

7 days Machine Learning Algorithms

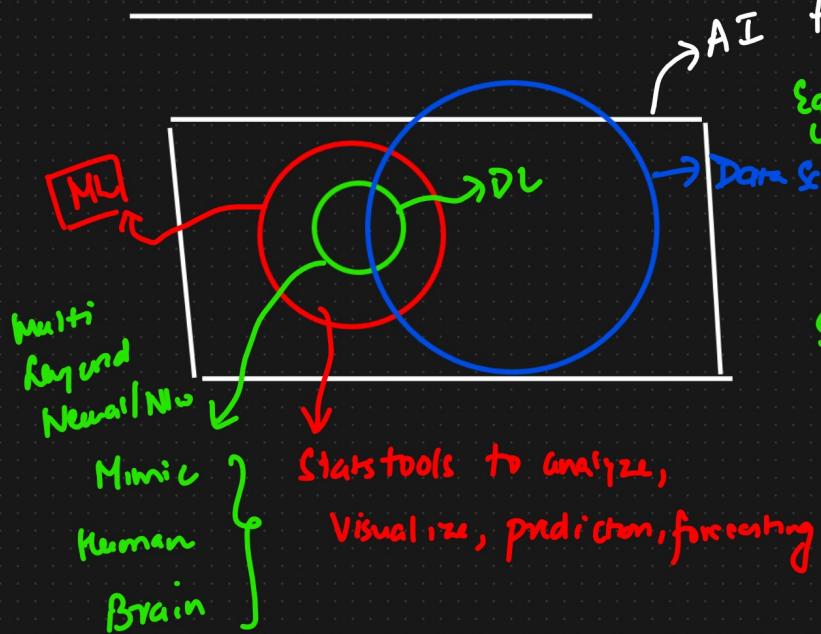
Purpose : Clear the Interviews

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 it own task without any human intervention

Eg: Netflix → Action → Recommendation
Domestic → Comedy → "

Amazon.in → iPhone → Headphones }

Sufi Driving Cars →

Machine & Deep learning

Reinforcement



Supervised ML

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



Independent features \rightarrow Age

Dependent feature \rightarrow Weight



① Regression Problem

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

② Classification

Binary classification
Multi-class classification
O/p

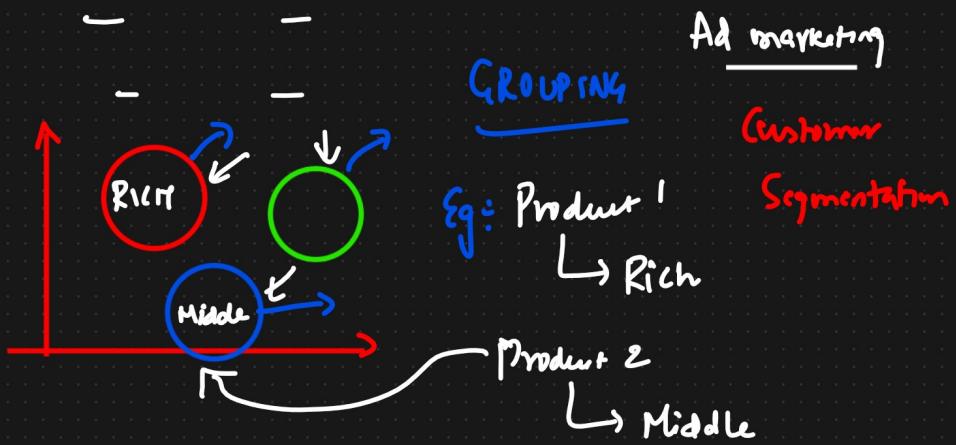
No. of hours	No. of play hours	No. of sleep	P/F
-	-	-	P
-	-	-	F
-	-	-	P
-	-	-	F

②

Unsupervised ML \rightarrow Clustering
Dimensionality Reduction

Salary	Age	$\rightarrow \{$ No Dependent variable $\}$
-	-	
-	-	

Clustering \rightarrow Customer Segmentation



② Dimensionality Reduction

1000 → lower dimension
 ↓
 100

PCA, LDA

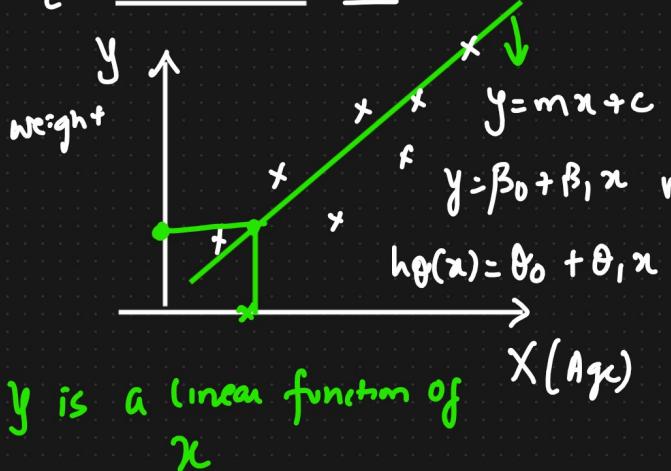
Supervised

- 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

Model

Hypothesis

O/P weight

Credits : Andrew Ng

Equation of a straight line

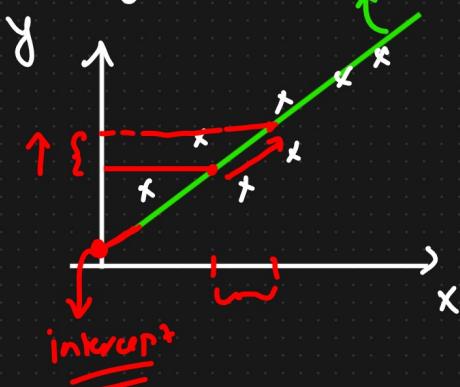
$$h_\theta(x) = \theta_0 + \theta_1 * x$$

When $x=0$

θ_0 = Intercept

θ_1 = Slope or Coefficient

x_i = data points



Linear Regression

Minimize

Bst fit line

Start at point \rightarrow best fit line



Hypothesis

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

Purpose
Derivation

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

Cost function

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

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

→ Cost function

↳ Squared Error Function

What we need to solve

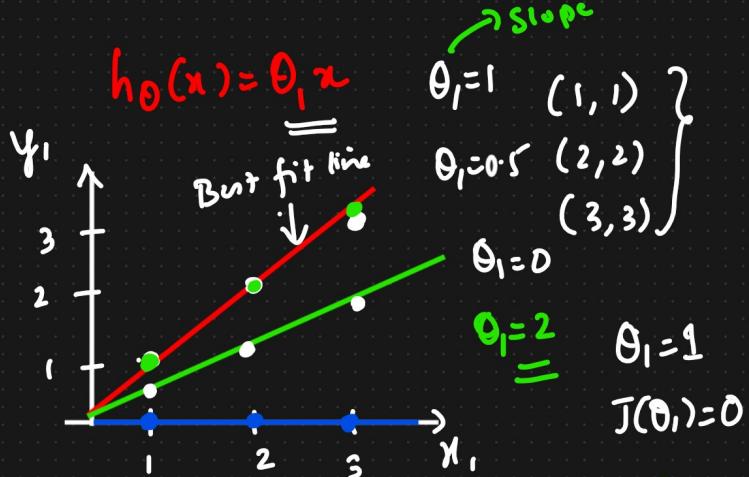
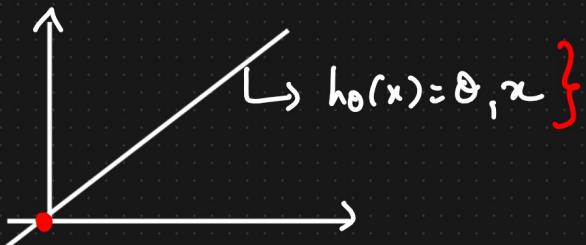
$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

\Downarrow

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

$$\theta_0, \theta_1$$

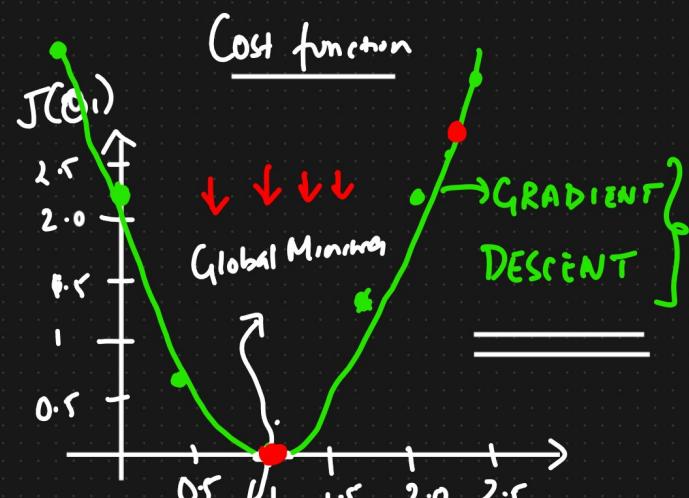
* $h_\theta(x) = \theta_0 + \theta_1 x \quad \text{If } \theta_0 = 0$



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^3 (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

$$J(\theta_1) = 0$$



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^3 (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

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

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

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

$$\approx 2.3$$

$\alpha \Rightarrow$ large

$$\alpha = 0.01$$

Convergence Algorithm

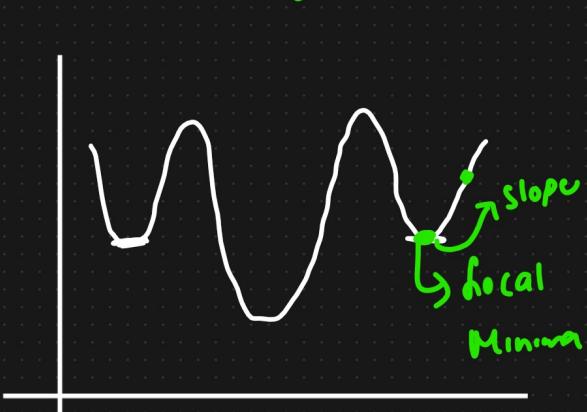
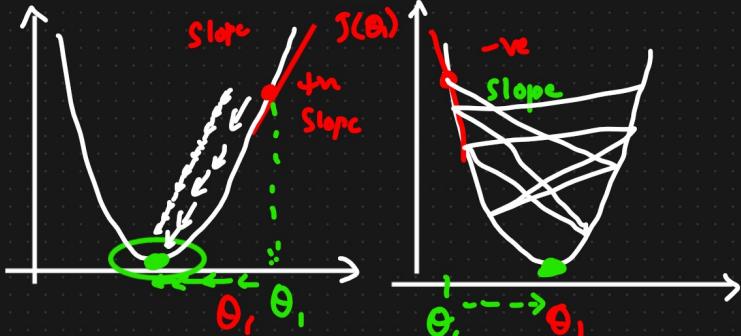
Repeat until convergence derivative (slope)

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \boxed{\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}} \end{array} \right.$$

