

# 7 days Machine Learning Algorithms

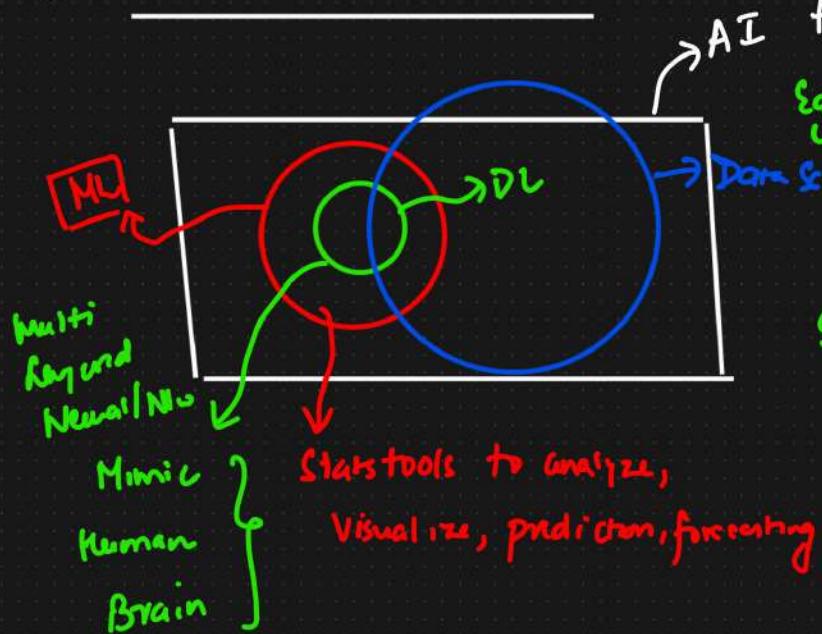
Purpose : Clear the Interviews

## Agenda

- ① Introduction to ML (AI Vs ML Vs DL Vs DS)
- ② Supervised ML and Unsupervised ML
- ③ Linear Regression (Maths & Geometric Intuition)
- ④  $R^2$  & Adjusted  $R^2$
- ⑤ Ridge and Lasso Regression

## AI application

### ① AI Vs ML Vs DL Vs DS



**AI application** is able to do its own task without any human intervention

Eg: Netflix → Action → Recommendation  
Dom Scientist → Comedy → "  
Amazon.in → iPhone → Headphones }  
Self Driving Cars → }

## Machine & Deep Learning

## Reinforcement



## Supervised ML

Age	Weight	O/p	hypothesis	O/p
24	62			weight
25	63		Independent features	$\rightarrow$ Age
21	72		Dependent feature	$\rightarrow$ weight
27	62			

y = mx + c

### ① Regression Problem

Age	Weight	O/p
24	72	
23	71	
25	71.5	
-	-	

Continuous variable

$\Downarrow$

Regression Problem

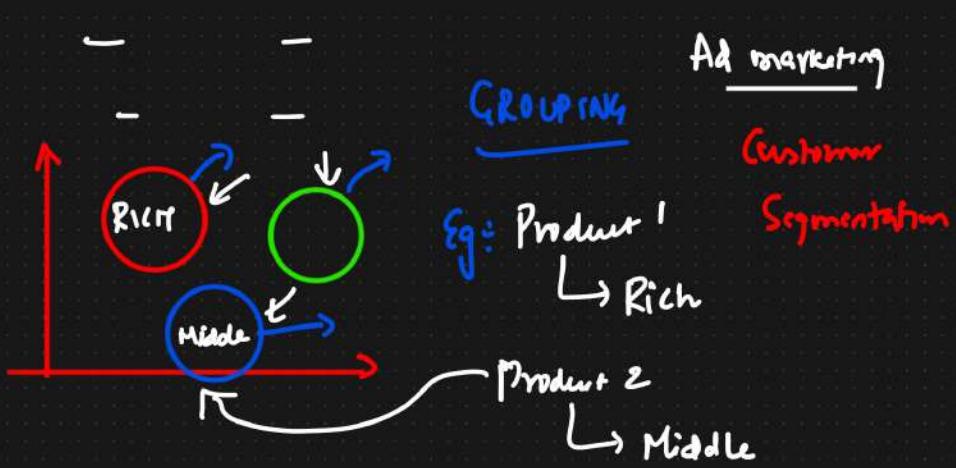
### ② CLASSIFICATION

No. of hours	No. of play hours	No. of sleep	P/F
-	-	-	P
-	-	-	F
-	-	-	P
-	-	-	F

### ③ Unsupervised ML

Salary	Age	→ {No Dependent variable}
-	-	
-	-	

Clustering → Customer Segmentation



## ② Dimensionality Reduction

1000 → lower Dimension  
 ↓  
 100

PCA, LDA

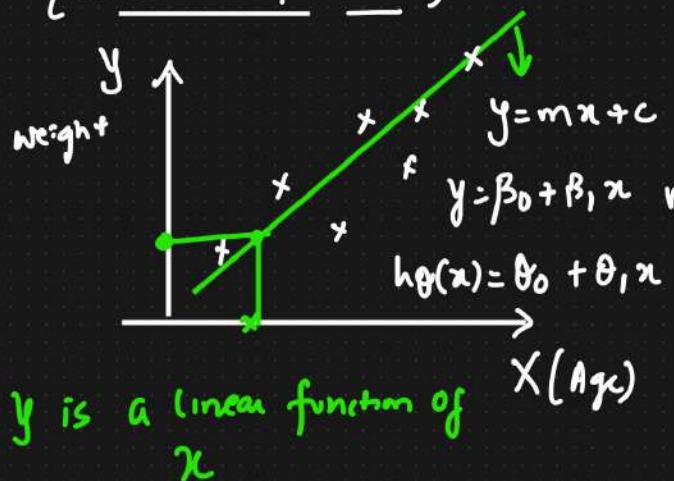
### Supervised

- ① Linear Regression
- ② Ridge & Lasso
- ③ Logistic Reg
- ④ Decision Tree
- ⑤ AdaBoost
- ⑥ Random Forest
- ⑦ Gradient Boosting
- ⑧ Xgboost
- ⑨ Naive Bayes
- ⑩ SVM
- ⑪ KNN

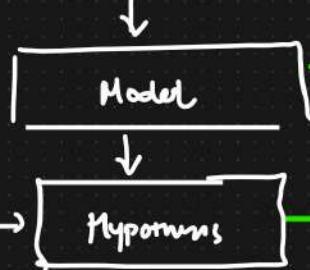
### Unsupervised

- ① K Means
- ② DBScan
- ③ Hierarchical
- ④ K Nearest Neighbor Cluster
- ⑤ PCA
- ⑥ LDA

# ① { Linear Regression }



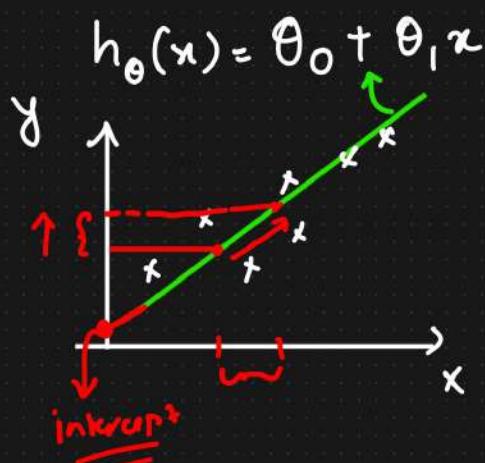
TRAIN DATASET



O/P weight

(redits : Andrew NG)

## Equation of a straight line



When  $x=0$

$\theta_0$  = Intercept  
 $\theta_1$  = Slope or Coefficient

$x_i$  = data points

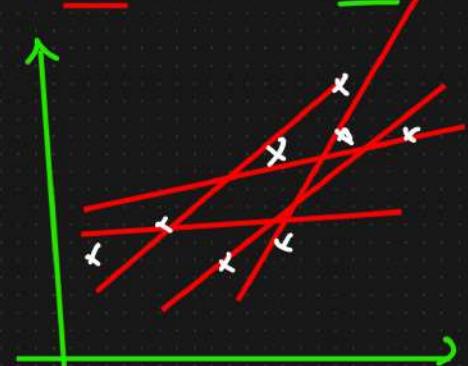
Start at point  $\rightarrow$  best fit line

## Linear Regression



Minimise

Best fit line



Hypothesis

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

Purpose  
Derivation

Cost function

$$x^n = n x^{n-1}$$

$$\frac{\partial J(x^L)}{\partial x} = \frac{x^L}{n}$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

→ Cost function

## ↳ Squared Error Function

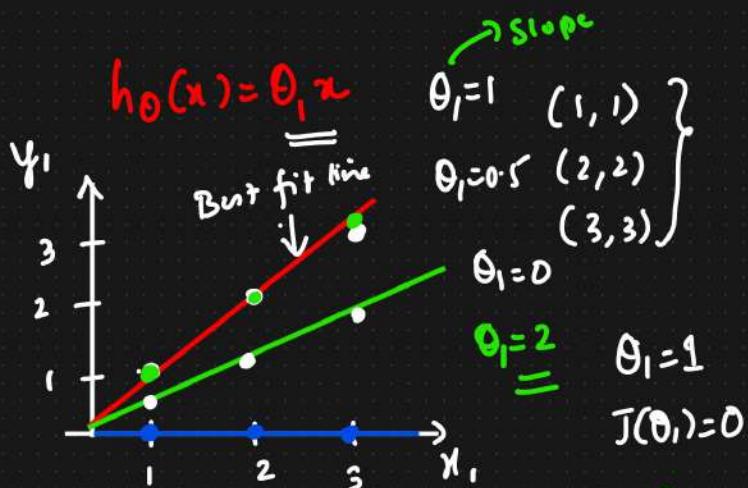
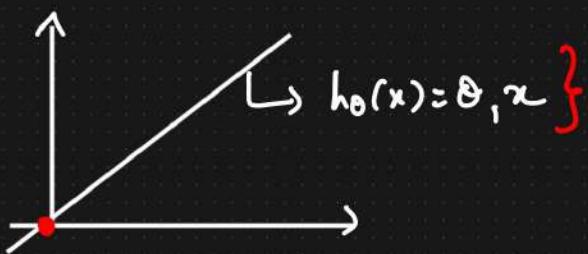
What we need to solve

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$\Downarrow$

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

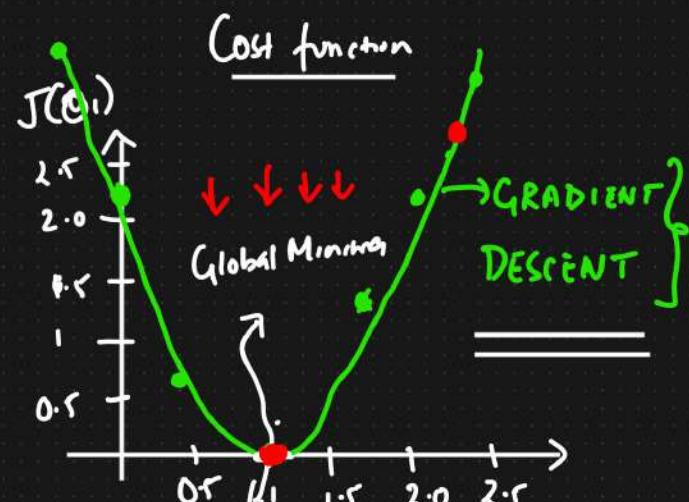
\*  $h_\theta(x) = \theta_0 + \theta_1 x \quad \text{If } \theta_0 = 0$



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

$$J(\theta_1) = 0$$



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \left[ (0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right]$$

$$= \frac{1}{2m} [0.25 + 1 + 2.25] \approx 0.58$$

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

$$= \frac{1}{6} [1+4+9]$$

$$\approx 2.3$$

$\alpha \Rightarrow$  huge

$$\alpha = 0.01$$

Convergence Algorithm

Repeat until convergence  $J(\theta_1)$

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

↳ Decaying Rate

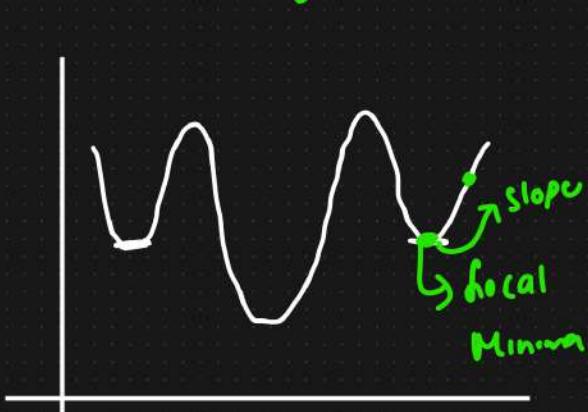
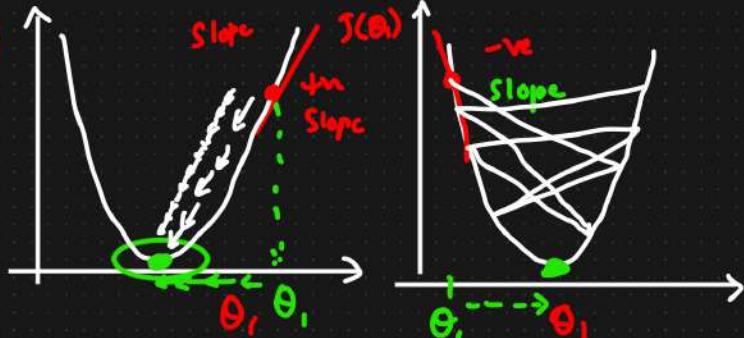
$$\theta_1 := \theta_1 - \alpha (+ve)$$

$$\theta_1 := \theta_1 - \alpha (0)$$

$$\theta_1 := \theta_1$$

$$\theta_1 := \theta_1 - \alpha (-ve)$$

$$\theta_1 := \theta_1 + \alpha (-ve)$$



## GRADIENT DESCENT Algorithm

Repeat until convergence

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

↳

$J=0$  and 1

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Convergence Algorithm :

$$\left\{ \begin{array}{l} j=0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ j+1 \Rightarrow \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right. \quad \begin{array}{l} h_\theta(x) = \theta_0 + \theta_1 x \\ \frac{\partial}{\partial \theta_0} (h_\theta(x)) = 1 \\ \frac{\partial}{\partial \theta_1} (h_\theta(x)) = x \end{array}$$

