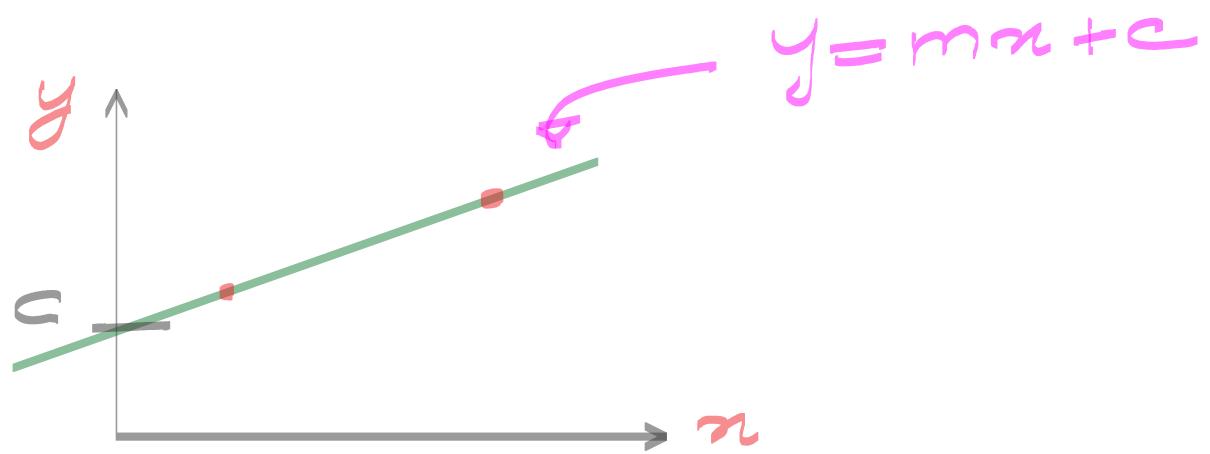


# \* Equation of Line

$$y = mx + c$$

$$\hat{y} = a + bx$$

$$\hat{y} = \theta_0 + \theta_1 x$$



$m$  → slope of line  
 $c$  → intercept

$$\hat{y} = a + bx$$

If  $x = 2$ ,  $\hat{y} = 13$   
 $x = 5$ ,  $\hat{y} = 28$

## \* Types of LR

1. Simple LR

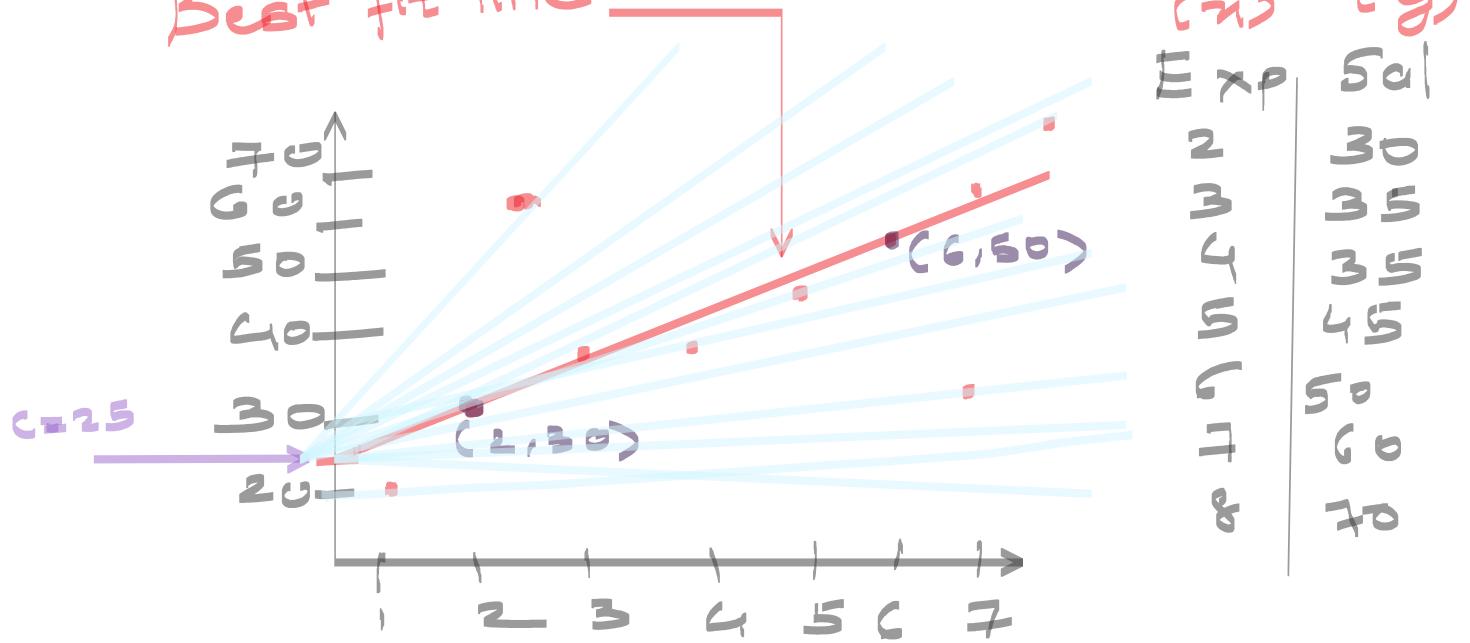
1. Simple LR

$$y = mx + c$$

2. Multiple LR

$$y = m_1x_1 + m_2x_2 + \dots + m_nx_n + c$$

Best fit line



$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{50 - 30}{6 - 2} = \frac{20}{4} = 5$$

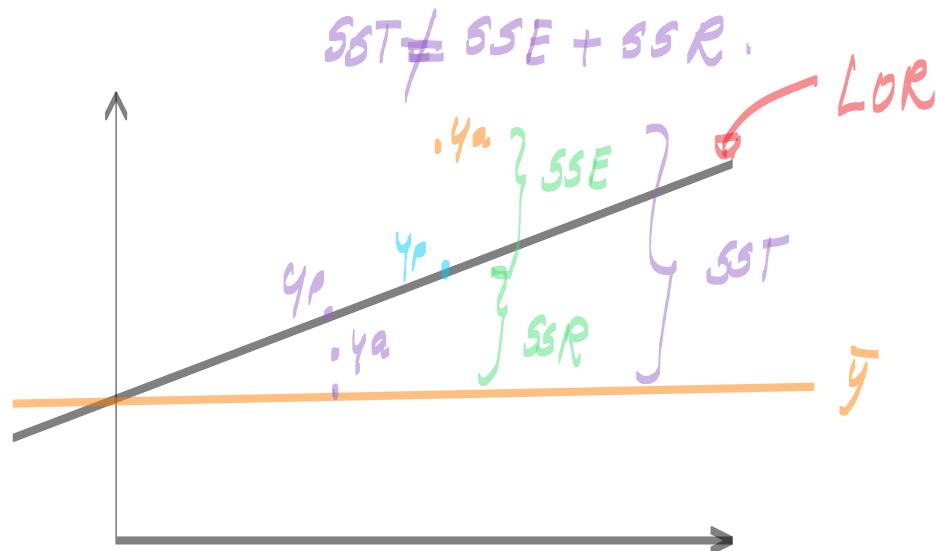
$$y = 5x + 25$$

Linear Model

\* Blue lines → infinite number of possibilities

## Linear Regression

08 December 2022 07:31



$MSE \rightarrow$  Scale Variant

Ht.	Temp.	MSE
5.1	150	
5.2	160	
5.4	170	
6.0	183	
6.1	185	

$25.36$

$28900$

$R^2$ -squared  $\rightarrow R^2 \rightarrow$  Scale Invariant

0 to 1, -ve

$$1. \quad R^2 = 1$$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$= 1 - \frac{0}{SST} = 1$$

$$2. \quad R^2 = 0$$

$$R^2 = 1 - n \quad SSE = SST$$

$$2. R^2 = 0 \quad R^2 = 1 - 0 \quad , \quad SSE = SST$$

$$= \frac{1}{1}$$

$$3. R^2 = +ve , \quad SSE < SST$$

$$R^2 = 1 - 0.1 = \underline{\underline{0.9}}$$

$$4. R^2 = -ve \quad , \quad SSE > SST$$

$$R^2 = 1 - 1.4 = -0.4$$

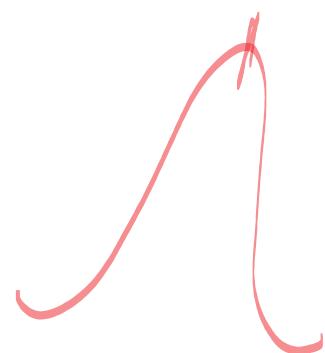
$$\underline{\underline{-0.4}}$$

Non-linear.