{ decreasing Rate }

$$\left\{ \begin{array}{l} \theta_1 := \theta_1 - \alpha (+ve) \\ \theta_1 := \theta_1 - \alpha (-ve) \end{array} \right.$$

$$\left. \begin{array}{l} \theta_1 := \theta_1 - \alpha (-ve) \\ \theta_1 := \theta_1 + \alpha (+ve) \end{array} \right\}$$



GRADIENT DESCENT Algorithm

Repeat until convergence

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \boxed{\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}} \end{array} \right.$$

{ $j=0$ and 1 }

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

$$\frac{1}{2m} x^2 \sqrt{\frac{2}{2^m}} x$$

Convergence Algorithm:

$$\left\{ \begin{array}{l} j=0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ j+1 \Rightarrow \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right. \quad \begin{array}{l} h_\theta(x) = \theta_0 + \theta_1 x \\ \frac{x^2}{2} \\ \frac{\partial}{\partial \theta_0} (x, \theta_0) = x \end{array}$$

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

Repeat until converge

{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

}



Performance Metrics

R^2 and Adjusted R^2

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

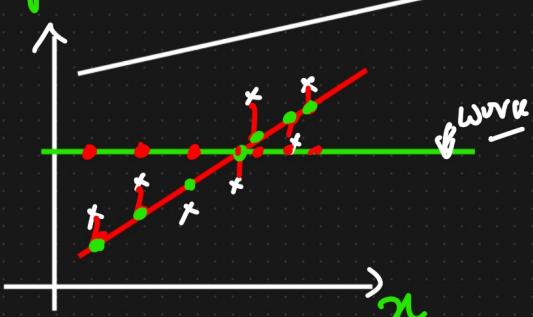
$$h_\theta(x)$$

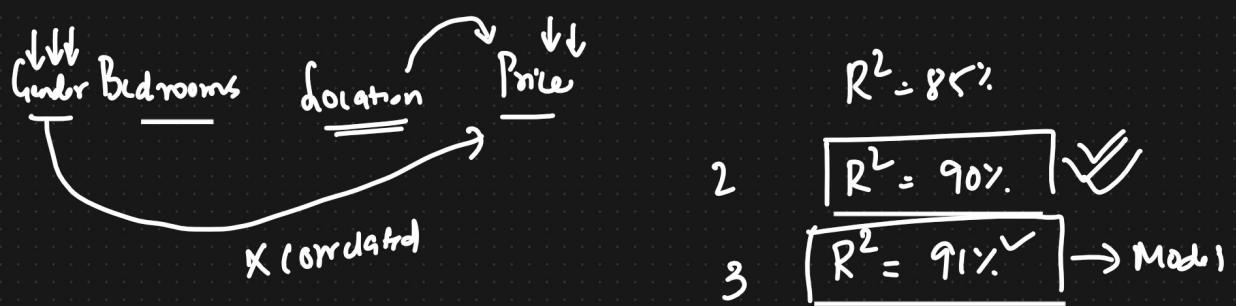
$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Higher

Small number
Big number
90%

$$1 - \frac{\text{low}}{\text{High}}$$





Adjusted R^2

$p = \text{features or predictors}$

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

$\checkmark \uparrow \text{Big} \checkmark$

$$\left\{ \begin{array}{l} p=2 \quad R^2=90\% \quad R^2_{\text{adjusted}}=86\% \\ p=3 \quad R^2=91\% \quad R^2_{\text{adjusted}}=82\% \end{array} \right.$$

$\downarrow \downarrow$

$p=2 \quad > \quad N-p-1 \quad >> \quad p=3$

$R^2 \uparrow \uparrow \uparrow$

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

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

$p >>>$

Day 2 - Linear Machine Learning Algorithm

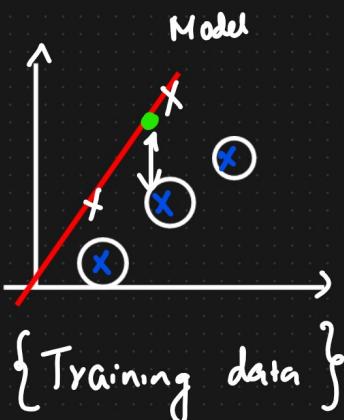
Agenda

- 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 }

Model performs well → Training data

Fails to perform well → Test Data ✓

(High Variance)

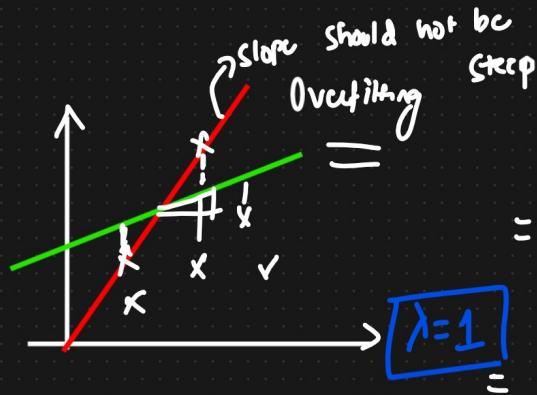
Model 1

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 }

$$J(\theta_1) = 0$$

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

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

Model 3

Training Acc = 70%.

Test Acc = 65%.



Underfitting

High Bias, High Variance

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

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

$$h_\theta(x) = \theta_1 x \quad \text{L} \rightarrow \underline{\text{slope}}$$

Ridge (L2 Regularization)

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

iterations { Hyperparameter }

$$= 4/4 \downarrow \downarrow \text{bb}$$

R², adjusted R²

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



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

Convergence

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



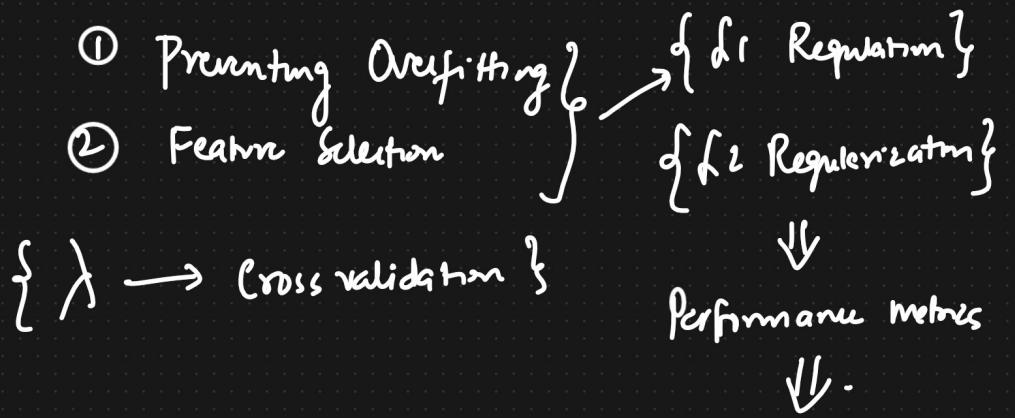
$$\approx 3 \downarrow \downarrow$$

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 (λ_2 Norm)

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

④

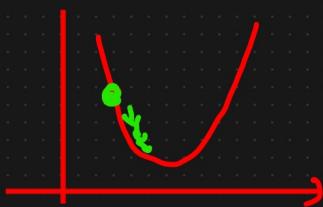
Purpose : Preventing Overfitting

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

Lasso Regression (λ_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

$$X_3 \quad \boxed{X_1 \quad \cancel{X_2}} \quad \boxed{Y}$$

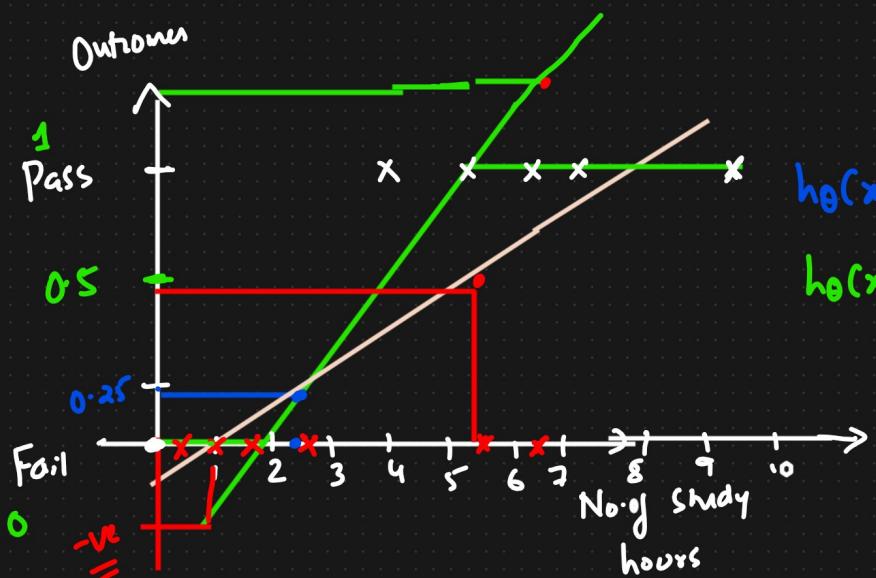
Variation Inflation factor?

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

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

$$\boxed{h_{\theta}(x) = \theta^T x}$$

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

Squash



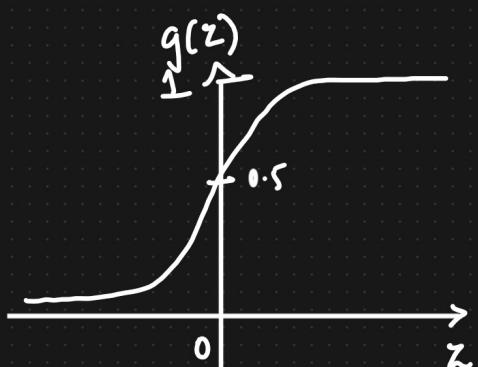
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1)$$

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

$$h_{\theta}(x) = g(z)$$

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

Sigmoid or Logistic function



$$g(z) \geq 0.5 \quad \left\{ \begin{array}{l} \checkmark \\ \cdot \end{array} \right.$$

When $z \geq 0$

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

Training Set

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

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

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

(change parameter θ_1 ?)

Cost function

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

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

Logistic Regression

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

Logistic Regr
Cost function $= \frac{1}{2} (h_{\theta}(x^{(n)}) - y^{(n)})^2 \quad \left. \begin{array}{l} \text{We cannot use this} \\ \text{cost function for logistic} \end{array} \right\}$

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

Gradient Descent

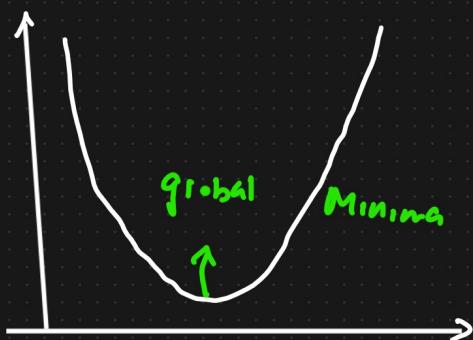
Non convex function



Local Minima
Problem

Gradient Descent

Convex function



global
Minima

Logistic Regression Cost function

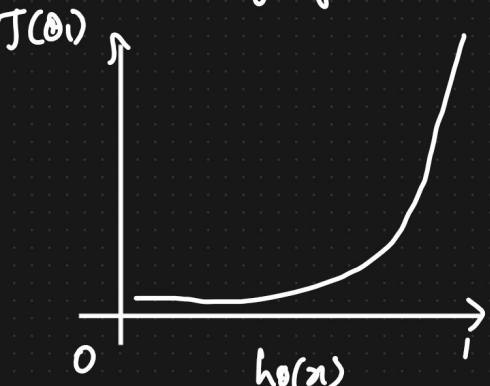
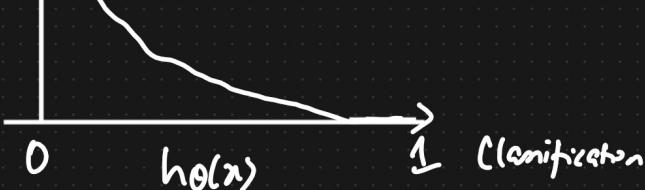
$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}}$$

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

if $y=0$

$$J(\theta_1) \text{ if } y=1$$

Cost = 0 if $y=1, h_{\theta}(x)=1$



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

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

if $y=1$ \Downarrow cost function.

