_3)}$$

$$\Pr(A/B) = \frac{\Pr(A) * \Pr(B/A)}{\Pr(B)}$$

$$\Pr(Y/(x_1, x_2, x_3)) = \frac{\Pr(y) * \Pr(x_1, x_2, x_3/y)}{\Pr(x_1, x_2, x_3)}$$

$$= \frac{\Pr(y) * \Pr(x_1/y) * \Pr(x_2/y) * \Pr(x_3/y)}{\Pr(x_1) * \Pr(x_2) * \Pr(x_3)}$$

<u>DATASET</u>			$\downarrow$ Pred.
$x_1$	$x_2$	$x_3$	0/P?
-	-	-	Yes
-	-	-	No
-	-	-	Yes
-	-	-	No
-	-	-	Yes

$$\Pr(Y=Yes/(x_1, x_2, x_3)) = \frac{\Pr(Yes) * \Pr(x_1/Yes) * \Pr(x_2/Yes) * \Pr(x_3/Yes)}{\Pr(x_1) * \Pr(x_2) * \Pr(x_3)} = 0.60$$

Remove       ~~$\Pr(x_1) * \Pr(x_2) * \Pr(x_3)$~~       1

$$\Pr(Y=No/(x_1, x_2, x_3)) = \frac{\Pr(No) * \Pr(x_1/No) * \Pr(x_2/No) * \Pr(x_3/No)}{\Pr(x_1) * \Pr(x_2) * \Pr(x_3)} = 0.40$$

~~$\Pr(x_1) * \Pr(x_2) * \Pr(x_3)$~~

Let's Solve This Problem

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
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	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Outlook

Sunny

Yes  
2  
No  
3

$P(E/\text{Yes})$   
 $P(E/\text{No})$

$2/9$   
 $3/5$

Overcast

Yes  
4  
No  
0

$4/9$   
 $0/5$

Rain

Yes  
3  
No  
2

$3/9$   
 $4/5$

Temperature

$$(\text{Sunny}, \text{Hot}) = 0/9$$

Yes/No

Yes  
No  
 $P(E/\text{Yes})$   
 $P(E/\text{No})$

Hot  
2  
2  
 $2/9$   
 $4/5$   
Yes  
9  
 $P(Y) = 9/14$

Mild  
4  
2  
 $4/9$   
 $2/5$   
No  
5  
 $P(N) = 5/14$

Cool  
3  
1  
 $3/9$   
 $1/5$

$$Pr(Y_{\text{Yes}} | (\text{Sunny}, \text{Hot})) = \frac{Pr(Y_{\text{Yes}}) * Pr(\text{Sunny}/Y_{\text{Yes}}) + Pr(\text{Hot}/Y_{\text{Yes}})}{Pr(\text{Sunny}) + Pr(\text{Hot})}$$

$$= 9/14 * 2/9 + 2/5$$

$$= \frac{2}{14} = \boxed{0.031}$$

$$Pr(N_{\text{Yes}} | (\text{Sunny}, \text{Hot})) = \frac{Pr(N_{\text{Yes}}) * Pr(\text{Sunny}/N_{\text{Yes}}) * Pr(\text{Hot}/N_{\text{Yes}})}{Pr(\text{Sunny}) + Pr(\text{Hot})}$$

Constant

$$= 5/14 * 3/5 * 2/5$$

$$= \frac{0.085}{0.031 + 0.085}$$

Finally

$$\Pr(\text{Yes} | (\text{Sunny}, \text{hot})) = \frac{0.031}{0.031 + 0.085} = 0.27 = 27\%$$

$$\Pr(\text{No} | (\text{Sunny}, \text{hot})) = \frac{0.085}{0.031 + 0.085} = 0.73 = 73\%$$

Now DATA  $\left[ \begin{matrix} \text{Sunny, Hot} \end{matrix} \right] \Rightarrow \boxed{73\%} \Rightarrow \text{No} \Rightarrow \circ$

$27\% \Rightarrow \text{Yes}$



$\boxed{\text{Person will Not play}}$

## Variants of Naive Bayes

① Bernoulli Naive Bayes

② Multinomial Naive Bayes

③ Gaussian Naive Bayes

### ① Bernoulli Naive Bayes

Whenever your features are following a Bernoulli Distribution, then we use Bernoulli Naive Bayes

Dataset:

$f_1$	$f_2$	$f_3$	O/P
Yes	Pass	Male	Yes
No	Fail	Female	No
Yes	Pass	Male	Yes
No	Pass	Female	No.
Yes	Pass	Female	No.

Bernoulli  $\rightarrow$  0, 1

$$P(\text{Success}) = 1 = P$$

$$P(\text{Fail}) = 0 = 1 - P$$

### ② Multinomial Naive Bayes $\Rightarrow$ I/P = Text

Dataset : Sentiment Analysis

$\pi_{IP}$	<u>O/P</u>
Review	Message
The product is really good	Positive
The product is bad	Negative

↓

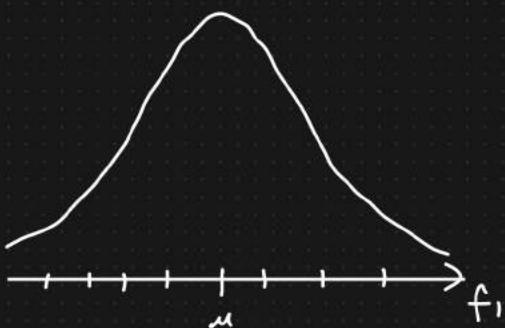
Numerical Values  $\Rightarrow$  Natural Language Processing

- ① BOW
- ② Tf-IDF
- ③ Word2Vec

### ③ Gaussian Naive Bayes

If the features are following Gaussian Distribution then we use Gaussian Naive Bayes Algorithm to solve classification problem

IRIS Dataset



[continuous features]

Age	Height	Weight	Yes/No
25	170	78	
28	160	75	
32	170		
34	140		

# Ensemble Technique

① What is Ensemble?