$$R^2 = 1 - \frac{SSE}{SST} \rightarrow \frac{(y_a - y_p)^2}{(y_a - \bar{y})^2}$$

Exp	$y_a$
2	30
3	40
4	50
5	60
6	70

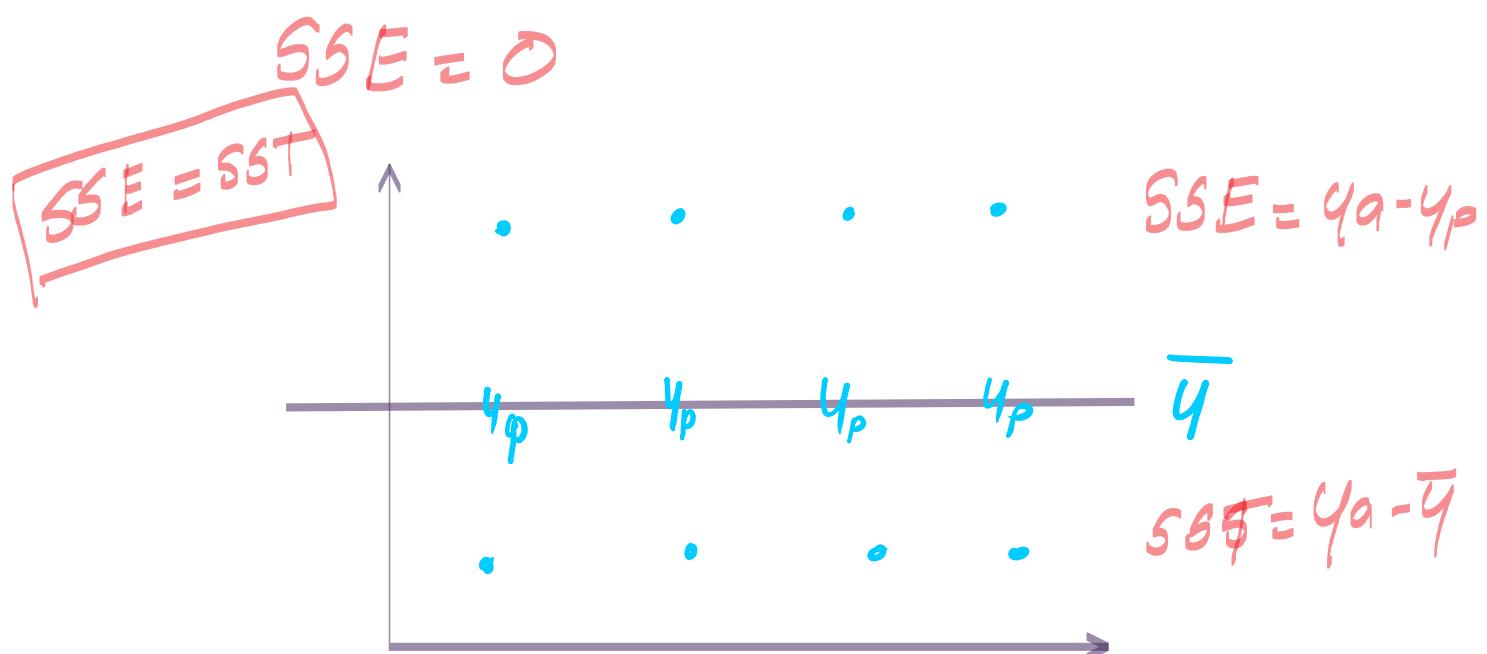
$\left\{ y_a , \bar{y} = 50 \right.$



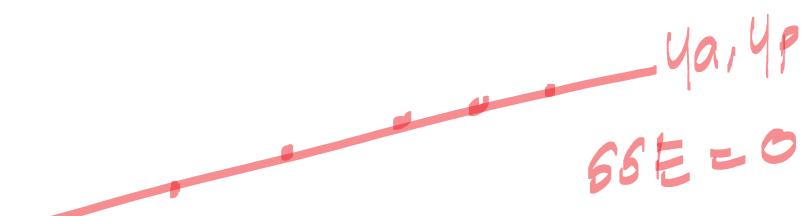
$$R^2 = 1 - \frac{SSE}{SST}$$

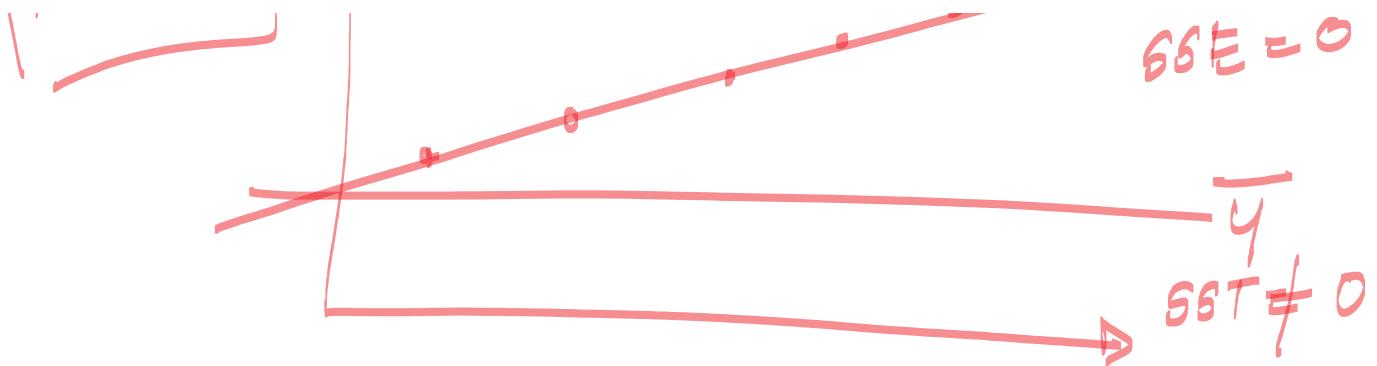
$$R^2 = 1 - \frac{\text{Unexplained Variation}}{\text{Total Variation}}$$

$$R^2 = \frac{SST - SSE}{SST} \rightarrow \frac{\text{Explained Var}}{\text{Total Var}}$$



$$R^2 = 1$$





R-squared

Coeff. of Corr.

$$R^2 = R \times R \rightarrow SLR$$

$$R^2 \neq R \times R \rightarrow MLR$$

$R^2 \rightarrow$  Goodness of Best fit line  
 (R<sub>2</sub>-score)

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$y$
0.9	0.4	0.8	0.7	-0.5	

3 features  $\rightarrow R^2 = 0.85$  } 0.852  
 $(x_1, x_3, x_4)$

4 features  $\rightarrow R^2 = 0.86$  }  
 $(x_1, x_2, x_3, x_4)$

Adjusted R-squared  $\rightarrow \bar{R}^2$

$$\bar{R}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

$n$  = No of samples (rows)