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

if $y=0$

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

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

↓
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}

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

Predicted	1	0	Actual
1	TP	FP ↓	Confusion matrix
0	FN ↓	TN	

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

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

$$\begin{aligned} & 0 \rightarrow 600 \\ & 1 \rightarrow 400 \end{aligned} \quad \begin{aligned} \text{Balanced} & \\ \text{Data} & \end{aligned} \quad \begin{aligned} 0 : 900 \\ 1 : 100 \end{aligned} \quad \begin{aligned} \text{Balanced} & \\ = & \end{aligned}$$

$$\left\{ \text{Model} \rightarrow 0 = \frac{900}{1000} = 90\% \right\}$$

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 CANCER OR NOT } → Recall

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

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

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

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

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

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

F2 Score

3rd Day → Machine Learning Algorithms

Agenda

- ① Practicals
- ② Naive Bayes Intuition
- ③ KNN algorithms

→ Simple Examples

Previous Session

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

⇒ Complex

① Naive Bayes Intuition {Classification}



{Baye's Theorem}

Rolling a Dice

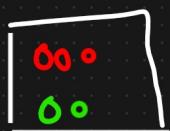
{1, 2, 3, 4, 5, 6}

{Independent Events}

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

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

Dependent Event



$$P(R) = \frac{3}{5} \rightarrow R \quad \text{Dependents} \quad P(G) = \frac{2}{5}$$

↙ Green Marble

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

$$P(R) = \frac{3}{4}$$

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
(RUx)

Naive Bayes

I/P

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad \dots \quad x_n$$

$$y \rightarrow \alpha_p$$

→ - - - - - - - - - ✓
→
→

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

$$= \frac{P(y) * P(x_1/y) * P(x_2/y) * P(x_3/y_3) * \dots * P(x_n/y_n)}{P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)}$$

DEFINITION

$$x_1 = x_2 = x_3 = x_4 = y$$

$$\rightarrow \boxed{x_1, x_2, x_3, x_4} \quad \text{Yes } \checkmark$$

No ✓

$$P(y=\text{Yes}/x_i) = \frac{P(y_{\text{Yes}}) * P(x_1/y_{\text{Yes}}) * P(x_2/y_{\text{Yes}}) * P(x_3/y_{\text{Yes}}) * P(x_4/y_{\text{Yes}})}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

Constant → $P(x_1) * P(x_2) * P(x_3) * P(x_4)$ # fixed
Ignore

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

Constant → $P(x_1) * P(x_2) * P(x_3) * P(x_4)$ # fixed

x_i ; $\begin{cases} \rightarrow \text{Yes} \\ \rightarrow \text{No} \end{cases}$

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

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

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

$$< 0.5 \Rightarrow 0$$

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

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

Dataset

Day	Outlook	Temperature	Humidity	Wind	Binary Item	
					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	2/9	3/5
Overcast	4	0	4/9	0/5
Rain	3	2	3/9	2/5

Total 9 5

PLAY

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

$$\begin{array}{c} \text{Yes} \quad 9 \\ \text{- No} \quad 5 \\ \hline \text{Total} \quad 14 \end{array} \quad \begin{array}{c} P(Y_{\text{Yes}}) \\ \boxed{\frac{9}{14}} \end{array} \quad \begin{array}{c} P(N) \\ \boxed{\frac{5}{14}} \end{array}$$

→ Test (Sunny, Hot) → O/P

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

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

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

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

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

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

$$= \cancel{\frac{1}{7} * \frac{3}{5}} * \frac{2}{5}$$

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

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

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

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

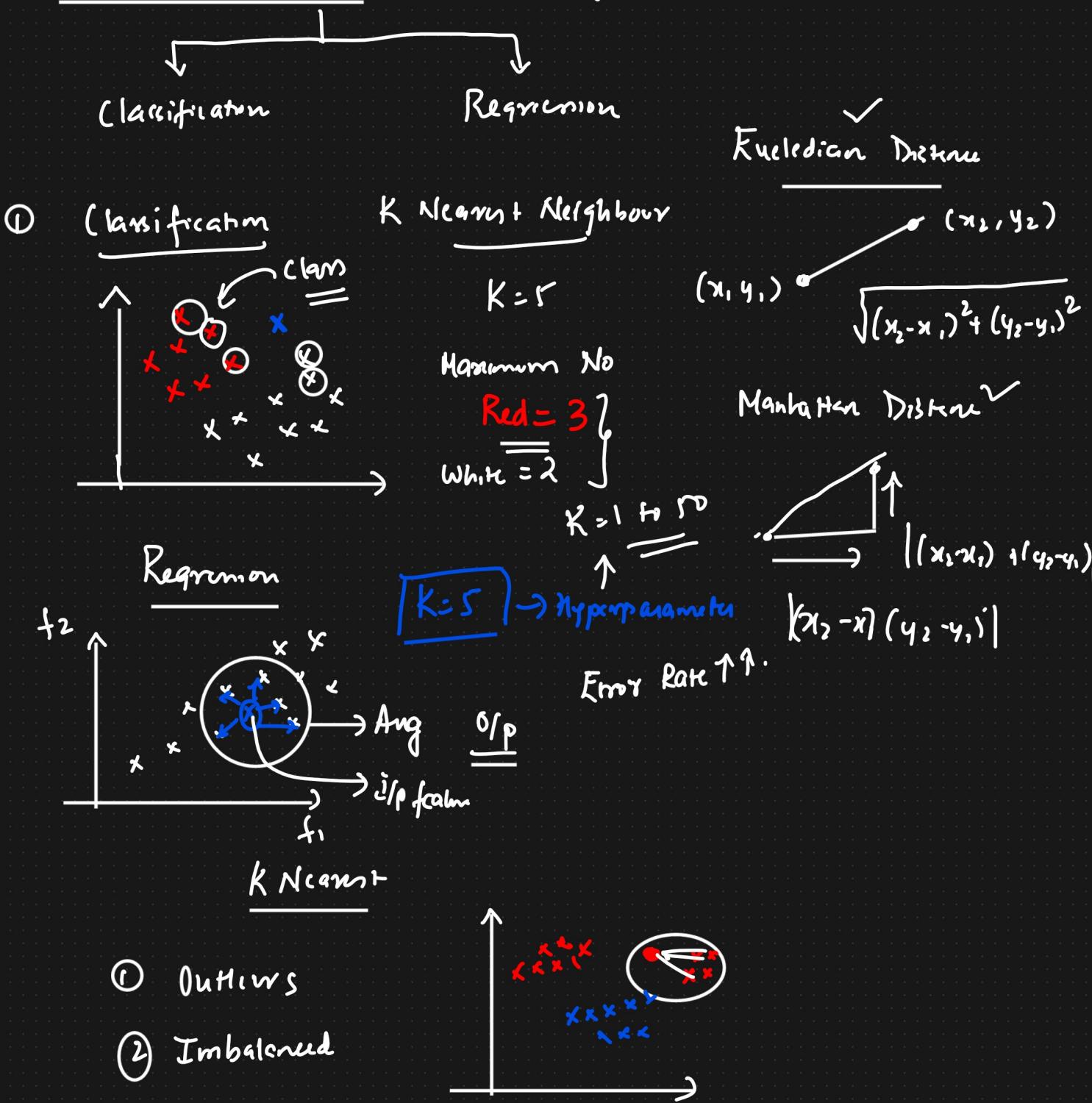
Always $\rightarrow \text{No} \checkmark$

Assignment

$(\text{Overcast, Mild}) \rightarrow \text{Naive Bayes?}$

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 \leq 18):

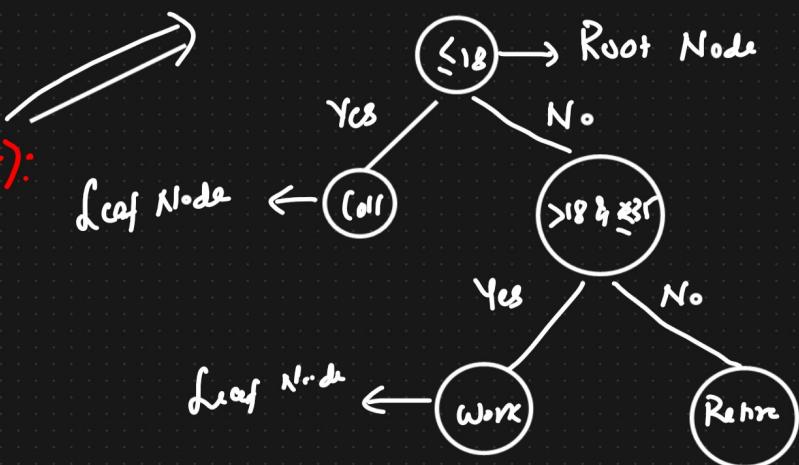
Print ("College")

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

Print ("Work")

else :

Print ("Retire")

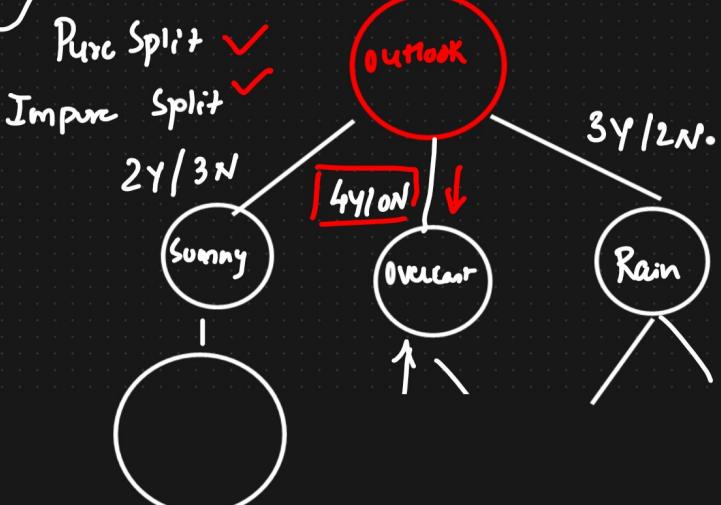


DECISION TREE

Nest if else \Rightarrow Decision Tree

CLASSIFICATION 9Y/5N

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

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

+ } Binary

\Rightarrow 50% for.