Combining Multiple Models  $\Rightarrow$  TRAIN  $\Rightarrow$  Predictions

Two Types

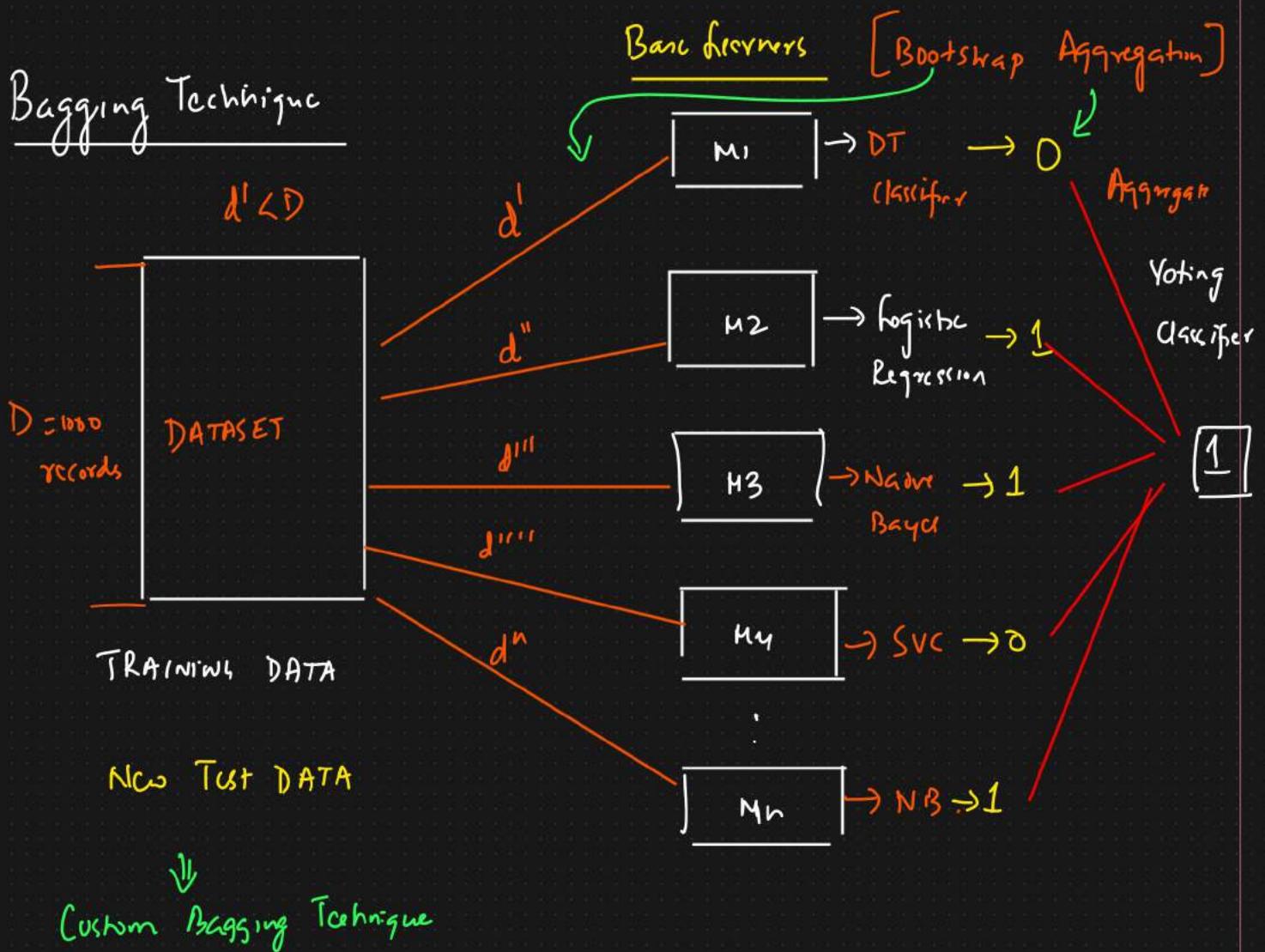
① Bagging

① Random Forest classifier  
And Regressor

② Boosting

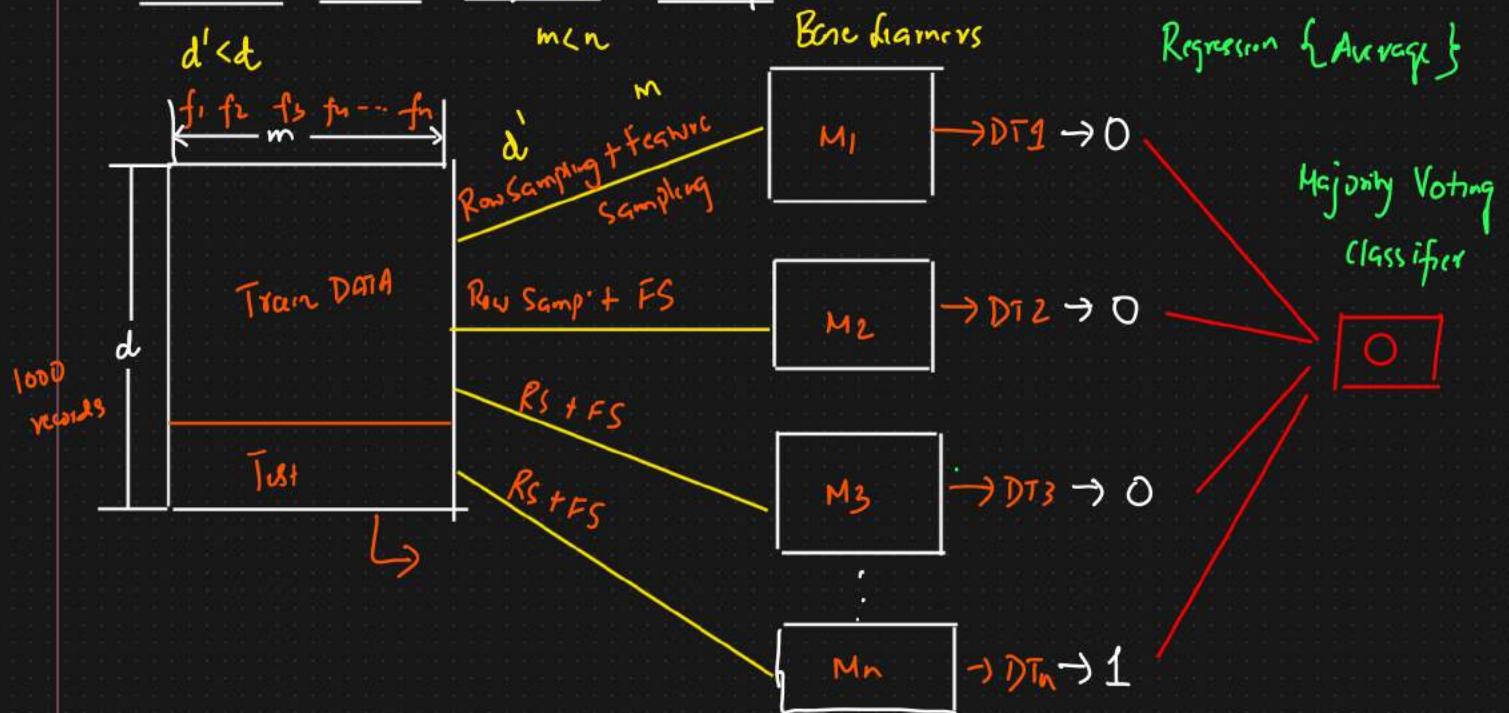
- ① AdaBoost
- ② GRADIENT BOOST
- ③ Xgboost

Bagging Technique



Reqnun → o/p → Average

# Random Forest Classifier And Regressor

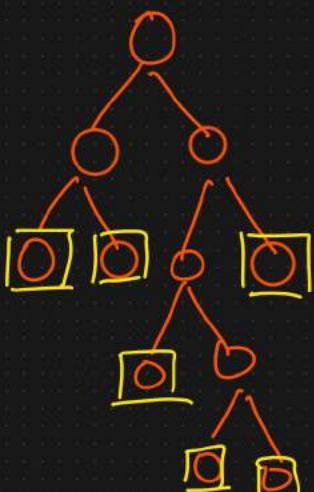


Note :

Classification  $\rightarrow$  Majority Voting classifier }  
 Regression  $\rightarrow$  Average O/P of the Model }  $N_{\text{no Data}}$

④ Why should we use Random Forest instead of DT?

## Decision Tree



## Overfitting

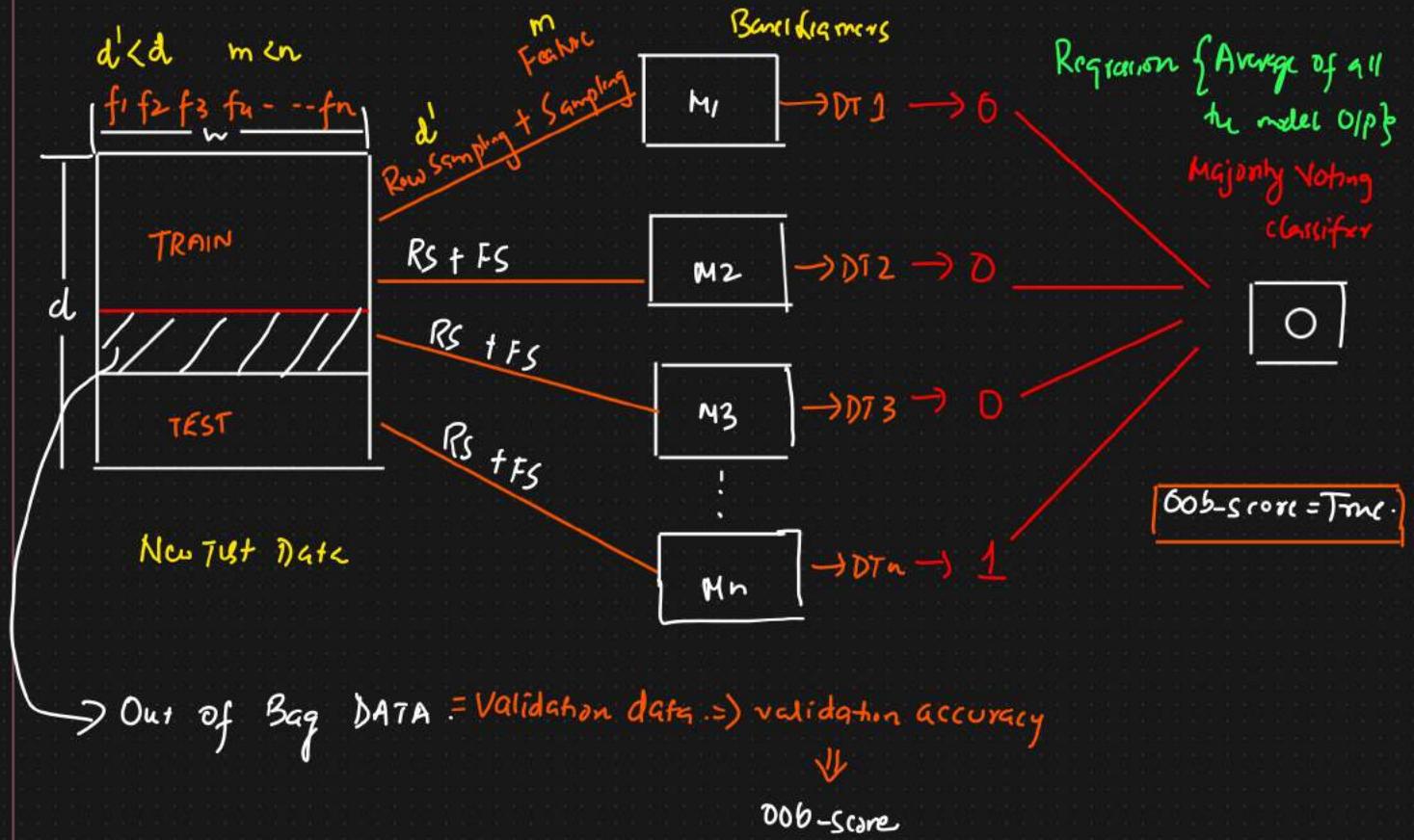
TRAINING ACC  $\uparrow \rightarrow$  Low Bias  $\rightarrow$  Low Bias

TEST DATA  $\downarrow \downarrow \rightarrow$  High Variance  $\rightarrow$  Low Variance

## Generalized Model

{ Random Forest }

## Out of Bag Score



# Boosting Algorithms

- ① Adaboost
- ② Gradient Boost
- ③ Xgboost.

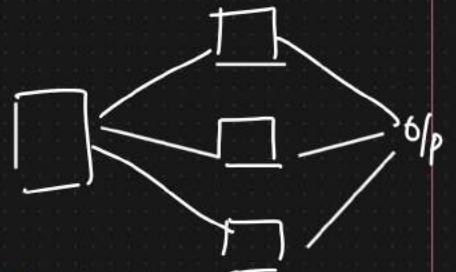
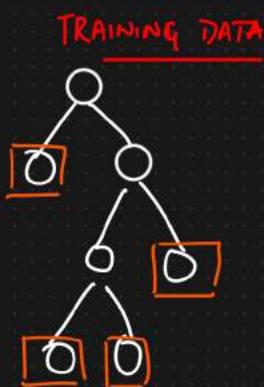
## Boosting

### Decision Tree

Overfitting :

TRAIN Acc ↑↑

TEST Acc ↓↓



### Generalized Model

Low Bias

High Variance

$\Rightarrow$  Random Forest

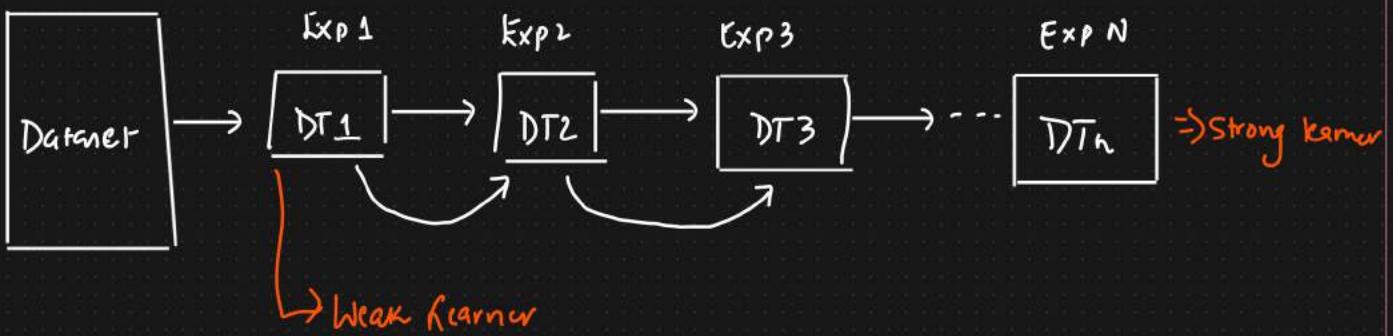
Low Bias

Low Variance

$\rightarrow \square \rightarrow \underline{\square} \rightarrow \underline{\square} \rightarrow \underline{\square} \rightarrow q_p$

### Boosting { Decision Trees Sequentially Connected }

KBC



Weak learners : Haven't learnt much from the Training Dataset

Bagging → Random Forest → Majority Voting [Average of all O/P]