$p$  = No of Predictors

$$R^2 \geq \bar{R}^2$$

$$R^2 \quad \bar{R}^2$$

$$0.85 \quad 0.85$$

$$0.84 \quad \downarrow$$

①  $R = 0.3$

$$0.86$$

$$0.84 \quad \downarrow$$

②  $R = 0.9$

$$0.87$$

$$0.86$$

$$R^2 \approx 1 - \frac{SSE}{SST}$$

$$\bar{R}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

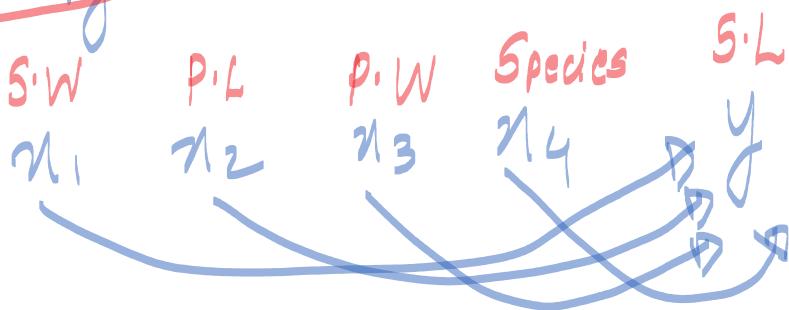
← 999  
← 994

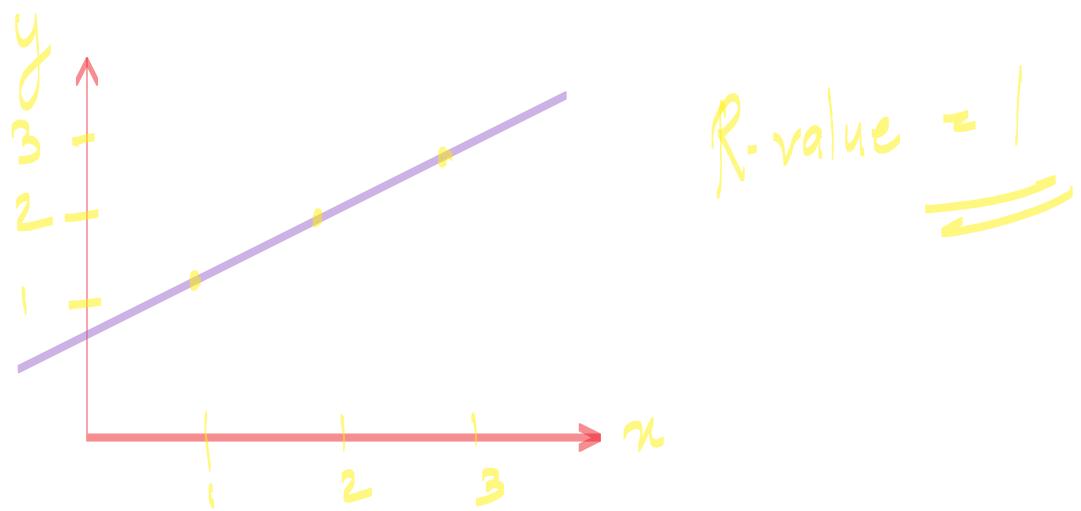
$$n = 1000, p = 5$$

$$\boxed{n > p}$$

5 rows  
6 columns X

\* Linearity



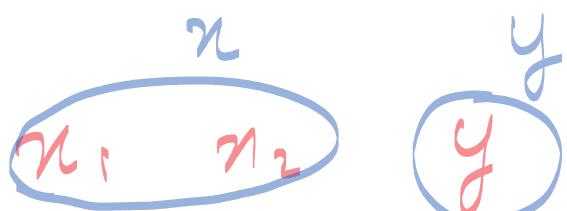


150 rows  
 $\text{df} \rightarrow n_1 \ n_2 \ n_3 \ n_4 \ y$

Tr  
Test  $\xrightarrow{\leftarrow} n \rightarrow Df \rightarrow n_1 \ n_2 \ n_3 \ n_4$   
 Tr  
Test  $\xrightarrow{\leftarrow} y \rightarrow \text{series} \rightarrow y$

Train  
 $n_{\text{-train}} \quad y_{\text{-train}}$   
 120 rows 120 rows

Test  
 $n_{\text{-test}} \quad y_{\text{-test}}$   
 30 rows 30 rows



Tr  $\leftarrow 1$   
 Tr  $\leftarrow 4$   
 Tr  $\leftarrow 5$

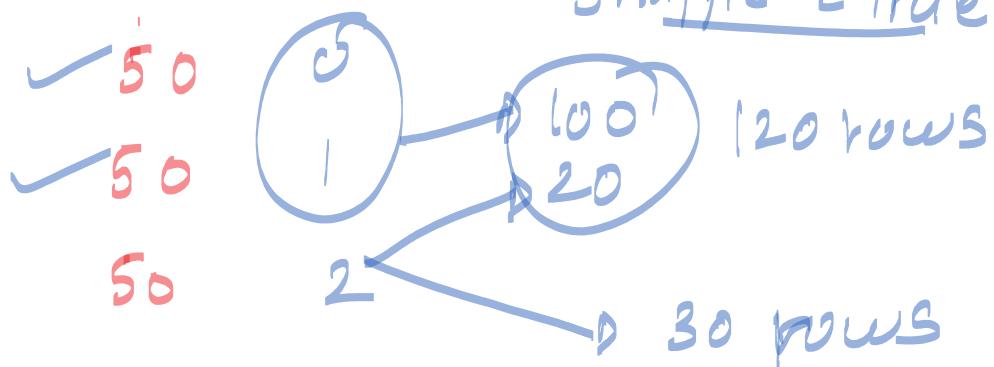
at Line 9.6.5

Tr ✓  
Tr ✓  
Tr ✓

✓ Train → 3 row  
✓ Test → 2 row

u-train, 2, 4, 5  
y-train, 2, 4, 5  
u-test 1, 3  
y-test 1, 3

shuffle = True



100 row      80 → Train ✓

20 → Test

\* What is Data Science ?

Train 4 lines → 3 correct 75%  
evaluate 1 incorrect 25%

\* What is KNN ?

Test 4 line → 2 correct 50%

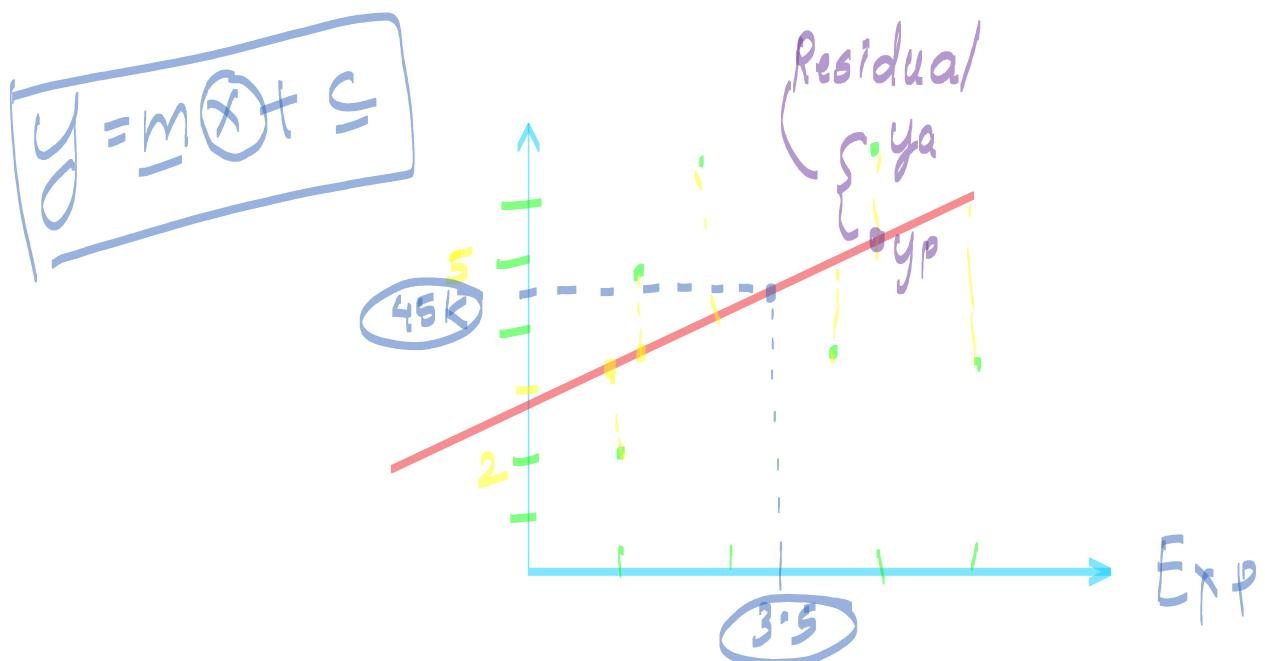
\* What is R<sup>2</sup>  $\rightarrow$  2 covered 50%