$6Y|3N$

f_1

P_+ = Probability of Yes

$3Y|1N$

c_1

c_2

\rightarrow pure split

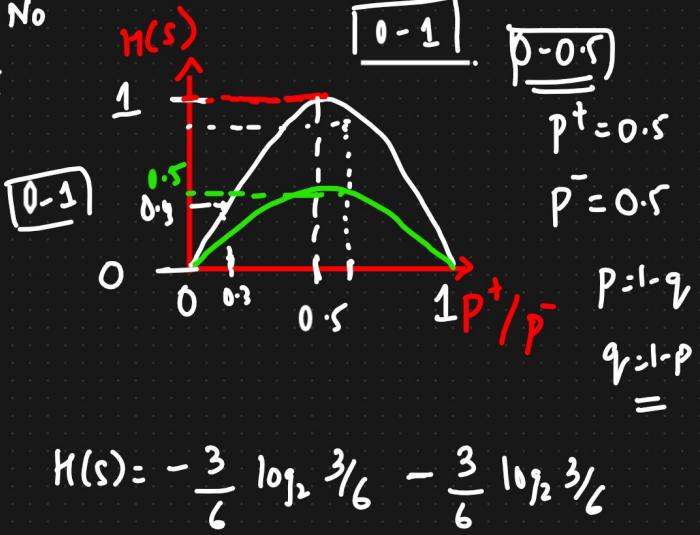
Entropy $H(S) = -\frac{3}{6} \log_2 \frac{3}{3} - \frac{3}{6} \log_2 \frac{3}{3}$

 $= -1 \log_2 1$
 $= 0 \rightarrow \text{Pure Split}$

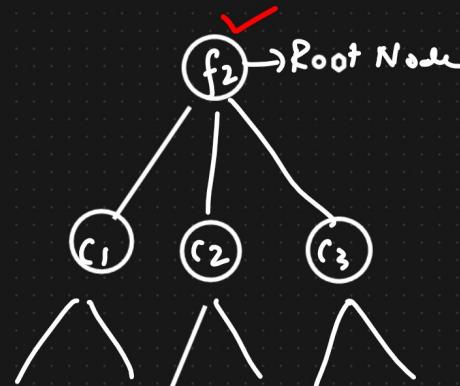
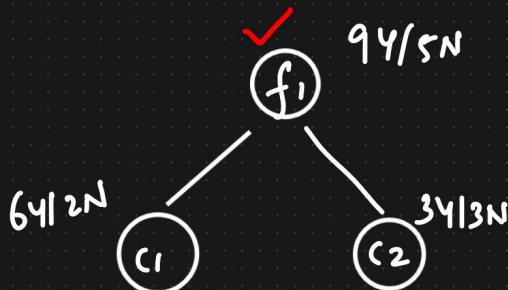
Purity Test \rightarrow Entropy

① Gini Impurity

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



- ② Which feature to take to split??



Information Gain

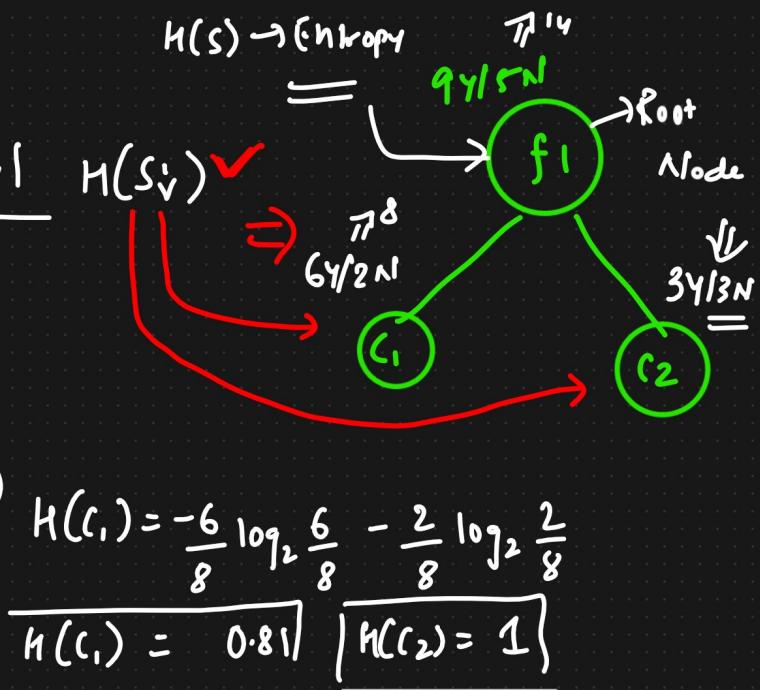
$$\text{Gain}(S, f_1) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

↗ Root Node

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

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

$$\approx = \boxed{0.94}$$



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

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

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 }

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

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

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

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



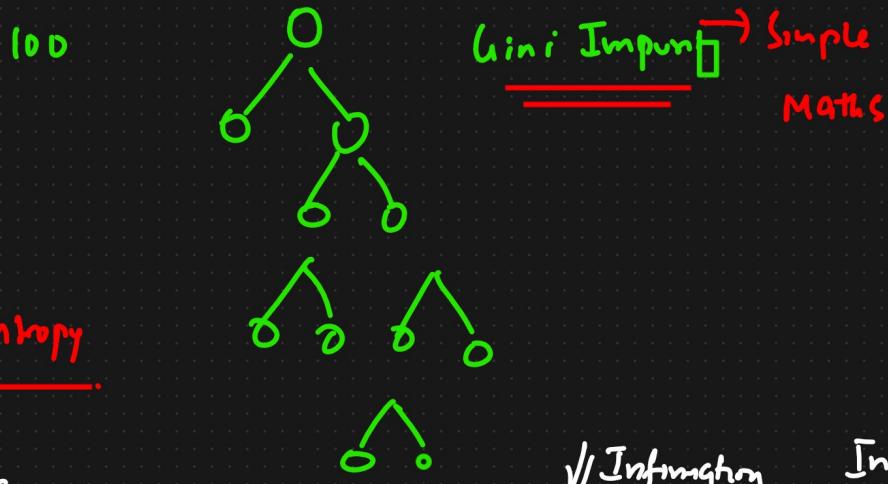
$$\text{Entropy} = 1$$

$$\text{Gini Impurity} = 0.5$$

Entropy \rightarrow fog

{Entropy}

Fast
Gini > Entropy



continuous

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

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

$$\frac{1.3}{1.3}$$

$$\frac{2.3}{3}$$

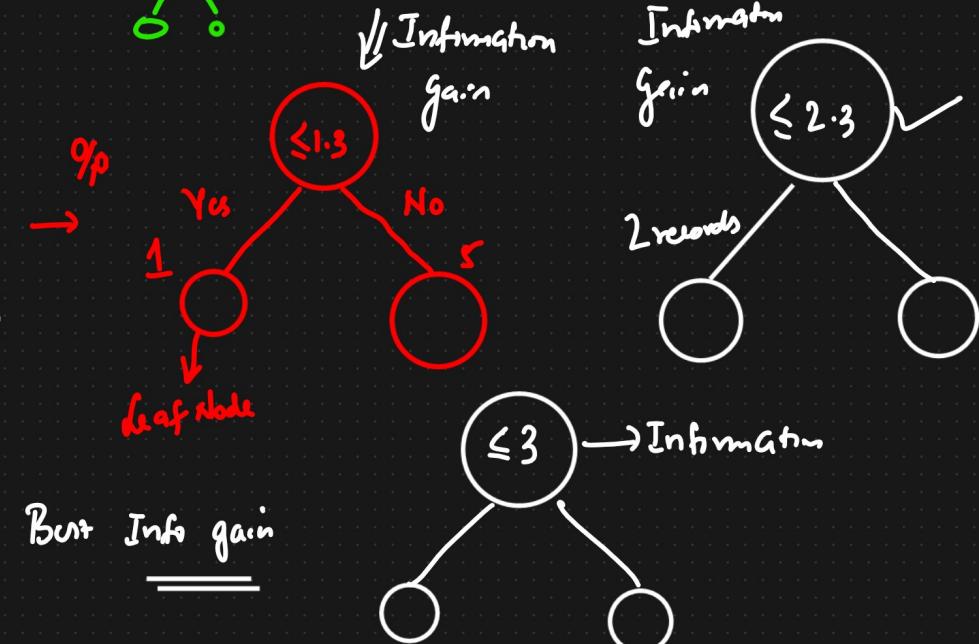
$$\frac{3}{4}$$

$$\frac{4}{5}$$

$$\frac{5}{7}$$

$$\frac{7}{3}$$

But Info gain



Decision Tree Regressor

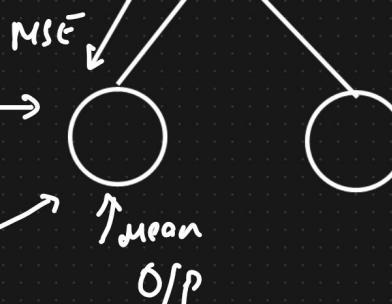
$f_1 \ f_2 \text{ O/P}$

(continuous)



f_1 Mean $\boxed{\text{MSE Or MAE}}$

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



Overfitting

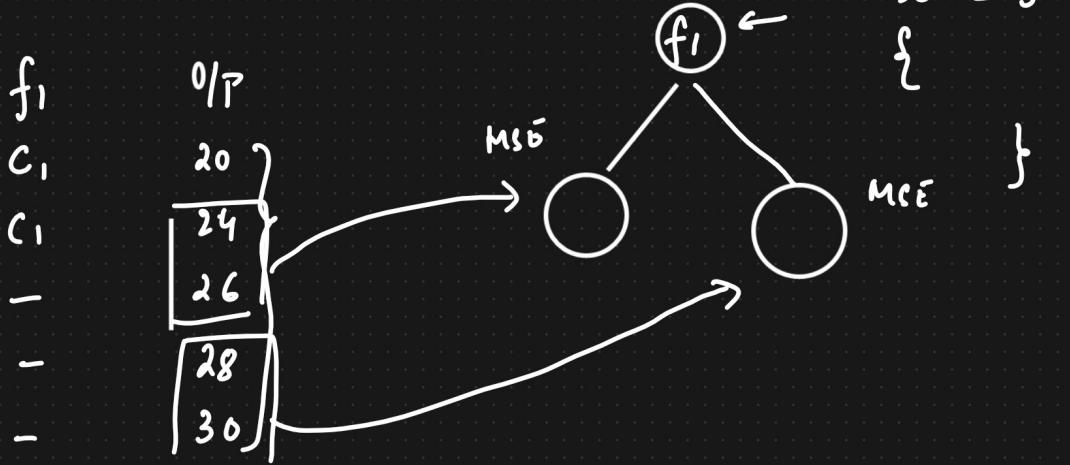
Hyp parameters

Decision \rightarrow Overfitting

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



Decision Tree Regression



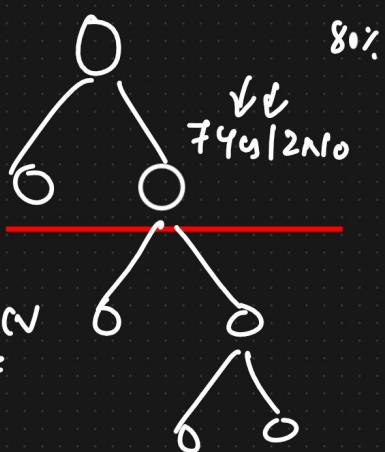
post pruning

pre pruning

Hypoparameter

max_depth, max_leaf

GridSearchCV



Day 5 — Machine Learning Algorithms

Agenda

- ✓ ① Ensemble Techniques
 - Bagging { DJANGO }
 - Boosting { FLASK }
- ✓ ② Random Forest { EDA }
- ✓ ③ AdaBoost { Deep Learning }
- ④ Xgboost → Youtube channel { NLP }