$\downarrow \alpha = 0.001 \quad \downarrow \alpha = \text{Learning Rate}$

Repeat until converge

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right.$$



## Performance Metrics

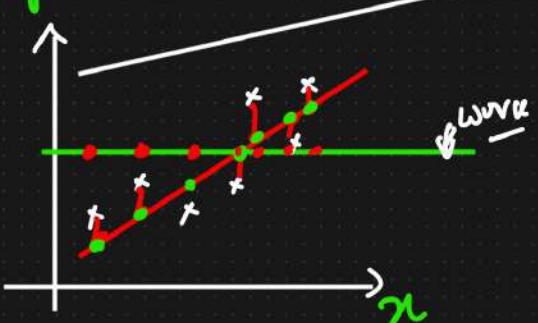
$R^2$  and Adjusted  $R^2$

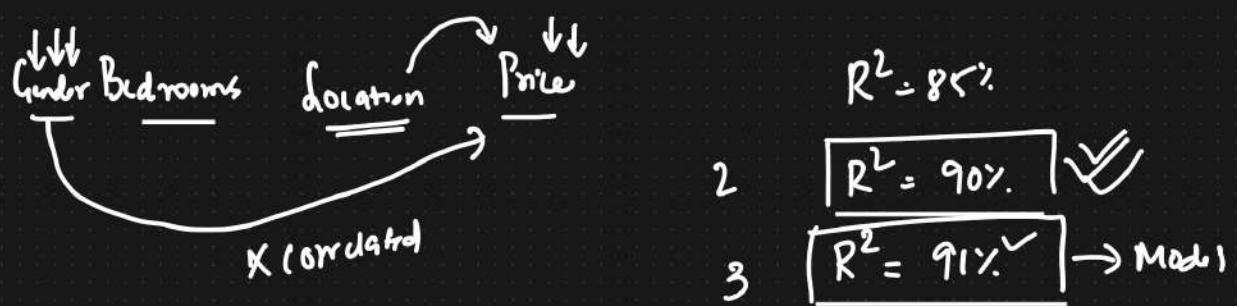
$$R^2 = 1 - \frac{SS_{Res}}{SS_{Total}}$$

$$h_\theta(x)$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Small number  
Big number  
= 90%





Adjusted  $R^2$

$p = \text{features or predictors}$

$$R^2_{\text{adjusted}} = 1 - \frac{\sqrt{N-1} R^2}{(N-1)} \quad \left\{ \begin{array}{l} p=2 \quad = R^2 = 90\% \quad R^2_{\text{adjusted}} = 86\% \\ p=3 \quad = R^2 = 91\% \quad R^2_{\text{adjusted}} = 82\% \end{array} \right.$$

$$p=2 > \frac{N-p-1}{N-p-1} \gg p=3$$

$= N = \text{No. of data points}$

$P = \text{No. of predictors}$

$R^2 \uparrow \uparrow$

$p \ggg$

# Day 2 - Linear Machine Learning Algorithm

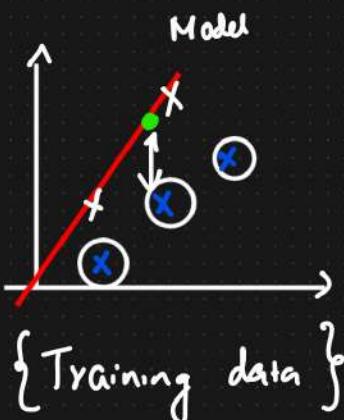
## Agenda

- ① Ridge and Lasso Regression
- ② Assumption of Linear Regression
- ③ Logistic Regression
- ④ Confusion Matrix
- ⑤ Practical Implementation

## ① Ridge And Lasso Regression

$$\text{Cost function} = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\theta_0 = 0$$



$$J(\theta_0, \theta_1) = 0 \quad \downarrow \downarrow \downarrow$$

Underfitting { High Bias  
High Variance }

- { ① Model Accuracy is bad with Training data  
② Model Accuracy is also bad with Test data }

Overfitting

✓  
(Low Bias)

Model performs well  $\rightarrow$  Training data

Fails to perform well  $\rightarrow$  Test Data ✓

(High Variance)

### Model 1

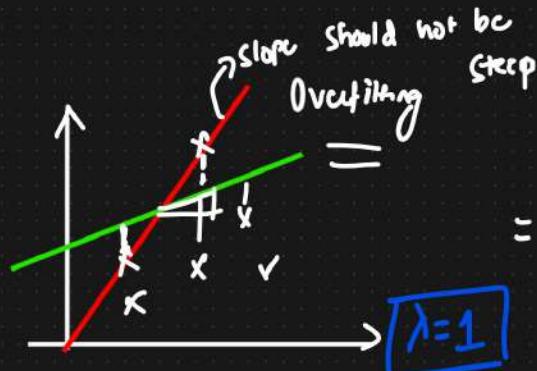
Training Acc = 90%.

Test Acc = 80%.



Overfitting

{ Low Bias, High Variance }



### Model 2

Training Acc = 92%.

Test Acc = 91%.



{ Generalized Model }

{ Low Bias }  
{ High Variance }

$$J(\theta_1) = 0$$

$$= \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y^{(i)})^2$$

$$= (\hat{y}_i - y^{(i)})^2 + \lambda (\text{slope})^2 \checkmark$$

### Model 3

Training Acc = 70%.

Test Acc = 65%.



Underfitting

High Bias, High Variance

$$h_{\theta}(x) = \hat{y} \quad \theta_1 = 2 \quad \theta_0 = 0$$

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

$$h_{\theta}(x) = \theta_1 x \quad \text{↓ slope}$$

### Ridge (L2 Regularization)

$$= 0 + 1(2)^2$$

iterations { Hyperparameter }

$$= 4/ \downarrow \downarrow \downarrow$$

R<sup>2</sup>, adjusted R<sup>2</sup>

$$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2 \quad \lambda \rightarrow \text{Hyperparameter} \checkmark$$

{ Prevent Overfitting }

$$\downarrow \quad (\text{Small value}) + 1(1.5)^2$$

Convergence

$$= (\text{Small value}) + 2.25$$

\\$

$$\approx 3 \downarrow \downarrow$$

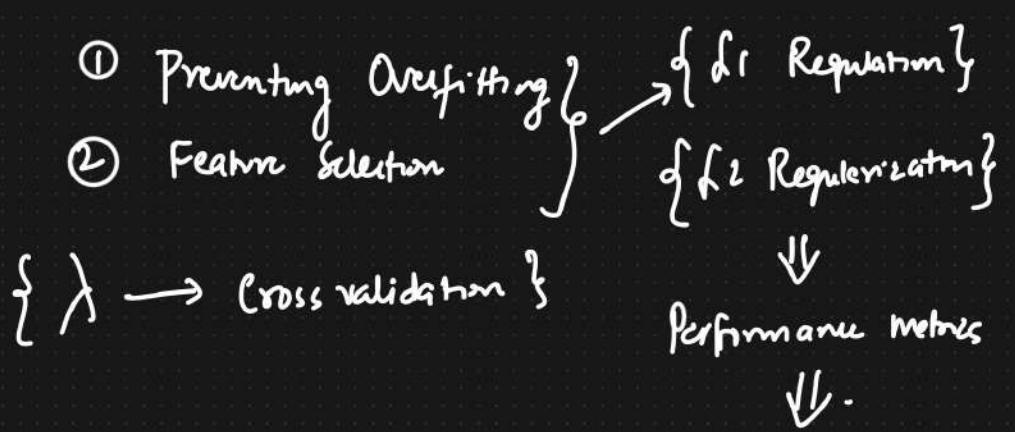
feature selection

### Lasso (L1 Regularization)

$$= (\hat{y} - y)^2 + \lambda |\text{slope}| \quad |\theta_0 + \theta_1 + \theta_2 + \theta_3 +$$

$$\theta_4 + \theta_5 + \dots + \theta_n|$$

$$h_{\theta}(x) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots + \theta_n x_n$$



## Ridge Regression ( $\lambda_2$ Norm)

$$\text{Cost function} = (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda (\text{slope})^2$$

④

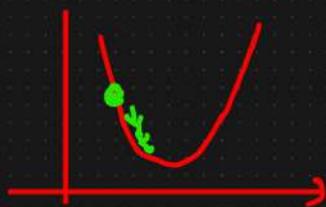
Purpose : Preventing Overfitting

$$|\theta_0 + \theta_1 + \theta_2 + \theta_3 + \dots + \theta_n|$$

## Lasso Regression ( $\lambda_1$ Reg)

$$\text{Cost function} = (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda |\text{slope}|$$

Purpose : 1) Prevent Overfitting  
2) Feature Selection



## Assumption of Linear Regression

① Normal / Gaussian Distribution  $\rightarrow$  Model will get trained well

✓ ② {Standardization {Scaling data}  $\rightarrow$  Z-score  $\mu=0, \sigma=1$ }

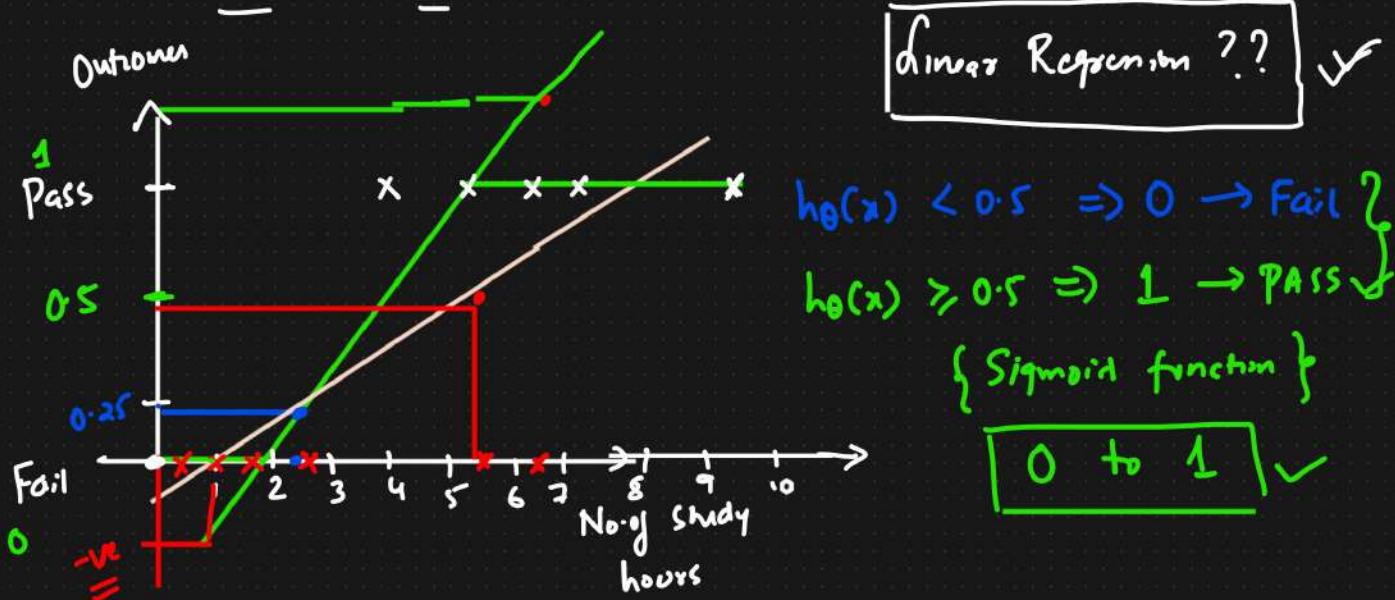
③ Linearity  $x_3$   $\boxed{x_1 \quad \cancel{x_2}}$   $\boxed{y}$

Variation Inflation factor?

