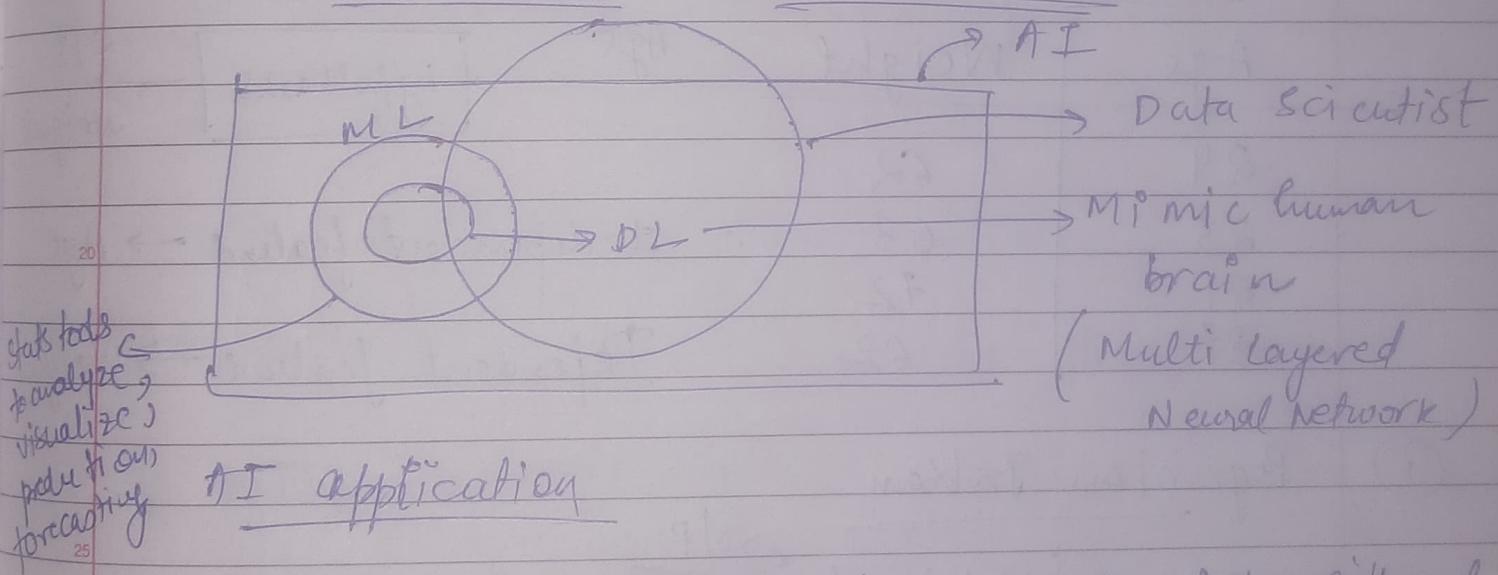


Agenda

- ① Introduction to ML (AI v/s ML v/s DL v/s 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 v/s ML v/s DL v/s DS



AI application is able to do its own task without any human interventions.

Eg - Netflix

```

    graph LR
      Netflix[Netflix] --> Action[Action]
      Netflix --> Comedy[comedy]
      Action --> Recommendation[Recommendation]
      Comedy --> Recommendation
  
```

Action → Recommendation

comedy → //

Amazon.in → iPhone → Headphones

Self Driving Cars → Tesla

Machine & Deep learning

[Reinforcement]

Supervised ML
Regression

Classification

Unsupervised ML → CLUSTERING

→ DIMENSIONALITY
REDUCTION

Supervised ML

Age	Weight
24	62
25	63
21	72
27	62

Age → [hypothesis] → old weight

Independent feature → Age

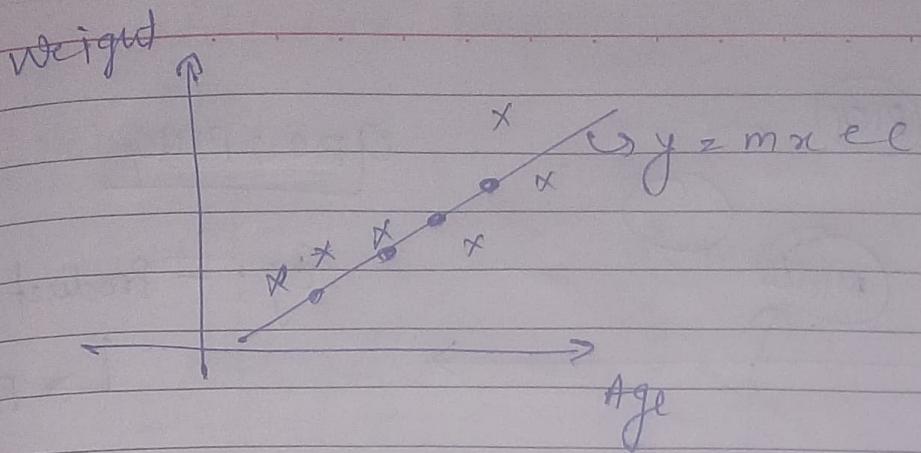
Dependent feature → Weight

①

Regression Problem

Age	Weight	continuous variable
24	72	
23	71	
25	71.5	
-	-	

Regression Problem



② CLASSIFICATION

	No. of hours	No. of play hours	No. of sleep	O/P
10	—	—	—	P
11	—	—	—	P
12	—	—	—	R
13	—	—	—	P
14	—	—	—	R

Two O/P \rightarrow Binary Classification

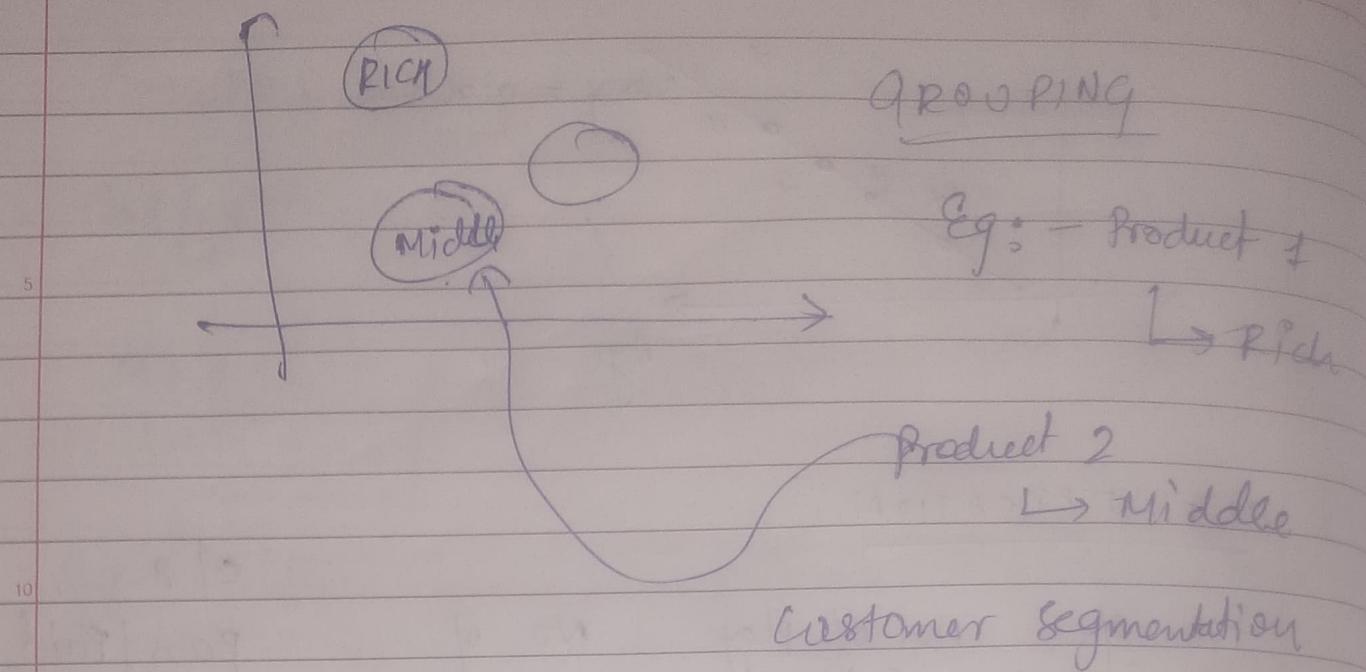
More than Two O/P \rightarrow Multi class classification

Unsupervised ML \rightarrow Clustering

\rightarrow Dimensionality Reduction

Salary Age } No Depend Variable }

Clustering \rightarrow Customer segmentation



Dimensionality Reduction

Two → lower Dimension
L → 1D

PCA, LDA

Supervised

① Linear Regression

② Ridge & Lasso

③ Logistic Reg.

④ Decision Tree

⑤ Adaboost

⑥ Random forest

⑧ Xg boost

⑪ KNN

⑩ SVM

⑨ Naive Bay's

Unsupervised

① K Means

② DB Scan

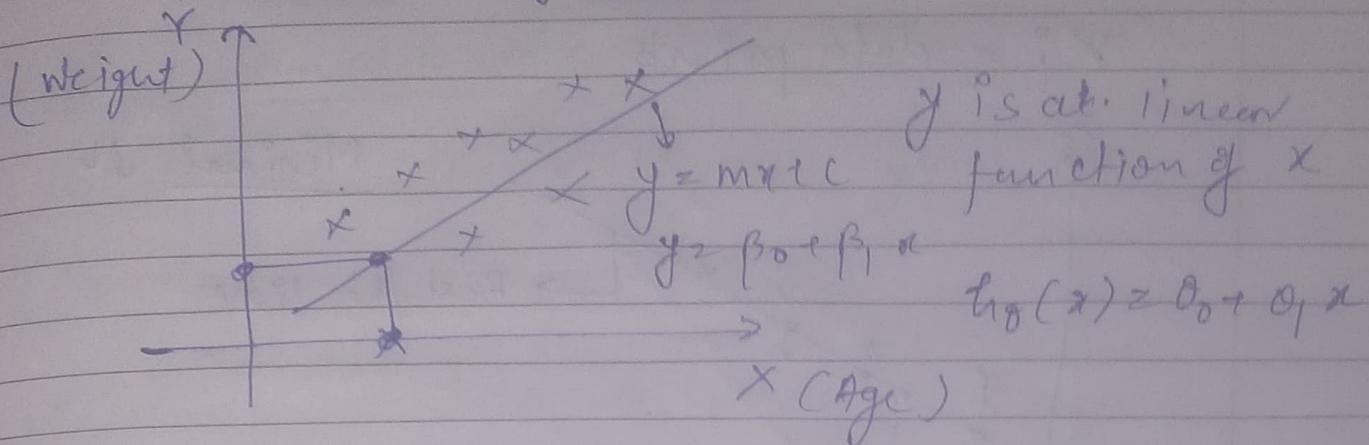
③ Hierarchical

④ K - Nearest Neighbor

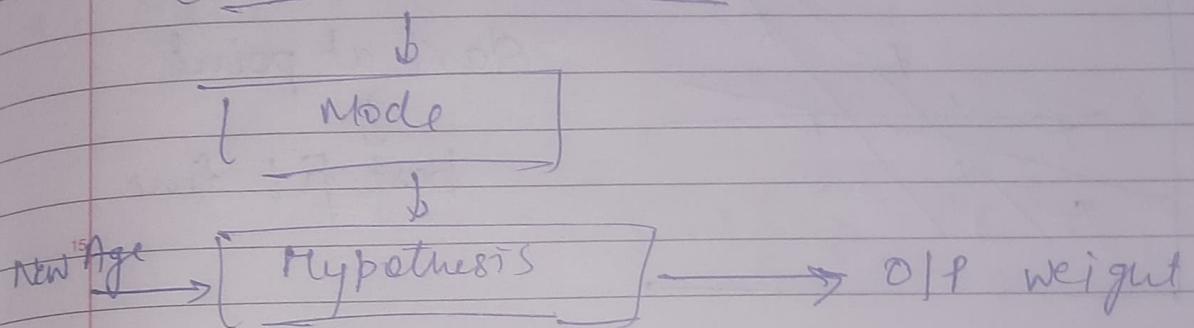
⑤ PCA

⑥ LDA

Linear Regression.