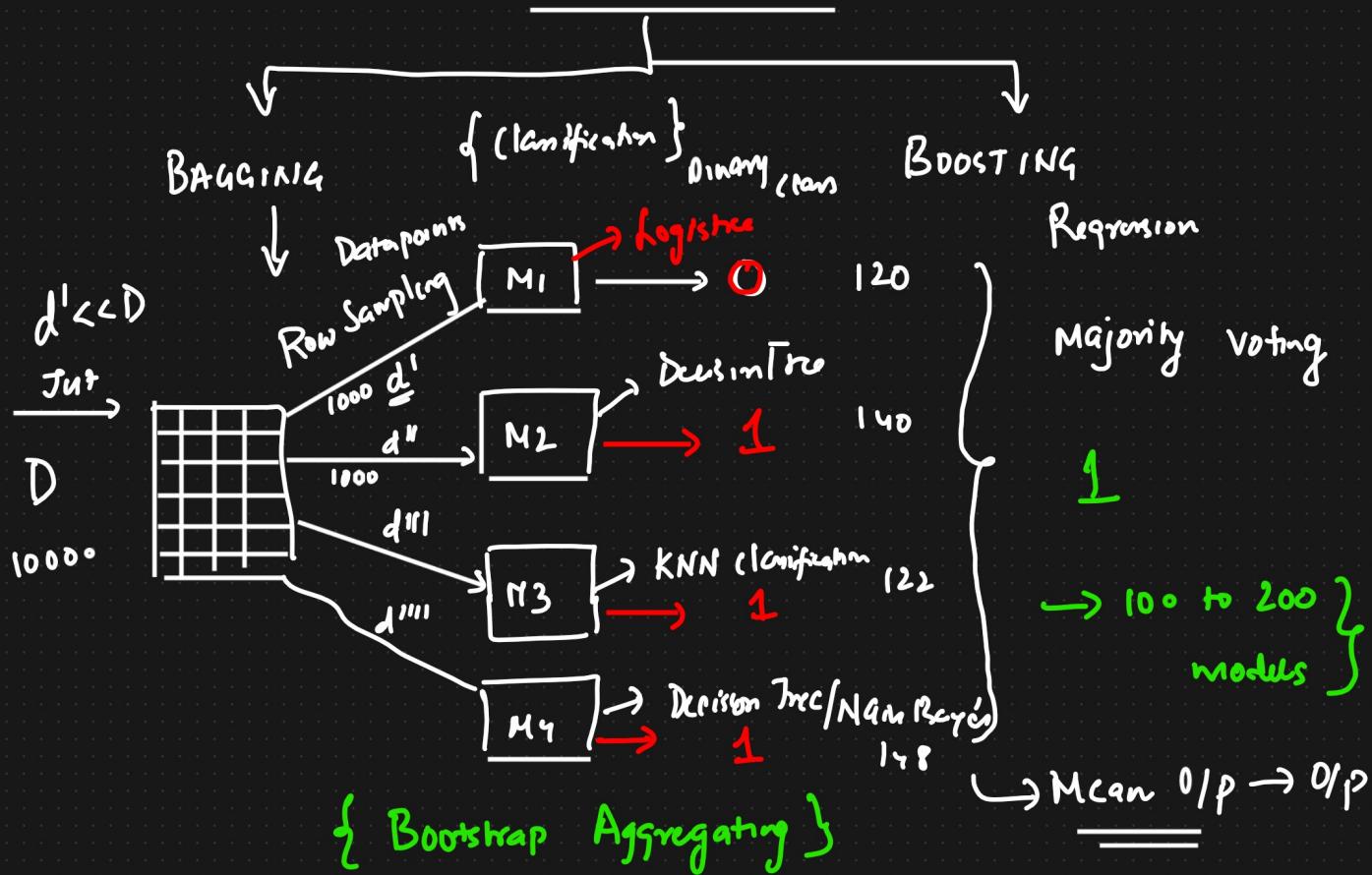
Ensemble Techniques ✓

① Classification & Regression

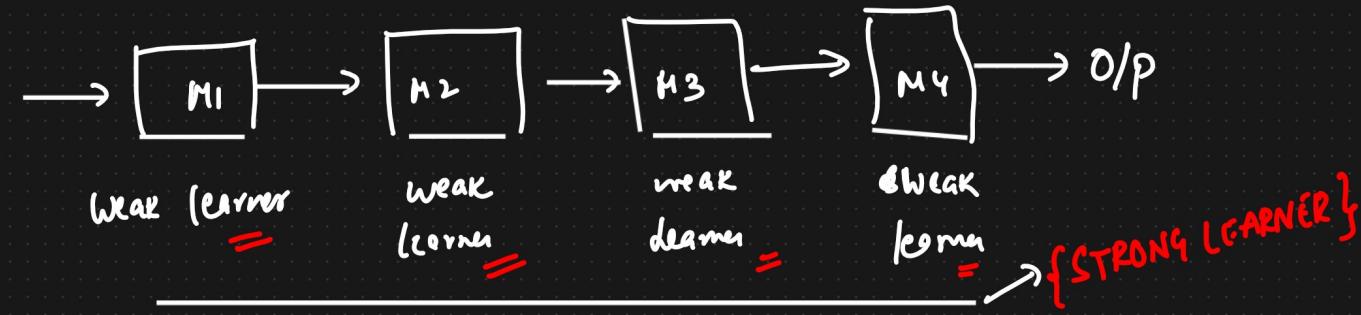
↳ 1 Algorithm $\xrightarrow{\text{TC}}$ Reg

Multiple Algorithms to solve a problem?

Ensemble Techniques



Boosting



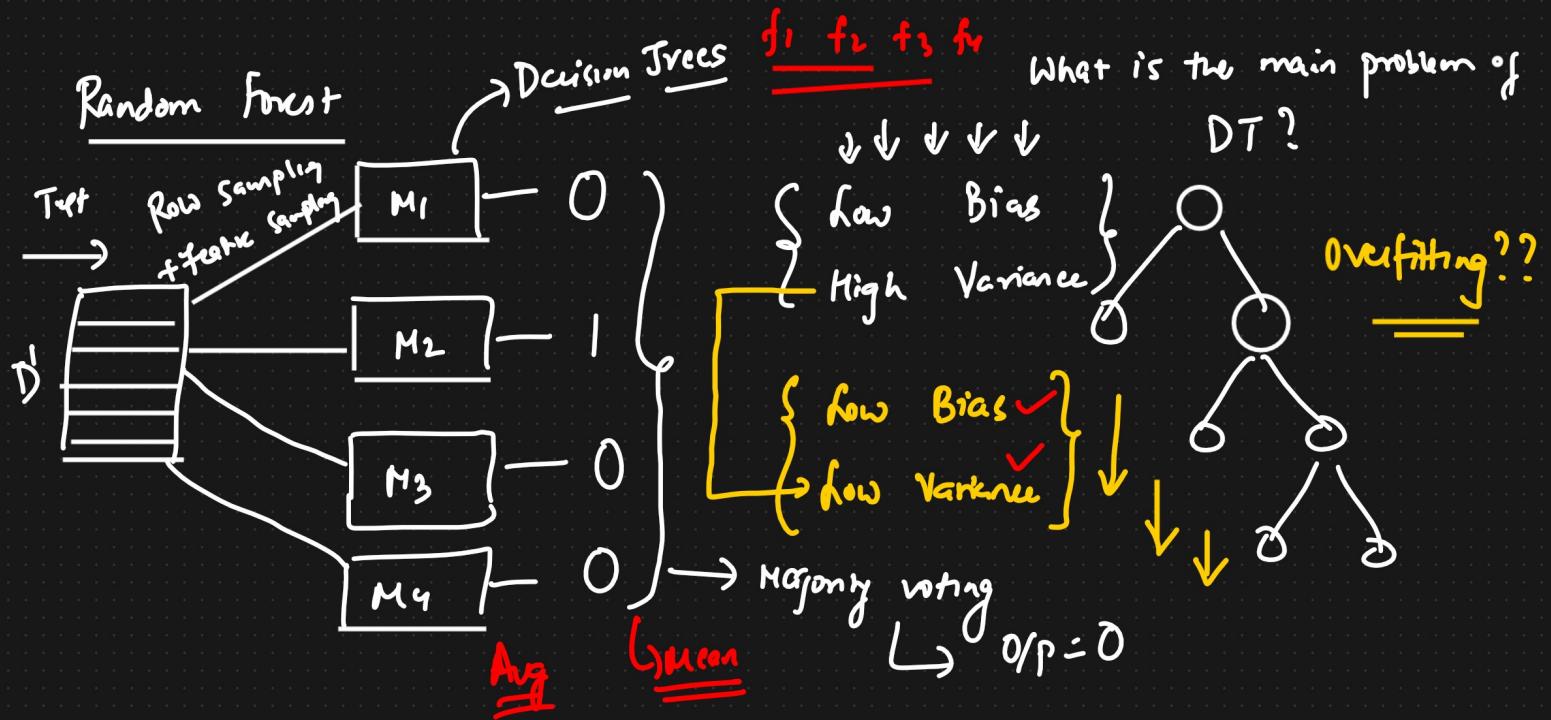
BAGGING

- ↓
- {
- 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

↓

↓

Yes

Impacted by

Yes

Yes

Outlier??