Test - 4 line → 2 covered

Evaluate

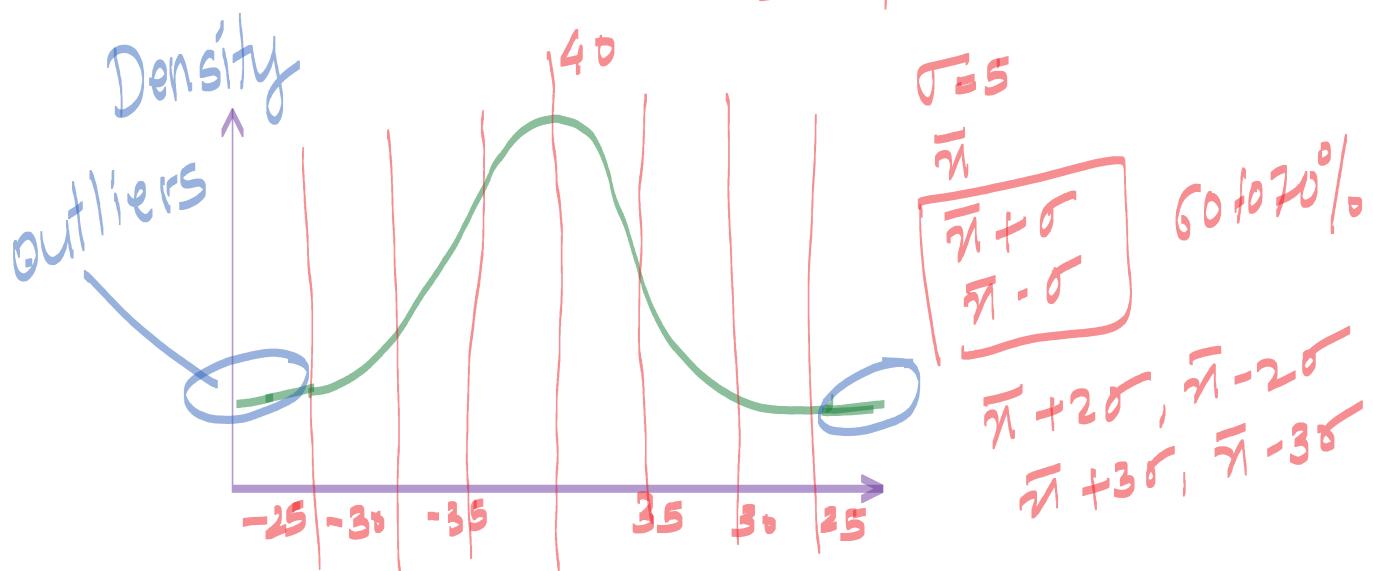
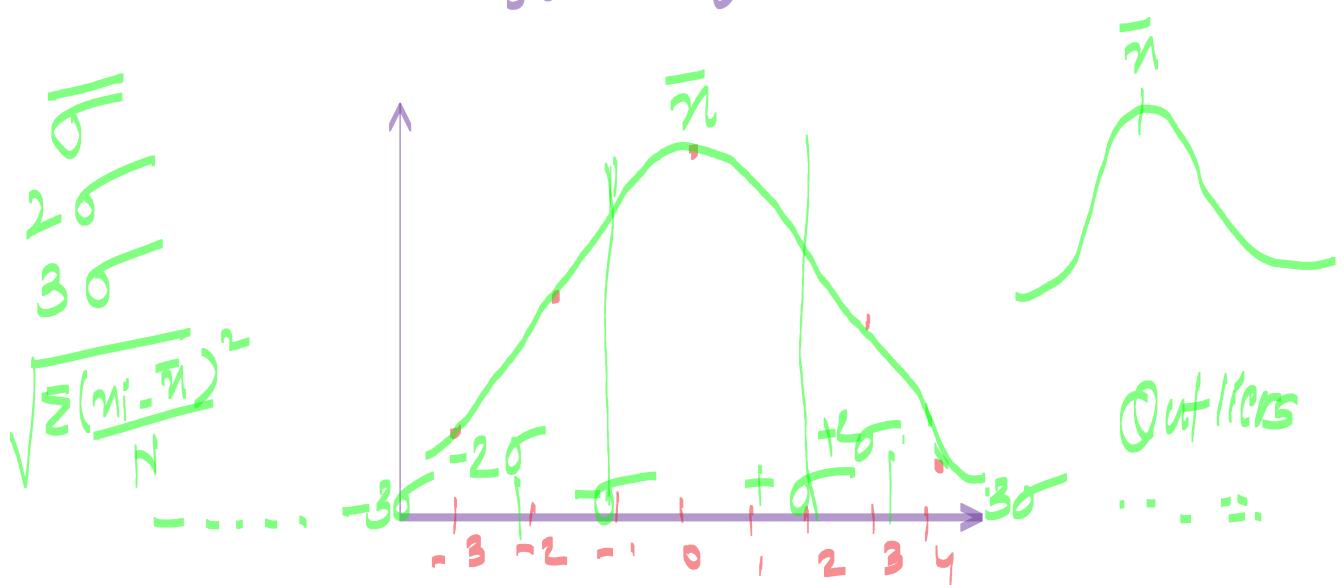
unseen

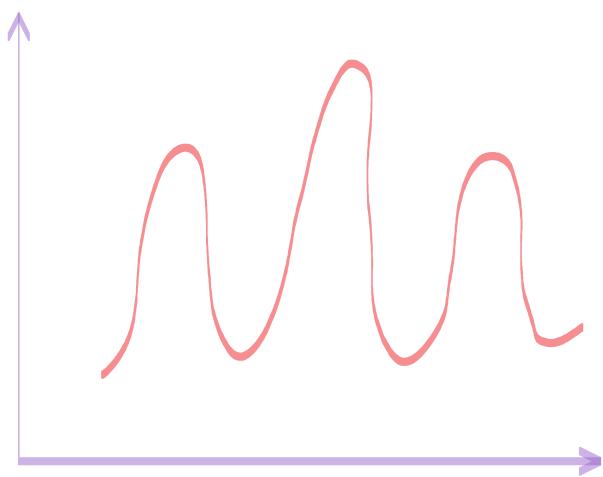
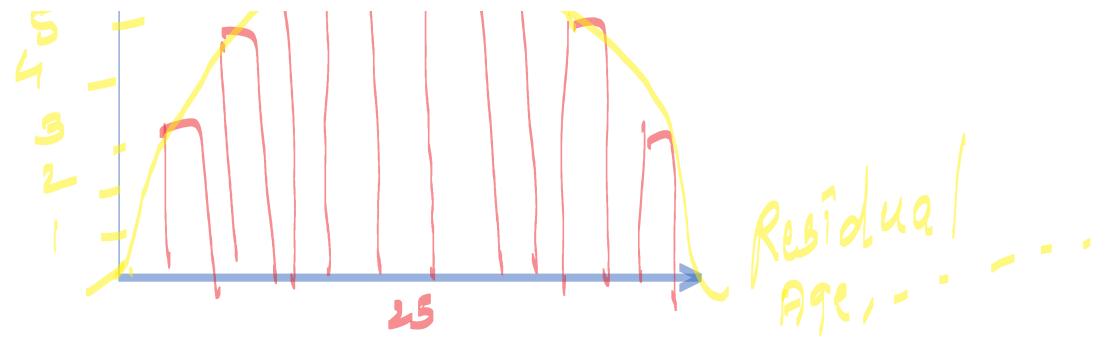
	Area	Bed	Bath	Location	Price
Train	1000	3	2	Baner	2 cr
	700	2	1	Walkad	80
	800	1	1	Bule .	75
Test Unseen	650	2	2	Baner	90 L