[TRAINING DATABASE]



Equation of a straight line

$$y = mx + c$$

$$y = \beta_0 + \beta_1 x$$

$$h_{\theta}(x_0) = \theta_0 + \theta_1 x_0$$

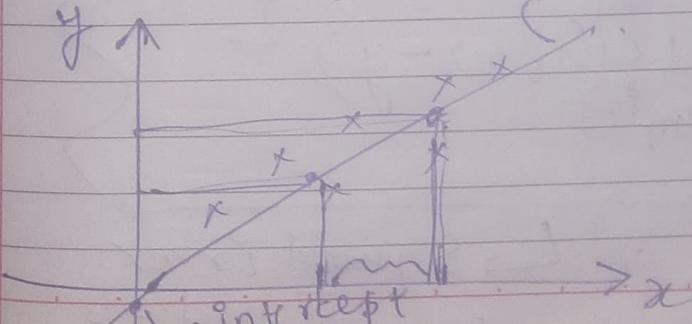
when $x_1 = 0$

$$h_{\theta}(0) = \theta_0$$

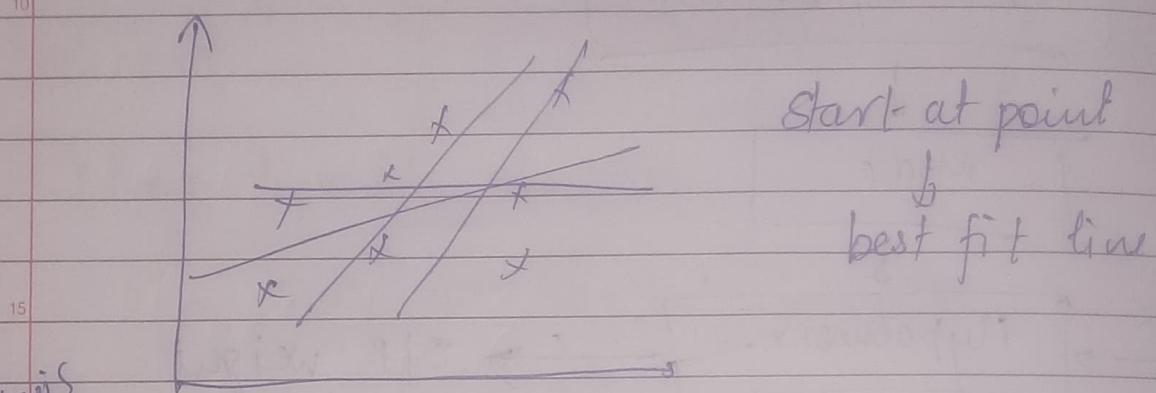
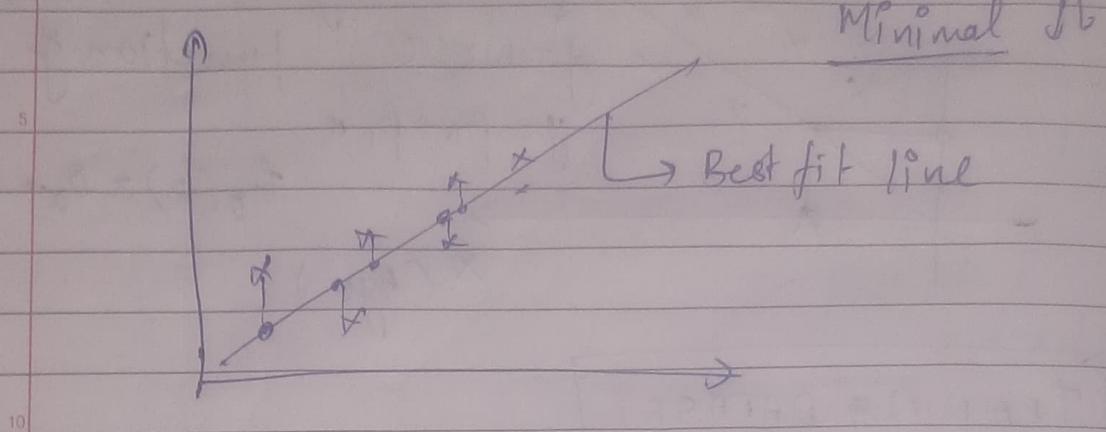
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

θ_0 = Intercept

θ_1 = Slope or coefficient



n_i = data points



Hypothesis

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

cost function

~~$$\sum_{i=1}^n (h_0(x_i) - y_i)^2$$~~

$$\boxed{\text{Cost func}^m = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2}$$

$$\boxed{J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2}$$

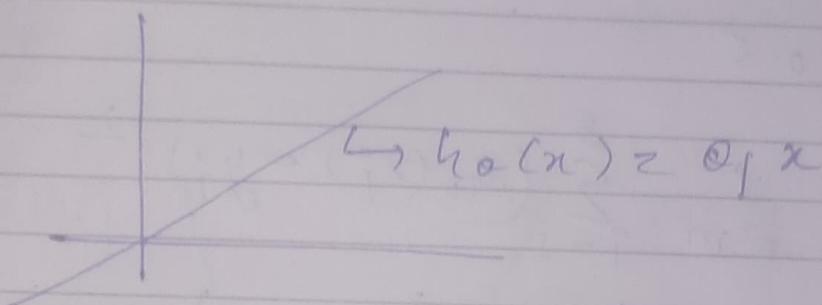
↳ Squared Error Function

What we need to solve

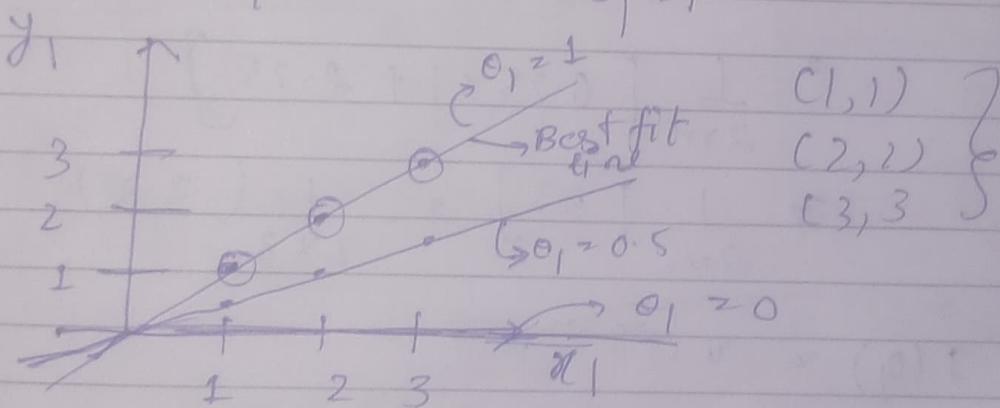
$$\underset{(\theta_0, \theta_1)}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

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

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



$$h_\theta(x) = \theta_1 x \quad \theta_1 \rightarrow 1$$



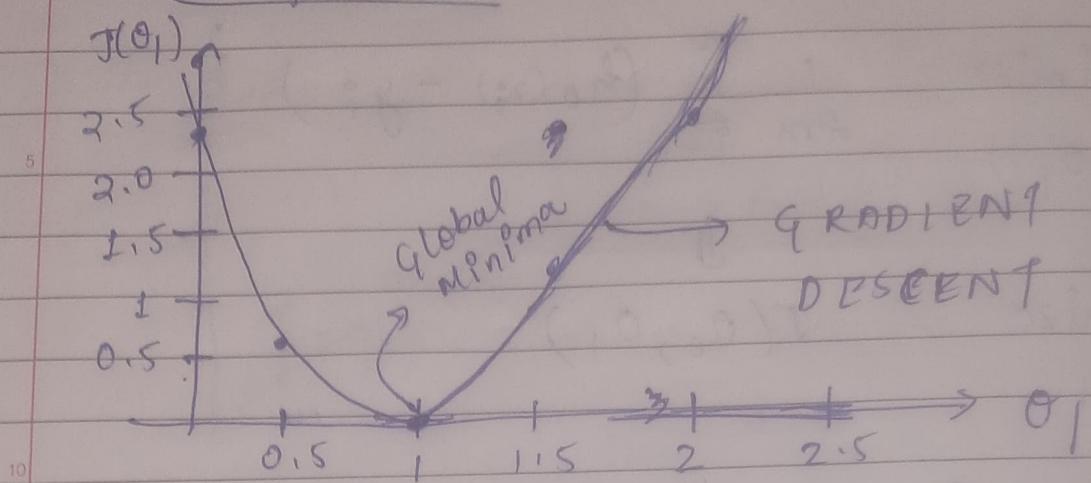
$$\theta_1 = 1$$

$$J(\theta_1) = 0$$

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

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

$$J(\theta_1) = 0$$

Cost function

let $\theta_1 = 0.5$

$$15 \quad J(\theta_1) = \frac{1}{2m} \sum_{i=1}^3 (h_\theta(x_i) - y_i)^2$$

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

$$20 \quad = \frac{1}{2m} [0.25 + 1 + 2.75]$$

$$= \frac{1}{2m} [3.5] = \frac{1}{2 \times 3} (3.5)$$

25 $J(\theta_1) = 0.58$

let $\theta_1 = 0$

$$30 \quad J(\theta_1) = \frac{1}{2(3)} [(0-1)^2 + (0-2)^2 + (0-3)^2]$$

$$= \frac{1}{6} [1+4+9] = \frac{14}{6} = 2.3$$

Convergence Algorithm

Repeat until convergence

$\alpha \rightarrow$ learning rate

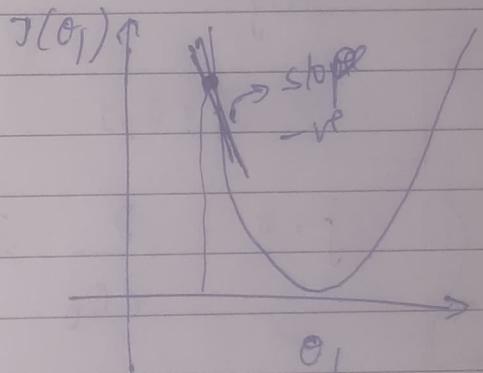
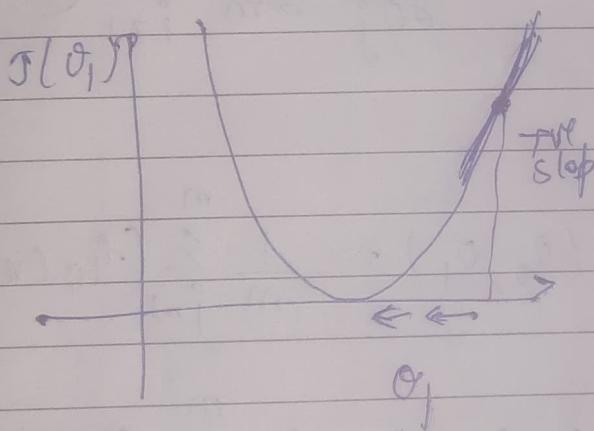
5
S

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

~~$\theta_j = \theta_j - \alpha$~~ $\theta_j = \theta_j - \alpha$

derivative

6
S

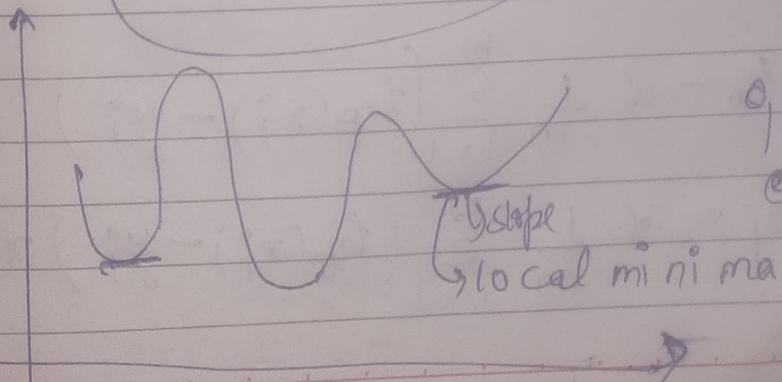


$$\theta_1 := \theta_1 - \alpha (+ve)$$

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

($\theta_1 + \alpha (+ve)$)