{ Euclidean, Manhattan } =

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

Bagging = Random Forest

Custom Bagging



②

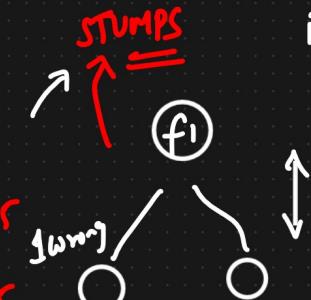
Boosting

i) Adaboost → Decision Tree

Overall = 1

information gain & entropy

f_1	f_2	f_3	f_4	O/p	<u>Weight</u>
-	-	-	-	Yes	$\checkmark \frac{1}{2} 0.05$
-	-	-	-	No	$\checkmark \frac{1}{7} 0.05$
-	-	-	-	-	$\checkmark \frac{1}{7} 0.05$
X	-	-	-	-	$\frac{1}{7} \uparrow 0.349$
-	-	-	-	-	$\checkmark \frac{1}{7} 0.05$
-	-	-	-	-	$\checkmark \frac{1}{7} 0.05$
-	-	-	-	-	$\checkmark \frac{1}{7} 0.05$



② Performance of Stump

$$= \frac{1}{2} \log_e \left(\frac{1 - TE}{TE} \right)$$

$$\text{Total Error} = \frac{1}{7} (TE) = \frac{1}{2} \log_e \left(\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right) = 0.895$$

Correct Records

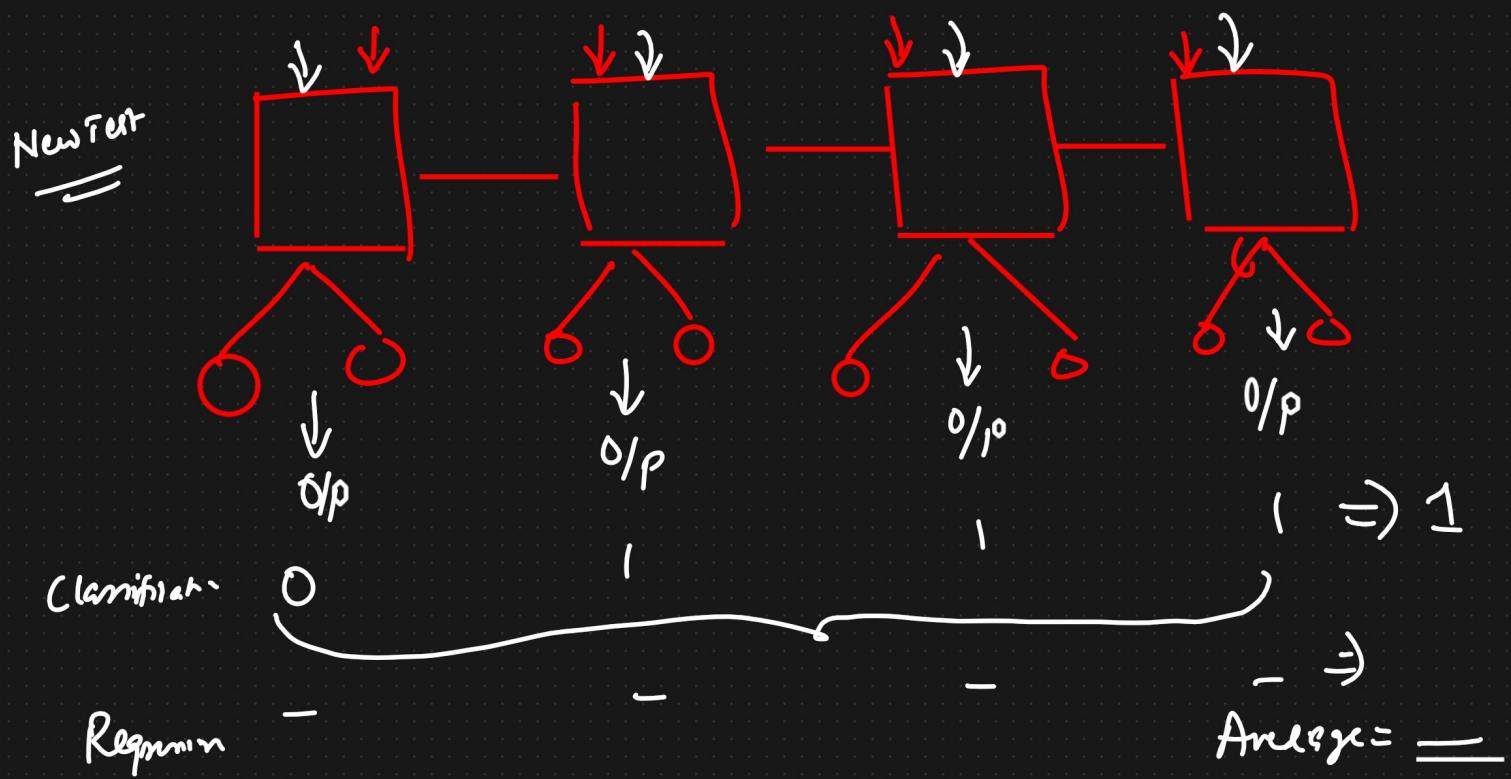
③ New Sample = Weight \times P_s

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

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

<u>New weight</u>	<u>Is it 1??</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	{ $[0.21 - 0.287]$ } \checkmark	
0.05		0.07	$[0.287 - 0.351]$	$0.537 [0 - 1]$
0.05		0.07	$[-]$	0.21
0.05		0.07	$[-]$	0.347
$\frac{0.649}{0.649}$	≈ 1			

Randomly
Create
some number
between



Black models VS White box Models

ANN \rightarrow Black Box

\rightarrow Linear Regression \rightarrow White 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

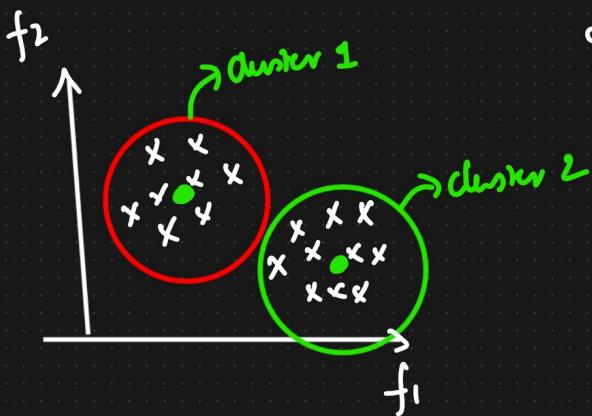
Agenda {
① SVM & SVR
② XGBoost
③ PCA}

Unsupervised ML

Op f_1 f_2

Clusters
↓
Similar kind of
data

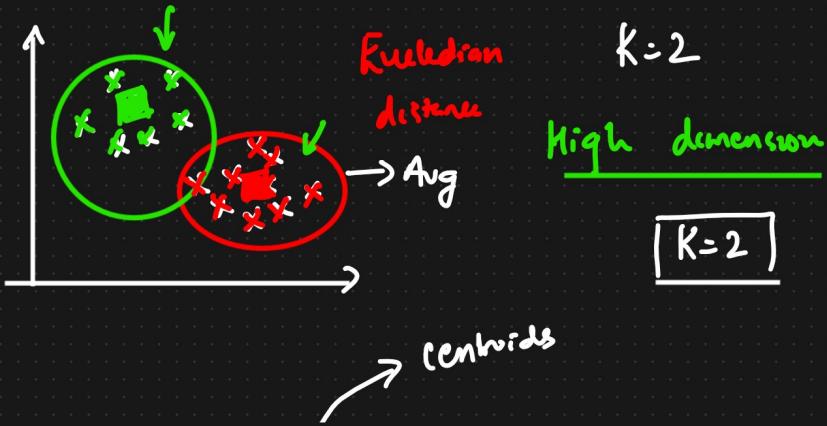
K Means Clustering



Custom Ensemble Technique



K Means $K = \underline{\text{centroids}}$



① We try K values \Rightarrow suitable $K=2$

② Initialize K number of centroids ✓

③ Compute the avg to update centroids ✓

Validating

Elbow method (K value)

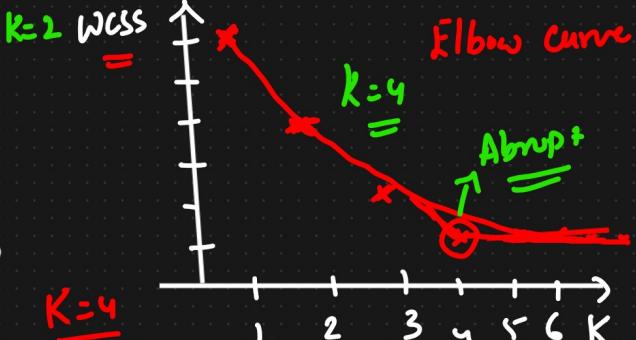
$K=$

within cluster sum of square

$K=1$

for $i=1, 10$

$K=2$ WCSS



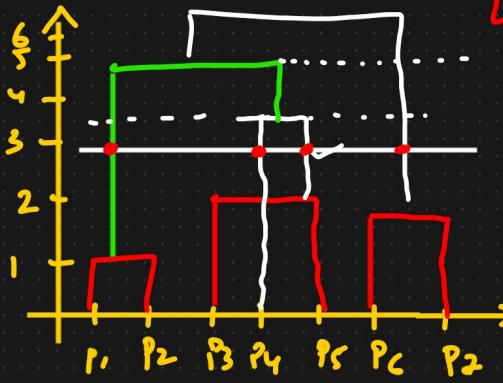
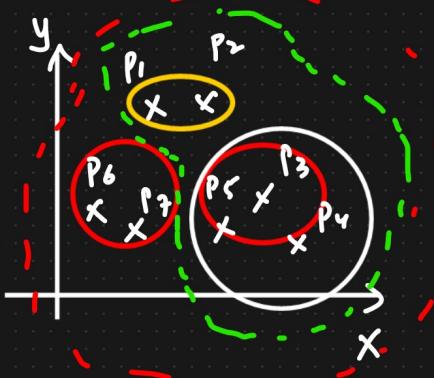
$K=4$



$K=4$

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

② Hierarchical clustering



Dataset is small }
Dataset is large }
Kmeans

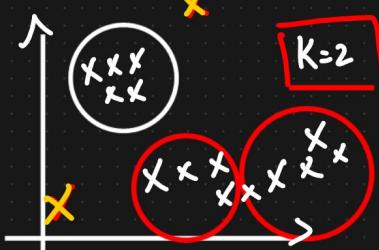
Max Time is taken by KMeans or

Hierarchical clustering ?? ✓

Max Time

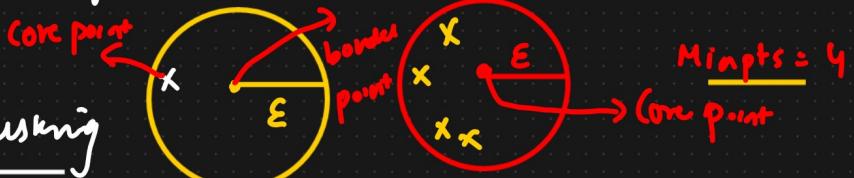
$K=2$

$\Rightarrow \{ K \text{ means} + t \}$



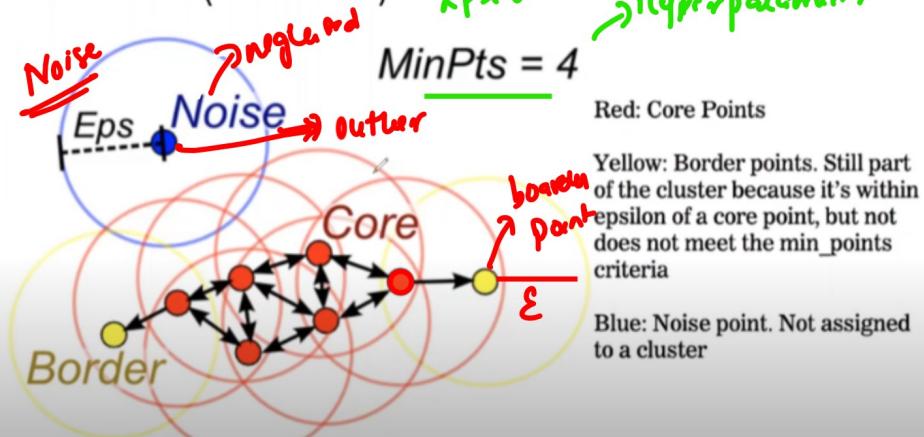
Validate Clustering Models

DBScan Clustering



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

Epsilon → Hyperparameter



① Epsilon

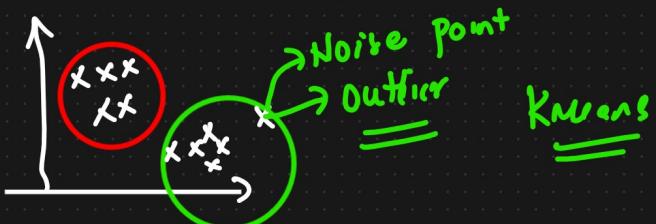
② Min pts

③ Core points

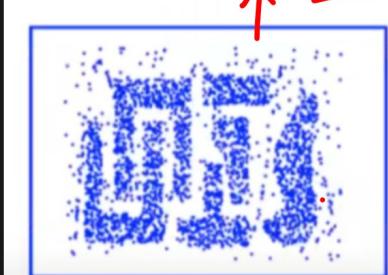
④ Border points

⑤ Noise point

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.