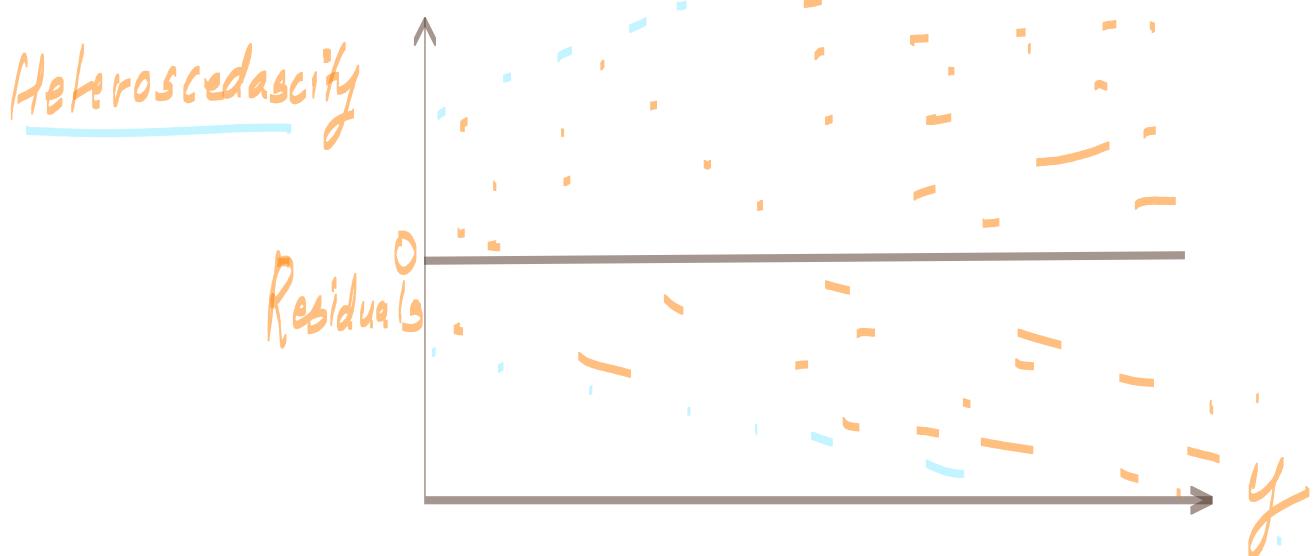
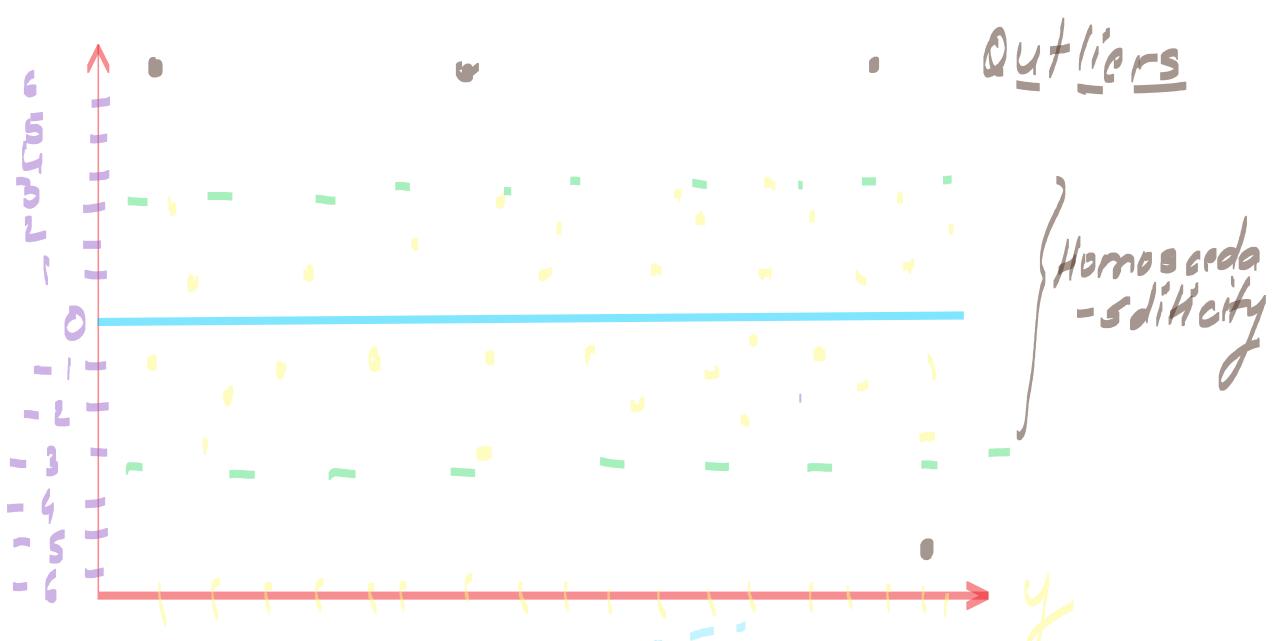
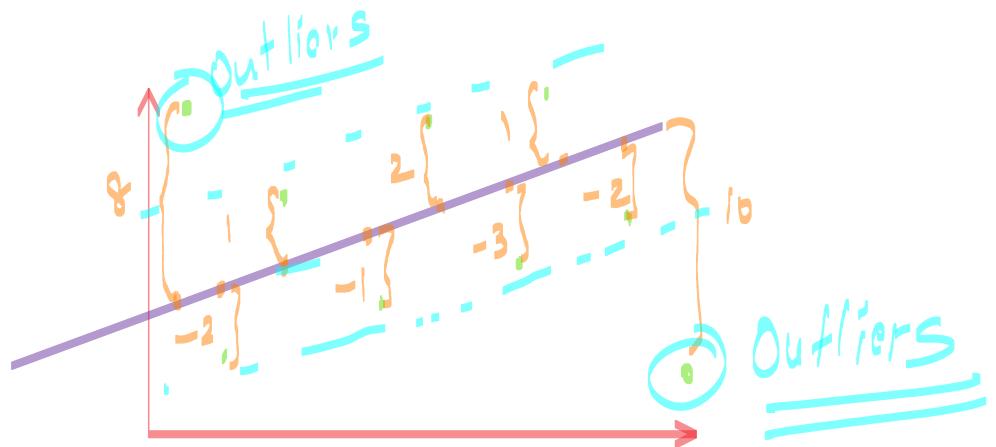


$x_1$	$y_a$	$y_p$	Residual
	10	12	-2

$X_1$	$y_a$	$y_p$	Residual
10	12	-2	
20	17	3	
30	34	-4	
25	28	-3	
15	12	3	
45	41	4	
50	50	0	







MSE → square (Differential Term)  
 $(y_a - y_p)^2$

$$\partial(\mathcal{L}) \rightarrow (\text{MSE})$$

$$\text{MAE} \rightarrow |y_a - y_p|$$

Outliers  
100, 2, 3, 4, 5, 6

$$\frac{\partial C(F)}{\partial m, c} \rightarrow (\text{MSE})$$

$$\text{MSE} \rightarrow \begin{cases} 16 \\ 60 \end{cases} \rightarrow \text{MSE}$$

MSE is less robust  $\rightarrow$  highly outliers

4	9	16	$\dots$	49	$\rightarrow$ MSE
2	3	4	$\dots$	7	$\rightarrow$ MAE

{MAE}  $\rightarrow$   
 {RMSE}  $\rightarrow$

Price  $\rightarrow$  X

✓ Train       $R^2 \rightarrow 0.9$

X Test       $R^2 \rightarrow 0.6$

} X

## 1. Simple LR

 $x \quad y$ 

$$y = mx + c$$

G.D.R

## 2. Multiple LR

 $x_1 \quad x_2 \quad x_3 \quad x_4 \quad y$ 

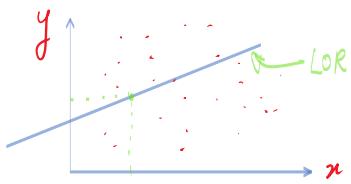
$$y = m_1x_1 + m_2x_2 + \dots + m_nx_n + c$$

$$m = 10, c = 2$$

$$y = 10x + 2$$

$$x = 2, y_p = ?$$

$$y_p = 10 \times 2 + 2 = 22$$

1. Regression  $\rightarrow$  Target (y) continuous.2. Classification  $\rightarrow$  Categorical

Independent	Target (dependent)
$x_1 \quad x_2 \quad x_3 \quad x_4$	y
Con Con Cat Cat	Con (Req)
Con cat Con cat	Cat (Class)

## \* Assumptions of LR.

## 1. Linearity

Linear Rel.  $x$  &  $y$ .Correlation  $\rightarrow$  Coeff or R value

df.corr()
-----------

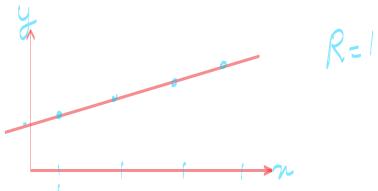
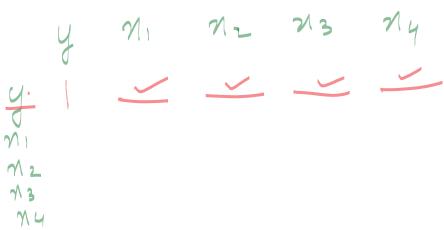
Positive Corr  
Good Predictors  $R \geq 0.7$ Negative Corr  
Good Predictors  $R < -0.7$  $n = n^2$

## Good Predictors

No corr -0.3 R 0.3

$R = 1 \rightarrow$  strong +ve  
 $R = -1 \rightarrow$  strong -ve  
 $R = 0$

$$R = \frac{\text{Cov}}{\sigma_x \sigma_y}$$



## 2. No Multicollinearity



$$\begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix} \quad \left. \quad \right\} VIF$$

1 to inf

$x_1$ , VIF = 1, No corr

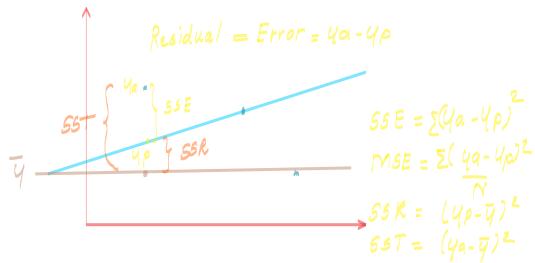
VIF = 1 to 5, Moderate

VIF > 5, Highly

VIF > 10, Significantly.