$\alpha \rightarrow$ small



GRADIENT DESCENT ALGORITHM

repeat until convergence

{

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

{

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

 $j = 0, 1$

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

Converge
algorithm

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

Repeat until convergence

{

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

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

{

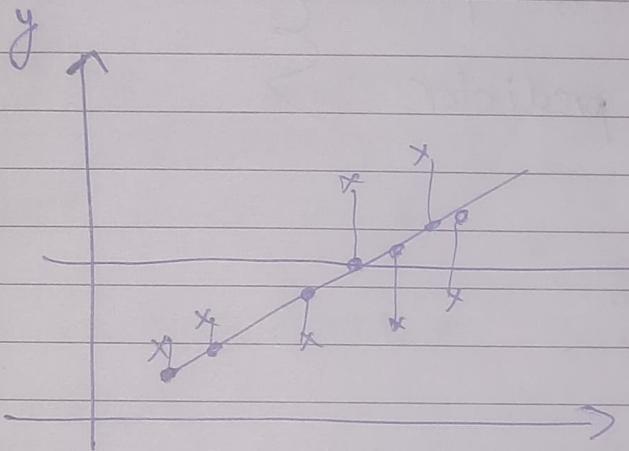


← multi variable

Performance metrics

R² and Adjusted R²

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



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

$\approx 1 - \frac{\text{low}}{\text{high}}$

small number

(Gender)

Bed rooms

Price Location R^2

85%

$R^2 = 90\% \checkmark$

$R^2 = 91\%$

no correlation

Adjusted R^2

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

p = features or predictors

$$p = 2 \Rightarrow R^2 = 90\% \quad R^2_{\text{adjusted}} = 86\%$$

$$p = 3 \Rightarrow R^2 = 91\% \quad R^2_{\text{adjusted}} = 82\%$$

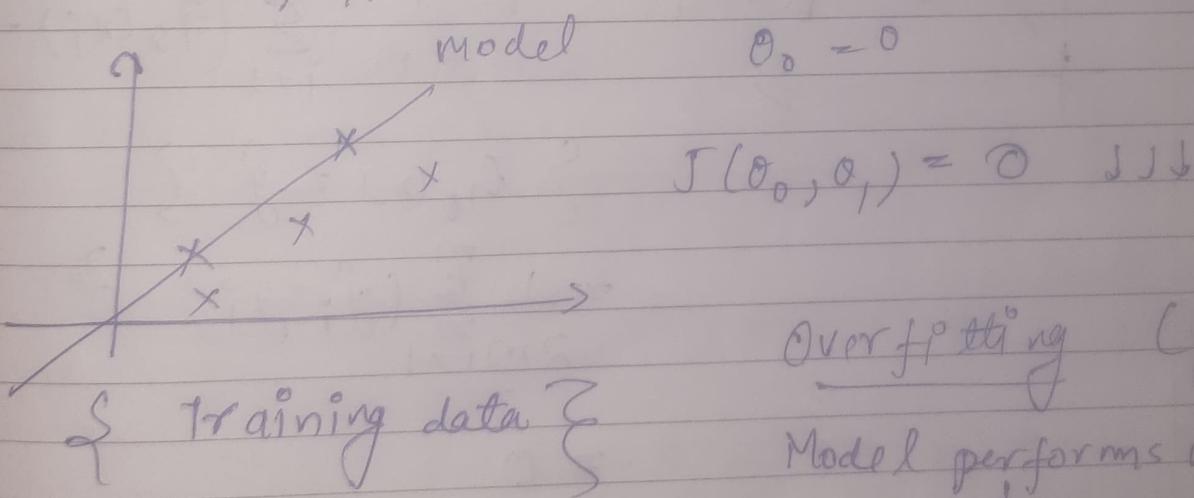
$$\left. \begin{array}{l} N = \text{No. of data points} \\ p = \text{No. of predictor} \end{array} \right\}$$

Agenda

- ① Ridge and Lasso Regression
- ② Assumption of Linear Regression
- ③ Logistic Regression
- ④ Confusion Matrix
- ⑤ Practicals for Linear, Ridge, Lasso & Logistics
Practicals Implementation

Ridge And Lasso Regression

Cost function of Linear Regression = $\frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$



Underfitting } { High bias
High variance }

- ① Model ~~fails~~ Accuracy is bad with training data
- ② Model Accuracy is also bad with Test data

Model 1

Model 2

Model 3

Training Acc = 90%

Test Acc = 80%

Overfitting

Low Bias

High Variance

Training Acc = 92%

Test Acc = 91%

Generalised Model

Low Bias

Low Variance

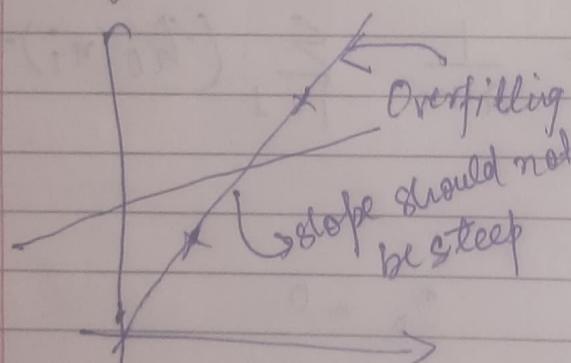
Training Acc = 75%

Test Acc = 6.5%

Underfitting

High Bias

High Variance



$$\hat{J}(0_1) = 0$$

$$= \frac{1}{m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

$$\left\{ h_0(x) = \hat{y} \right\}$$

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

$$= 0$$

Ridge (L2 Regularization)

$$J(\theta_1) = (\hat{y}_i - y_i)^2 + \lambda (\text{slope})^2$$

Let $\lambda = 1$ $\theta_1 = 2$

$$= 0 + 1(2)^2 = 4$$

$$J(\theta_1) = (\hat{y}_i - y_i)^2 + \lambda (\text{slope})^2$$

(small value) $\leftarrow 1(1.5) + 1(1.5)^2$

$\approx \underline{\underline{3}}$

~~\$\nabla\$~~ prevent overfitting }

$\lambda \rightarrow$ hyperparameter

Lasso (L1 Regularization)

$$J(\theta_1) = (\hat{y}_i - y_i)^2 + \lambda |\text{slope}|$$

\hookrightarrow feature selection

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

$$|\text{slope}| = |\theta_1 + \theta_2 + \theta_3 + \dots + \theta_n|$$

- ① Preventing Overfitting } → Lasso
 ② Feature selection } (L1 regularization)
 { $\lambda \rightarrow$ found out through cross-validation

Ridge Regression (L2 Norm)

$$\text{Cost function} = (h_0(x_i) - y_i)^2 + \lambda (\text{slope})^2$$

Purpose: Preventing overfitting

Lasso Regression (L1 Reg)

$$\text{Cost function} = (h_0(x_i) - y_i)^2 + \lambda |\text{slope}|$$

Purpose: 1) Prevent overfitting
 2) Feature selection

Assumption of Linear Regression

- ① Normal / Gaussian Distribution → Model will get trained well

- ② Standardization } Scaling data } → z-score
 Standardization

③ Linearity

95% correlated

④ Multi collinearity

$$x_1 \begin{bmatrix} x_2 & x_3 \end{bmatrix} \gamma$$

Variation In plation factor?

Logistic Regression (Classification)

No. of study No. of Play

P / P

→ Binary classification

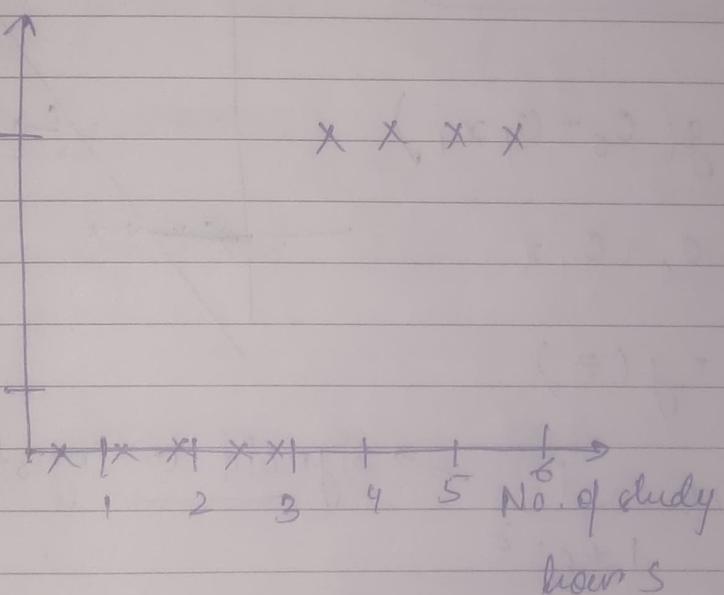
P

→ Multi class

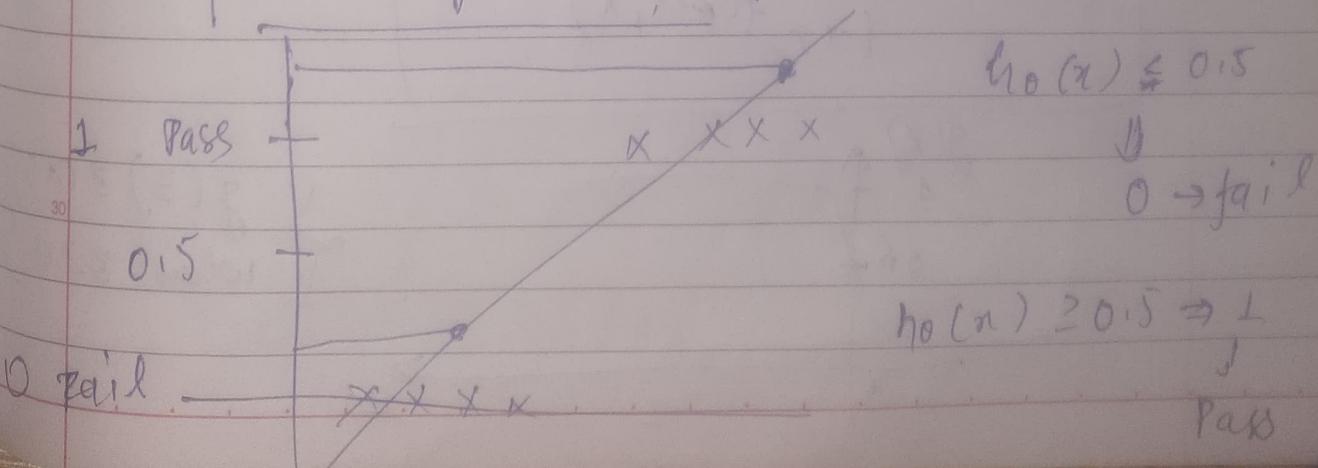
R

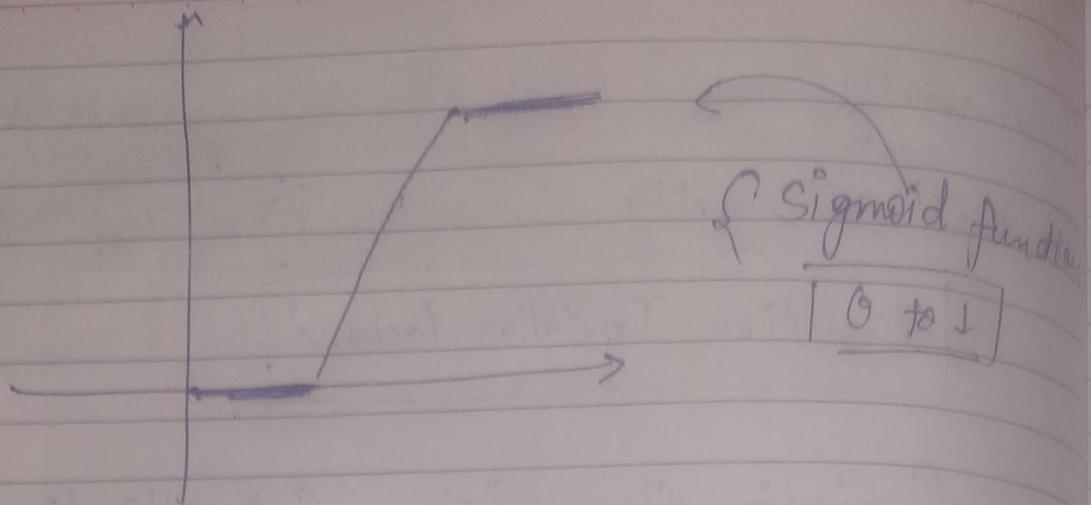
classification

↓
Fail Pass



[Linear Regression??]