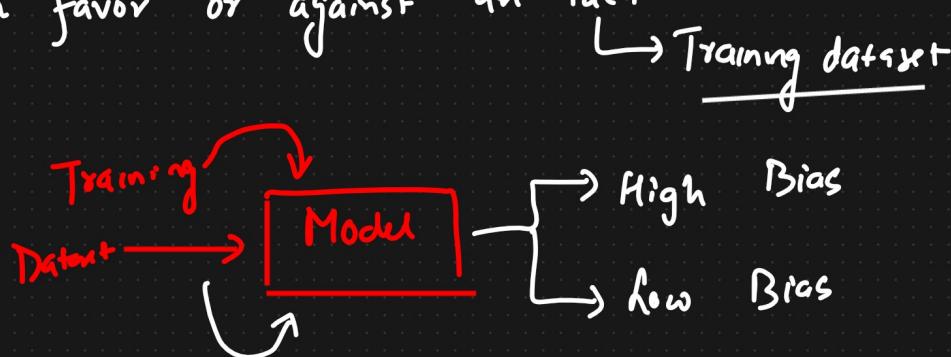
Defn of Bias And Variance

$$\left. \begin{array}{l} \text{Training Data} = 90\% \\ \text{Test Data} = 10\% \end{array} \right\} \Rightarrow \text{Overfitting}$$

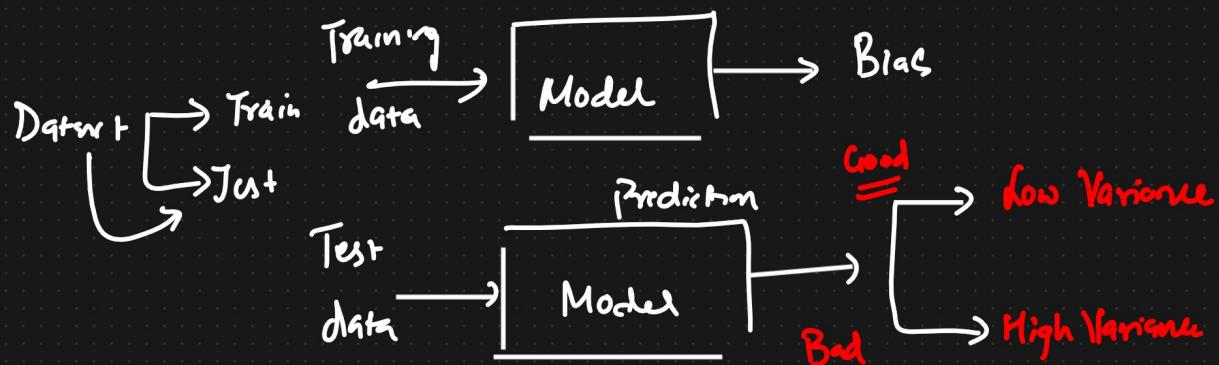
↓

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

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 = 75%

Model 2

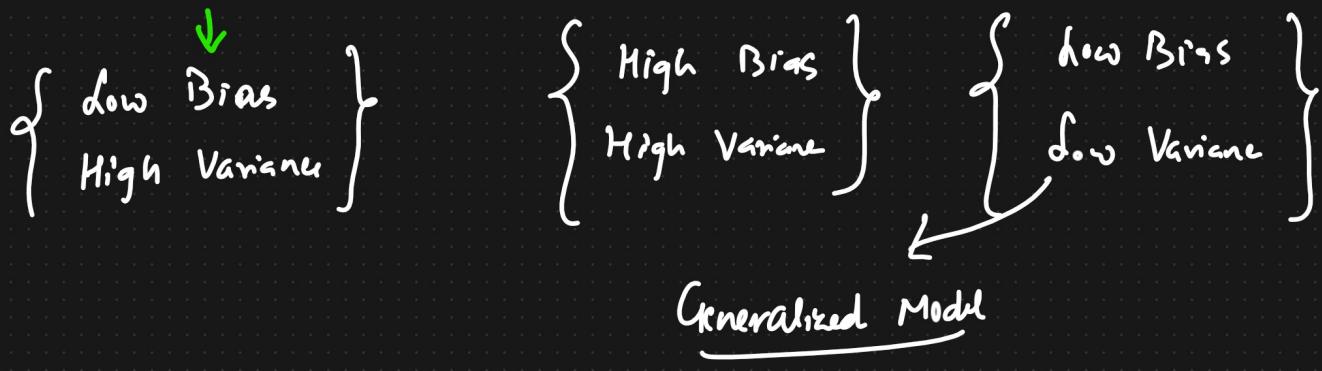
Train Acc = 60%

Test Acc = 55%

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

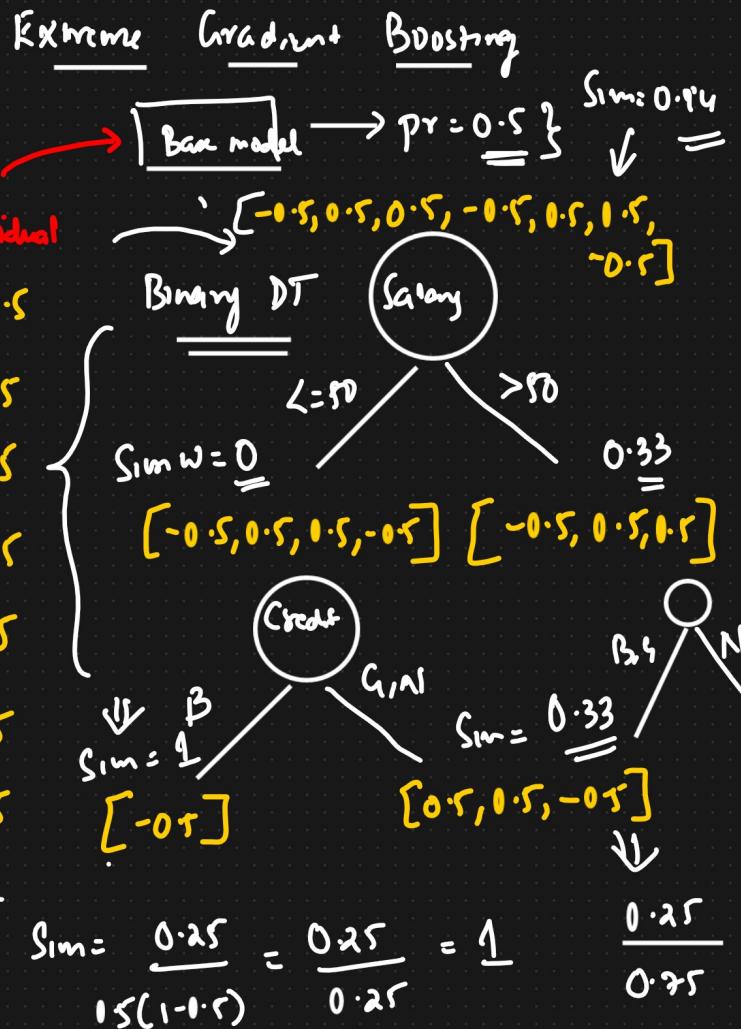
$$\log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.5}{0.5}\right) = 0$$

① Xgboost Classifier

{ Dataset }

Salary	Credit	Approval	Residual
≤ 50	B	0	-0.5
≤ 50	G	1	0.5
≤ 50	G	1	0.5
> 50	B	0	-0.5
> 50	G	1	0.5
> 50	N	1	0.5
≥ 50	N	0	-1.5

$$[\lambda=0] \quad \lambda ??$$



- ① Create a Binary Decision Tree using the feature

- ② Calculate Similarity weight

$$= \frac{\sum (\text{Residual})^2}{\sum (\Pr(1-\Pr) + \lambda)}$$

$$= \frac{0 + \alpha(1)}{0 + \alpha(1)}$$

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

$$1 + 0.33 - 0 = 1.33$$

$\boxed{\alpha = 0.01}$

Sigmoid

$$\textcircled{3} \text{ Information gain}$$

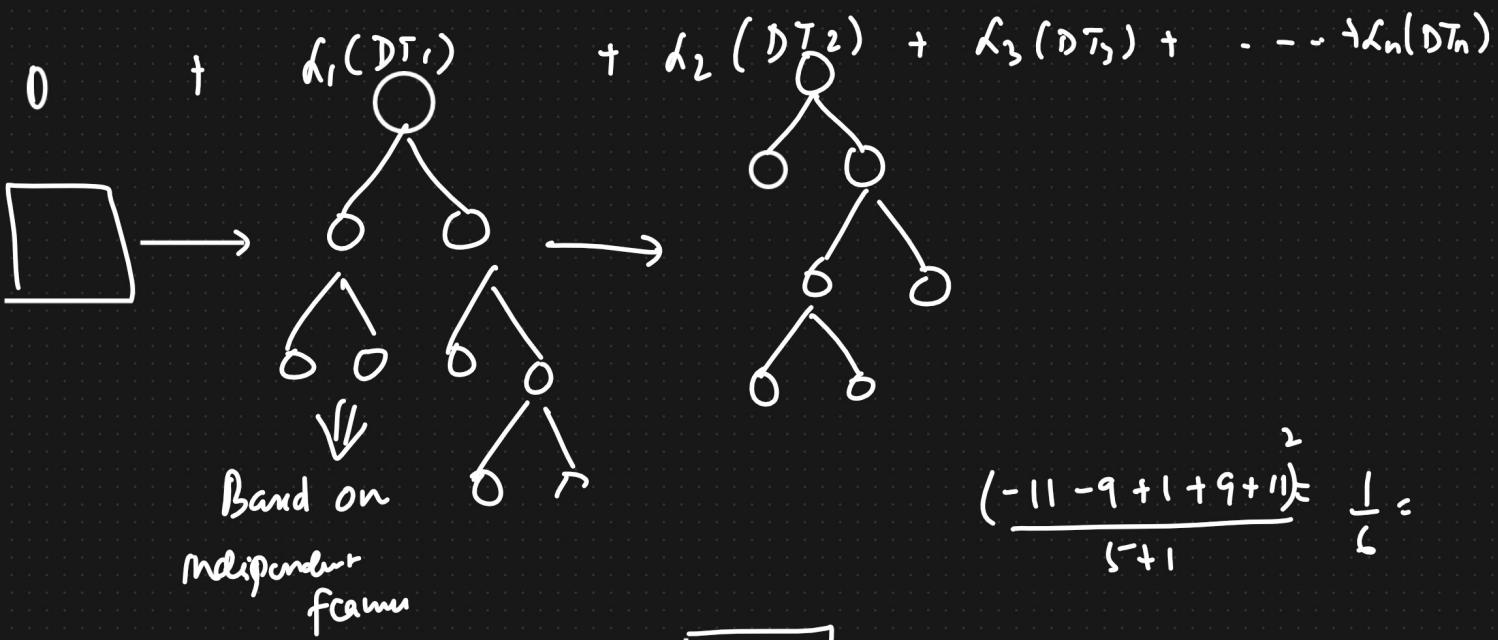
New O/P

$$\rightarrow \left[\sigma \left[0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n) \right] \right]$$