Linearity  $x_1, x_2 \rightarrow y$   
 No Multico



$R^2 \rightarrow R$  squared (Coeff. of Determination)

$$R^2 = 1 - \frac{SSE}{SST}$$

$$\boxed{SST = SSE + SSR}$$

$$= \frac{SST - SSE}{SST} = \frac{SSR}{SST}$$

$R^2 \rightarrow 0$  to 1

$R^2 \rightarrow -ve$  (Non linear)

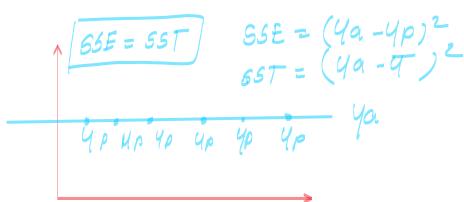
$$1. R^2 = 1, SSE = 0$$

$$R^2 = 1 - \frac{SSE}{SST} = 1$$



$$2. R^2 = 0 \quad SSE = SST$$

$$R^2 = 1 - \frac{SSE}{SST}$$



$$3) R^2 = +ve \quad SSE < SST$$

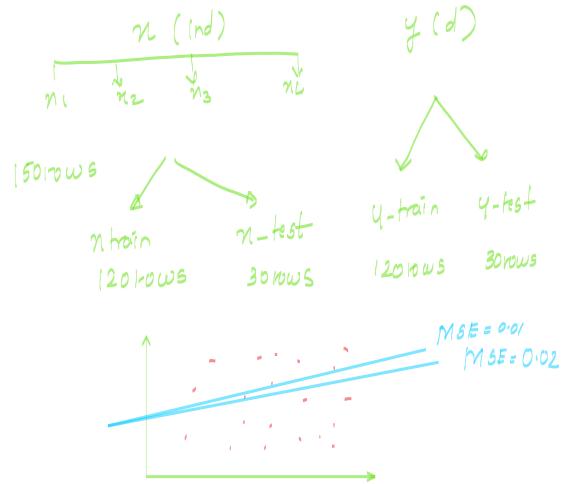
$$R^2 = 1 - \frac{SSE}{SST}$$

$$4) R^2 = -ve, \boxed{SSE > SST}$$

Non Linear

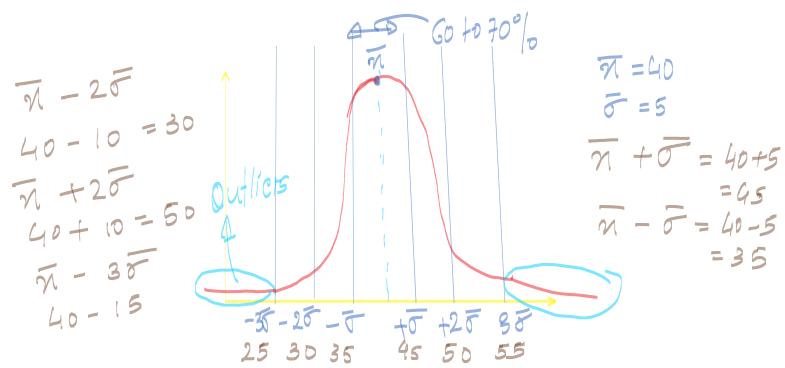
Species	$x_1$	$y$
0	0	-
1	1	-
2	2	-

shuffle = True



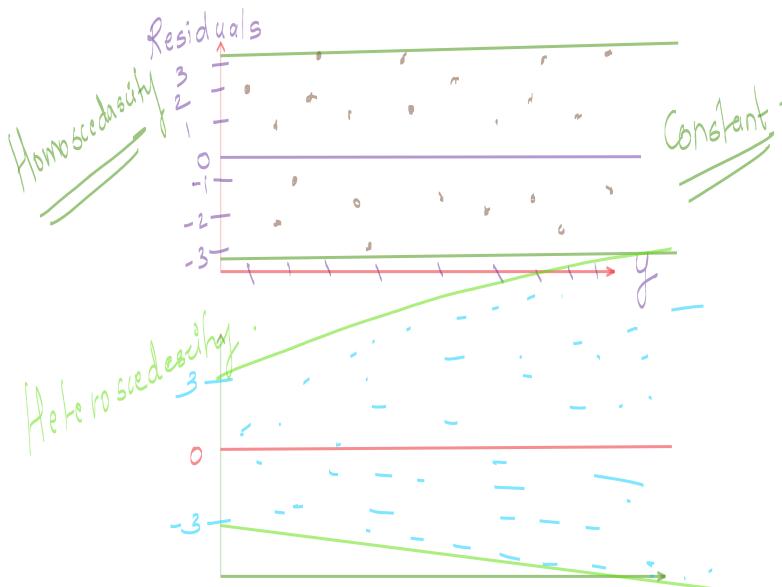
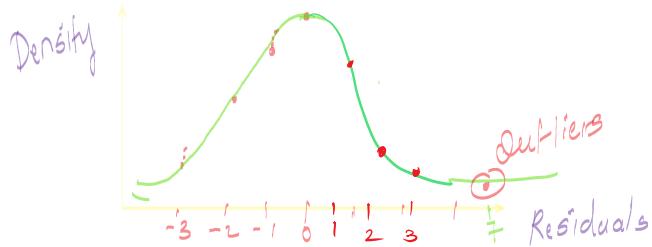
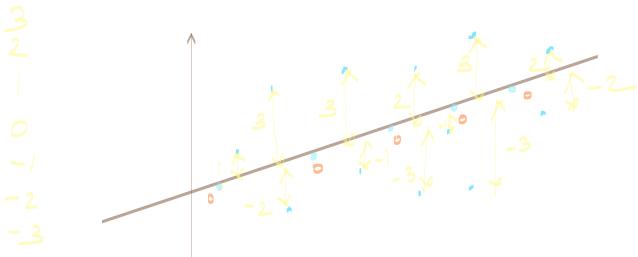
1. What is Data Science? (Train)
- 5 line → 4 lines → 80%
2. What is kNN?
- 5 line → 2 line → Test (Unseen)

	$x_1$	$x_2$	$x_3$	$x_4$	$y_a$	$y_p$
Train	1	2	3	4	20	22
Test	2	3	4	5	24	25



$$\begin{aligned}\bar{x} &= 40 \\ \bar{\sigma} &= 5 \\ \bar{x} + \bar{\sigma} &= 40 + 5 \\ &= 45 \\ \bar{x} - \bar{\sigma} &= 40 - 5 \\ &= 35\end{aligned}$$





(CF)  $MSE \rightarrow$  Differential  $f^2(y_a - y_p)$

$MAE \rightarrow |y_a - y_p|$  cm, Price variant

$\sqrt{RMSE} \rightarrow \sqrt{MSE}$

$MSE$

length cm  $\frac{\text{cm}^2}{\text{cm}^2}$   
price Rs  $\frac{\text{Rs}^2}{\text{Rs}^2}$

$\checkmark MAE$

$MSE$

2 3 ... 5 25 21 100

$4 \quad 9 \quad \rightarrow 3 \rightarrow 10$

less robust  
(Highly sensitive)

\* Outliers

C.F  $\rightarrow$  MSE

$\frac{2}{2}$

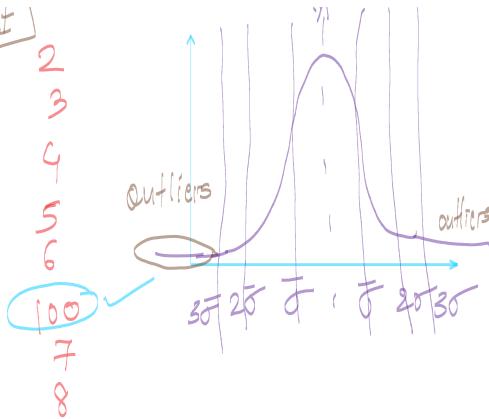


## \* Outliers

$$\text{Mean without-out} = 4.5$$

$$\text{Mean with-out} = 50$$

$\rightarrow \text{C.F} \rightarrow \text{MST}$



## \* Missing

Continuous  $\rightarrow$  mean ✓

Continuous  $\rightarrow$  median

categorical  $\rightarrow$  mode .