④ Multi Collinearity

# Logistic Regression (Classification) → Binary Classification

No. of study      No. of play      P/F  
 —                    —                    P  
 —                    —                    F



Linear Regression ?? ✓

$h_\theta(x) < 0.5 \Rightarrow 0 \rightarrow \text{Fail}$

$h_\theta(x) \geq 0.5 \Rightarrow 1 \rightarrow \text{PASS}$  ✓

{ Sigmoid function }

0 to 1 ✓

## Decision Boundary Logistic Regression

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

$$h_\theta(x) = \theta^T x$$

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

Squash

x x x ↗

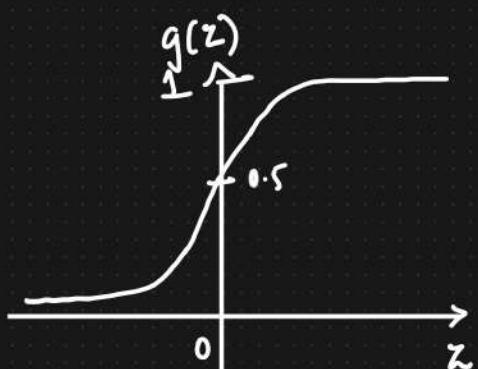


$$h_\theta(x) = g(\theta_0 + \theta_1 x_1)$$

$$\text{let } z = \theta_0 + \theta_1 x$$

$h_\theta(x) = g(z)$  Sigmoid or Logistic function

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



$g(z) \geq 0.5 \quad \{ \quad \checkmark$   
 When  $z \geq 0$

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

## Training Set

$$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^n, y^n)\}$$

$y \in \{0, 1\} \rightarrow 2 \text{ o/p}$

$$h_{\theta}(z) = \frac{1}{1 + e^{-z}}$$

$$z = \theta_0 + \theta_1 x$$

(change parameter  $\theta_1$  ?)

## Cost function

Linear Regression  $J(\theta_0) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^i) - y^i)^2$

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

Logistic Regression

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

Logistic Regr  
Cost function  $= \frac{1}{2} (h_{\theta}(x^{(n)}) - y^{(n)})^2$

We cannot use this

cost function for logistic

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

## Gradient Descent

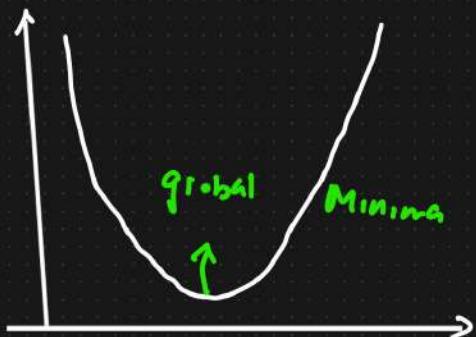
### Non convex function



Local Minima  
Problem

## Gradient Descent

### Convex function



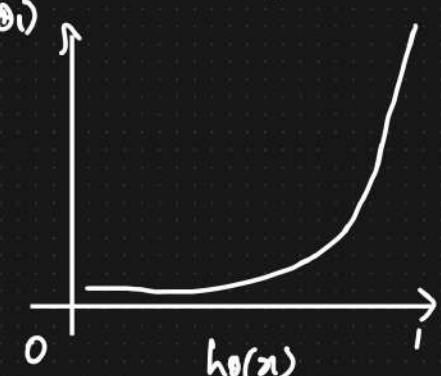
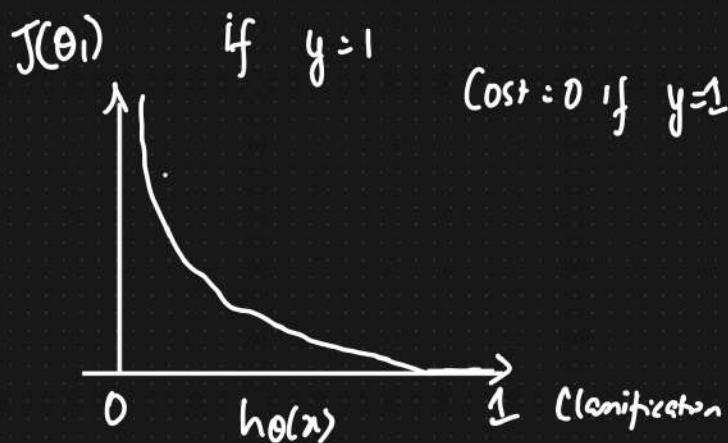
Global  
Minima

## Logistic Regression Cost function

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

$$J(\theta_1) = \begin{cases} -\log(h_{\theta}(x^i)) & y=1 \\ -\log(1-h_{\theta}(x^i)) & y=0 \end{cases}$$

if  $y=0$



$$\text{Cost}(h_{\theta}(x^i), y) = \begin{cases} -\log(h_{\theta}(x^i)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x^i)) & \text{if } y=0 \end{cases}$$

$$\boxed{\text{Cost}(h_{\theta}(x^i), y) = -y \log(h_{\theta}(x^i)) - (1-y) \log(1-h_{\theta}(x^i))}$$

if  $y=1$        $\Downarrow$   
cost function.

$$\text{Cost}(h_{\theta}(x^i), y) = -\log(h_{\theta}(x^i)) \quad \left. \right\}$$

if  $y=0$

$$\text{Cost}(h_{\theta}(x^i), y) = -\log(1-h_{\theta}(x^i)) \quad \left. \right\}$$

$$J(\theta_0) = -\frac{1}{2m} \sum_{i=1}^m \left[ (y^i \log(h_\theta(x^i)) + (1-y^i) \log(1-h_\theta(x^i))) \right]$$

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

Repet until convergence

→ {  
 $\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$ .  
} } } }

### Performance Metrics {Classification Problem}

		Actual		Pred	Actual
$x_1$	$x_2$	$y$	$\hat{y}$		
-	-	0	1	1	3
-	-	1	1	0	1
-	-	0	0	1	2
-	-	1	1	0	1
-	-	0	1		
-	-	1	0		

Predicitn      Confusion matrix

	1	0	Actual
Pred	TP	FP	↓
	FN	TN	Confusion matrix

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\begin{aligned} \textcircled{1} \quad & 0 \rightarrow 900 \\ & 1 \rightarrow 100 \end{aligned} \quad \begin{aligned} \text{Imbalance} &= \frac{3+1}{3+2+1+1} = \frac{4}{7} \\ \text{DATASPLIT} &= 0.57 = 57\% \end{aligned}$$

$$\begin{aligned} & 0 \rightarrow 600 \\ & 1 \rightarrow 400 \end{aligned} \quad \begin{aligned} \text{Balanced} & \quad 0 : 900 \\ & \quad 1 : 100 \end{aligned} \quad \begin{aligned} \text{Model} \rightarrow 0 &= \frac{900}{1000} = 90\% \end{aligned}$$

• TPR, Sensitivity

① Precision

$$\frac{TP}{TP + FP}$$

② Recall

$$\left\{ \frac{TP}{TP + FN} \right\}$$

③ F-Score.

		Actual	
		1	0
Pred	1	TP	FP
	0	FN	TN

{ Tom Stock market  
is going to crash } → Precision  
{ Spam classification } → Recall  
{ Has CONFE or NOT }

$$\underline{\underline{F-\text{Beta}}} = (1+\beta^2) \frac{\text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

$$\beta=1 \quad \approx (1+1) \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{F1-Score} \quad \approx \frac{2(\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

$$\frac{\text{Harmonic Mean}}{=} \frac{2xy}{x+y}$$

$$\beta=0.5 \quad (1+(0.5)^2) \frac{P \times R}{(0.25) P + R}$$

F0.5 Score

$$\beta=2 \quad FN \gg FP$$

F2 Score

# 3<sup>rd</sup> Day → Machine Learning Algorithms

## Agenda

- ① Practicals
- ② Naive Baye's Intuition
- ③ KNN algorithms

→ Simple Examples

## Previous Session

- ① Linear Regression
- ② Ridge & Lasso
- ③ Logistic Regression

⇒ Complex

## ① Naive Baye's Intuition {Classification}



{Baye's THEOREM}

Rolling a Dice

{1, 2, 3, 4, 5, 6}

{Independent Events}

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

$$P(2) = \frac{1}{6}$$

Dependent Event

First Event ✓

$P(R) = \frac{3}{5} \rightarrow R$  ✓

Dependents =  $P(G) = \frac{2}{5}$

↙ Green Marble

$P(G) = \frac{2}{4} = \frac{1}{2} \rightarrow G$

Conditional probability

$$P(R \text{ and } G) = P(R) * P(G|R)$$

$$P(A \text{ and } B) = P(A) * P(B|A)$$

$$\Rightarrow P(A \text{ and } B) = P(B \text{ and } A) \quad \{ \text{Yes} \}$$

$$P(A) * P(B/A) = P(B) * P(A/B)$$

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

## Bayes Theorem

CRUX

## Naïve Bayes

$$x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ \cdots x_n \ \boxed{y} \rightarrow y^p$$

→ - - - - - - - - ✓  
→  
→

$$P(y/x_1, x_2, x_3, \dots, x_n) = \frac{P(y) * P(x_1, x_2, \dots, x_n | y)}{P(x_1, x_2, \dots, x_n)}$$

$$= P(y) * P(x_1|y) * P(x_2|y) * P(x_3|y) \dots P(x_n|y)$$

$$P(x_1) \neq P(x_2) * P(x_3) * \dots * P(x_n)$$

DARMET

$$x_1 = \underline{\hspace{2cm}} \quad x_2 = \underline{\hspace{2cm}} \quad x_3 = \underline{\hspace{2cm}} \quad x_4 = \underline{\hspace{2cm}} \quad y = \underline{\hspace{2cm}}$$

→ | \_\_\_\_\_ | Yes ✓  
| \_\_\_\_\_ | No ✓

$$P(y_{\text{obs}} | x_i) = \frac{P(y_{\text{obs}}) * P(x_1 | y_{\text{obs}}) * P(x_2 | y_{\text{obs}}) * P(x_3 | y_{\text{obs}}) * P(x_4 | y_{\text{obs}})}{\sum_{y_{\text{pred}}} P(y_{\text{pred}}) * P(x_1 | y_{\text{pred}}) * P(x_2 | y_{\text{pred}}) * P(x_3 | y_{\text{pred}}) * P(x_4 | y_{\text{pred}})}$$

Constant  $\rightarrow P(x_1) \neq P(x_2) \neq P(x_3) \neq P(x_4)$  #fixed  
Ignore

$$P(y = \text{No} / x_i) = P(\text{No}) * P(x_1 / \text{No}) * P(x_2 / \text{No}) * P(x_3 / \text{No}) * P(x_4 / \text{No})$$

**constant**  $\longrightarrow P(x_1) \wedge P(x_2) \wedge P(x_3) \wedge P(x_4)$  #fixed

$x_i$ :  $\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

$$P(\text{Yes} | x_i) = \boxed{0.13}$$

$$P(\text{No} | x_i) = \boxed{0.05}$$

$$\Downarrow \quad \geq 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

$$P(\text{Yes} | x_i) = \frac{0.13}{0.13 + 0.05} = 0.72 = 72\% \quad \boxed{\quad}$$

$$P(\text{No} | x_i) = 1 - 0.72 = 0.28 = 28\%$$

### DATA SET

Day	Outlook	Temperature	Humidity	Wind	Binary Class	
					Play Tennis	
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	

$x_i$  Outlook  $P(\text{Sunny} / \text{Yes})$

	Yes	No	$P(Y)$	$P(N)$
Sunny	2	3	2/5	3/5
Overcast	4	0	4/4	0/4
Rain	3	2	3/5	2/5
<u>Total</u>	<u>9</u>	<u>5</u>		

### Temperature

	Yes	No	$P(Y)$	$P(N)$
Hot	2	2	2/4	2/4
Mild	4	2	4/6	2/6
Cold	3	1	3/9	1/9
<u>Total</u>	<u>9</u>	<u>5</u>		

PLAY

	Yes	No	$P(\text{Yes})$	$P(N)$
Yes	9			
No		5		
Total		14	$\frac{9}{14}$	$\frac{5}{14}$

→ Test (Sunny, Hot) → O/P

$$P(\text{Yes} | (\text{Sunny}, \text{Hot})) = P(\text{Yes}) * P(\text{Sunny} / \text{Yes}) * P(\text{Hot} / \text{Yes})$$

$$\cancel{P(\text{Sunny}) * P(\text{Not})}$$

$$= \frac{1}{14} * \frac{2}{7} * \frac{2}{9}$$

$$= \frac{2}{63} = 0.031$$

$$P(\text{No} | \text{Sunny, Not}) = P(\text{No}) * P(\text{Sunny} / \text{No}) * P(\text{Not} / \text{No})$$

$$\cancel{P(\text{Sunny}) * P(\text{Not})} \rightarrow \text{constant}$$

$$= \frac{8}{14} * \frac{3}{7} * \frac{2}{5}$$

$$= \frac{3}{35} = 0.085$$

$$P(\text{Yes} | \text{Sunny, Not}) = 0.031 = 1 - 0.73 = 0.27 = 27\%$$

$$P(\text{No} | \text{Sunny, Not}) = 0.085 = \frac{0.085}{0.031 + 0.085} = 0.73 = 73\%$$