Boosting → Assignments of weights to the weak learners

Boosting

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

Classification  
Regression

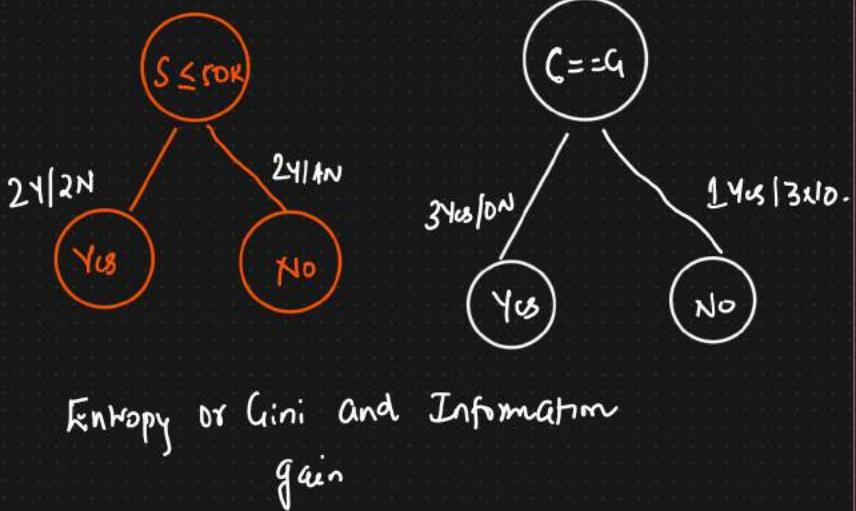
$$\{\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_n\} \Rightarrow \underline{\text{Weights}}$$

AdaBoost Classifier → Boosting Technique

Datgar

<u>Salary</u>	<u>Credit</u>	<u>Approval</u>	<u>SW</u>
$\leq 50K$	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	$\frac{1}{7}$

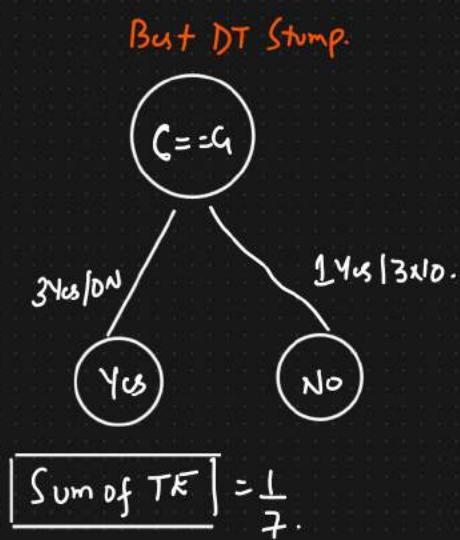
① We create Decision Tree Stump and we select the best Stump



### ③ Sum of Total Errors And performance of Stamp

Datgaf

<u>Salary</u>	<u>Credit</u>	<u>Approval</u>	<u>SW</u>
<=50K	B	No	1/7
<=50K	G	Yes	1/7
<=50K	G	Yes	1/7
>50K	B	No	1/7
>50K	G	Yes	1/7
>50K	N	Yes	1/7
<=50K	N	No	1/7



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

$$= \frac{1}{2} \ln \left[ \frac{1 - \gamma_7}{\gamma_7} \right]$$

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

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

$$f = \alpha_1(M_1) + \alpha_2(M_2) + \dots + \alpha_n(M_n) = \frac{1}{2} \ln[6] \approx 0.896$$

$$f_1 = 0.896$$

$\Rightarrow$  weight.

- ④ Update the weights for correctly and Incorrectly classified points.

Salary	Credit	Approval	SW	Updated weights
$\leq 50K$	B	No	$\frac{1}{7} \downarrow$	0.058
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	0.058
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	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
$\leq 50K$	N	No	$\frac{1}{7} \downarrow$	0.058

For correctly  
classified points

$$\begin{aligned} &= \text{Weight} * e^{-\text{SW}} \\ &= \frac{1}{7} * e^{-(0.896)} \\ &= 0.058 \\ &= \end{aligned}$$

For Incorrect Classified  
points

$$\begin{aligned} &= \text{Weight} * e^{\text{SW}} \\ &= \frac{1}{7} * e^{(0.896)} \\ &= 0.349 \end{aligned}$$

- ⑤ Normalized Weights Computation And Assigning Bins

0.41

0.17

Salary	Credit	Approval	SW	Updated weights	Normalized Wt	Bin Assignment
$\leq 50K$	B	No	$\frac{1}{7} \downarrow$	0.058	0.013	$0 - 0.08$
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	0.058	0.013	$0.08 - 0.16$ $\leftarrow 0.095$
$\leq 50K$	G	Yes	$\frac{1}{7} \downarrow$	0.058	0.013	$0.16 - 0.24$
$> 50K$	B	No	$\frac{1}{7} \downarrow$	0.058	0.013	$0.24 - 0.32$
$> 50K$	G	Yes	$\frac{1}{7} \downarrow$	0.058	0.013	$0.32 - 0.40$
$> 50K$	N	Yes	$\frac{1}{7} \uparrow$	0.349	0.500	$0.40 - 0.90$
$\leq 50K$	N	No	$\frac{1}{7} \downarrow$	0.058	0.013	$0.90 - 0.98$

## ⑥ Select data points to form Next Stump

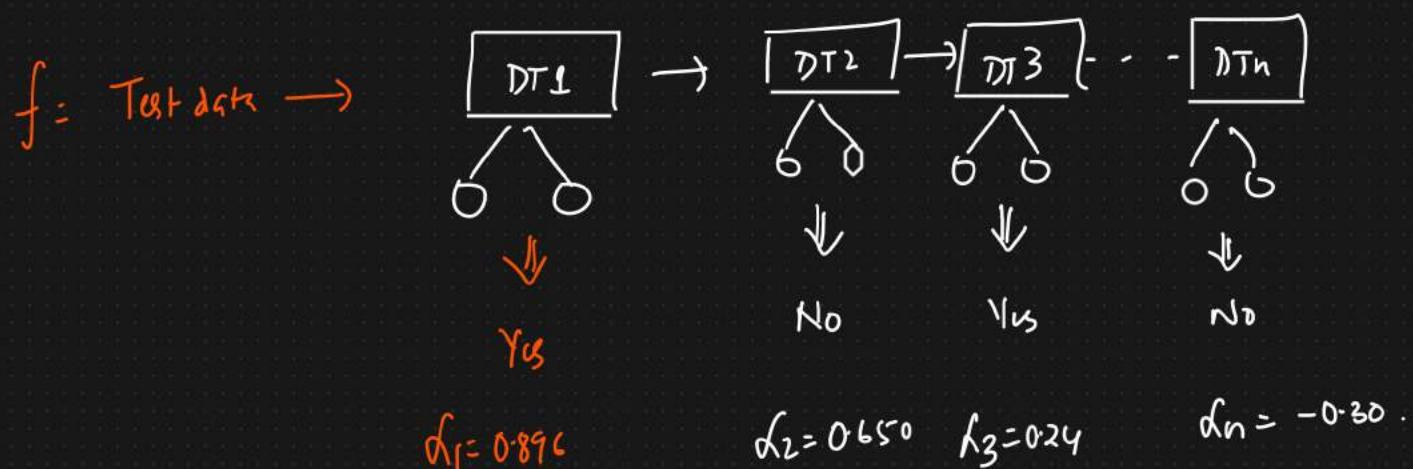
Salary	Credit	Approval	Bin Assignment
		No.	
$\leq 50K$	B	Yes	$0 - 0.08$
$\leq 50K$	G	Yes	$0.08 - 0.16 \leftarrow 0$
$\leq 50K$	G	Yes	$0.16 - 0.24$
$> 50K$	B	No	$0.24 - 0.32$
$> 50K$	G	Yes	$0.32 - 0.40$
$> 50K$	N	Yes	$0.40 - 0.90$
$\leq 50K$	N	No.	$0.90 - 0.98$

① Iteration process selecting random values 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

## ⑦ Final Prediction

Test data  $(\leq 50K, G) \leftarrow$



$$f = 0.896(\text{Yes}) + 0.650(\text{No}) + 0.24(\text{Yes}) - 0.30(\text{No})$$

$$= 1.136(\text{Yes}) + 0.350(\text{No}) \Rightarrow \text{O/P} \Rightarrow \text{Yes}$$

Performance of say(Yes) = 1.136 >

Performance of say(No) = 0.350

# Gradient Boosting Algorithms

① Regression

② Classification

Dataset		$y$	$(y - \hat{y})$	Predicted $R_2$	$\hat{y}$	$R_3$
<u>Exp</u>	<u>Degree</u>	Salary	$R_1$			
→ 2	B.E	50K	-25K	-23K	74.77	-24.77
→ 3	Masters	70K	-5K	-3K	74.97	-4.97
5	Masters	80K	5K	3	-	-
6	PhD	100K	25K	20K	-	-
		75K				

## Steps

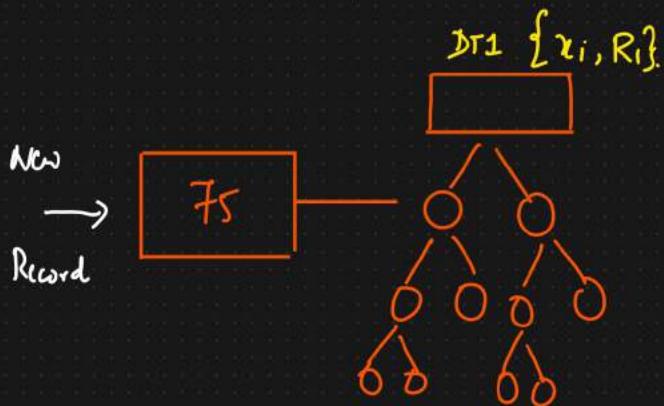
① Create a Base Model

$$\text{I/P} \rightarrow \boxed{\quad} \rightarrow \text{O/P} = 75K$$

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

② Compute Residuals, Error

③ Construct a Decision Tree consider inputs  $x_i$  and o/p  $R_i$



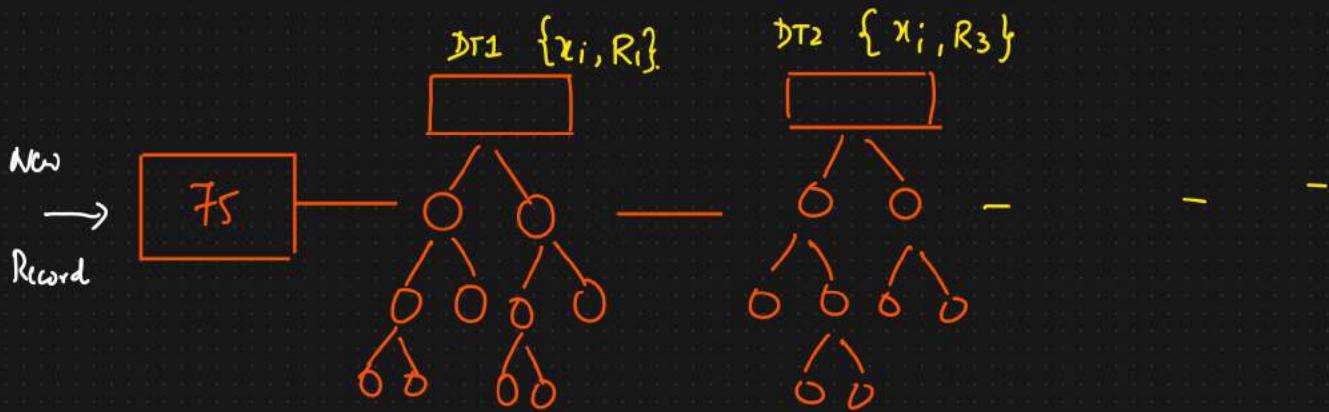
Basic Model

$$\text{Predicted O/P} = \underset{\downarrow}{75} + \underset{\downarrow \text{DT1}}{\alpha} (-23) = 75 - 23 = 52 \{ \text{Overfitting} \}$$

$$\begin{aligned}\text{Predicted O/P} &= \underset{\downarrow}{75} + \underset{\downarrow}{\alpha} (-23) & \alpha = \text{Learning Rate} \\ &= 75 + (0.01)(-23) & \alpha = 0.01 \rightarrow \{0 \text{ to } 1\} \\ &= 75 - 0.23 \\ &= 74.77\end{aligned}$$

$$\textcircled{1} \quad \text{Predicted O/P} = \underset{\downarrow}{75} + \underset{\downarrow}{\alpha} (-3)$$

$$\begin{aligned}&= 75 + (0.01)(-3) \\ &= 75 - 0.03 \\ &= 74.97\end{aligned}$$



Base Model       $\text{O/P} = \underset{\downarrow}{75} + \alpha_1(\text{DT}_1) + \alpha_2(\text{DT}_2) + \dots + \alpha_n(\text{DT}_n)$ .