Decision Boundary Logistic Regression

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

$\boxed{h_0(x) = \theta^T x}$

squash
 ~~x~~

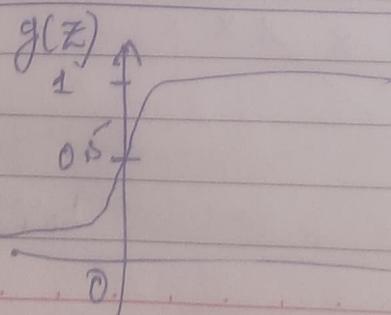
$$h_0(x) = g(\theta_0 + \theta_1 x_1)$$

$$\text{Let } z = \theta_0 + \theta_1 x_1$$

$$h_0(x) = g(z)$$

$$h_0(x) = \frac{1}{1 + e^{-z}} \rightarrow \text{Sigmoid or Logistic function}$$

$$\boxed{h_0(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}}$$



$\left\{ \begin{array}{l} g(z) \geq 0.5 \\ \text{when } z \geq 0 \end{array} \right.$

Training set

$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\}$$

$y \in \{0, 1\} \rightarrow \underline{\text{0 or 1}}$

$$h_0(z) = \frac{1}{1+e^{-z}} \quad [z = \theta_0 + \theta_1 x]$$

Let $\theta_0 = 0$

$$\rightarrow [z = \theta_1 x]$$

Change parameter θ_1 ?

Cost function

Linear Regression $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$

Logistic Regression

$$h_0(x) = \frac{1}{1+e^{-(\theta_1 x)}}$$

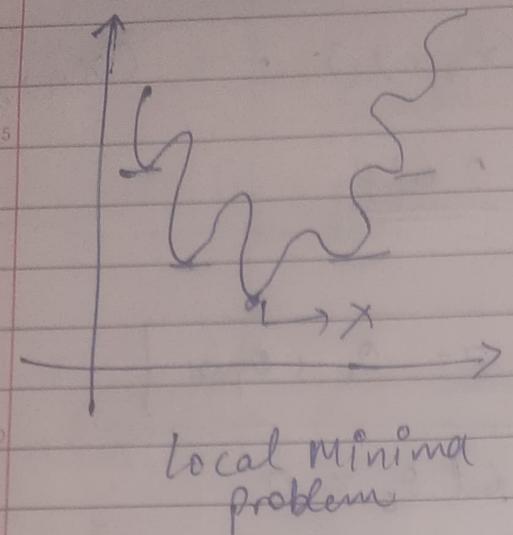
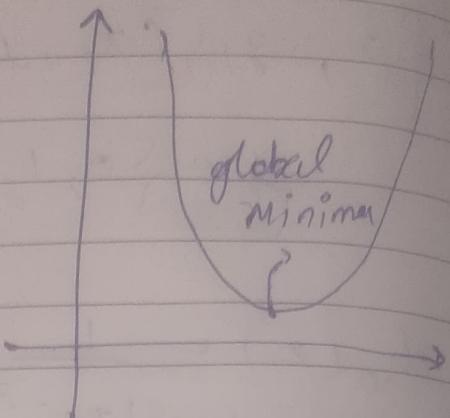
logistic

Regression Cost function = $\frac{1}{2} \sum (h_0(x_i) - y_i)^2$ ∇ we cannot use this

$$h_0(x) = \frac{1}{1+e^{-\theta_1 x}}$$

cost function
for logistic

Non Convex function

Non Convex functionConvex functionLogistic Regression Cost function

$$\text{cost}(h_0(x_i) \neq y_i)$$

$$h_0(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$$

$$J(\theta_1) = \begin{cases} -\log(h_0(x)) & y=1 \\ -\log(1-h_0(x)) & y=0 \end{cases}$$

if $y=1$

$$J(\theta_1)$$

cost = 0 if $y=1$ $h_0(x)=1$

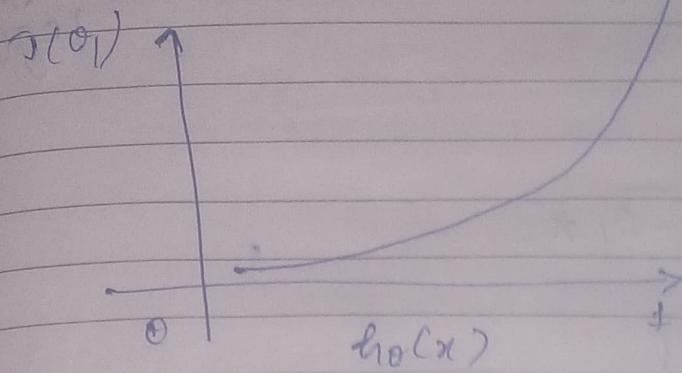
30

0

 $h_0(x)$

classification

if $y = 1$



$$\text{Then } \text{cost}(h_0(x_i), y) = \begin{cases} -\log(h_0(x_i)), & \text{if } y = 1 \\ -\log(1-h_0(x_i)), & \text{if } y = 0 \end{cases}$$

$$\left[\text{cost}(h_0(x_i), y) = -y \log(h_0(x_i)) - (1-y) \log(1-h_0(x_i)) \right]$$

if $y = 1$ cost function

$$\text{cost}(h_0(x_i), y) = -\log(h_0(x_i))$$

if $y = 0$

$$\text{cost}(h_0(x_i), y) = -\log(1-h_0(x_i))$$

$$J(\theta_1) = -\frac{1}{m} \sum_{i=1}^m \left(-y_i \log(h_0(x_i)) - (1-y_i) \log(1-h_0(x_i)) \right)$$

$$h_0(x_i) = \frac{1}{1 + e^{-\theta_1 x_i}}$$

converge repeat until

δ

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

δ

Performance Metrics of Classification Problems

		Actual		Prediction		Actual	
		y	\hat{y}	1	0	1	0
x_1	x_2	-	-	0	1	1	0
-	-	-	-	1	1	2	2
-	-	-	-	0	0	1	1
-	-	-	-	1	1	-	-
-	-	-	-	1	1	-	-
-	-	-	-	0	1	-	-
-	-	-	-	1	0	-	-
-	-	-	-	1	0	-	-
		Predicted		Actual			
		T_P	F_P	T_N	F_N		

6 confusion matrix

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

$$\approx \frac{3+1}{3+2+1+1} = \frac{4}{7} = 0.57 \approx 0.57 \text{ or } 57\%$$

$0 \rightarrow \text{Normal}$ $1 \rightarrow \text{Abnormal}$ } Imbalanced data | Biased

$0 \rightarrow \text{Good}$ } Balanced data
 $1 \rightarrow \text{Bad}$

(Sensitivity)

① Precision

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

② Recall

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

③ F-score

		Actual	
		1	0
Pred	1	TP	FP
	0	FN	TN

{ Spam Classification } → Precision

{ HAS CANCER OR NOT } → Recall

{ Tomm, Stock Market } → Recall / Precision
 is going to crash }

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

$$\beta = 1 \Rightarrow \frac{(1 + 1)}{2} \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$
$$= \frac{2(\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

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

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

$$\beta = 2 \quad \frac{(1 + 2^2)}{4P + R} \frac{P \times R}{P + R}$$

20 Beta score

25

30

Agenda

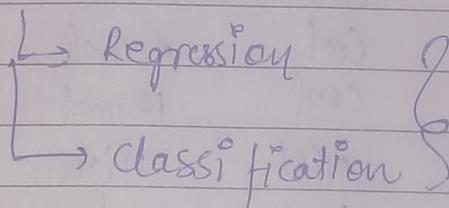
① Decision Tree Classification

② Decision Tree Regression

③ Practical Implementation

④ Ensemble Techniques

Decision Tree



if ($age \leq 18$):

 print ("College")

elif ($age > 18$) and

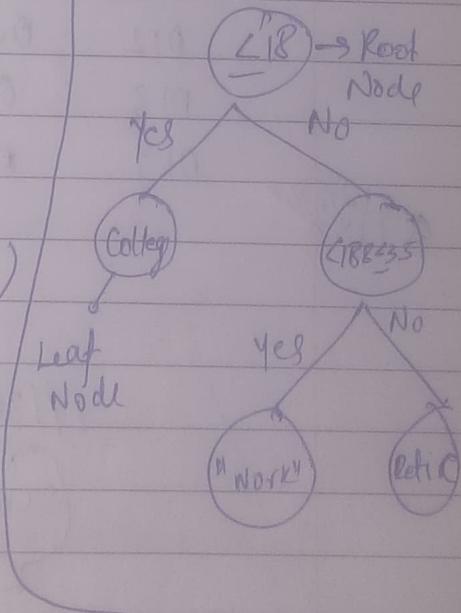
$(age \leq 35)$ and $(age > 18)$

 print ("work")

else

 print ("Retire")

Decision Tree

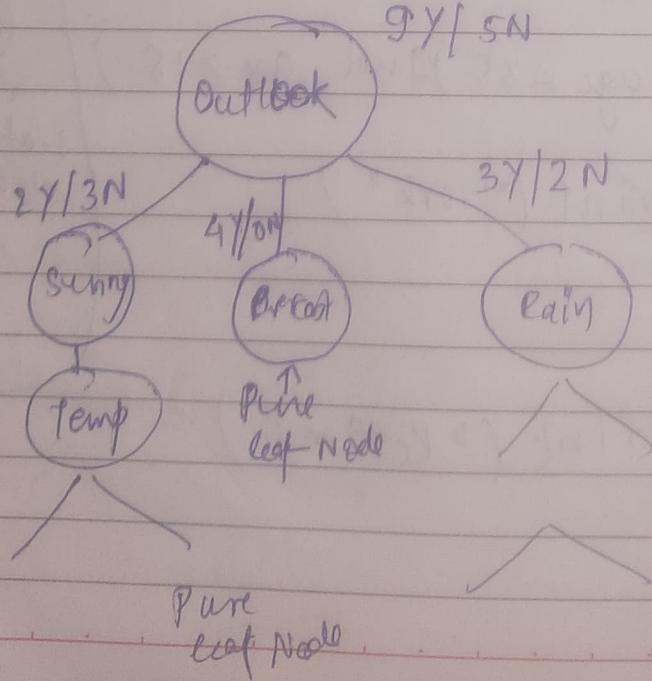


Regression & Classification ??

DECISION TREE

Nest if else → Decision tree

Day	<u>Outlook</u>	<u>Temperature</u>	<u>Humidity</u>	<u>Wind</u>	<u>Play</u>
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



Pure split ✓

Impure split ✓

① Purity \rightarrow Pure split ??

\hookrightarrow Entropy $\quad \mathfrak{S}$
 \hookrightarrow Gini coefficient/
impurity

② How the features are selected

\hookrightarrow Information Gain ?