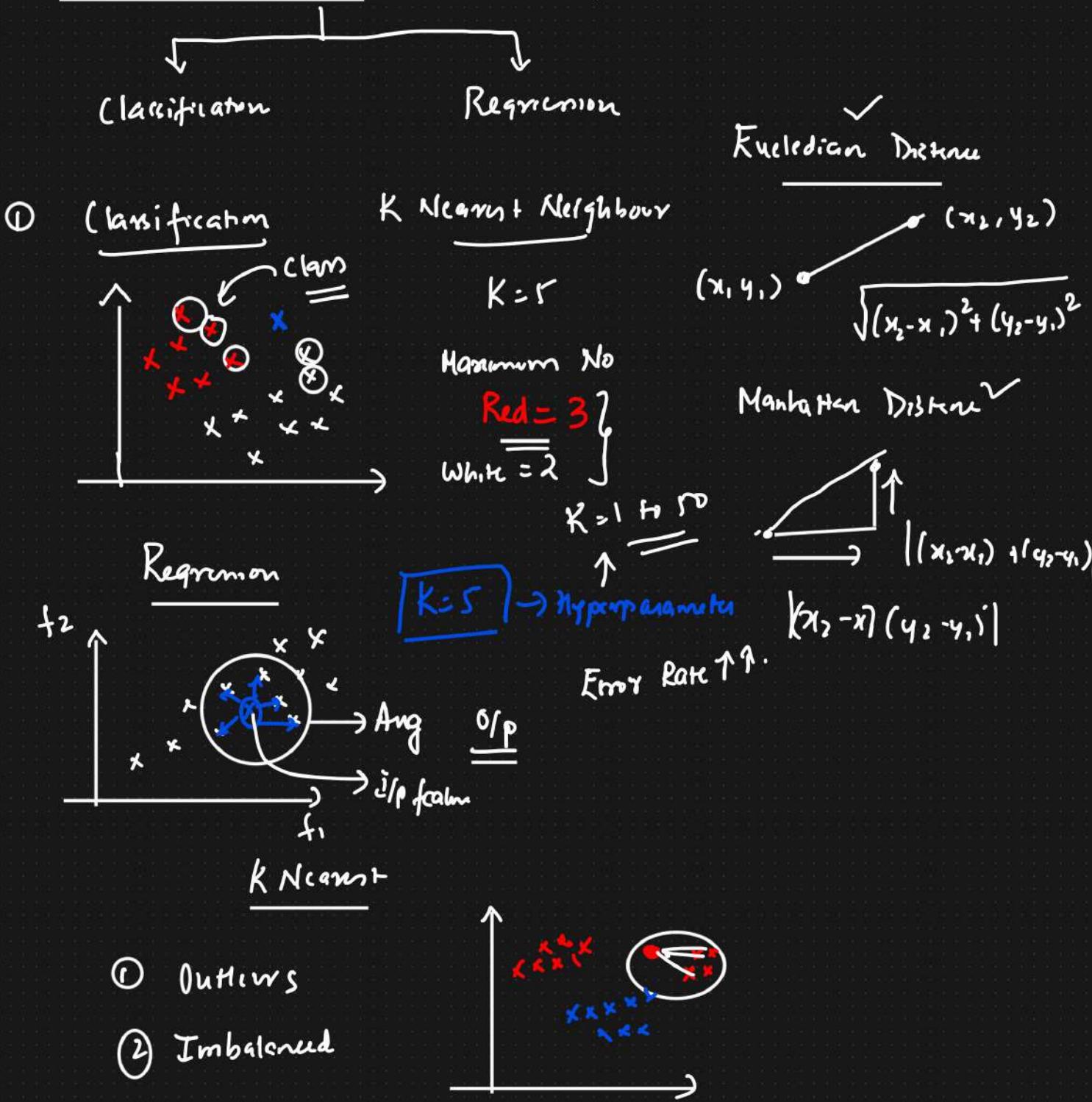
$\rightarrow (\text{Sunny, Not}) \rightarrow \text{Yes or No}$

Always  $\rightarrow \text{No}$  ✓

Assignment

(Overcast, Mild)  $\rightarrow$  Naive Bayes?

## ② KNN Algorithm {K Nearest Neighbour}



## Day 4 - Machine Learning Algorithms

## Agenda

- ① Decision Tree CLASSIFICATION
  - ② DECISION TREE REGRESSION
  - ③ PRACTICAL IMPLEMENTATION
  - ④ Ensemble Techniques

## Agenda

{ DAY 1, DAY 2, DAY 3 }

## Experience

Decision Tree {Solving many usecases}



if (age <= 18):

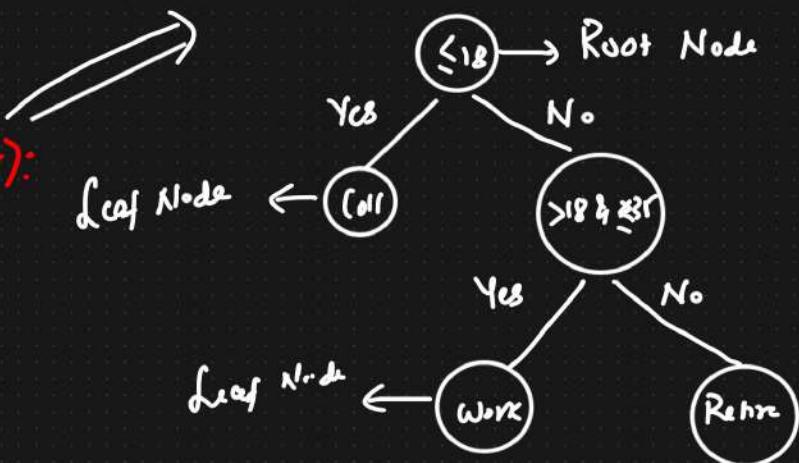
Print ("College")

if ( $age > 18$  and  $age \leq 35$ ):

### Print (work)

*else :*

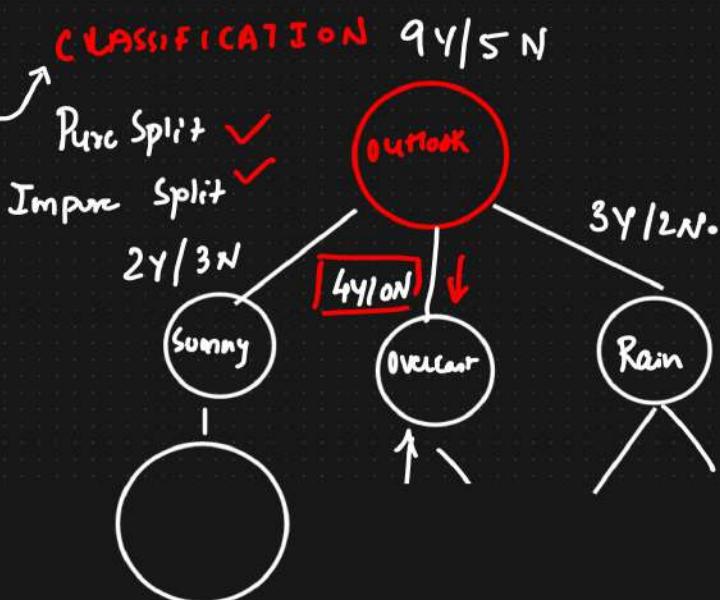
Print ("Retire")



## DECISION TREE

Nest if else  $\Rightarrow$  Decision Tree

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
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

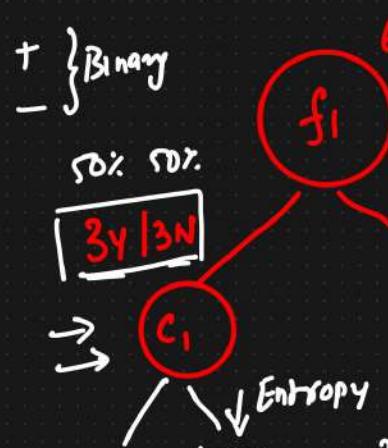




- ② How the features are selected
- ↳ Information Gain ??

### ① Entropy

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

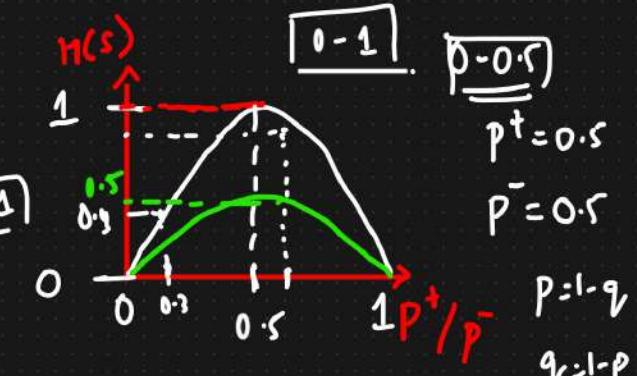


$$\begin{aligned} \text{Entropy } H(S) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= -\frac{1}{2} \log_2 \frac{1}{2} \\ &= \boxed{0} \rightarrow \text{Pure Split} \end{aligned}$$

Purity Test → Entropy

### ① Gini Impurity

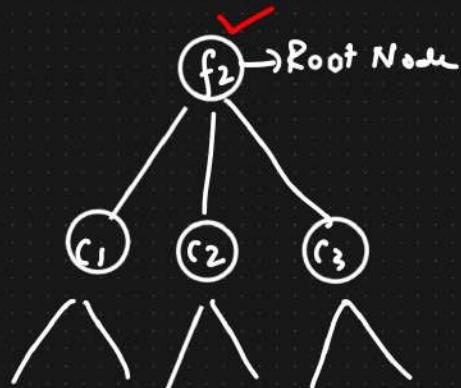
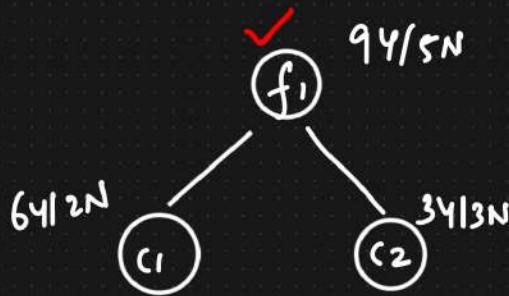
$$G.I. = 1 - \sum_{i=1}^n (P_i)^2$$



$$\begin{aligned} H(S) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= \boxed{1} \quad \checkmark \end{aligned}$$

↳ Impure Split

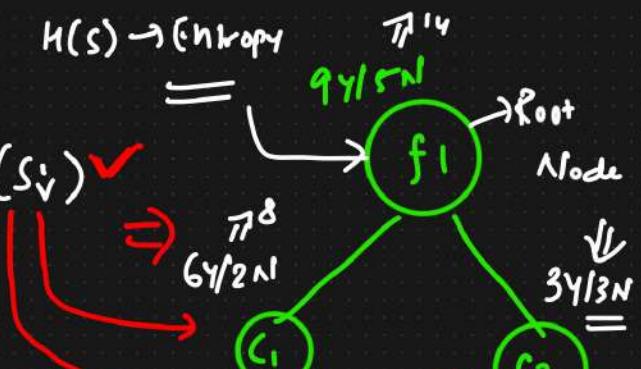
- ② Which feature to take to split??



## Information Gain

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

Root Node



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

$$= -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$\approx = 0.94$$

$$H(c_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$\frac{H(c_1) = 0.81}{H(c_2) = 1}$$

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

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

Using which feature  
Should I start splitting  
first

$$\text{Gain}(S, f_2) = 0.051$$

$$\text{Gain}(S, f_2) \gg \text{Gain}(S, f_1)$$

## Gini Impurity

$n=2$  output {Yes, No}

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

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

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

$$= 1 - \left[ \frac{1}{2} \right] = 0.5$$



Entropy = 1

Gini Impurity = 0.5

Entropy  $\rightarrow$  flag

Fast  
Gini >> Entropy

continuous

$$f_1 \text{ O/P} \Rightarrow \boxed{f_1}$$

$$\frac{2.3}{1.3} \rightarrow \frac{2.3}{2.3}$$

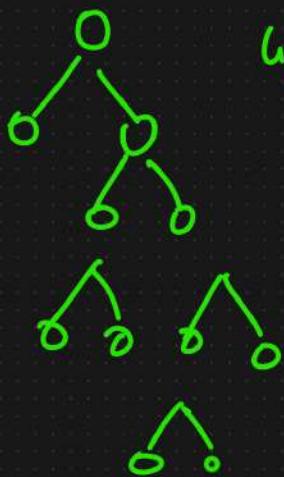
$$\frac{1.3}{1.3} \rightarrow \frac{3}{3}$$

$$\frac{4}{4} \rightarrow \frac{4}{4}$$

$$\frac{5}{5} \rightarrow \frac{5}{5}$$

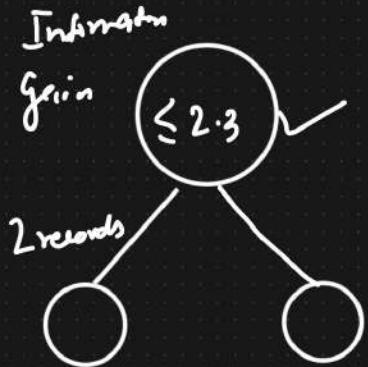
$$\frac{7}{7} \rightarrow \frac{7}{7}$$

$$\frac{3}{3} \rightarrow \frac{3}{3}$$

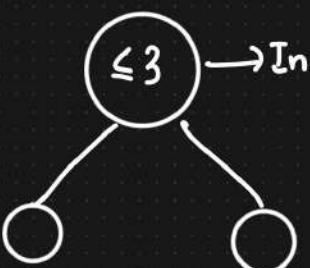


Gini Impurity  $\rightarrow$  Simple Maths

Information Gain



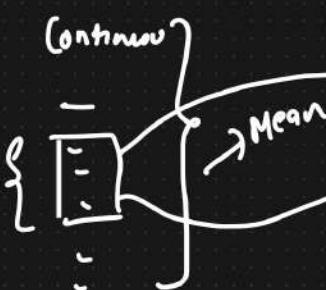
Information Gain



But Info gain

Decision Tree Regressor

$f_1 \ f_2 \text{ O/P}$



$f_1$  Mean

MSE

$\boxed{\text{MSE Or MAE}}$

$$\frac{1}{2m} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

O/P

Overfitting

Hypoparameters

Decision  $\rightarrow$  Overfitting

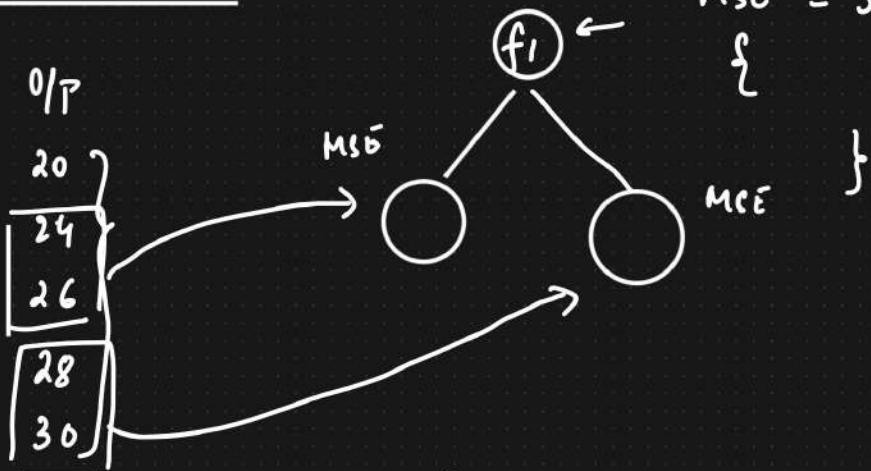
- { ① Post Pruning }
- { ② Pre Pruning }