Mathematical       $\alpha_0 = 1$

$$F(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_n h_n(x)$$

$\{\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_n\} \rightarrow \text{Learning Rate } \{0 \text{ to } 1\}$

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

# Xgboost ML Algorithm Classification

## Dataset

$x_1$ Salary	$x_2$ Credit	$y$ Approval	Error( $y - \hat{y}$ )
$\leq 50K$	B	0	-0.5
$\leq 50K$	G	1	0.5
$\leq 50K$	G	1	0.5
$> 50K$	B	0	-0.5
$> 50K$	G	1	0.5
$> 50K$	N	1	0.5
$\leq 50K$	N	0	-0.5

## Steps

① Construct a base Model

② Construct a Decision Tree with root node.

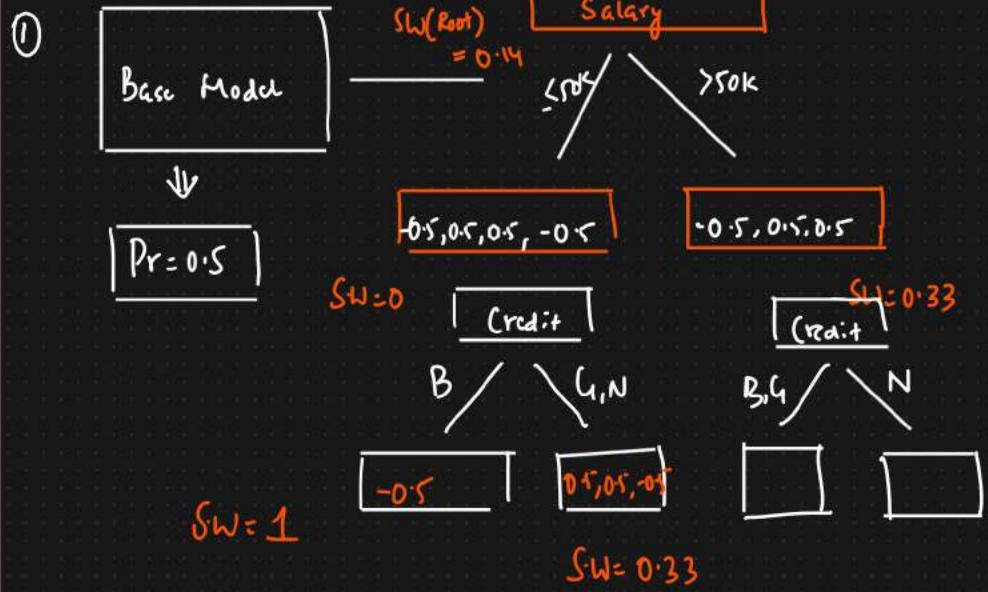
③ Calculate Similarity Weight

$$S\cdot W = \frac{(\sum \text{Residuals})^2}{\sum P_r(1-P_r)}$$

④ Calculate Gain

$$\frac{10 \cdot \frac{3}{4}}{1 \cdot \frac{7}{4}} = \frac{3}{7}$$

$$S\cdot W(\text{Root}) = \frac{0.25}{1.75} = 0.14$$

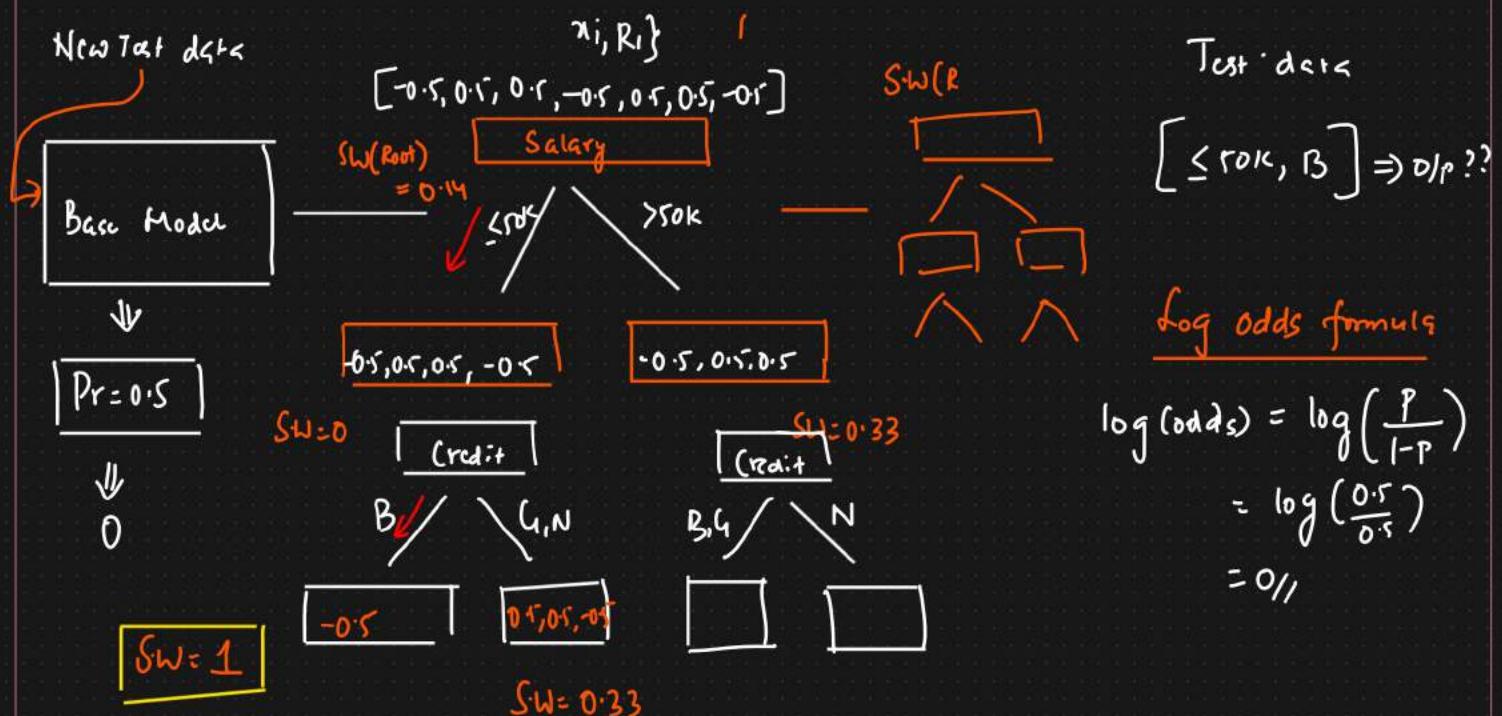


$$S\cdot W(LC) = \frac{(\sum \text{Residuals})^2}{\sum P_r(1-P_r)} = \frac{[-0.5 + 0.5 + 0.5 - 0.5]}{0.5(0.5) + 0.5(0.5) + 0.5(0.5) + 0.5(0.5)} = 0$$

$$S\cdot W(RC) = \frac{[-0.5 + 0.5 + 0.5]^2}{0.5(0.5) + 0.5(0.5) + 0.5(0.5)} = \frac{0.25}{0.75} = 0.33$$

$$\text{Gain} = 0 + 0.33 - 0.14 = 0.19\%$$

## Final O/P Classification



$$\text{Test Data O/p} = \sigma(0 + \alpha(1))$$

$\alpha$  = learning Rate

$$\alpha = 0.1$$

$$= \sigma(0 + 0.1)$$

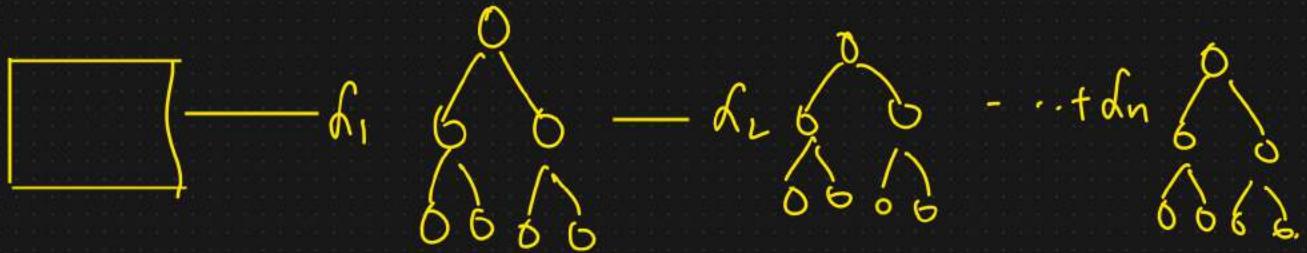
$$\sigma = \frac{1}{1+e^{-z}}$$

$$z = \frac{1}{1+e^{-0.1}}$$

$$= 0.52 // \Rightarrow \text{Threshold} \Rightarrow \boxed{0.6} \Rightarrow \text{Domain}$$

$$\boxed{0.52 < 0.6} \Rightarrow 0 \quad \text{Exptn -}$$

## Xgboost classifier



$$Op = \sum \left( \text{Base Learner} + \lambda_1(DT_1) + \lambda_2(DT_2) + \dots + \lambda_n(DT_n) \right)$$

XgBoost Classifier

# Xgboost Regressor MH Algorithm

<u>Dataset</u>	{Regression}		$\rightarrow$	$\boxed{\quad}$	$\downarrow n_K$	$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \vdots \\ \textcircled{n} \end{array}$	$\rightarrow$	$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \vdots \\ \textcircled{n} \end{array}$	$-$	$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \vdots \\ \textcircled{n} \end{array}$
	<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>				
$\rightarrow 2$	Yes		40K	-11	49.9	-9.9				$[51 + (0.1)(-10)] = 51 - 0.1 = 49.9$
$\rightarrow 2.5$	Yes		42K	-9	49.9	-7.9				
$\rightarrow 3$	No		52K	1	51.5	0.5				
4	No		60K	9	51.5	8.5				$[51 + (0.1)(5)] = 51 + 0.1 = 51.5$
4.5	Yes		62K	11	52.1	9.9				$[51 + 0.1(11)] = 51 + 1.1 = 52.1$
					$\approx 51K$					