$$R^2 = 1 - \frac{SSE}{SST}$$

$\overline{R^2}$   $\rightarrow$  Adjusted  $R^2$

$$\checkmark n_1 \rightarrow R = 0.85$$

$$\checkmark n_2 \rightarrow R = -0.3 \quad \checkmark$$

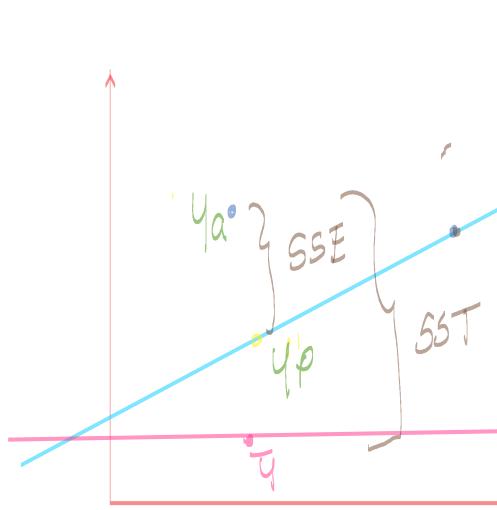
$$n_1 \rightarrow R^2 = 0.85$$

$$n_2 \rightarrow R^2 = 0.86$$

Adjusted  $\overline{R^2}$

$$\overline{R^2} = 0.82 \downarrow$$

$$\overline{R^2} = 0.81$$



$$\checkmark SSE = (y_a - y_p)^2$$

$$SST = (y_a - \bar{y})^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

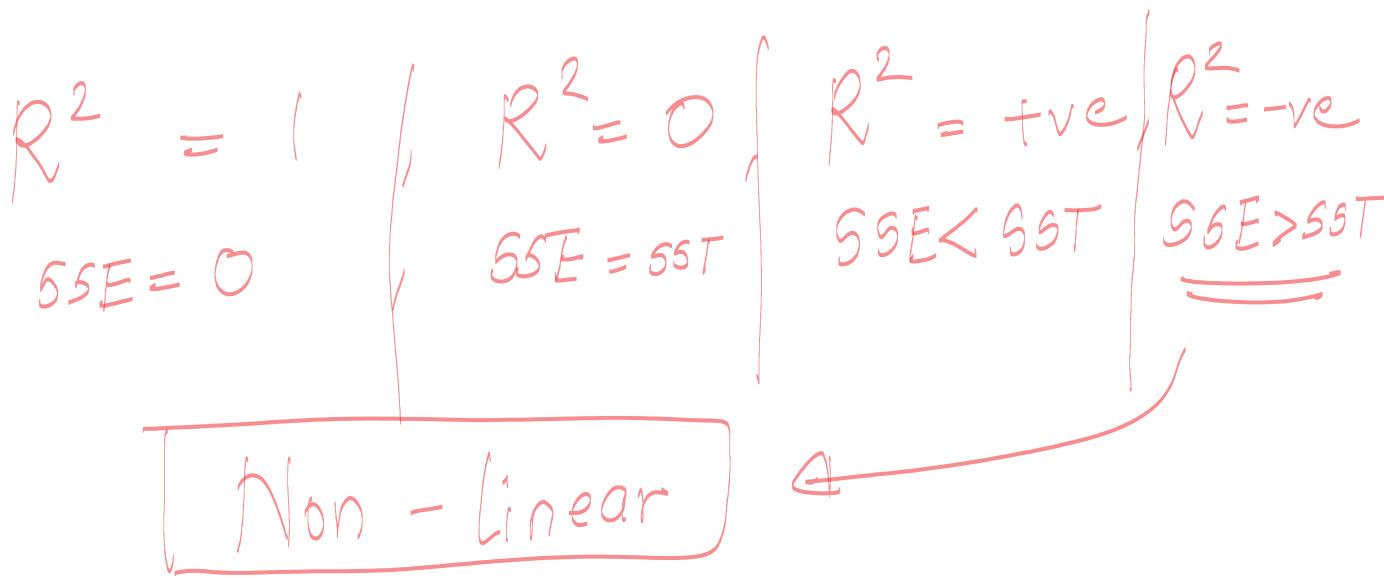
$$R^2 = 1 - \frac{SSE}{SST}$$

$R^2 = 1 - \frac{\text{Unexplained Variation}}{\text{Total Variation}}$

$$= \frac{SST - SSE}{SST}$$

$$= \frac{\text{Total var} - \frac{\text{Unexplained var}}{\text{Total}}}{\text{Total}}$$

$R^2 = \frac{\text{Explained Var}}{\text{Total Var.}}$



\*  $\overline{R^2} \rightarrow$  Adjusted.

$$\overline{R^2} = 1 - \frac{(1 - R^2)(n-1)}{n-p-1}$$

$n > p$

Simple  $\rightarrow R^2 = R \times R$

Multiple  $\rightarrow R^2 \neq R \times R$

$n =$  No. of samples (rows)

$p =$  No. of predictors (features)

$$n_1 = 0.85$$

$$(n_2 = -0.2)$$

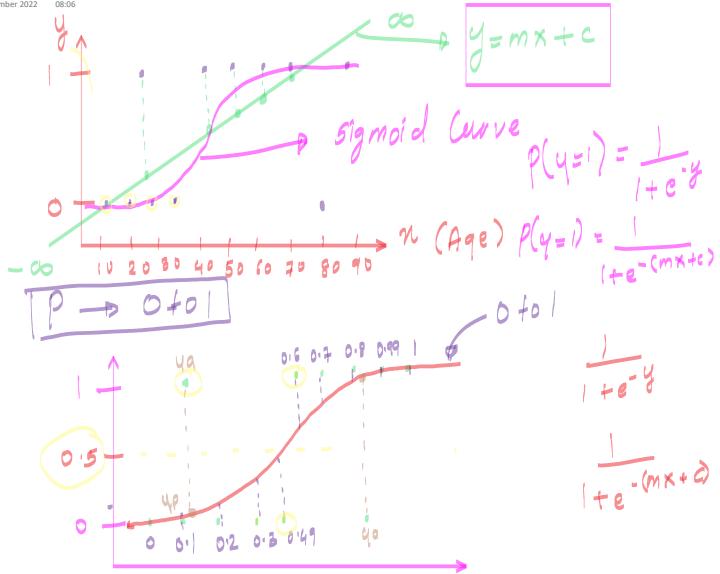
$$R^2 = 0.85$$

$$\overline{R^2} = 0.83$$

$$\overline{R^2} = 0.82$$

$$R^2 = 0.85$$
$$R^2 = 0.86$$

$$\bar{R}^2 = 0.82$$



$$p(y=1) = 0.8$$

$$p(y=0) = 1 - 0.8 = 0.2$$

✓  $p \geq 0.5 \rightarrow \text{Class 1}$   
 $p < 0.5 \rightarrow \text{Class 0}$

Logistic  $\xleftarrow{\text{logit}} \xrightarrow{\log\left(\frac{p}{1-p}\right)}$  Odds

x	$y_a$	$y_p$		
10	0	0.1	0 ✓	$p \geq 0.5 \rightarrow 1$
20	0	0.2	0 ✓	$p < 0.5 \rightarrow 0$
30	1	0.8	1 ✓	
40	0	0.9	0 ✗	misclassified
50	1	0.7	1 ✓	
60	1	0.3	0 ✗	misclassified
70	0	0.4	0 ✓	

$$\frac{x}{1-x} \rightarrow \text{odds}$$

$$\log\left(\frac{p}{1-p}\right) \Rightarrow \text{log odds}$$

Cost Function  $\rightarrow$  MSE  $\rightarrow \frac{(y_a - y_p)^2}{N} \rightarrow$  Linear Reg