## Decision Tree Regressor

$$MSE = 37$$

$f_1$	0/P
$C_1$	20
$C_2$	24
-	26
-	28
-	30



post pruning

pre pruning

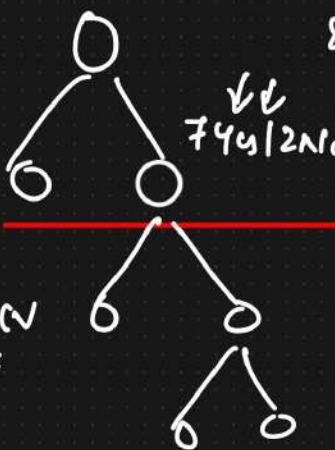
Hypoparameter

max\_depth, max\_leaf

GridSearchCV

80%

44/210



# Day 5 — Machine Learning Algorithms

## Agenda

- ✓ ① Ensemble Techniques
  - Bagging
  - Boosting
- ✓ ② Random Forest
- ✓ ③ AdaBoost
- ④ Xgboost → Youtube channel

{ DJANGO }

FLASK

Stats

ML

NLP

EDA

Deep learning

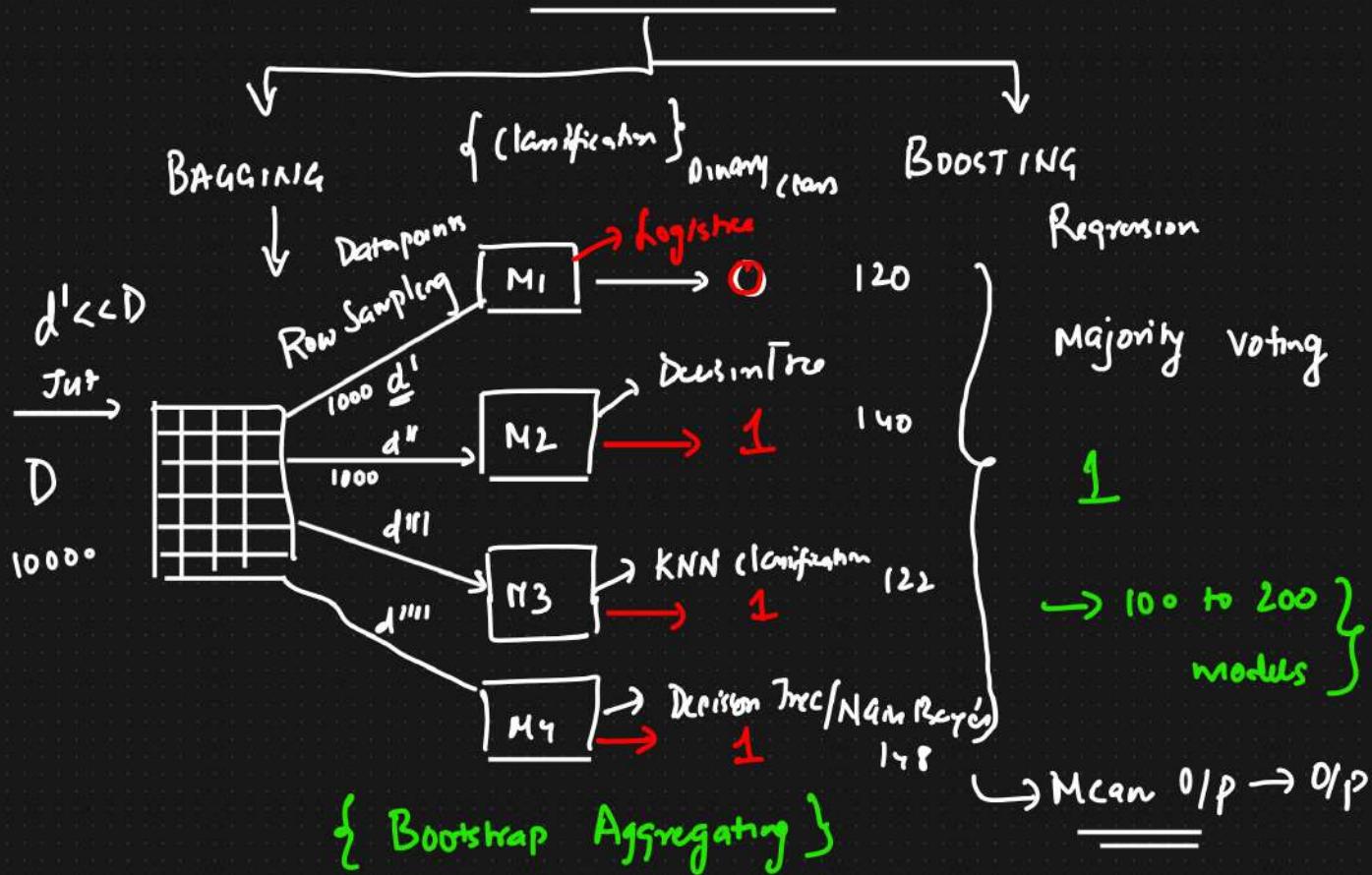
## Ensemble Techniques ✓

### ① Classification & Regression

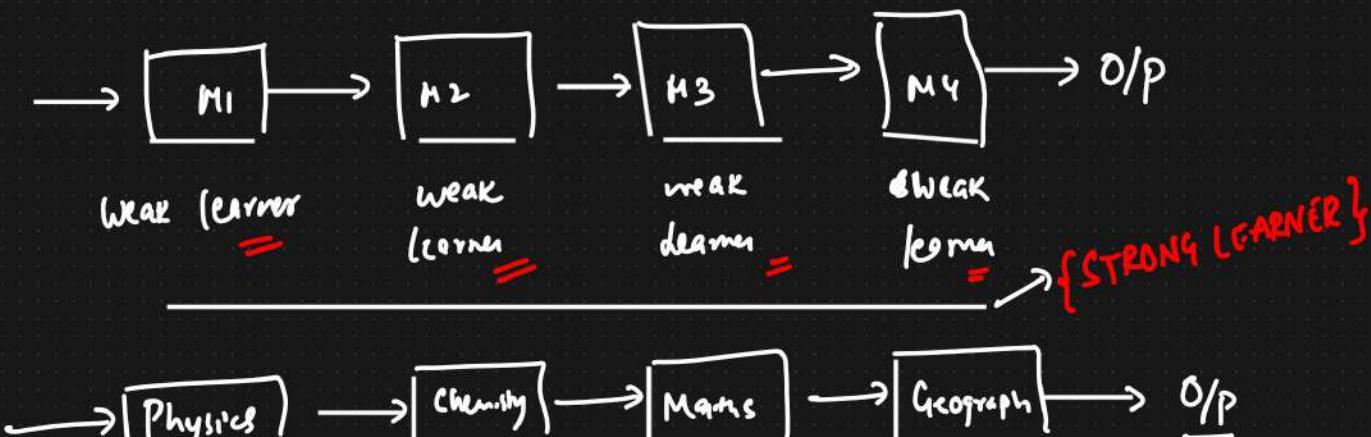
↳ 1 Algorithm  $\xrightarrow{\text{TC}}$  Reg

Multiple Algorithms to solve a problem?

## Ensemble Techniques



## Boosting



## BAGGING

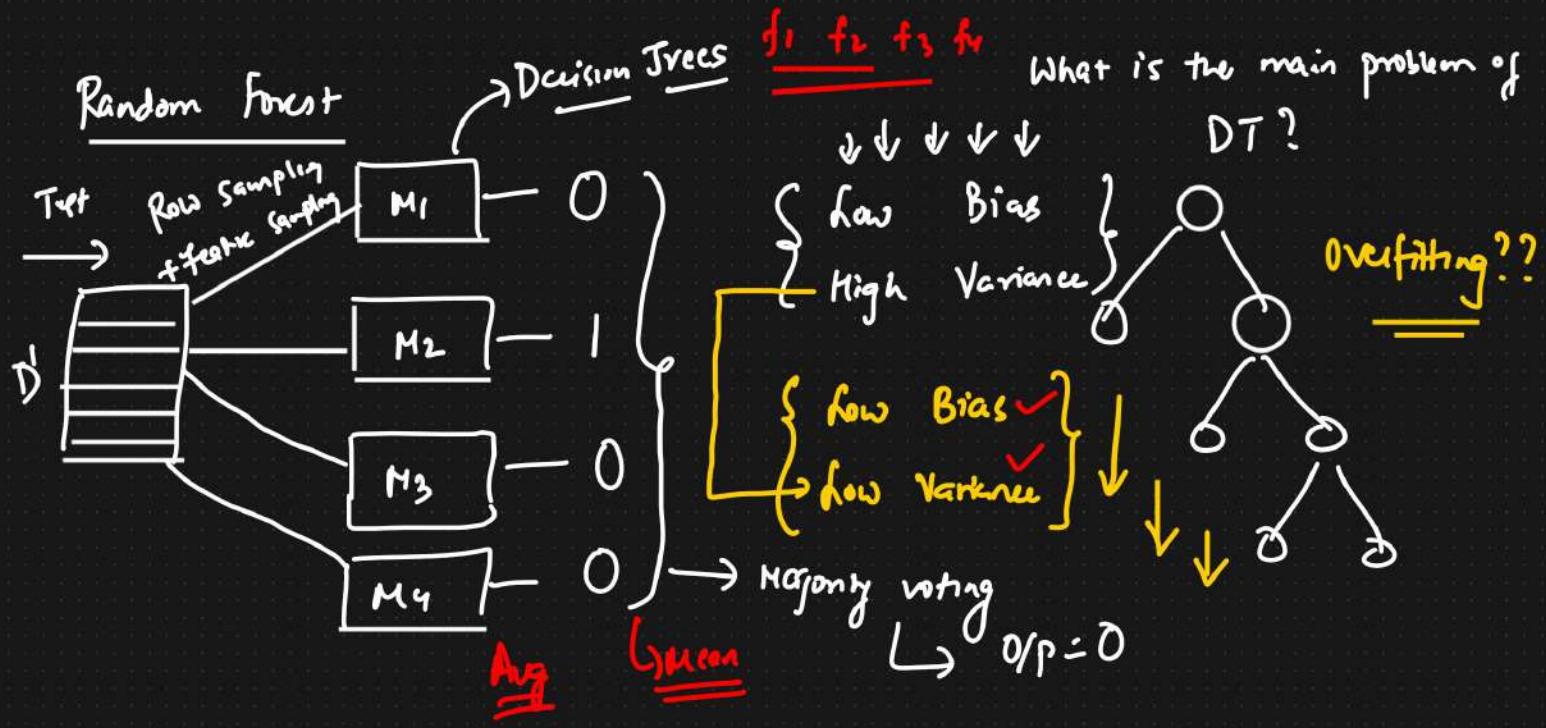
↓

{ ① RANDOM FOREST CLASSIFIER  
② Random Forest Regression

## BOOSTING

- ↓
- { ① Adaboost  
② Gradient  
③ Xgboost }

## ① Random Forest classifier And Regressor



① Normalization ??

or Decision Tree

No.



② KNN {Standardization} ??

Yes

↓  
Impacted

Yes

by  
Outlier??