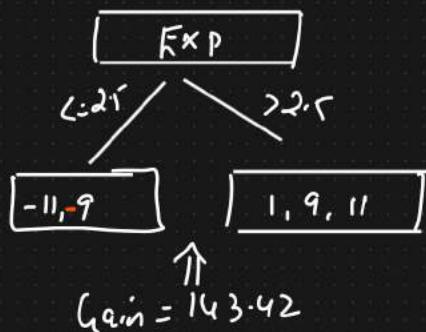
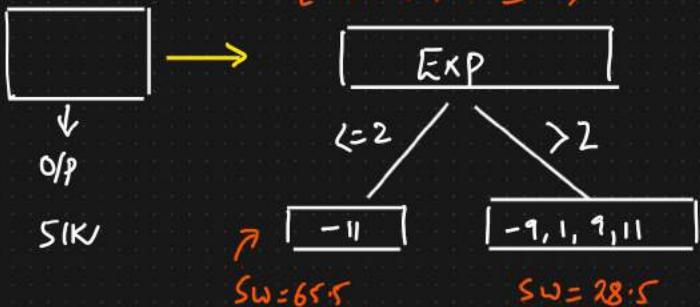
$$\text{Similarity weight} = \frac{\sum (\text{Residual})^2}{\sum p_r(1-p_r)}$$

Gain

## Steps

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

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



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

$\lambda = 1$   $\rightarrow$  No. of Residuals +  $\lambda \rightarrow$  Hyperparameter

$$SW(\text{regression}) = \frac{121}{1+1} = 121/2 = 65.5 \text{ // } \boxed{\lambda \uparrow SW \downarrow}$$

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

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

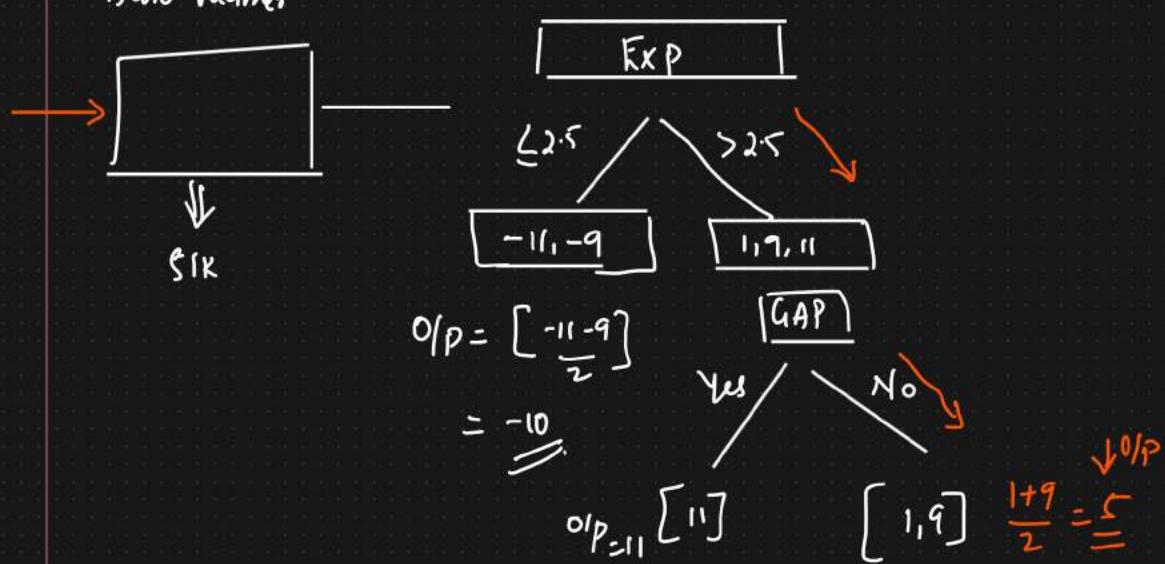
### ⑤ Calculate Gain

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

$\equiv$

Bare frame

DT 1



$\alpha$ : learning Rate     $\alpha = 0.1 \Rightarrow$  hyperparameter

$$XGB \text{ Classifier} = \text{Bare frame} + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n)$$

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

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

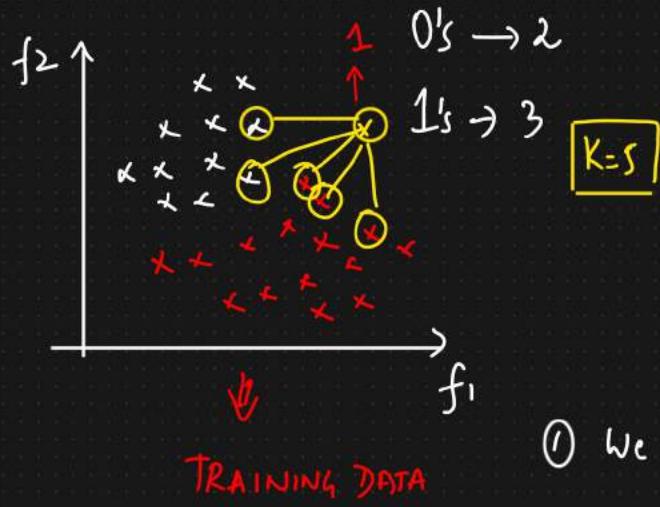


# K Nearest Neighbour (KNN)

① Classification

② Regression

① Classification



{Binary Categories}

{Multi-class}.

	$f_1$	$f_2$	$y$
-	-	0	
1			
1			
0			

① We have to initialize the  $K$  value

$$K \geq 0$$

$$K=5$$

$K=1, 2, 3, 4, 5, 6 \Rightarrow$  Hyperparameter

② Find the  $K$  Nearest Neighbour from the Test DATA

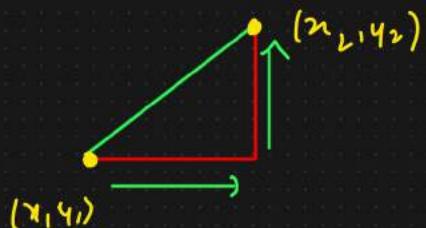
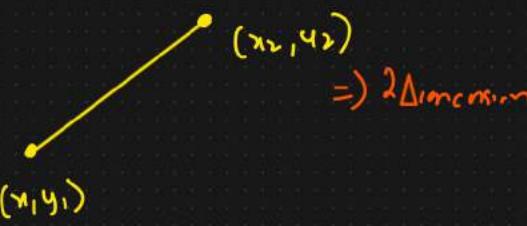
③ From those  $K=5$  how many neighbours belongs to 0 category and 1 category -

## Distance Metrics

① Euclidean Distance

② Manhattan Distance

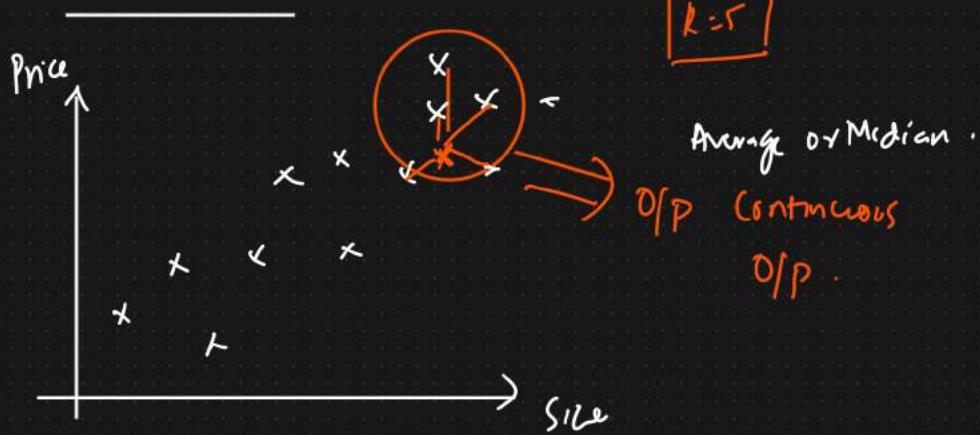




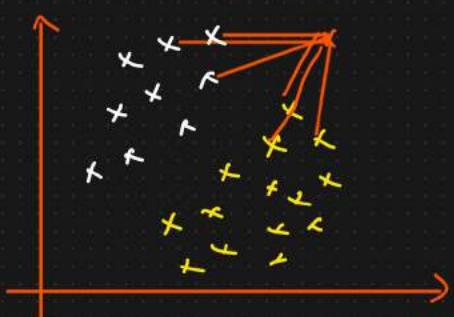
$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

## ② Regression



## Variant of KNN



Time Complexity  $\uparrow\uparrow\uparrow$

$O(n)$

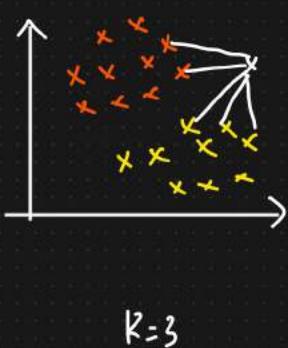
$\left. \begin{array}{l} \textcircled{1} \text{ K-D Tree} \\ \textcircled{2} \text{ Ball Tree} \end{array} \right\} \text{Binary Tree}$