Entropy

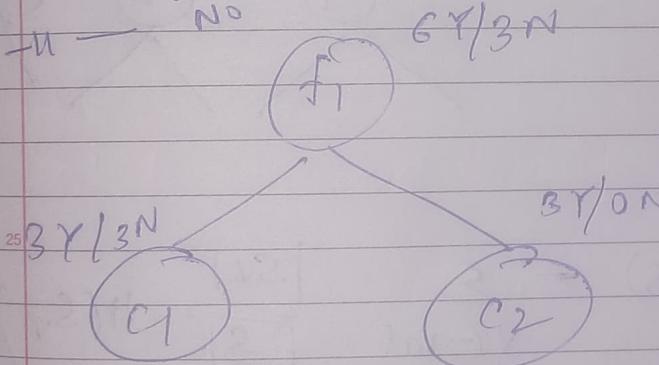
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

P_+ = probability of yes

$P_- = 1 - P_+$

Gini = Impurity

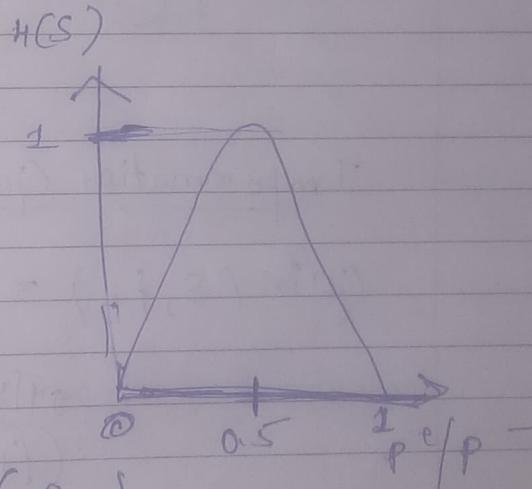
$$G^I = 1 - \sum_{j=1}^n (P_j)^2$$



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

$$-\frac{0}{3} \log_2 \frac{0}{3}$$

$$[H(S) = 0]$$



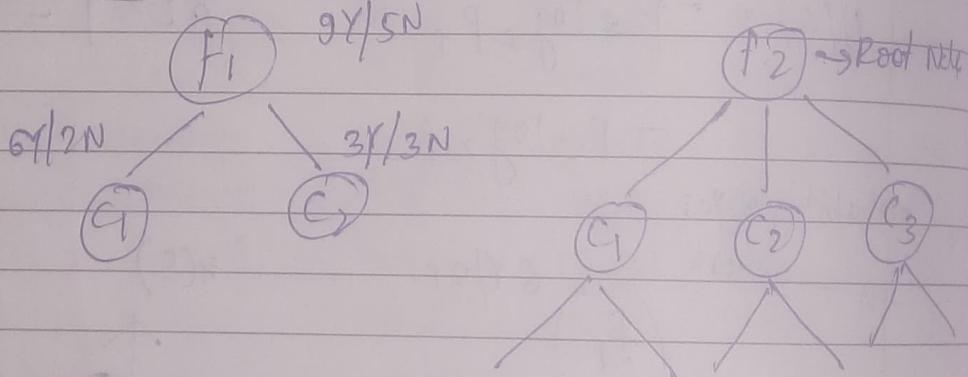
$$Q: H(S) = -\frac{3}{6} \log_2 \frac{3}{6} = \frac{3}{6} \log_2 \frac{3}{6}$$

$$\boxed{H(S) = 1}$$

Purity Test \rightarrow Entropy

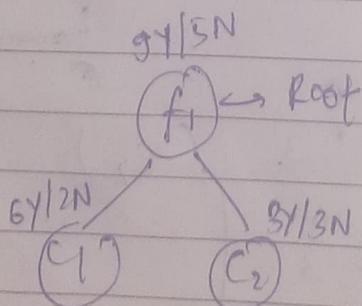
$H(S) = 1$ Impure split $\rightarrow H(S) \neq 0$
 $H(S) = 0$ pure split

→ Which feature do take to split? f_1, f_2, f_3



Information Gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{V(\text{Val})} \frac{|S_V|}{|S|} H(S_V)$$



$H(S) \rightarrow \underline{\text{entropy}}$

$$H(S) = -P_1 \log_2 P_1 - P_2 \log_2 (P_2)$$

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

$$H(S) = 0.94$$

$$\text{Gain } H(S) = H(S_{VC_1}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$\boxed{H(S_{VC_1}) = 0.8}$$

$$\boxed{H(S_{VC_2}) = 1}$$

$$\text{Gain}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.8 + \frac{6}{14} \times 1 \right]$$

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

$$\text{Gain}(S, f_2) = 0.05$$

Using which feature
should I start splitting
first?

f_1

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

Gini Impurity

$n = 2$ output {yes, no}

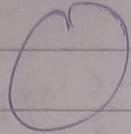
$$GI = 1 - \sum_{i=1}^n (P_i)^2$$

$$= 1 - [(P_1)^2 + (P_2)^2]$$

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

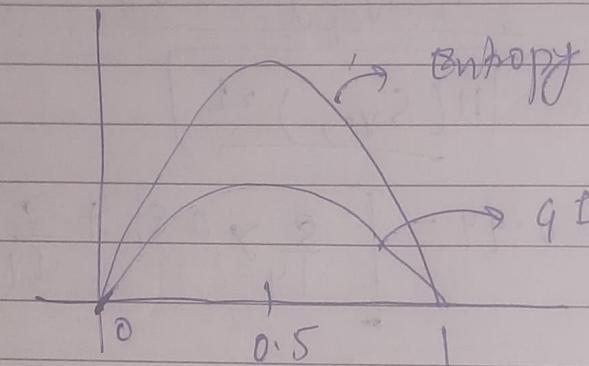
$$= 1 - \left[\frac{1}{2} \right] = 0.5$$

$2Y/2N$



Entropy \rightarrow

$$G_I = 0.5$$

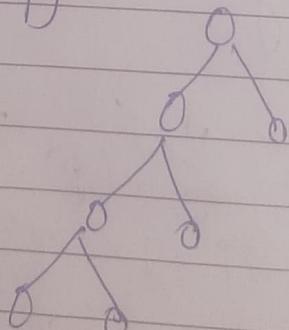


Entropy \rightarrow log \rightarrow More time

~~G_L~~ \rightarrow Simple maths \rightarrow Less time

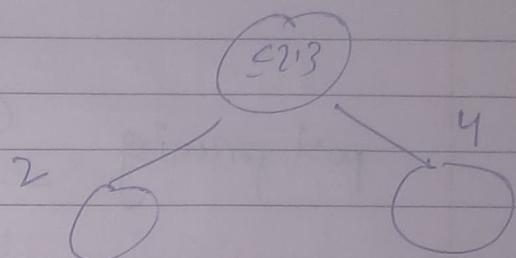
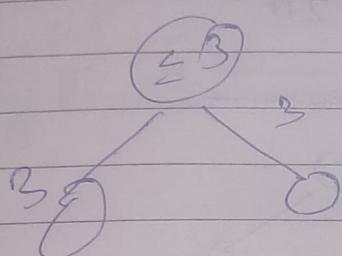
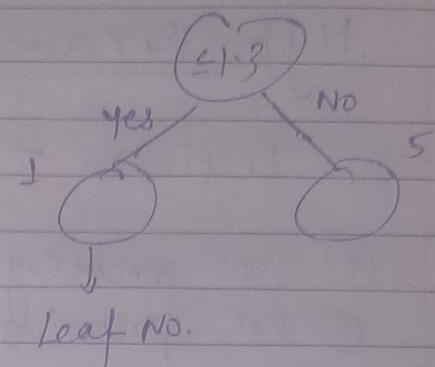
Fast

$Gini >> Entropy$



continues.

f_1	O/P	\Rightarrow	f_1
2.3			P
1.3			1.3
4			2.3
5			3
7			4
3			5



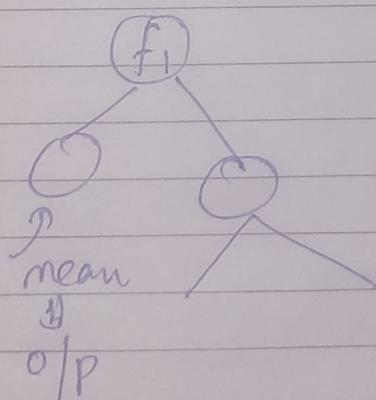
Decision tree Regressor

$f_1 \ f_2 \ O/P$

Continuous
O/P's
→ Mean

$$\text{MSE or MAE}$$

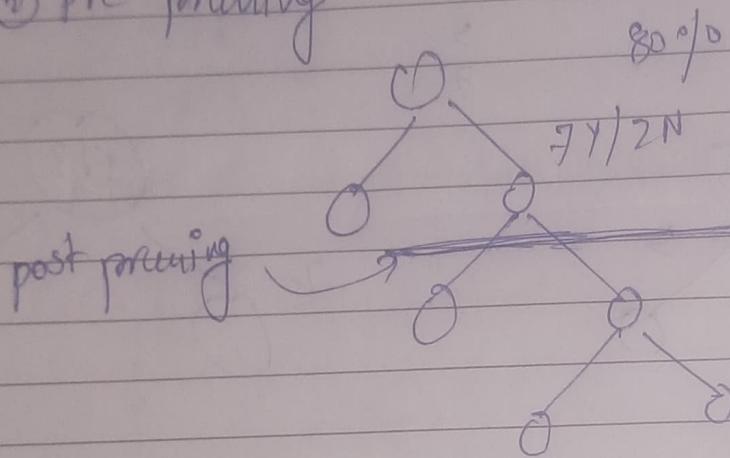
$$\frac{1}{2m} \sum_{i=1}^m (g_i - y_i)^2$$



Hyper Parameters

Decision \rightarrow Overfitting

- ① Post pruning
- ② pre pruning



prepruning

↳ hyper parameter

Grid search - CV

max-depth, max. leaf

Agenda

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

Simple Example

Previous Session

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

Naive Baye's Intuition & ClassificationRolling a Dice
 $\{1, 2, 3, 4, 5, 6\}$

BAYE'S THEOREM

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

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

Dependent Event0 - Red
1 - Green

0	0	0
0	0	

$$P(R) = \frac{3}{5} \rightarrow \text{Red}$$

Green Marble

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

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

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

$P(A \text{ and } B) = P(B \text{ and } A)$?? Yes ✓

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

$$\checkmark \quad \boxed{P\left(\frac{B}{A}\right) = \frac{P(B) * P(A|B)}{P(A)}}$$

→ Baye's theorem

$$\underbrace{x_1 \ x_2 \ x_3 \ x_4 \ \dots \ x_n}_{\text{Independent feature}} \quad \checkmark$$

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

$$P(x_1, x_2, x_3, \dots, x_n)$$

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

$$P(x_1) P(x_2) \dots P(x_n)$$

DATA [A SP]

$$n_1 \ n_2 \ n_3 \ n_4 \quad y$$

Yes

No

$$P(y=yes | no) = P(yes) * P(n_1/yes) * P(n_2/yes)$$

$$* P(n_3/yes) * P(n_4/yes)$$

$$P(n_1) * P(n_2) * P(n_3) * P(n_4)$$

$$P(y = \text{No} | x_i) = P(\text{No}) + P(x_1 | \text{No}) \times P(x_2 | \text{No}) + P(x_3 | \text{No})$$

$\approx P(x_4 | \text{No})$

$$P(x_1) \approx P(x_2) \approx P(x_3) \approx P(x_4)$$

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

$$P(y = \text{Yes} | x_i) = 0.13 \quad P(\text{No} | x_i) = 0.05$$

\downarrow $\geq 0.5 \rightarrow 1$
 $< 0.5 \rightarrow 0$
 Normalization

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

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

DATASES

Day	outlook	Temp.	Humidity	Wind	Play
D1	Sunny	Hot	High	No	
D2	Sunny				No
D3	Rain	Cool			Yes
D4	?	?	?		
D5	?	?	?		

I/P
Outlook

O/P
Play Tennis

Sunny
Overcast
Cloudy
Rain

Outlook

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		

Temperature 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
Total 9 5

PLAY

Yes 9
No 5
Total 14

P(Yes)
[9/14]

P(No)
[5/14]

$\rightarrow \text{Test}(\text{sunny}, \text{hot}) \rightarrow 0/P$

$$P(\text{Yes} | (\text{sunny}, \text{hot})) \rightarrow P(\text{Yes}) \times P(\frac{\text{sunny}}{\text{Yes}}) \times P(\frac{\text{hot}}{\text{Yes}})$$

Constant $\rightarrow P(\text{sunny}) \times P(\text{not})$

$$= \frac{9}{14} \times \frac{2}{9} \times \frac{2}{9}$$

(Normalize)

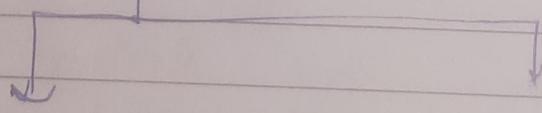
$$P(\text{Yes} | (\text{sunny}, \text{hot})) = 2/63 = 0.03 \rightarrow 0.27$$

$$P(\text{No} | (\text{sunny}, \text{hot})) = \frac{3}{35} = 0.085 \rightarrow 0.73$$

(Sunny, hot) \rightarrow Yes or No

Answer \rightarrow No

KNN Algorithm & K Nearest Neighbor

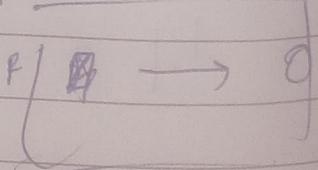
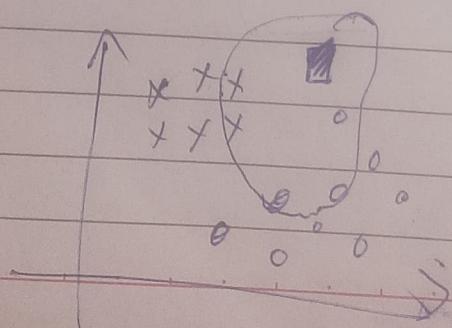