{Euclidean, Manhattan}

=

③ Random Forest → Outliers  $\Rightarrow \underline{\text{No}}$  → {check it google}

Bagging = Random Forest

Custom Bagging



②

Boosting

i) Adaboost → Decision Tree

Overall = 1

$f_1$	$f_2$	$f_3$	$f_4$	O/p	<u>Weight</u>	<u>STUMPS</u>	Information gain & Entropy
-	-	-	-	Yes	$\checkmark \frac{1}{2}$	$0.05$	
-	-	-	-	No	$\checkmark \frac{1}{7}$	$0.05$	
-	-	-	-	-	$\checkmark \frac{1}{7}$	$0.05$	
X	-	-	-	-	$\frac{1}{7}$	$0.349$	① {weak learner}
-	-	-	-	-	$\checkmark \frac{1}{7}$	$0.05$	② Performance of Stump
-	-	-	-	-	$\checkmark \frac{1}{7}$	$0.05$	$= \frac{1}{2} \log_e \left( \frac{1 - TE}{TE} \right)$
-	-	-	-	-	$\checkmark \frac{1}{7}$	$0.05$	$= \frac{1}{2} \log_e \left( \frac{1 - \frac{1}{7}}{\frac{1}{7}} \right)$
-	-	-	-	-	$\checkmark \frac{1}{7}$	$0.05$	$= 0.895$

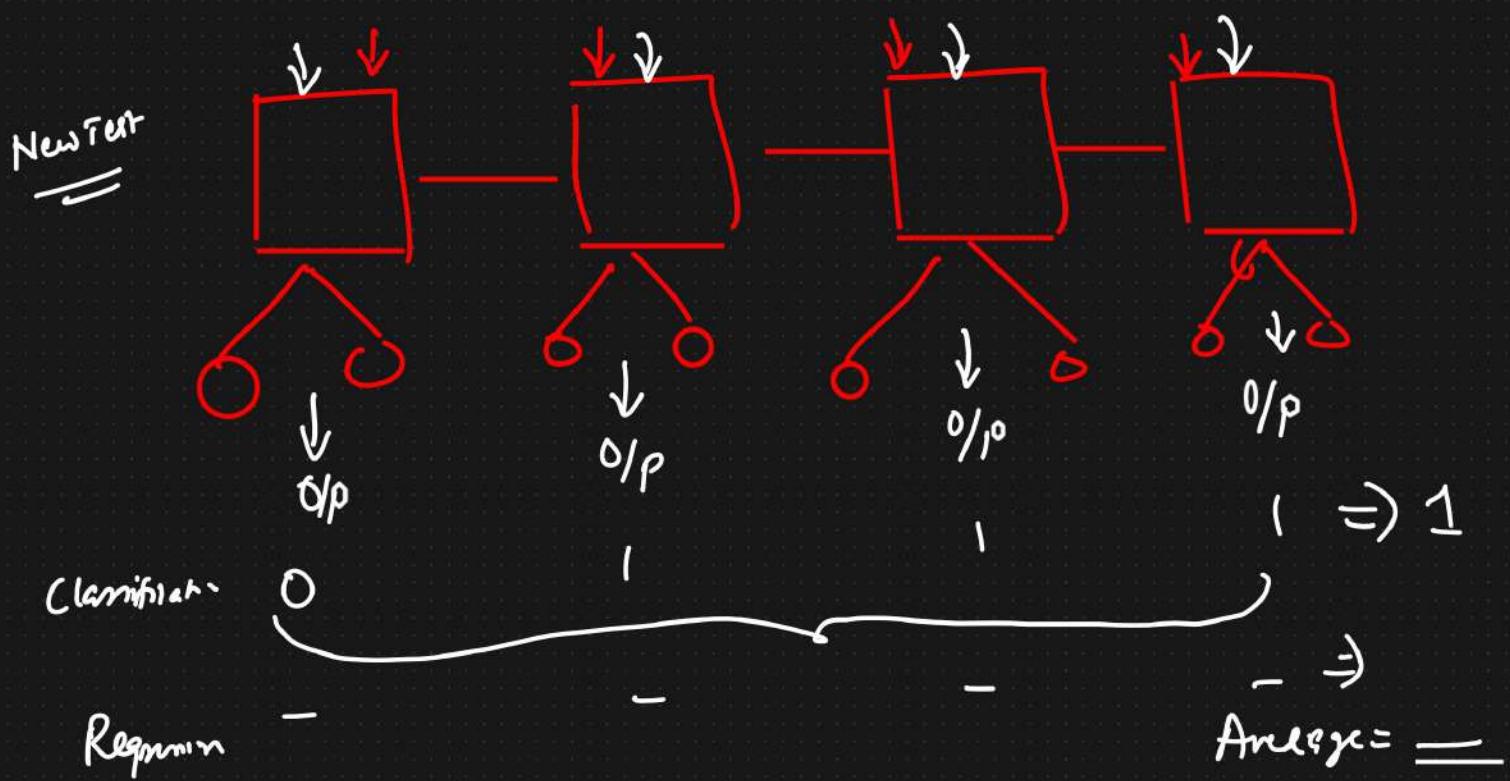
③ New Sample = Correct Records

$$\text{Weight} = \text{Weight} \times e^{-Ps} = \frac{1}{7} \times e^{-0.895} = 0.05$$

$$\text{Incorrect Record} = \text{Weight} \times e^{Ps} = \frac{1}{7} \times e^{0.895} = 0.349$$

<u>New Weight</u>	<u>Normalized weight</u>	<u>Buckets</u>	
$0.05 \div 0.649$	0.07	$[0 - 0.07]$	
$0.05 \div 0.649$	0.07	$[0.07 - 0.14]$	
$0.05 \div 0.649$	0.07	$[0.14 - 0.21]$	
$\rightarrow 0.349$	0.537	{ $[0.21 - 0.247]$ } ✓	
$0.05$	0.07	$[0.247 - 0.751]$	$1.537 [0 - 1]$
$0.05$	0.07	$[—]$	$0.21$
$0.05$	0.07	$[—]$	$0.717$
$\frac{0.05}{0.649}$	$\approx 1$		

Randomly  
Create  
of some number  
between



### Black models VS White box Models

$\rightarrow$  Linear Regression  $\rightarrow$  White box

ANN  $\rightarrow$  Black Box

Random Forest  $\rightarrow$  Black box

Decision Tree  $\rightarrow$  White box

# Day 6 – Machine Learning Algorithms

## Unsupervised ML

- ① K Means clustering
- ② Hierarchical clustering
- ③ Silhouette Score
- ④ DBScan clustering

Agenda

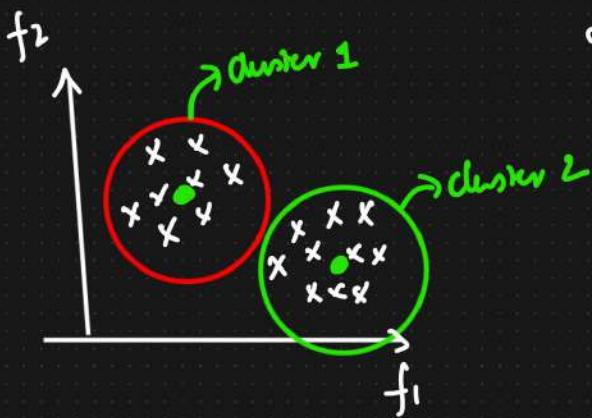
{ ① SVM & SVR  
② XGBoost  
③ PCA }

## Unsupervised ML

Op       $f_1$      $f_2$

Clusters  
↓  
{ Similar kind of }  
data      -    -  
-    -  
-    -

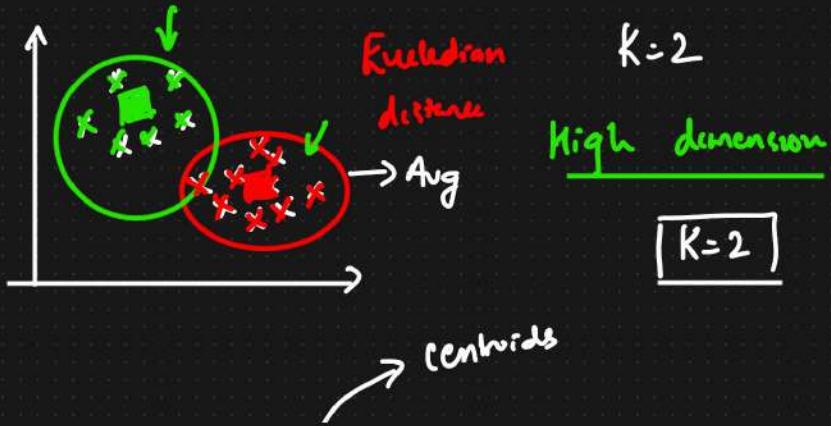
## K Means Clustering



## Custom Ensemble Technique



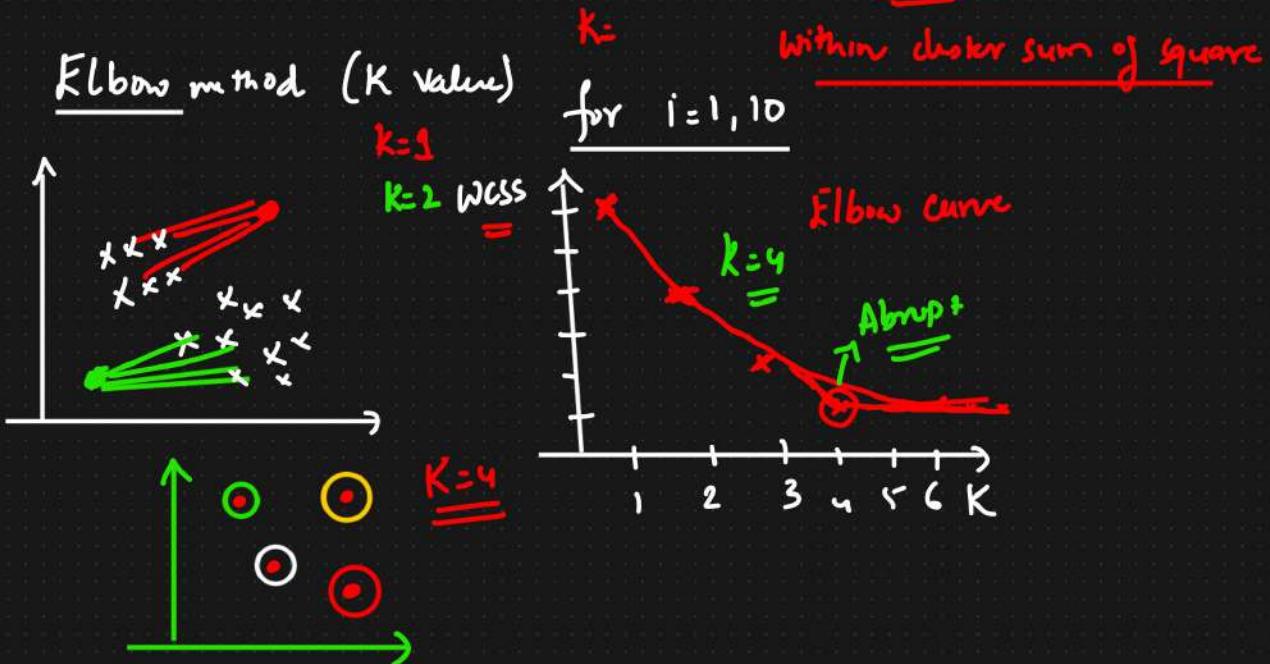
K Means      K = centroids



- ① We try K values  $\Rightarrow$  suitable  $K=2$
  - ② Initialise K number of centroids ✓
  - ③ Compute the avg to update centroids ✓

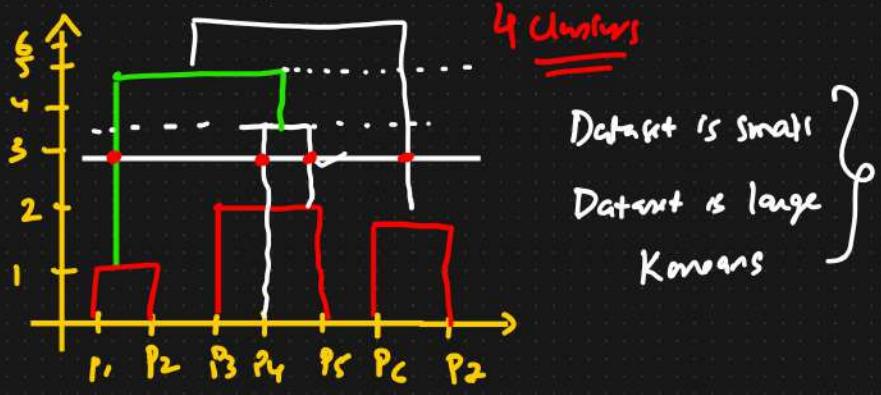
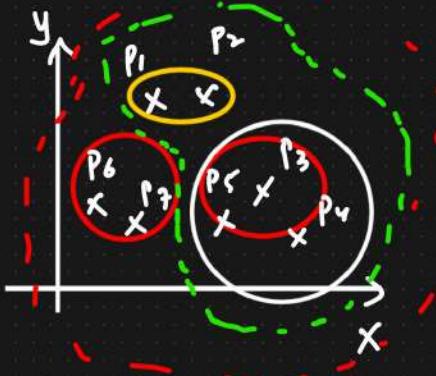
Validating

## Validating



You need find the longest vertical line that has no horizontal line passed through it.  
→ Dendrogram