$\Downarrow$   
Time complexity  $\downarrow\downarrow$

## Variants of KNN

Two Variants

- ① K D Tree ✓
  - ② Ball Tree ✓
- Binary Tree

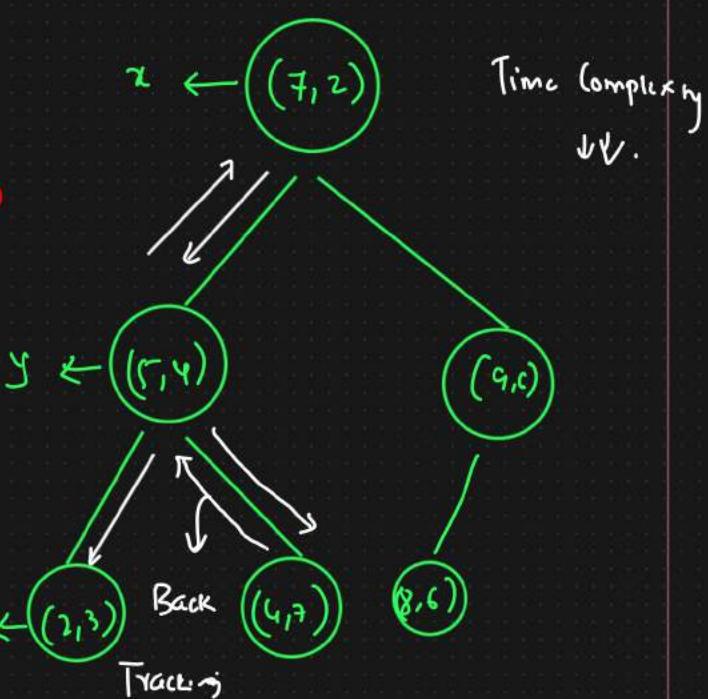
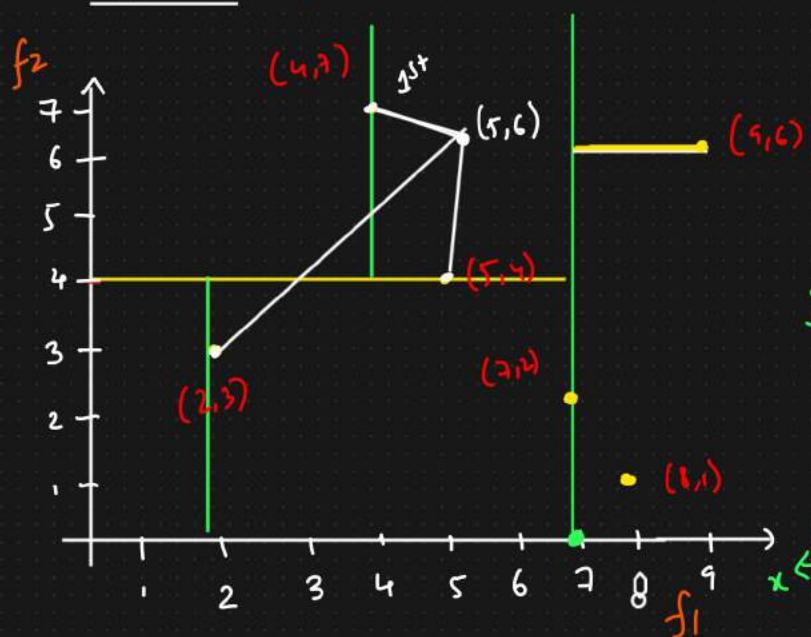


Time Complexity

Auto

Brute Search

① KD Tree  $\Rightarrow$  K-Dimensions Tree



① Median of the  $n$  coordinates

$$2, 4, \boxed{5, 7}, 8, 9$$

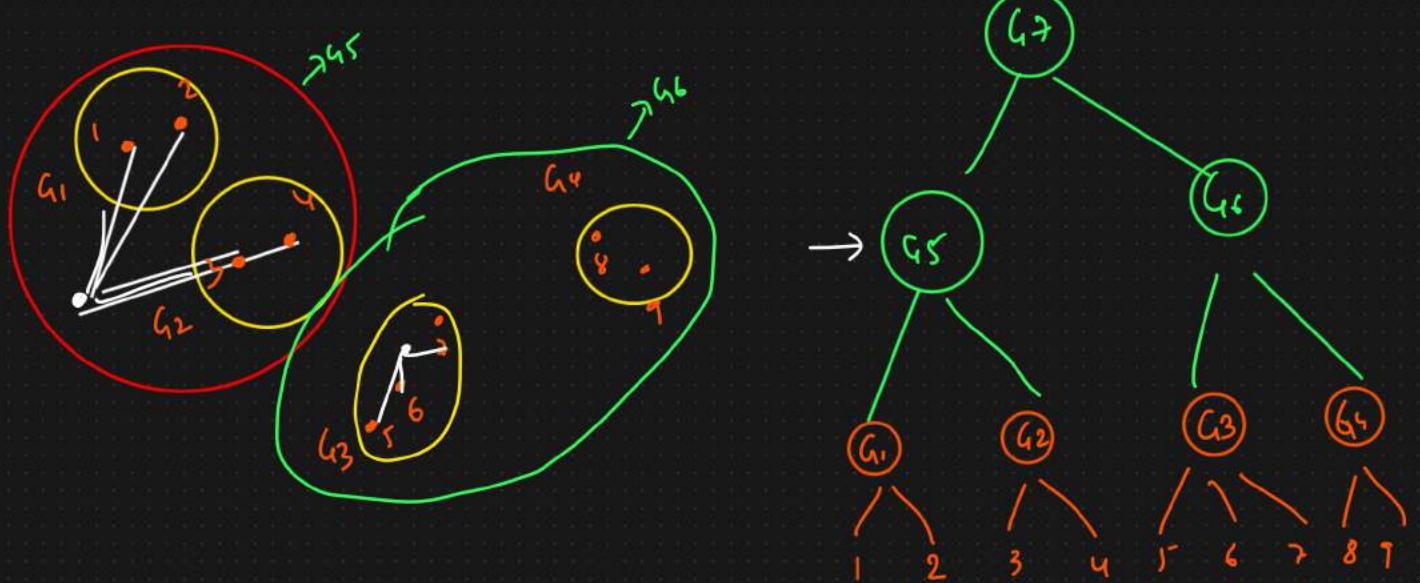
$$\downarrow$$

$$\frac{5+7}{2} = 6 \Rightarrow 7$$

② Median of  $y$  axis

$$1, 2, \boxed{3, 4}, 6, 7$$

## ② Ball Tree.



# Principal Component Analysis (PCA) [Dimensionality Reduction]

## ① Curse Of Dimensionality.

Dataset = 500 features

⇒ Price of the house ←

① House size

② No. of bedrooms

③ No. of bathrooms

3 features  
↓



Acc 1

6 features  
↓



Acc 2 ↑↑

15 features  
↓



M3

Acc 3 ↑↑

50 features  
↓



M4

Acc 4 ↓

→ Model is overfitted {

100 features →



Acc 5 ↓↓



House Price



↑  
Confused

↓↓

Price Rate

Acc ↓↓

Mumbai

1bhk



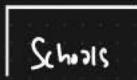
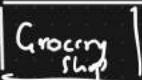
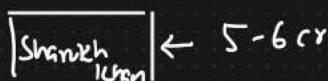
← 25-35 lakhs



← 80-100



← 2-5 cr



f1 f2 ... f<sub>n</sub>

↓  
PCA

f1 f2 ... f<sub>n</sub>

Two different ways to remove Curse of Dimensionality

100 D  
↓  
20 D  $\Rightarrow$  Problem

① Feature Selection



Imp features

② Dimensionality Reduction (PCA)



[100 features]

Feature Extraction



[20 features]

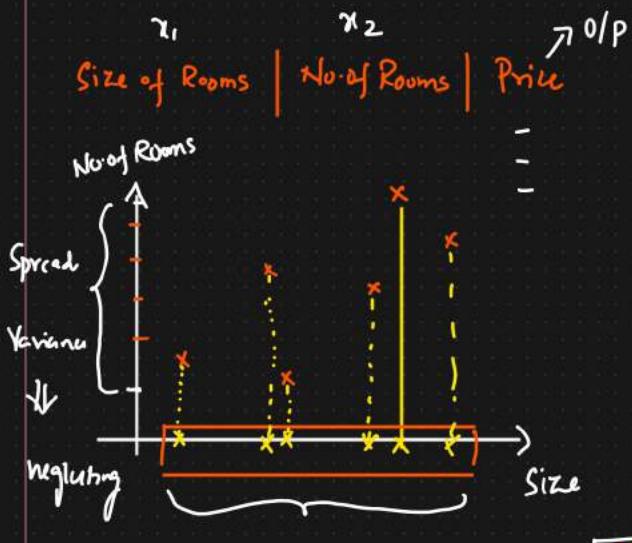
Principal Component Analysis



Eigen Value & Vectors

PCA Geometric Intuition

[Dimensionality Reduction]