Classification

Regression

① Classification

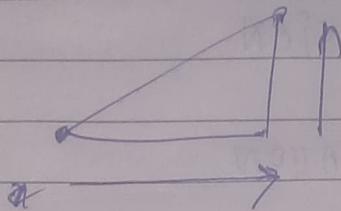
R2S



✓ Euclidean distance

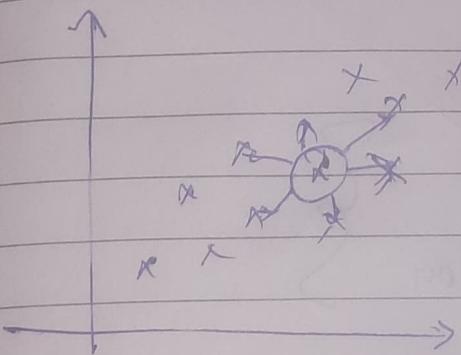
$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

✓ Manhattan distance



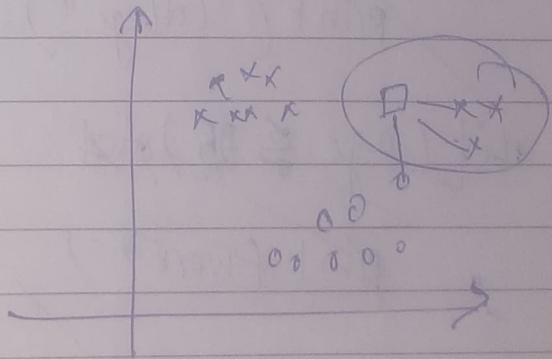
Regression

$$| \underbrace{k=5} \rightarrow \text{Hyperparameter} |$$



K - Nearest

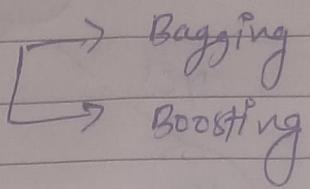
① Outliers



② Imbalanced

Agenda

(1) Ensemble Techniques



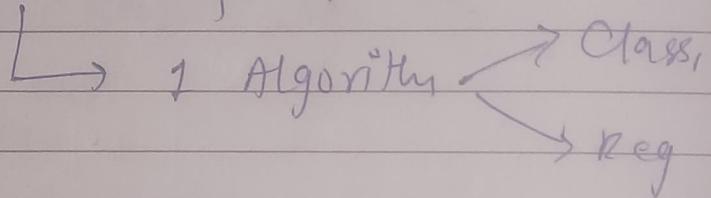
(2) Random Forest

(3) Ada Boost

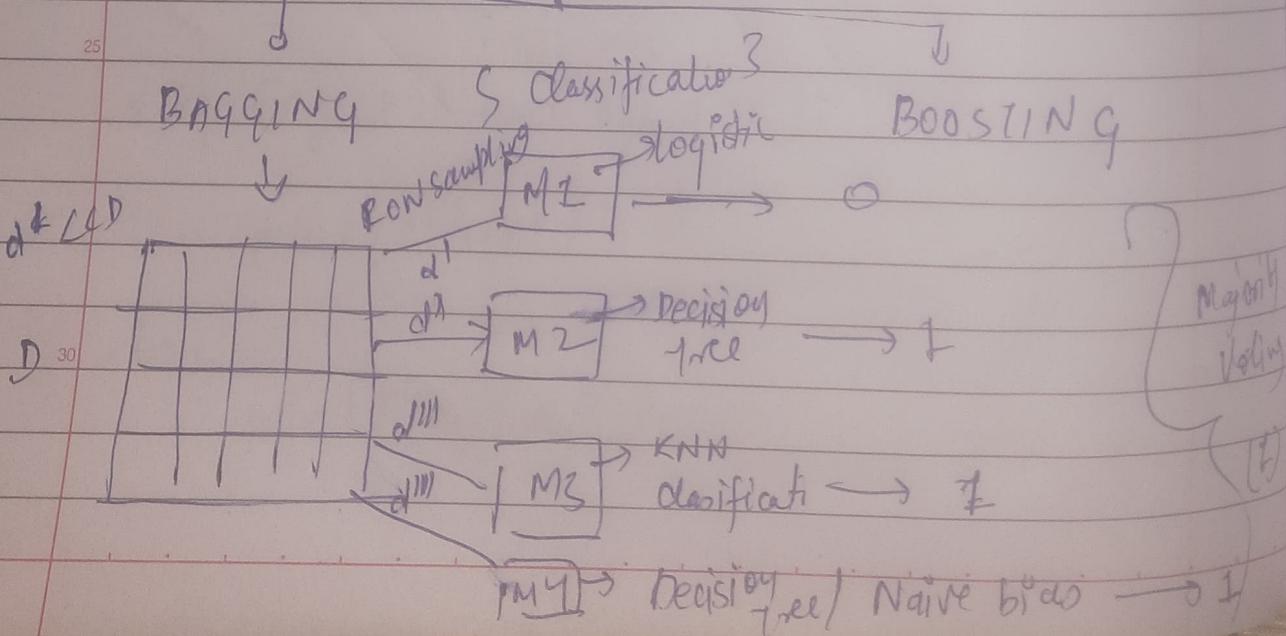
(4) Xg Boost

Ensemble Techniques

(1) Classification & Regression



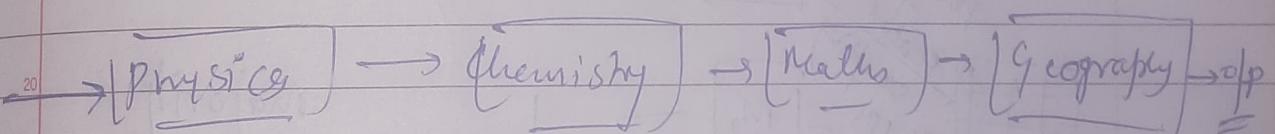
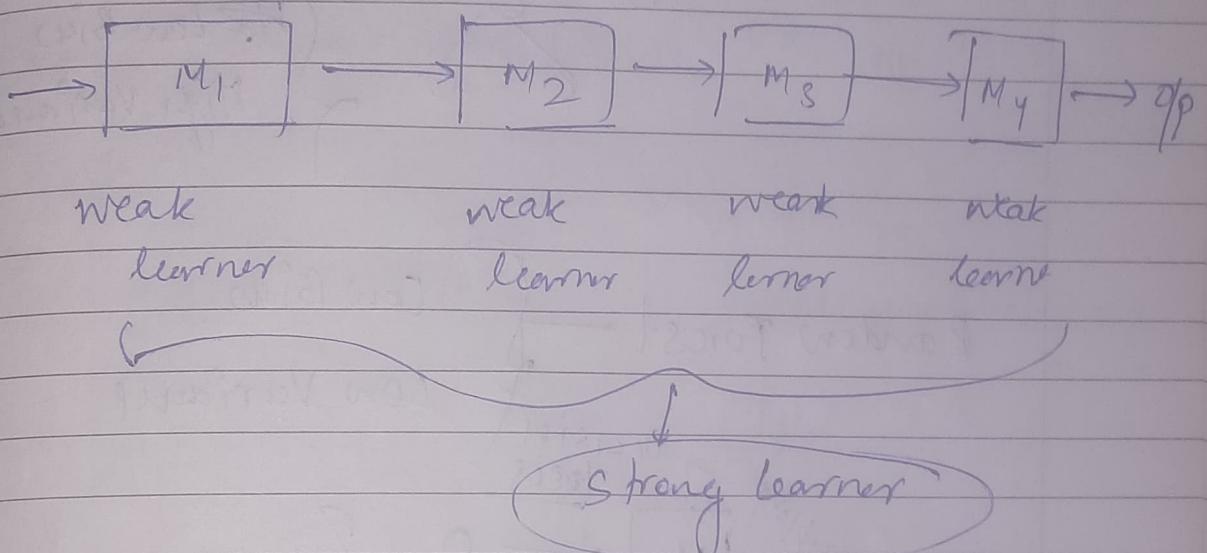
Multiple Algorithms to solve a problem ?

Ensemble Techniques

{ Bootstrap Aggregating }

Regression ? { Mean will be taken ? }

Boosting



BAGGING

- ① RANDOM FOREST CLASSIFIER

- ② Random forest Regressor

BOOSTING

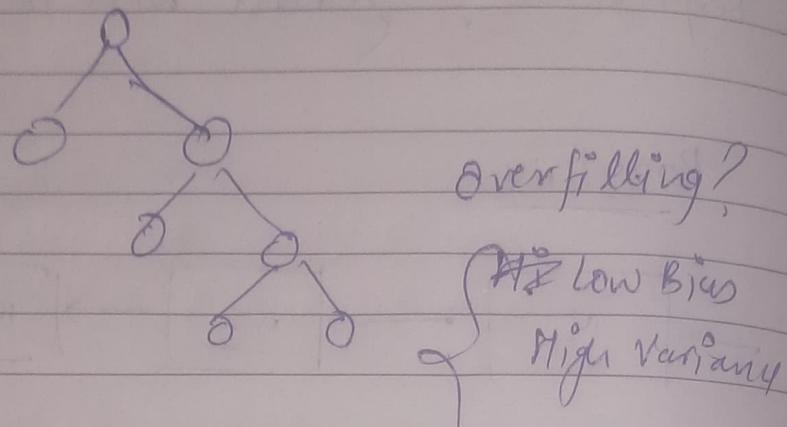
- ① Ada boost

- ② gradient boost

- ③ Xg boost

① Random Forest Classifier & Regressor

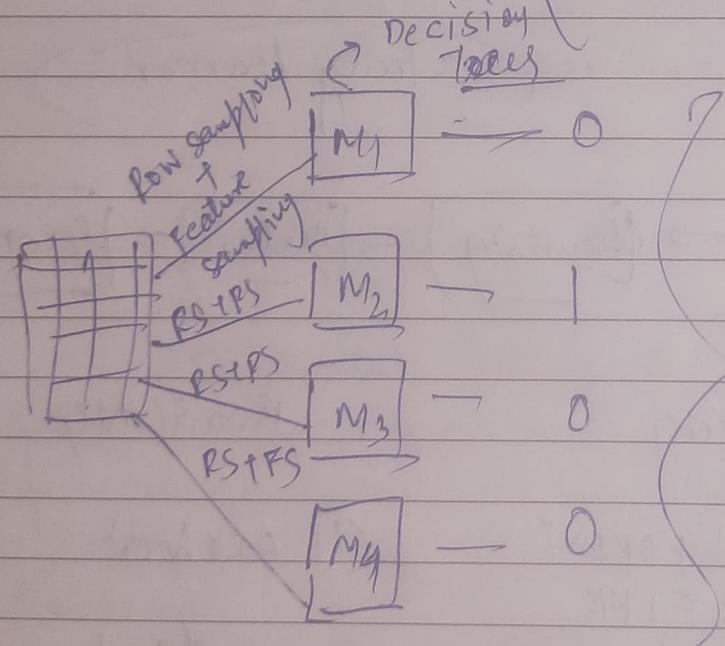
what is the main problem of DT?



Random Forest

low Bias

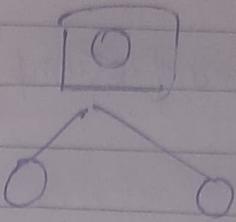
low Variance



↳ Regression

} Mean }

① Normalization?? NO



② KNN & Standardization?? Yes

↳ Euclidean, Manhattan }

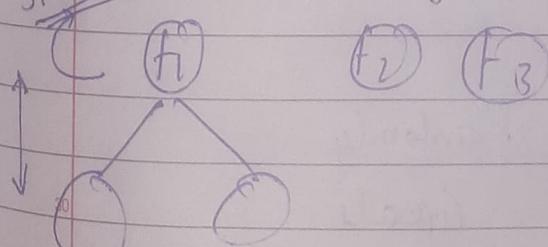
③ Random Forest → Outliers ?? NO

Boosting

1) Adaboost → Decision Tree overall weight = 1

f ₁	f ₂	f ₃	f ₄	o/p	weight
-	-	-	-	Yes	1/7
-	-	-	-	No	1/7
-	-	-	-	1	1/7
-	-	-	-	1	1/7
X	-	-	-	1	1/7
-	-	-	-	1	1/7
-	-	-	-	1	1/7

STUMPS Information gain, entropy



F₂ F₃

$$\text{total error} = \frac{1}{3}$$

↳ Weak learners }

① Performance of stump = $\frac{1}{2} \log e \left(\frac{1 - T^2}{T^2} \right)$

$$= \frac{1}{2} \log e \left(\frac{1 - 1/7}{1/7} \right) \Rightarrow 0.1895$$

$\sum_{i=1}^7 \text{total error}_i^2$

② New sample weight \Rightarrow Correct Records

$$\text{Weight} \propto e^{-P_s} = \frac{1}{7} \times e^{-0.1895}$$

$$\text{Incorrect record} = \text{Weight} \propto e^{P_s} = \frac{1}{7} \times e^{0.1895} = 0.349$$

New Weight

Is it + ??