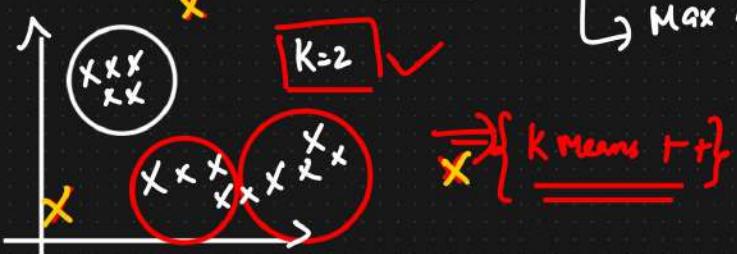
## ② Hierarchical Clustering



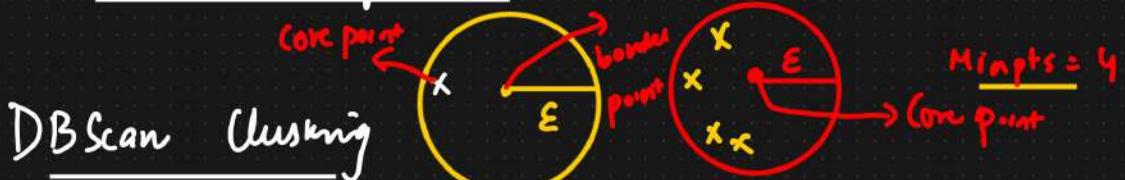
Max Time is taken by KMeans or

Hierarchical clustering ?? ✓

Max Time

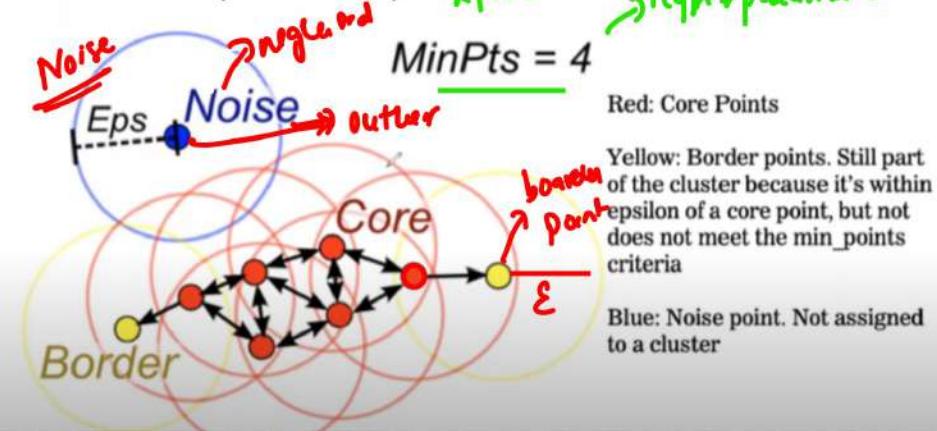


# Validate Clustering Models



Density-Based Spatial Clustering of Applications with Noise(DBSCAN)

Epsilon  $\rightarrow$  Hyperparameter



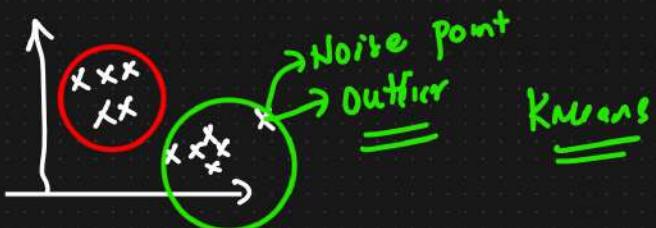
① Epsilon

② Min pts

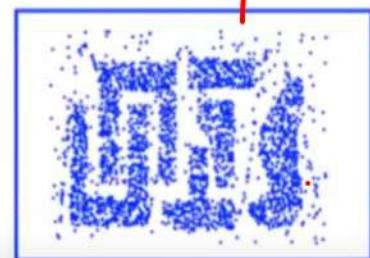
③ Core points

④ Border points

⑤ Noise point



KMeans



DBScan clustering

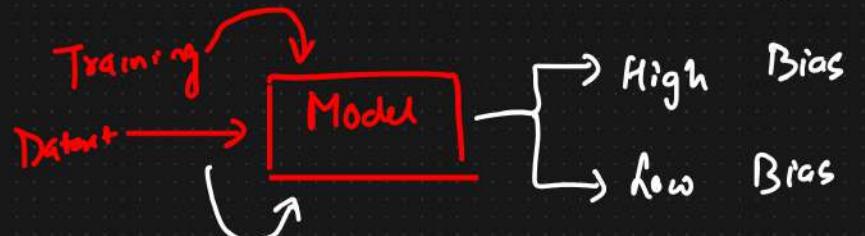
DBScan

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.

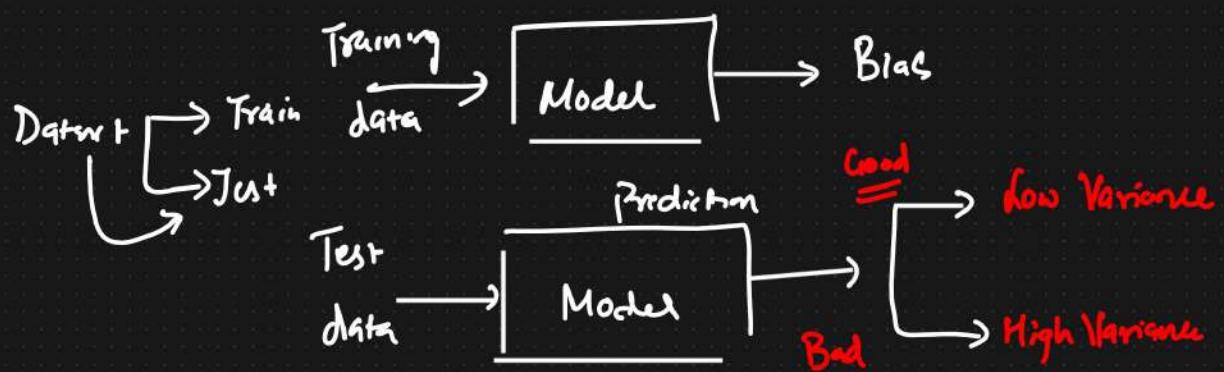
# Defn of Bias And Variance

Training Dataset = 90% }  
 Test Dataset = 70% }  $\Rightarrow$  Overfitting  
 ↓  
 { Low Bias ✓ }  
 { High Variance ✓ }

Bias : It is a phenomenon that skews the result of an algorithm in favor or against an idea.



Variance : Variance refers to the changes in the model when using different portions of the training or test data



## Model 1

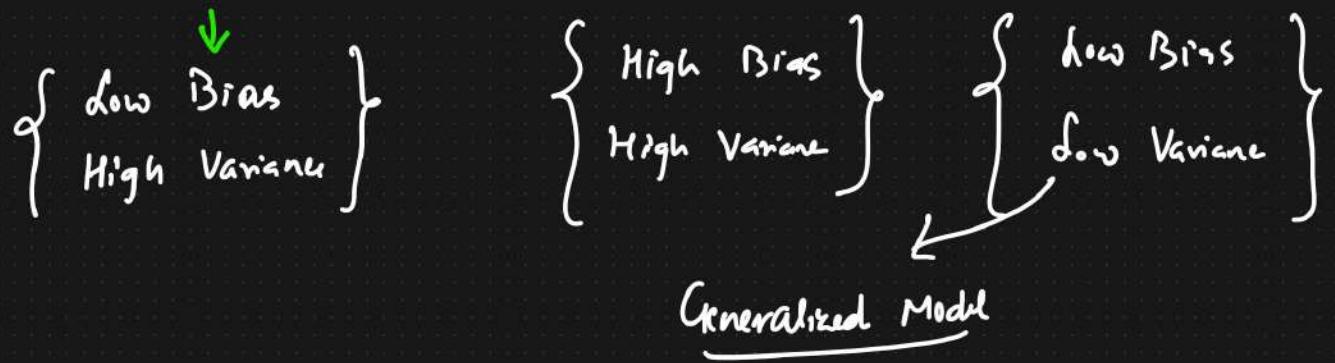
Train Acc = 90%  
 Test Acc = 78%

## Model 2

Train Acc = 60%  
 Test Acc = 85%

## Model 3

Train Acc = 90%  
 Test Acc = 92%



# Day 7 : Xgboost Classifier And Regressor

## Agenda

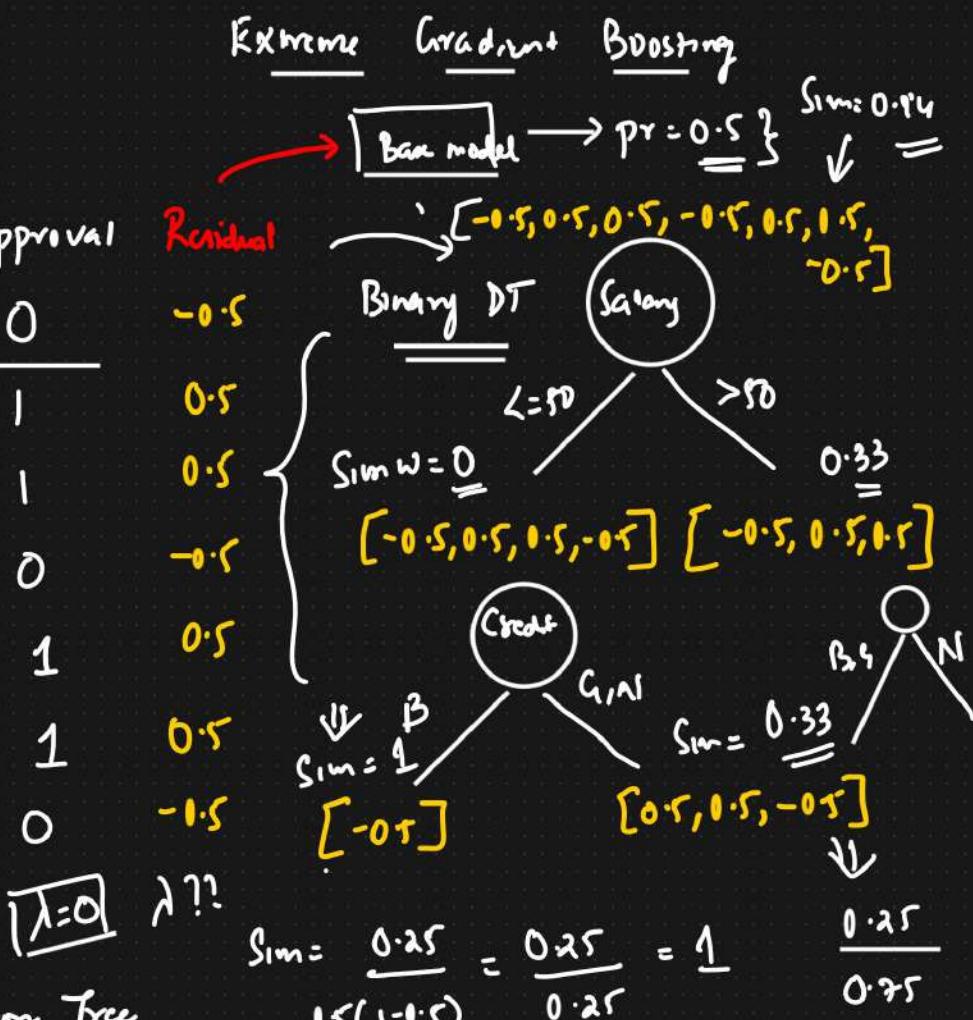
- ① Xgboost classifier
- ② Xgboost Regressor
- ③ SVM
- ④ SVR

$$\log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.5}{0.5}\right) = 0$$

## ① Xgboost Classifier

{ Dataset }

Salary	Credit	Approval	Residual
$\leq 50$	B	0	-0.5
$\leq 50$	G	1	0.5
$\leq 50$	G	1	0.5
$> 50$	B	0	-0.5
$> 50$	G	1	0.5
$> 50K$	N	1	0.5
$\leq 50K$	N	0	-1.5