① Feature Selection



2 D  $\rightarrow$  1 D

$3D - 2D$   
 $PC_1, PC_2, PC_3$

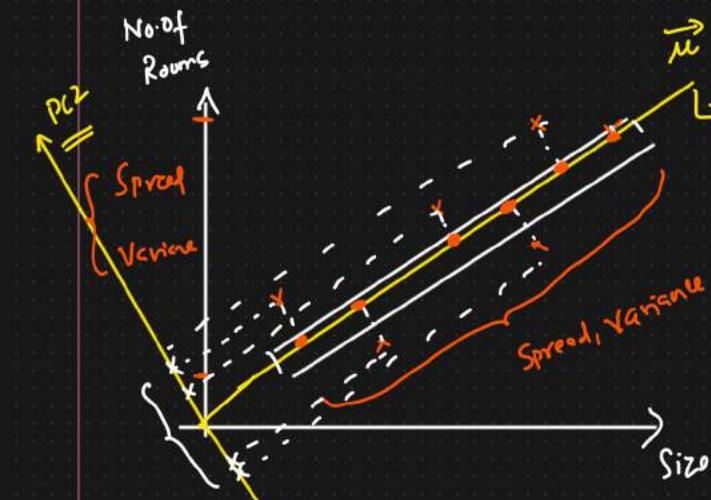
Feature Extraction



② PCA  $\rightarrow 2D \rightarrow 1D$

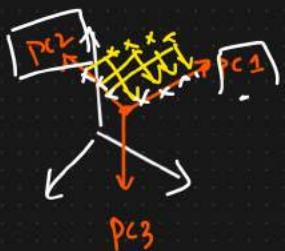
Variance, Spread  
 $\Downarrow$

Information of the DATA



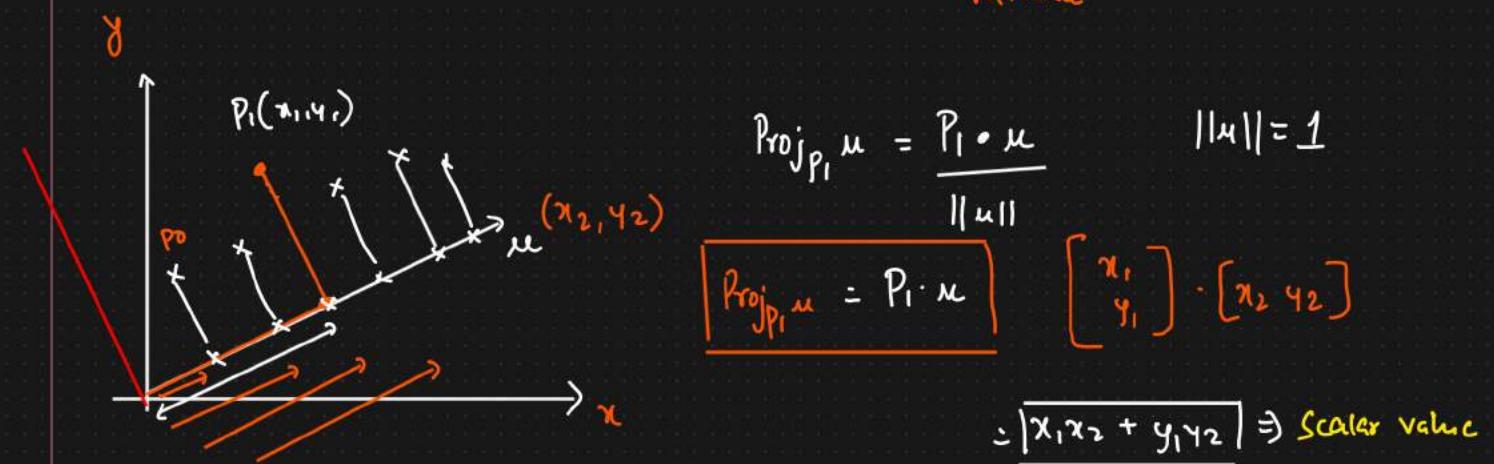
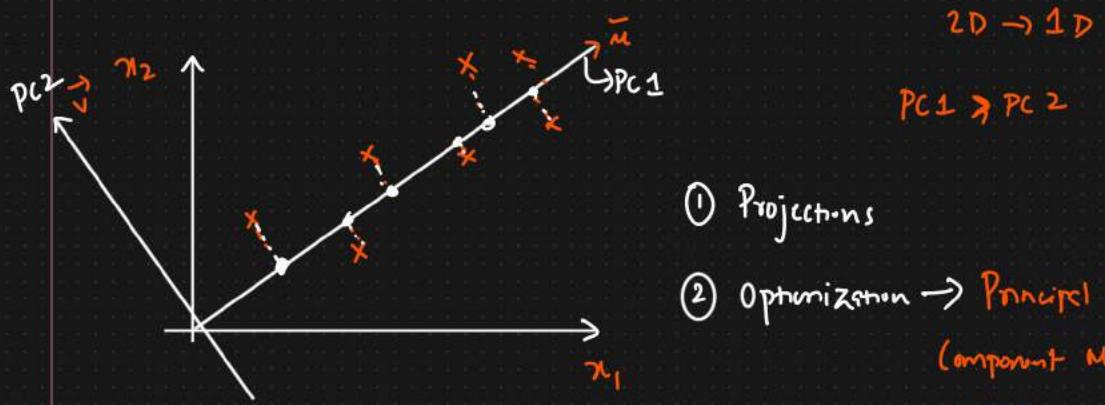
2 features  $\rightarrow$  PC1, PC2

3 feature  $\rightarrow$  PC1, PC2, PC3



3D - 2D

## Maths Intuition behind PCA Algorithms



$$\boxed{P'_0, P'_1, P'_2, P'_3, \dots, P'_n}$$

$\Downarrow$   
 Scalar Value.  
 $\Downarrow$   
 Variance

$P_0^1, P_1^1, P_2^1, P_3^1 \dots P_n^1$

$x_0, x_1, x_2, x_3 \dots x_n$

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

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

(Q) → How to find the vectors?

Eigen Value Decomposition → Eigen values and Eigen Vectors

$f_1 \ f_2 \text{ O/P}$   $dD \rightarrow 2D$

① Covariance Matrix between features  $Cov[f_1, f_2]$ .

② Eigen value and Eigen vector will be found out

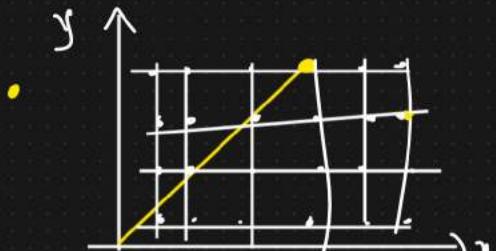
using this Covariance Matrix

$$A v = \lambda v$$

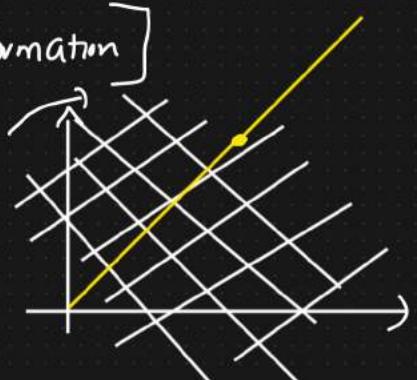
↓  
Eigen Value

③ Eigen Vector → Eigen Value → Capturing the maximum variance.

Eigen Vectors And Eigen Values  $\left[ \text{Linear Transformation} \right]$



$$Cov[x, y] = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$



$$\begin{bmatrix} \quad \end{bmatrix} * \begin{bmatrix} v \end{bmatrix} = \lambda \downarrow * v$$

↓  
Eigenvalue  
↓  
Magnitude

Eigen vector  $\rightarrow$  Max Magnitude



Max Eigen Value



PC1

Steps to calculate Eigen Value and Eigen Vector  $[2d \rightarrow 1d]$

① Covariance of features.

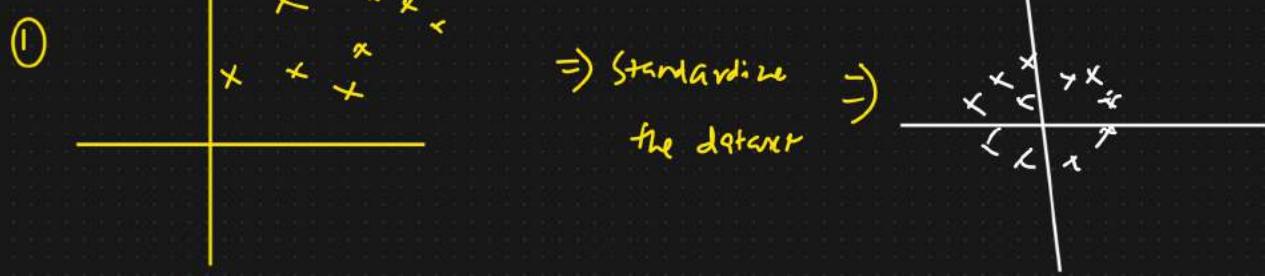
$$x, y, z \xrightarrow{\text{Cov}}$$

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

$$A = \begin{array}{c|cc} & x & y \\ \hline x & \text{Var}(x) & \text{Cov}(x,y) \\ y & \text{Cov}(y,x) & \text{Var}(y) \end{array} \quad \text{Cov}(x;x) = \text{Var}(x)$$

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

$3d - 2d \quad 3d - 1d$   
 $\Downarrow$   
 $\lambda_1 \quad \lambda_2 \quad \lambda_1$   
 $\Downarrow \quad \Downarrow \quad \Downarrow$   
 $\boxed{\lambda_1 \quad \lambda_2} \Rightarrow \text{Eigen Values} \quad \text{PC1} \quad \text{PC2} \quad \text{PC1}$



② Cov ( $x, y$ )

③ Find out the Eigen value & Eigen vector

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

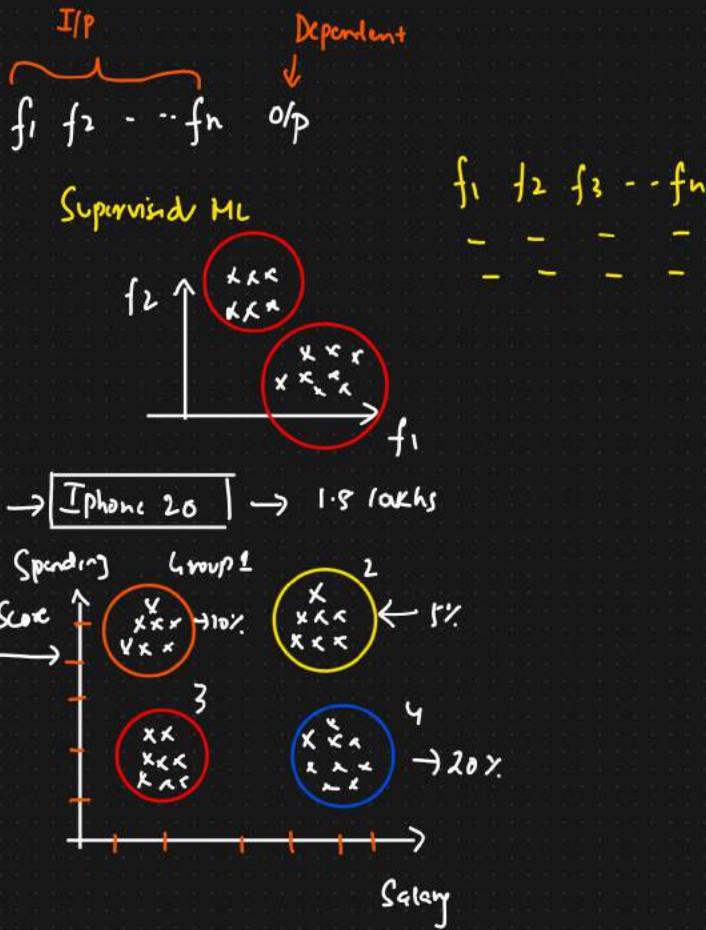
2 3 ... n

Variance ↑  $\Rightarrow$  PC1, PC2, PC3 ...

# Unsupervised Machine Learning

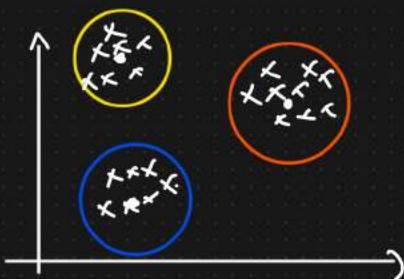
## Clustering Algorithms

- ① K Means Clustering
- ② Hierarchical Clustering.
- ③ DBScan Clustering.



## ① K Means clustering

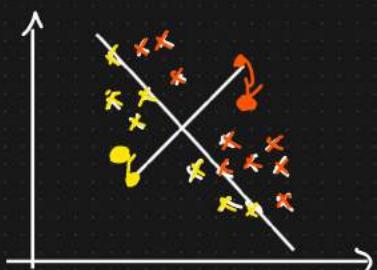
### Geometric Intuition



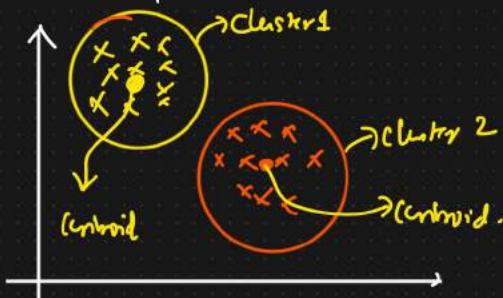
## K-Means Mathematical Intuition

Steps:

$k=2$



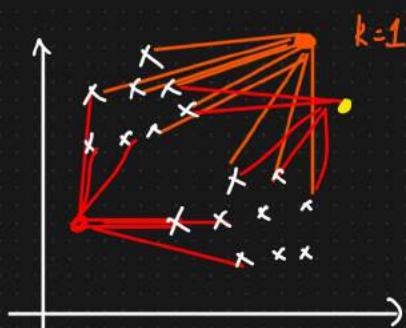
- ① Initialize some  $K \rightarrow$  centroids
- ② Points that are nearest to the centroid  $\rightarrow$  Group
- ③ Move the centroids  $\rightarrow$  Mean



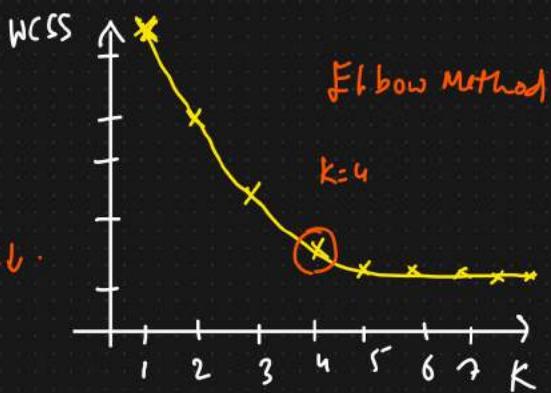
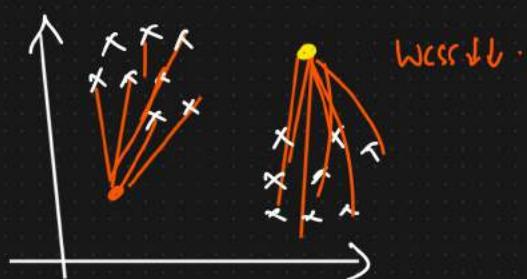
How do we select the K Value?