Normalized weight

Buckets

0.05

0.07

[0-0.07]

0.05

0.07

(0.07-0.14)

0.05

0.07

(0.14-0.21)

0.349

0.537

(0.21-0.74)

0.85

0.07

(0.747-0.85)

0.05

0.07

()

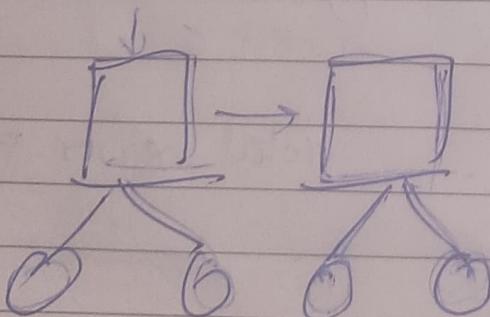
0.105

0.07

()

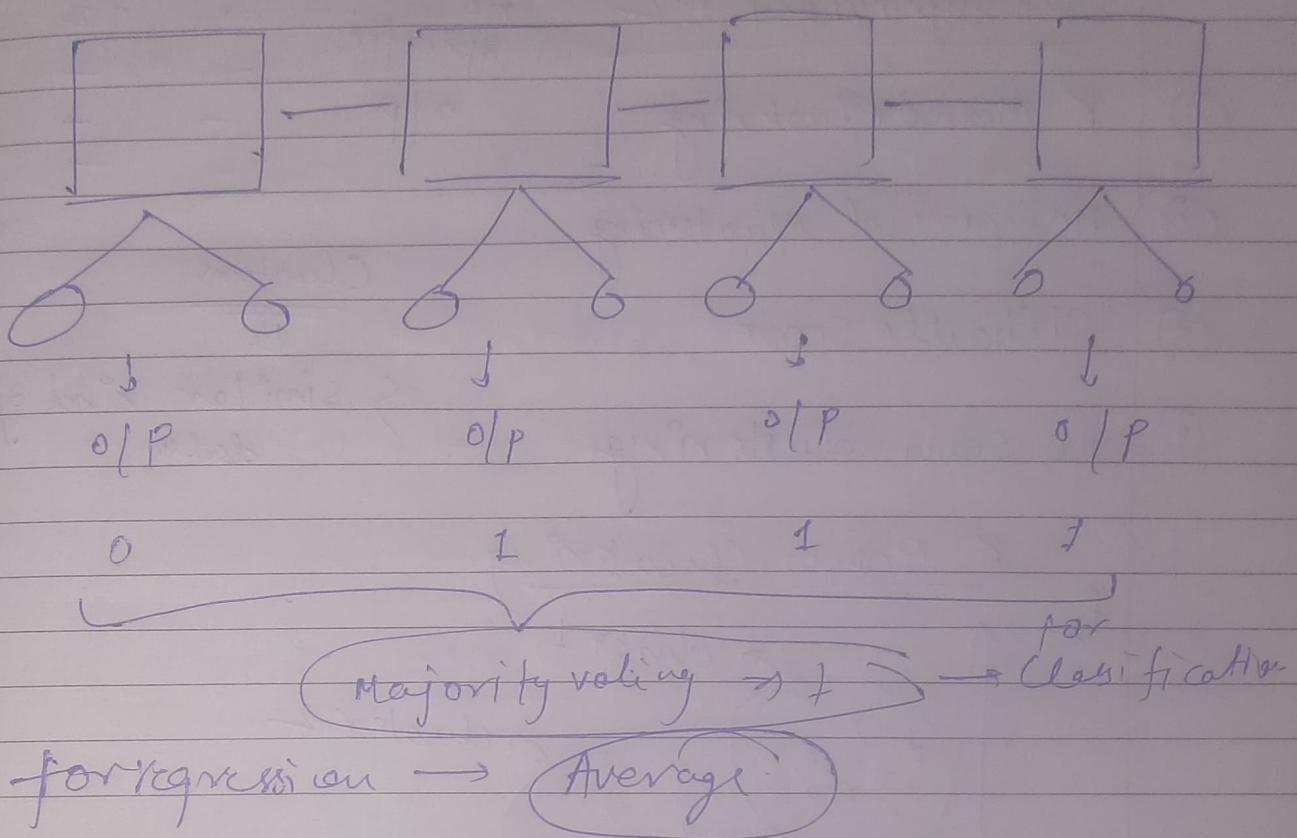
0.649

≈ 1



Randomly
Create
some number

b/w
0-1



Black box models V/s. white box Models

Linear Regression \rightarrow White box ANN \rightarrow Black Box

Random Forest \rightarrow Black box

Decision Tree \rightarrow white box

Unsupervised ML

no specific f_1 f_2

① K Means Clustering

of p

② Hierarchical clustering

Clusters

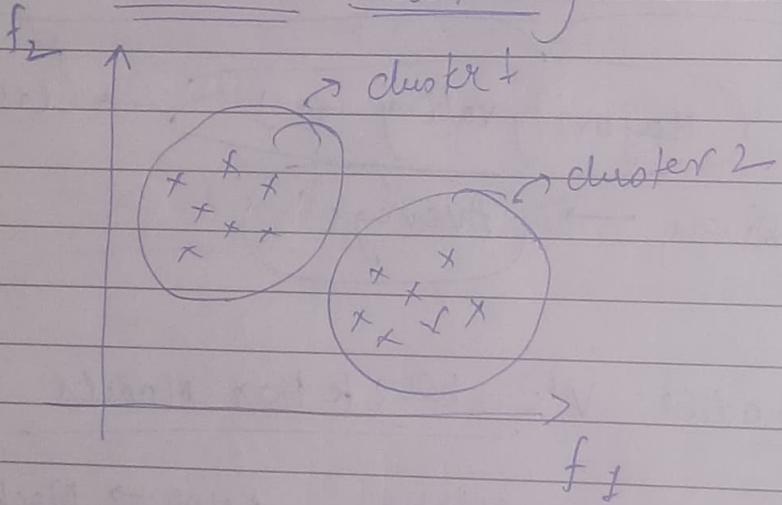
③ Silhouette score

$\frac{1}{15}$

④ DB scan clustering

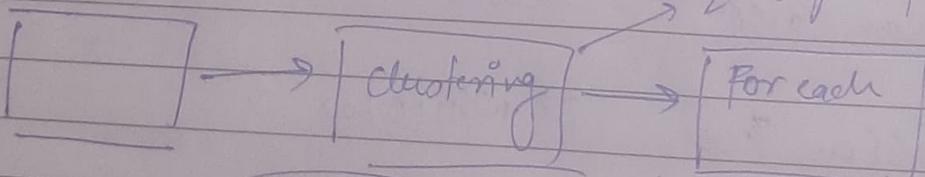
similar kind of data

K Means Clustering



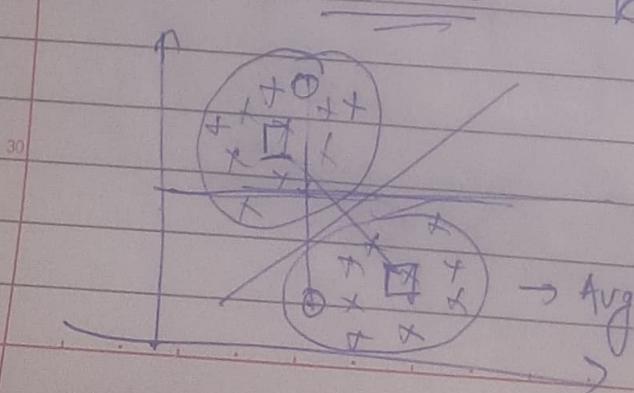
Custom Ensemble Technique

2-3 groups



K Means

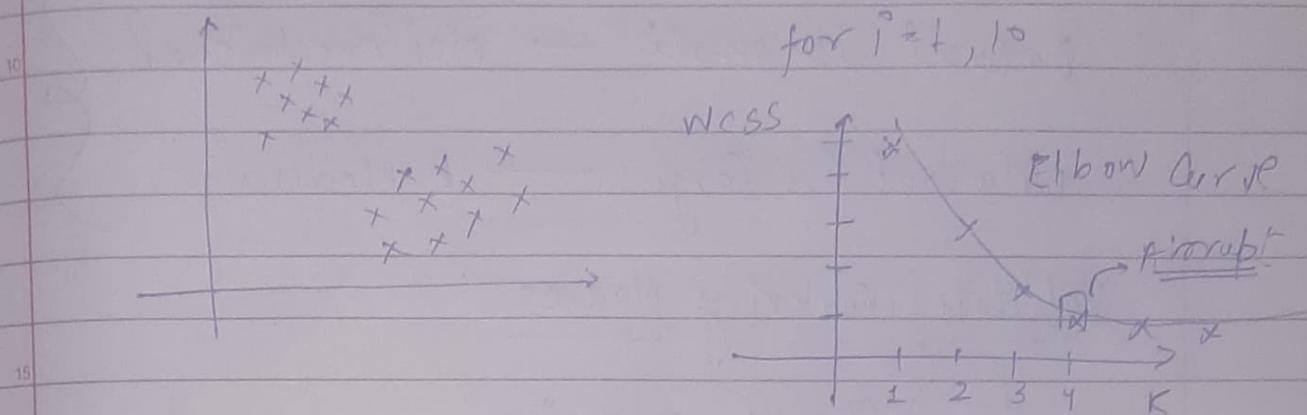
$K = 2$ centroids



High dimension

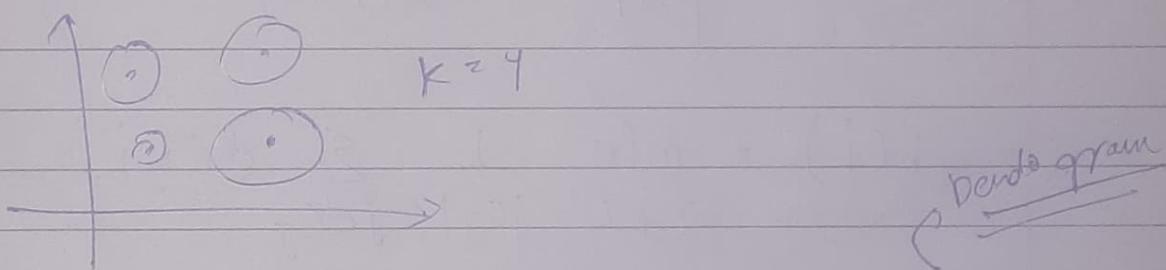
- ① we try with different K values \rightarrow suitable
- ② Initialize K number of centroids
- ③ Compute mean average to update centroid

Elbow method (K value ??)