Xgboost → BLACK BOX Model

Pruning

$\lambda \rightarrow$ Cross Validation



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

Base Model → 51
Sum = 16

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

Exp

Exp	Gap	Salary	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	

$$S_{\text{new}} = 68.5$$

$$[-11] \quad [-9, 1, 9, 11]$$

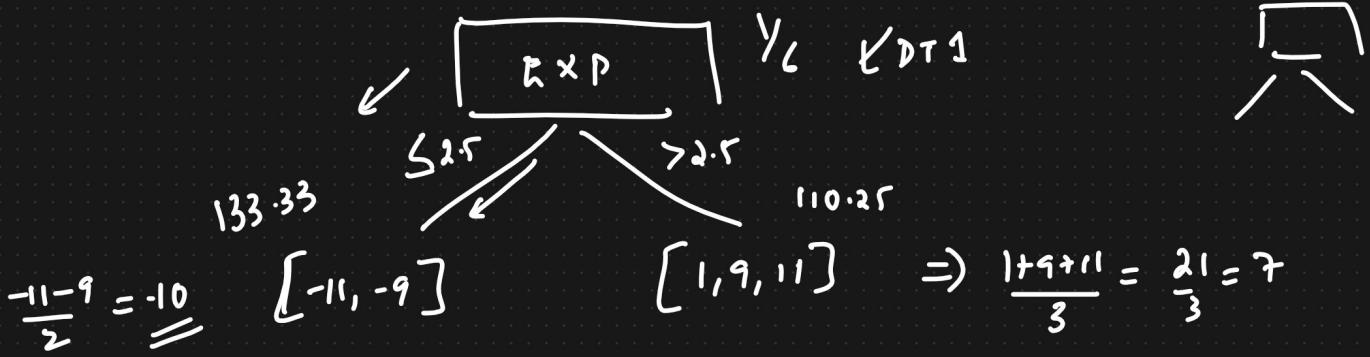
$$\frac{121}{1+1} = 60.5$$

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

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

Information

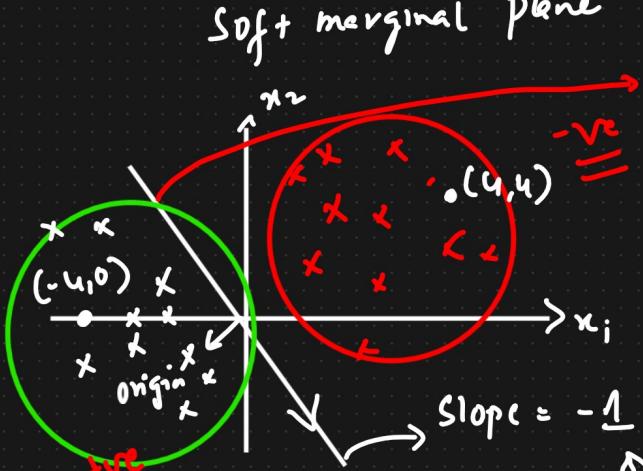
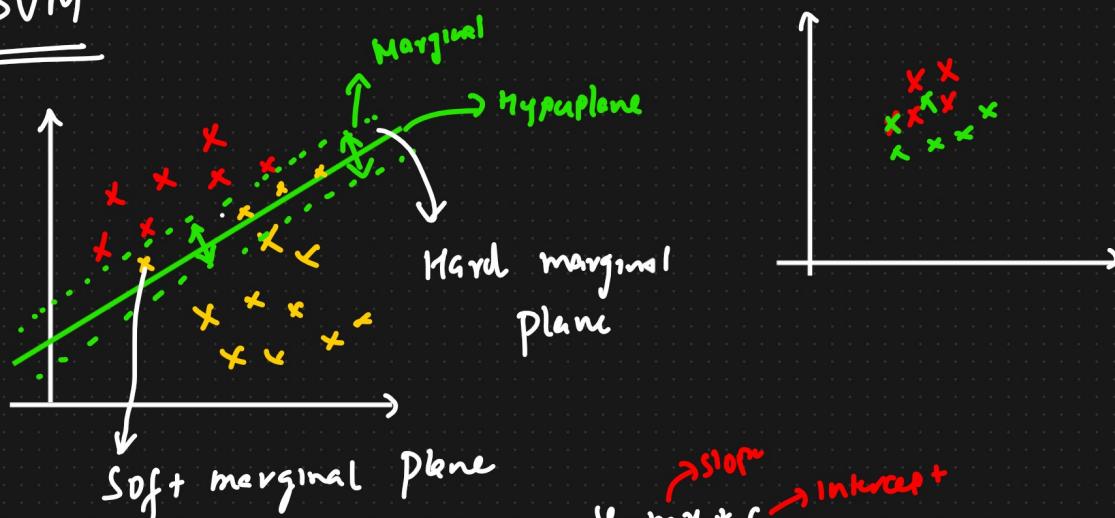
$$= 60.5 + 28.8 - \frac{1}{6} = 89.13$$



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

① {EDA & Feature Engineering} ✓

(3) SVM



$$y = mx + c \quad \text{Slope } m \quad \text{Intercept } c$$

$$\Downarrow \Downarrow$$

$$ax + by + c = 0 \quad \text{Equation of a straight line}$$

$$m = -a/b \quad c = -c/b$$

$$y = -\frac{c}{b} - \frac{ax}{b}$$

$$y = mx + c$$

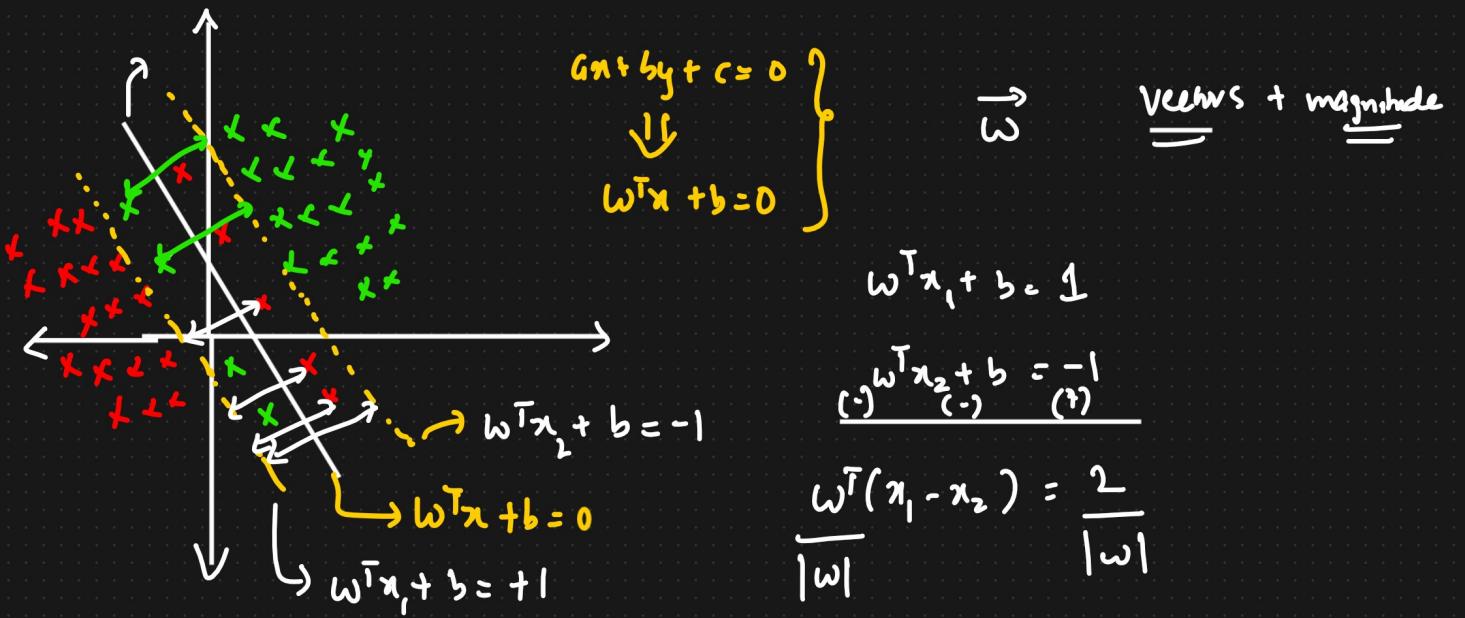
$$y = w_1x_1 + w_2x_2 + \dots + b$$

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 & 0 \end{bmatrix} \quad w \rightarrow -1 \rightarrow \boxed{y = w^T x + b}$$

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

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

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



Maximize (w, b) $\frac{2}{\|w\|} \Rightarrow \text{Marginal} =$
 Such that $y_i \begin{cases} +1 & w^T x_i + b > 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$

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

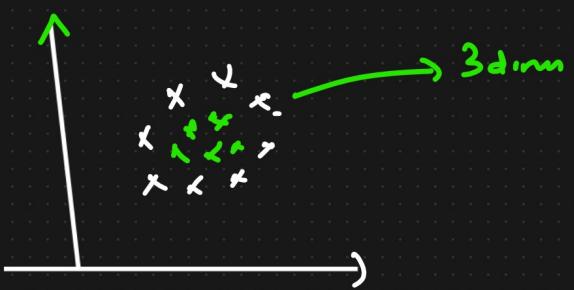
Maximize (w, b) $\frac{2}{\|w\|} \Leftrightarrow \text{Min}_{(w, b)} \frac{\|w\|}{2}$

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



of ξ $\left\{ \begin{array}{l} \text{How many} \\ \text{Errors we can} \\ \text{have} \end{array} \right\}$ SVR = Explore

SVM Kernel



Geometrically
=

A geometric diagram showing a parallelogram with vertices labeled x , x' , x'' , and x''' . The top edge of the parallelogram has arrows pointing to the left, indicating a transformation or mapping.