- ① Create a Binary Decision Tree using the feature

- ② Calculate Similarity weight

$$= \frac{\sum (\text{Residual})^2}{\sum [0 + \alpha(\text{Learned})]} = \frac{\sum [0 + \alpha(1)]}{\sum [0 + \alpha(1)]} = \frac{1}{1} = 1$$

$$1 + 0.33 - 0 = 1.33$$

$\boxed{\lambda = 0.01}$

Sigmoid

$\boxed{0 \text{ to } 1} \checkmark$

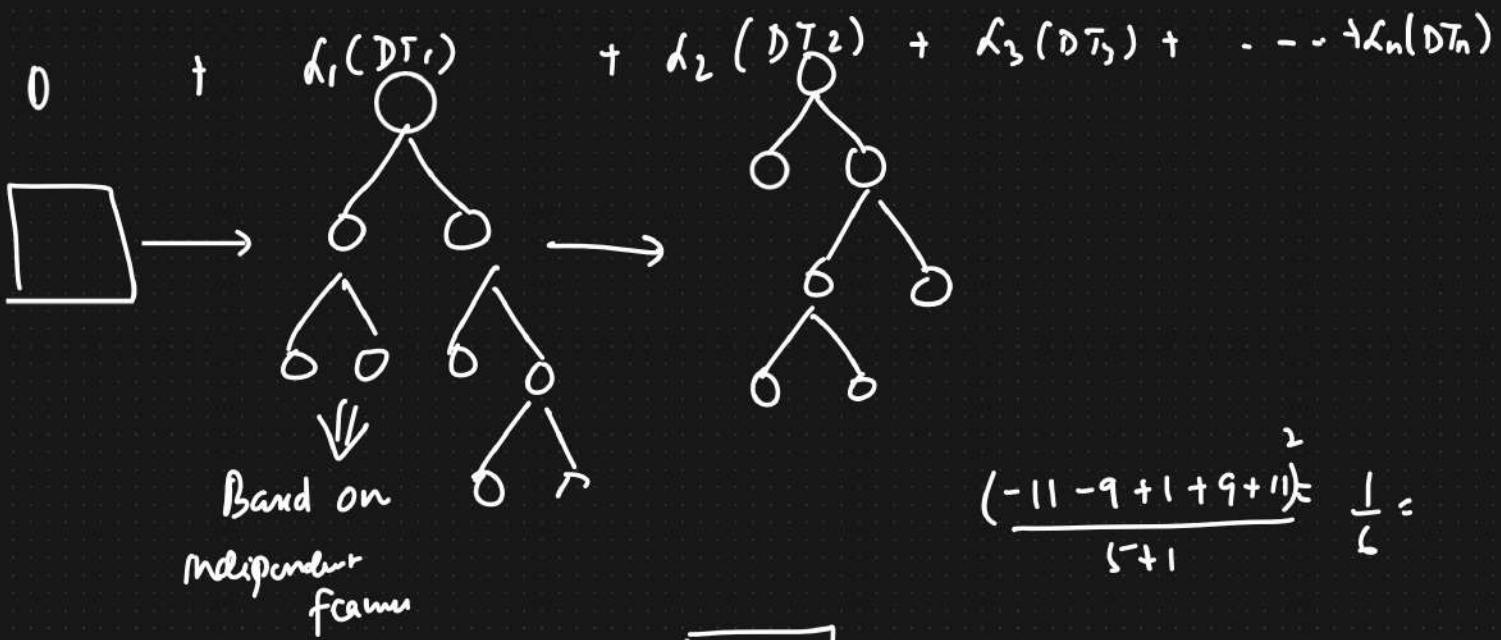
$$\textcircled{3} \text{ Information gain} \rightarrow \sigma \left[ 0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n) \right]$$

New O/P  
Reward

Xgboost → BLACK BOX Model

Pruning

$\lambda \rightarrow$  Cross Validation



$$\frac{(-11 - 9 + 1 + 9 + 1)^2}{5+1} = \frac{1}{6}$$

Bernoulli → 51  
Sum = 16

[-11, -9, 1, 9, 1]

Exp

② Xgboost Regressor

Exp	Gap	Salary	O/P	Res	$\lambda = 1$
2	Yes	40K	-11K		
2.5	Yes	42K	-9K		
3	No	52K	1K		
4	No	60K	9K		
4.5	Yes	62K	11K		

$$\text{Sum} = 60 \cdot 5 = 300$$

$$\frac{121}{1+1} = 60.5$$

$$\frac{(-9+1+9+11)^2}{4+1} = \frac{144}{5} = 28.8$$

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

$$\text{Information} = 60 \cdot 5 + 28.8 - \frac{1}{6} = 89.13$$

只

$\xi \times p$	$y_C$	$\xi_{DT1}$
----------------	-------	-------------

$\leq 2.5$        $> 2.5$

$133.33$        $110.25$

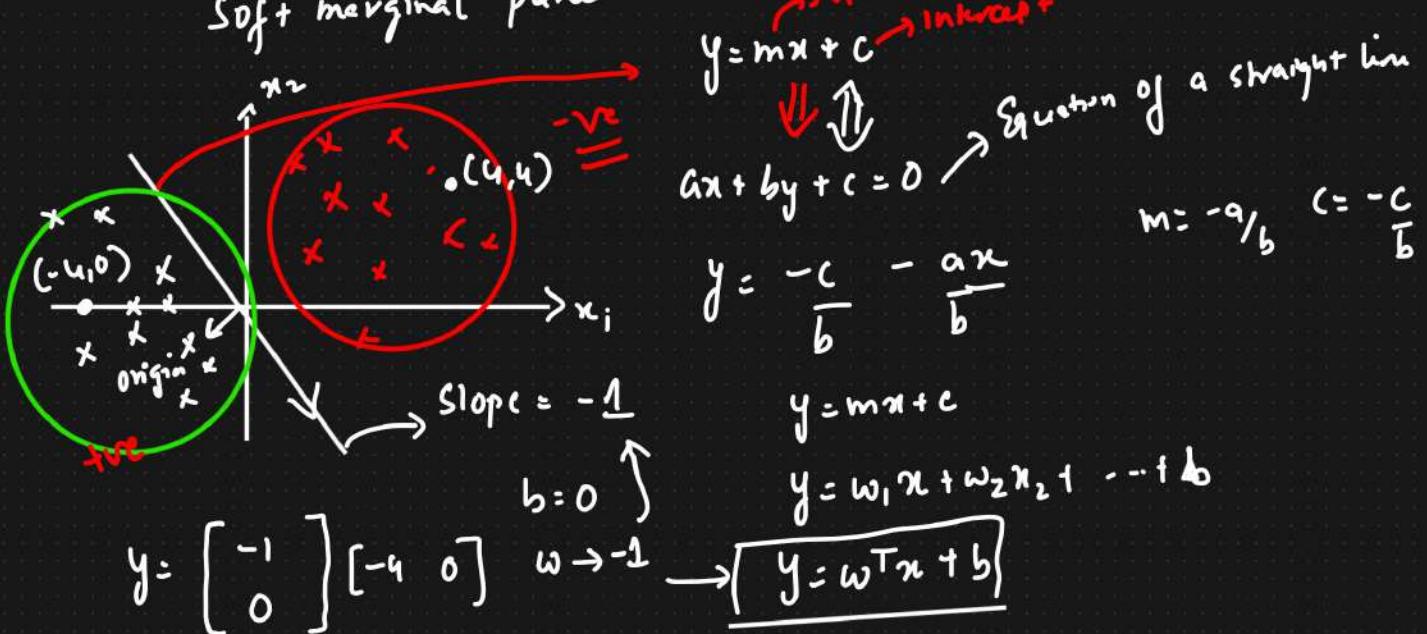
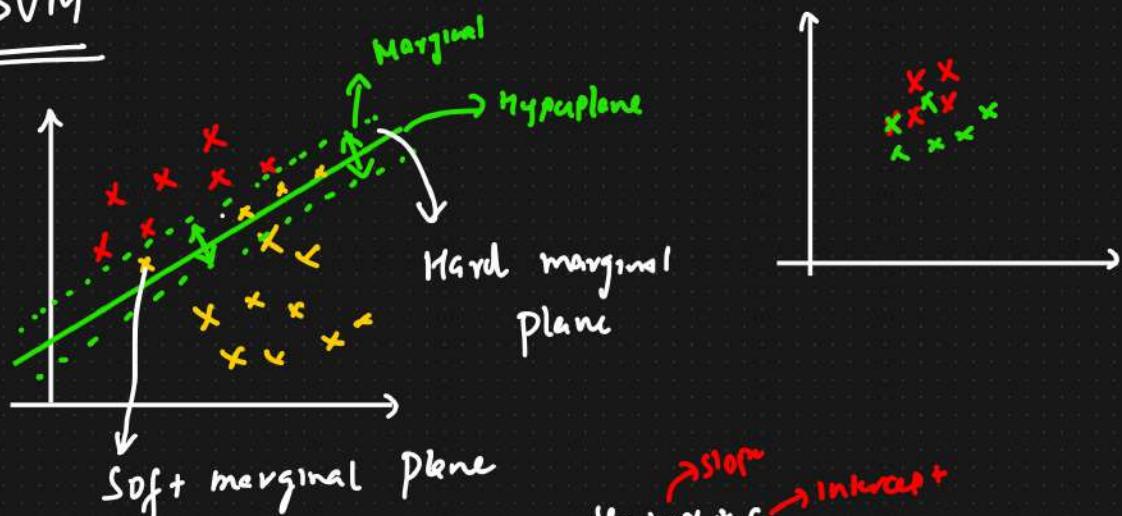
$$\frac{-11-9}{2} = -10 \quad [-11, -9]$$

$$[1, 9, 11] \Rightarrow \frac{1+9+11}{3} = \frac{21}{3} = 7$$

$$O/P = 51 + \alpha_1(-10) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)$$

① {EDA & Feature Engineering} ✓

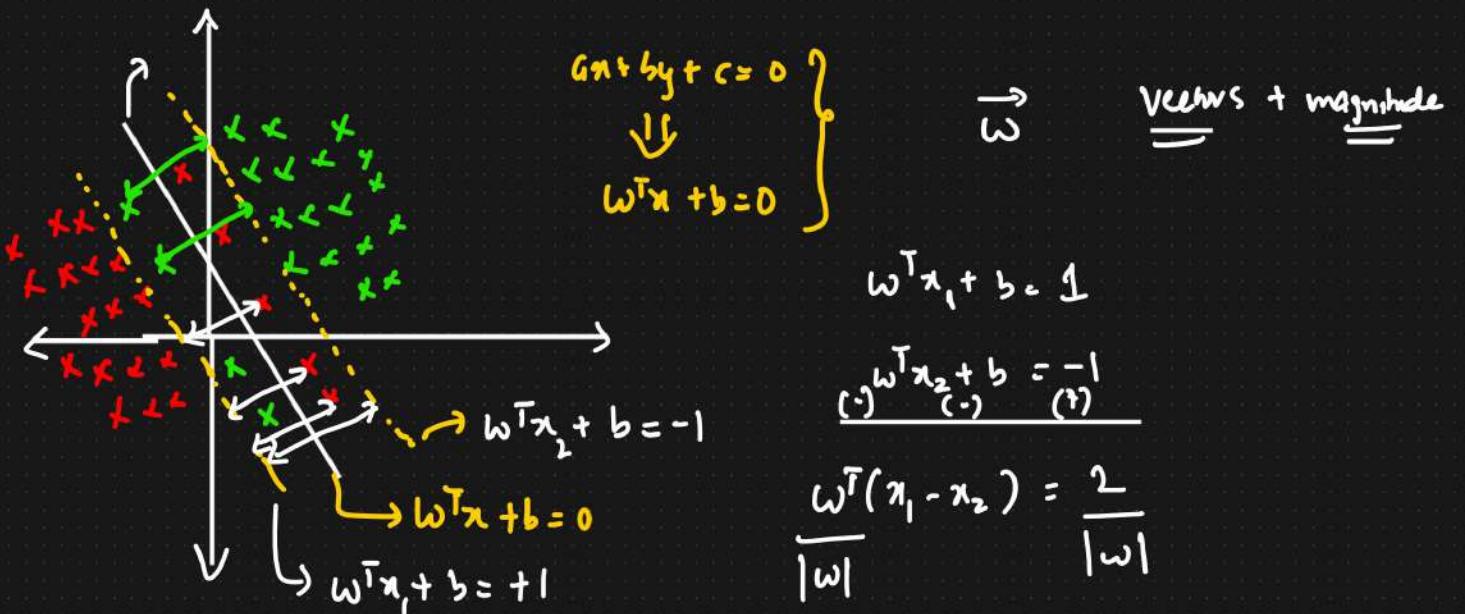
### ③ SVM



$$= 4 \Rightarrow +ve \text{ value}$$

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [4, 4]$$

$$= -4 + 0 = -4$$

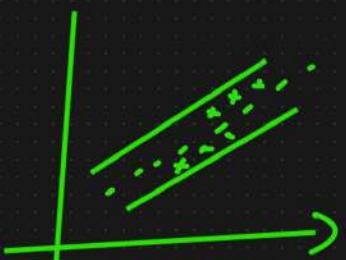


$\left. \begin{array}{l} \text{Maximize}_{(w,b)} \frac{2}{\|w\|} \Rightarrow \text{Marginal} \\ \text{such that } y_i \begin{cases} +1 & w^T x_i + b > 1 \\ -1 & w^T x_i + b \leq -1 \end{cases} \end{array} \right\}$

Major Aim  $y_i * (w^T x_i + b) \geq 1$   
 for correct point

Maximize  $(w,b)$   $\frac{2}{\|w\|} \Leftrightarrow$  Min  $(w,b)$   $\frac{\|w\|}{2}$

Min  $(w,b)$   $\frac{\|w\|}{2} + C_i \sum_{i=1}^n \xi_i$  → Summation of the distances of the wrong datapoints

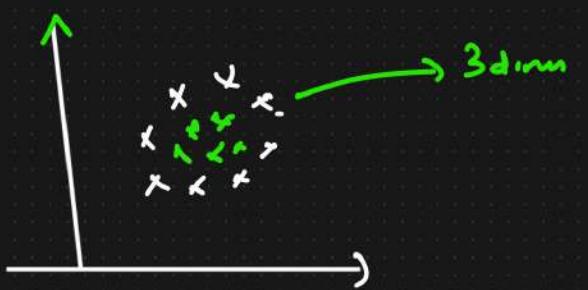


96%

$\left\{ \begin{array}{l} \text{How many} \\ \text{Errors we can} \\ \text{have} \end{array} \right\}$

SVR  
 $\underline{\underline{F}} \underline{\underline{X}} \underline{\underline{P}} \underline{\underline{o}} \underline{\underline{n}}$

## SVM Kernel



Probability

$$\frac{1}{1 + e^{-f(x)}}$$