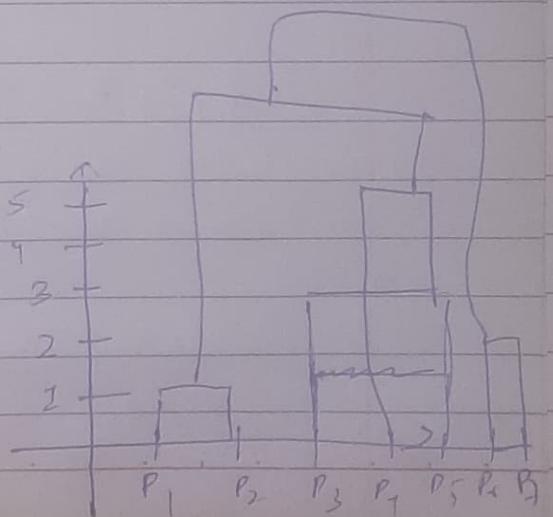
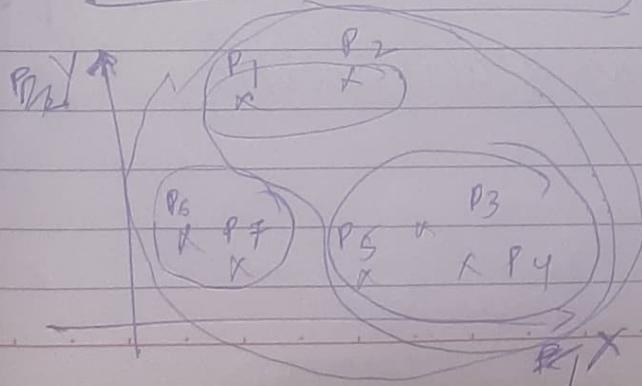


WCSS \rightarrow Within cluster sum of squares

Validating \rightarrow Silhouette score



Hierarchical Clustering



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

Q. Max line is taken by K Means or Hierarchical clustering?

Ans) Hierarchical Clustering

Dataset is small \rightarrow Hierarchical
Data set is large \rightarrow K Means

Validate Clustering Models

\rightarrow Silhouette Clustering

$$a(c_i) = \frac{1}{|C_i|-1} \sum_{j \in C_i, j \neq i} d(c_i, j) \rightarrow \text{Avg. distance}$$

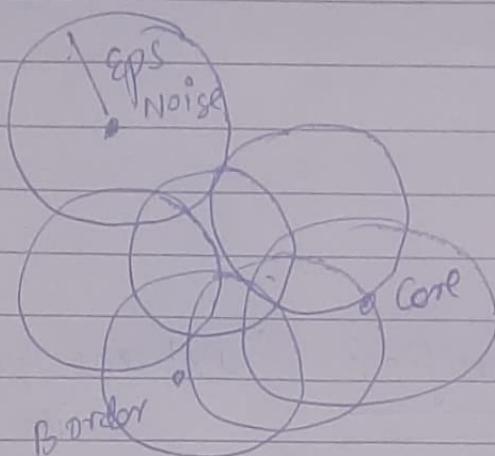
$$b(c_i) = \min_{j \neq i} \sum_{j \in C_j} d(c_i, j)$$

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

DB Scan Clustering

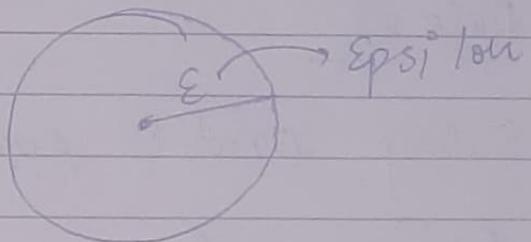
Density - Based Spatial Clustering of Application with Noise (DBSCAN)

C hyperparameter



minpts = 4 \rightarrow eps

- ① Min pts
- ② Core points
- ③ Border points
- ④ Noise points



Agenda

① Xgboost classifier

② Xgboost Regressor

③ SVM

④ SVR

Xgboost classifier (Extreme Gradient Boosting)

Data set

	Salary	Credit	Approval	Residual
15	<=80	B	0	-0.5
	<=10	g	1	0.5
20	<=250	g	0	0.5
	>50	B	0	-0.5
	>50	g	1	0.5
	>500	N	1	0.5
	<=50	N	0	-0.5

	Salary	Credit	Approval	Residual
15	<=80	B	0	-0.5
	<=10	g	1	0.5
20	<=250	g	0	0.5
	>50	B	0	-0.5
	>50	g	1	0.5
	>500	N	1	0.5
	<=50	N	0	-0.5

Box Model

 $\rightarrow \text{par} = \underline{0.5}$

① Create a Binary Decision Tree using two features

② Calculate Similarity weight (SW)

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

③ Information gain

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

Binary DT

Salary

$$SW = 0.14$$

$$c = 50$$

$$> 50$$

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

U

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

D

$$SW_U = f(-0.5)^2 f(0.5)$$

$$SW_D = 0.33$$

$$SW = (-0.5 + 0.5 + 0.5 + 0.5)$$

$$0.5(1-0.5) + 0.5(1-0.5)$$

$$+ 0.5(1-0.5) + 0.5(1-0.5)$$

$$SW = 0$$

$$\begin{aligned} \text{Information} &= 0 + 0.33 - 0.14 \\ \text{gain} &= 0.19 \end{aligned}$$

Credit

G, N

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

$$[SW = 0.33]$$

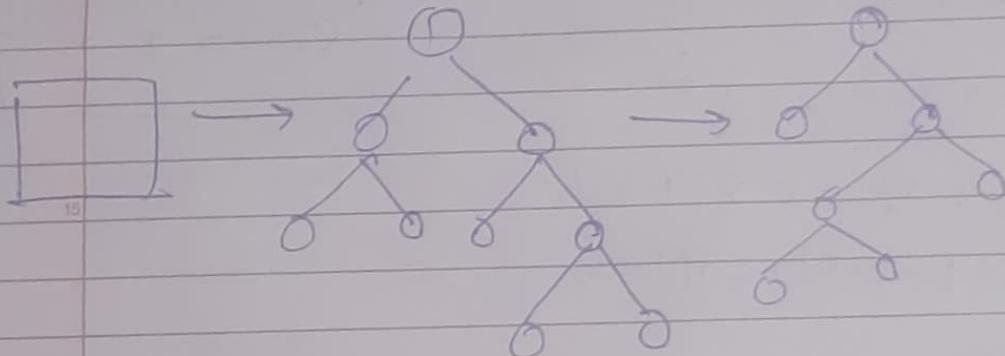
$$[SW = 1]$$

$$[I_g = 1 - 0.33 - 0 = 1.33]$$

$$O/P \rightarrow \left[\alpha_0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n) \right]$$

Δg Boost \rightarrow Black Box Model

pre-pruning



$$O/P \rightarrow \alpha_0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n)$$

$\lambda \rightarrow$ Cross Validation
(Hyperparameter)

XgBoost Regressor

Bsp	Gap	Salary	R _i	Bad Model
2	Yes	40K	-11K	↓
2.5	Yes	41K	10K	Avg. of all values
3	No	42K	0K	
9	No	60K	-9K	
4.5	Yes	62K	11K	(51K)

$(-4, 9, 1, 9, 11)$

Exp

$\text{SW} = 16$

≤ 2

> 2

(-11)

$(-9, 1, 9, 11)$

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

No. of Residuals \rightarrow

$\text{SW} = 60.55$

$\text{SW} = 28.8$

Information Gain = $60.5 + 28.8 - 1/6 = 89.13$

B_{np}

≤ 2.5

> 2.5

$(-4, 9)$

$(1, 9, 11)$

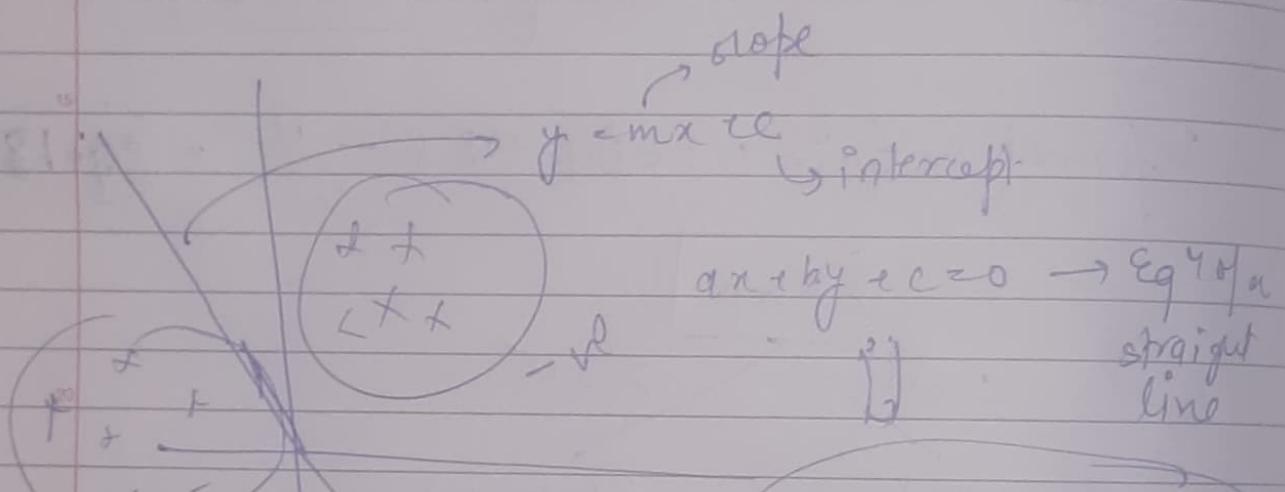
$\text{SW} = 133.33$

$\text{SW} = 110.25$

$\text{IG} = 133.33 + 110.25 - 1/6$

$SIZ \propto_1 (-10) + \propto_2 (D1_2) + \propto_3 (D1_3) - \dots$

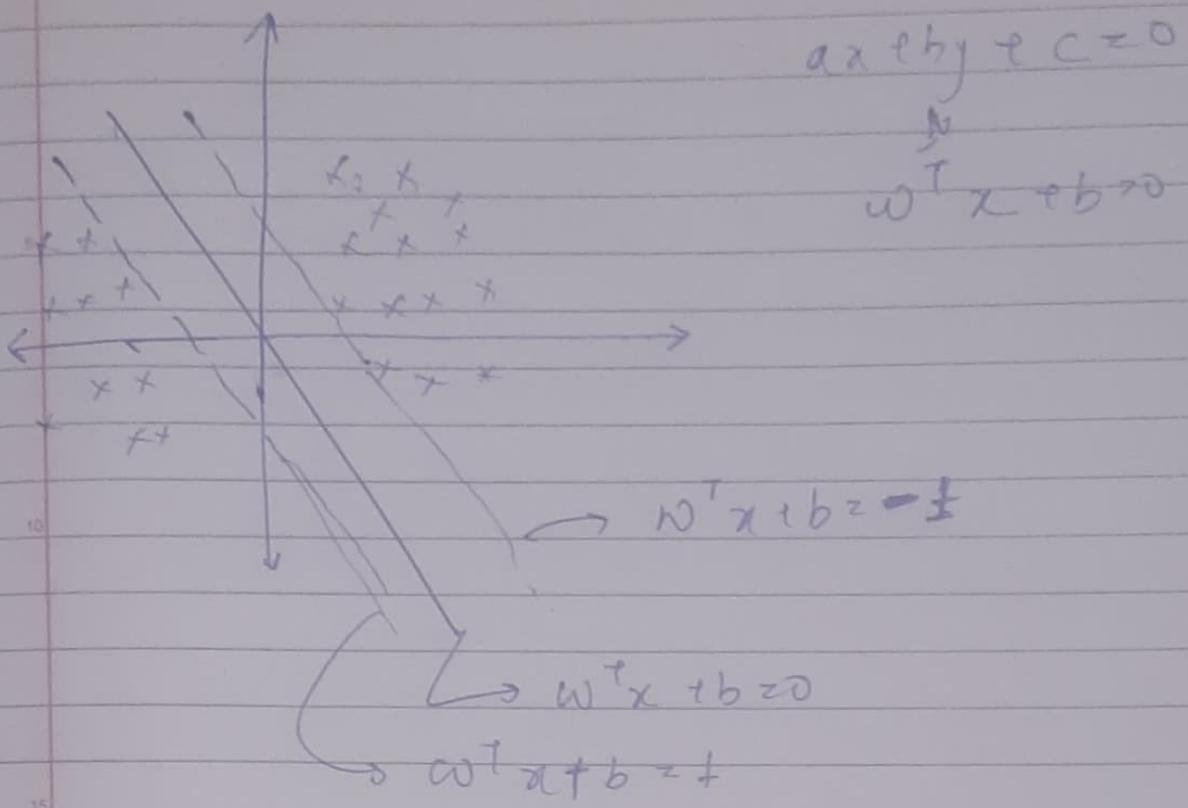
$d_n (D1_n)$

SVM


$$G_m = \frac{-a}{b}, \quad C^2 = \frac{-c}{b}$$

$$y = m x + C$$

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



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

$w^T x_2 + b = -1$

20. $\frac{w^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$

25. ~~Maximize~~ Maximise $\frac{2}{\|w\|}$
(w, b)

30. $y_i \left\{ \begin{array}{l} 1 \quad w^T x + b \geq 1 \\ -1 \quad w^T x + b \leq -1 \end{array} \right.$