WCSS = Within Cluster Sum of Squares

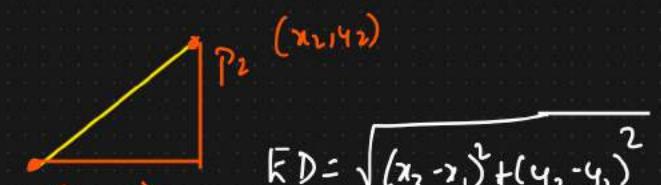
Initialize  $K=1$  to 20



$$WCSS = \sum_{i=1}^K \left( \text{Distance between point to the nearest centroid} \right)^2$$

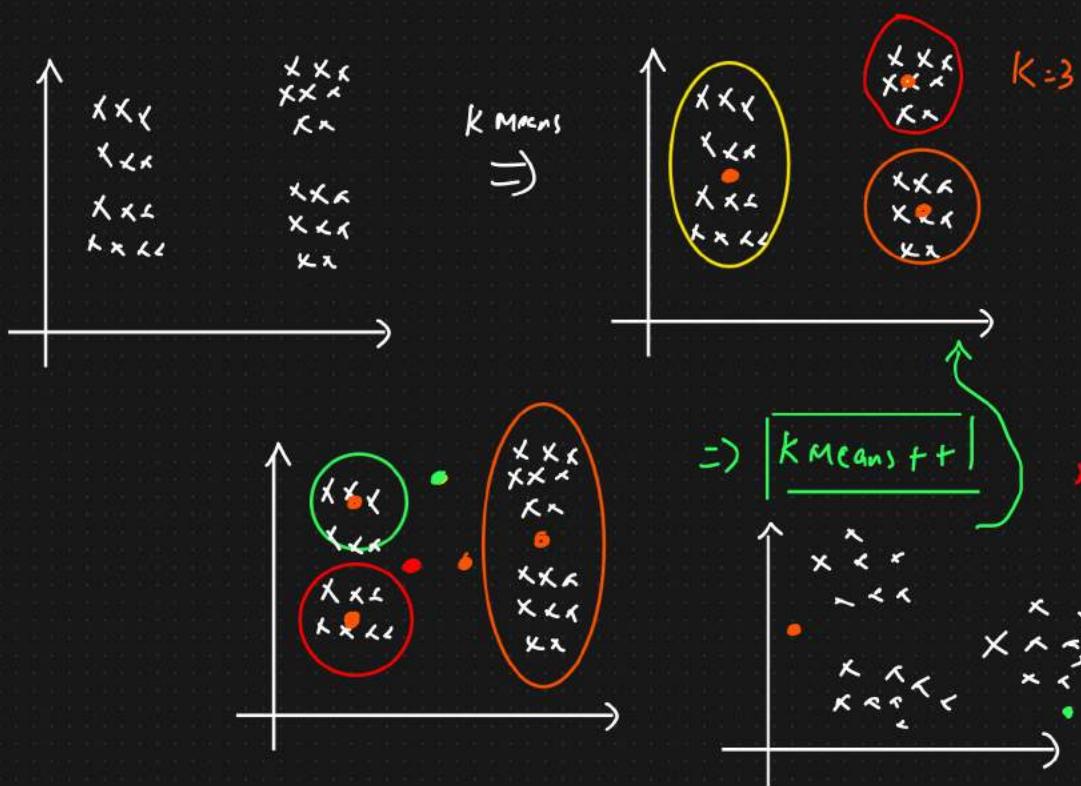


Euclidean Distance or Manhattan Distance



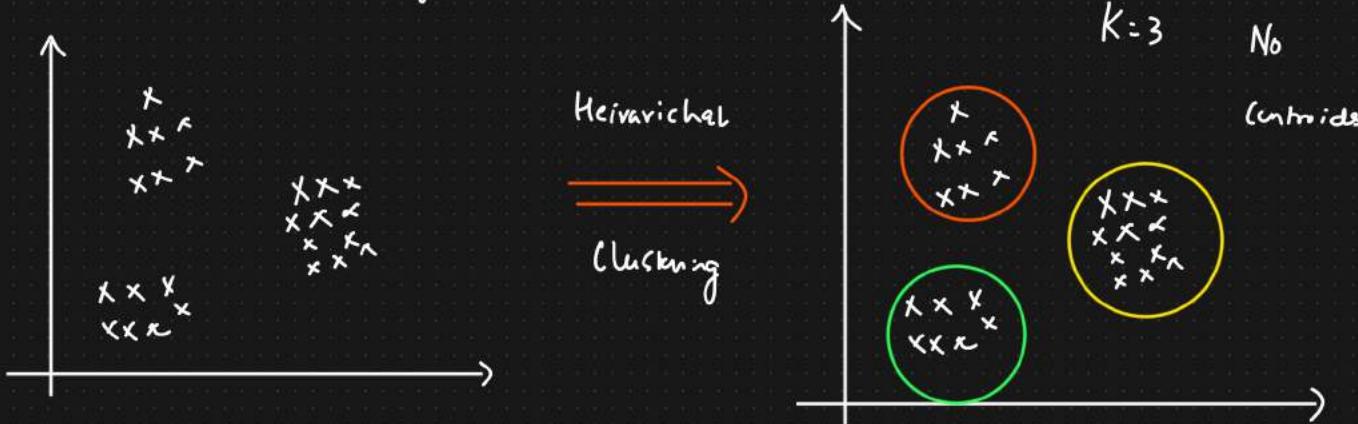
$$MD = |x_2 - x_1| + |y_2 - y_1|$$

## Random Initialization Trap (K-Means ++)



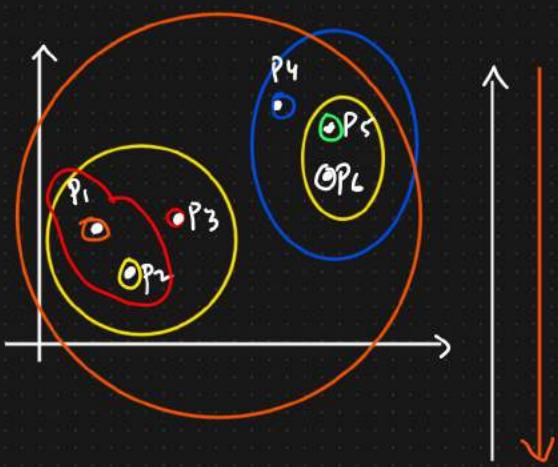
## K-Means++ Initialization Technique

## Hierarchical Clustering



## KC

- ① Agglomerative  
② Divisive }  $\Rightarrow$  Geometric Intuition

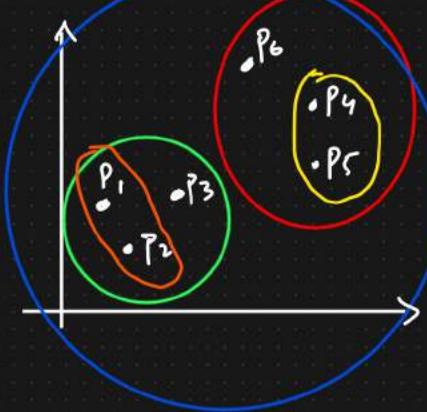


## Steps

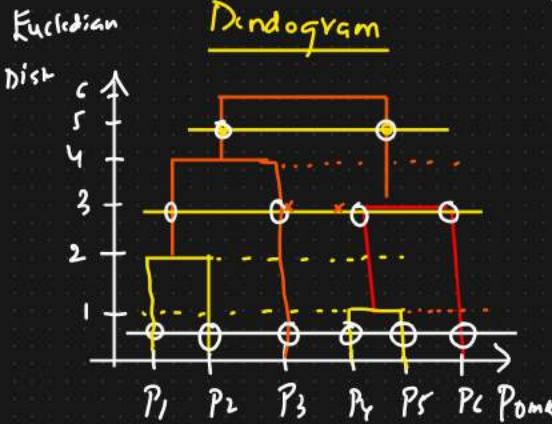
- ① For each point initially we will consider it as a separate cluster
- ② Find the nearest point and create a new cluster
- ③ Keep on doing the same process [Step 2] until we get a single cluster.

How many clusters?

$K=2$



## Threshold



- ⑧ Select the longest vertical line such that no horizontal line passed through it.

# K Means Vs Hierarchical Clustering

## Scalability And Flexibility

- ① Dataset size
  - Huge  $\Rightarrow$  K Means
  - Small  $\Rightarrow$  Hierarchical clustering
- ② Types of Data
  - Numerical data  $\rightarrow$  K Means or Hierarchical
  - Variety of data  $\rightarrow$  Hierarchical



# Silhouette Clustering

(1)

$$[a(i)]$$



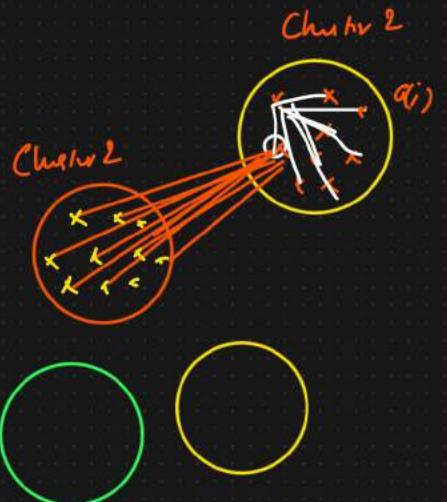
For data point  $i \in C_I$  (data point i in the cluster  $C_I$ ), let

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

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).

(2)

$$b(i)$$



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

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

to be the *smallest* (hence the **min** operator in the formula) mean distance of i to all points in any other cluster (i.e., in any 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.

③ Silhouette Score  $\text{small} \leftarrow \boxed{a_i/b_i}$   $\boxed{a_i > b_i} \Rightarrow \text{Good cluster}$

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

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

and

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

Which can be also written as:

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

From the above definition it is clear that

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

$\boxed{a_i < b_i} \Rightarrow \text{bad cluster}$