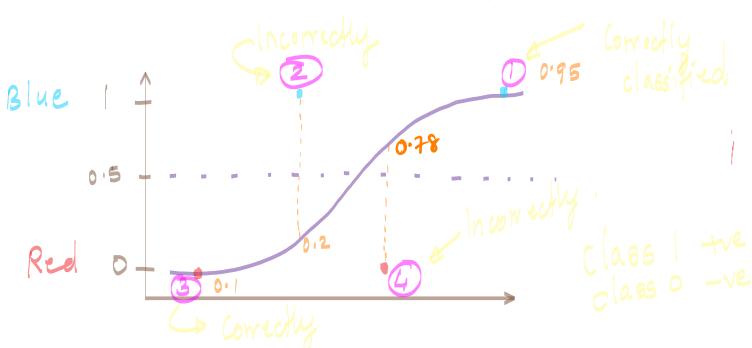
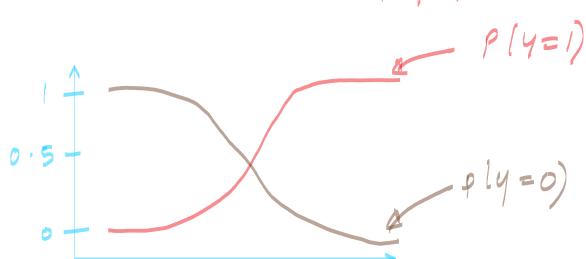
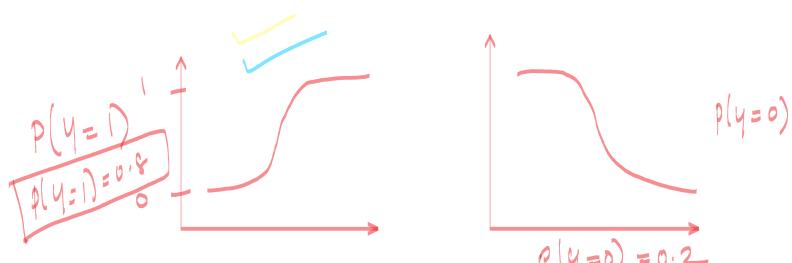
C.F  $\rightarrow$  Log loss  $\rightarrow$  Logistic Regression

$$\text{Log Loss} = -\frac{1}{N} \left[ \sum y_a \log(p) + (1-y_a) \log(1-p) \right]$$

$$\frac{1}{1+e^{-y}} \Rightarrow 0 \text{ to } 1$$

$$e = 2.78$$

$$\frac{1}{1+e^{0.8}}, \frac{1}{1+e^4}, \frac{1}{1+e^{-0.4}}, \frac{1}{1+e^{10}}$$



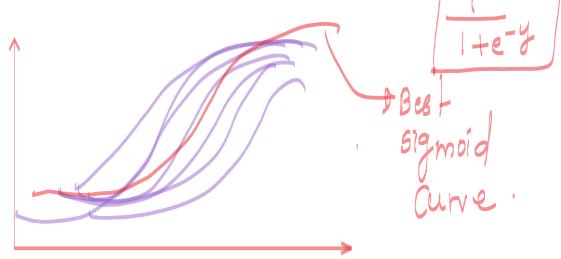
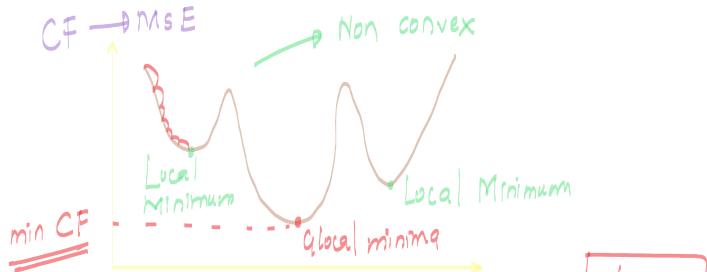
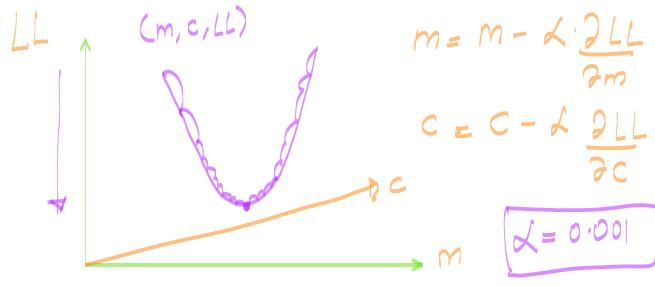
$$y_a = 1$$

$$-\frac{1}{N} \left[ \sum_{a=1}^N y_a \log P + (1-y_a) \log (1-P) \right]$$

$$\boxed{\text{Class 1} \rightarrow LL \rightarrow -\log P}$$

$$-\frac{1}{N} \left( \underbrace{y_a \log P}_{0} + (1-y_a) \underbrace{\log (1-P)}_{0} \right)$$

$$\boxed{\text{Class 0} \rightarrow LL = -\log (1-P)}$$



		Actual	
	1	1	0
1	TP	FP	Incorrect
0	FN	TN	Correct
	2.		
	0	1	3

Type I Error

Type II Error

Class 1 → +ve  
Class 0 → -ve

Correct =  $\frac{5}{7} \times 100 = 70\%$

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

$= \frac{2+3}{2+1+5+1} = \frac{5}{7} = 0.7 \times 100 = 70\%$

	Actual	
Pred	1	0
1	TP	FP
0	FN	TN
	0	1

	Pred	
0	0	1
1	TN	FP
0	FN	TP
	1	0

Learn

Sklearn

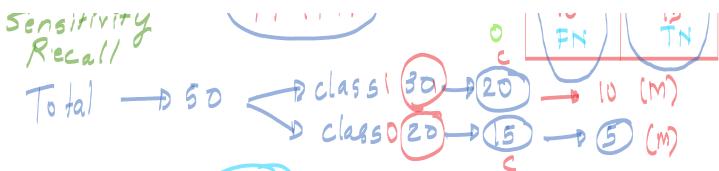
$$\text{① } \underline{\text{TPR}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Sensitivity  
Recall

$T_{\text{Total}} \rightarrow 50 \rightarrow \text{class 1} 30 \rightarrow 20 \rightarrow \text{Pred}$

$c \rightarrow 10 \rightarrow \text{Act} \rightarrow 10 \rightarrow \text{M}$

1	20	5
0	10	15
	FN	TN



$$\textcircled{2} \quad FNR = \frac{FN}{TP+FN}$$

$$\textcircled{3} \quad \begin{matrix} \text{specificity} \\ TNR = \frac{TN}{TN+FP} \end{matrix}$$

$$\textcircled{4} \quad FPR = \frac{FP}{TN+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

		Actual		Precision
		1	0	
Pred	1	TP	FP	
	0	FN	TN	
Actual	1			Recall
	0			

### Classification Report

	Precision	Recall	f1 score
Class 0	-	-	-
Class 1	-	-	-

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

		Act	
		1	0
Pred	1	25	4
	0	5	16

$$\text{Precision} = \frac{25}{25+4} = \underline{\underline{25\%}}$$

$$\text{Recall} = \frac{25}{25+5} = \underline{\underline{25\%}}$$

$$\textcircled{1} \quad \begin{matrix} 900 & 100 \\ 100 & 0 \end{matrix} \quad \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} \quad \text{Acc} = \frac{100+900}{100+900} = \underline{\underline{100\%}}$$

$$\textcircled{2} \quad \begin{matrix} 40 & 50 \\ 60 & 850 \end{matrix} \quad \text{Acc} = \frac{40+850}{6000} = \frac{890}{1000} = \underline{\underline{89\%}}$$

$$\textcircled{3} \quad \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} \quad \begin{matrix} 900 & 100 \\ 100 & 0 \end{matrix} \quad \text{Acc} = \frac{900}{100} = \underline{\underline{90\%}} \checkmark$$

50 - 50  
70 - 30  
60 - 40  
65 - 35

90 - 10  
80 - 20  
Imbalanced data

$$\begin{matrix} 500 & \rightarrow & 100\% \\ 50 & \rightarrow & 0\% \end{matrix}$$

50 → 50

10 → 1

Cancer Pred	
0	1
40	10 (FP)
5 (FN)	35

$$\text{Recall} \rightarrow \text{Increase} = \frac{TP}{TP+FN} \uparrow$$

## ② Spam - Filter

1 → Spam  
0 → Not Spams

Precision

		1	0
1	40	15 (FP)	55
0	10	35	45

Offer letter  
Not spam → Spam

Adv  
Spam → Not spam

## ③ f1-score

$$f_\beta\text{-score} = (1+\beta^2) \frac{P \times R}{(\beta^2 \times P) + R}$$

$$\beta=1, f_1\text{-score} = \frac{2PR}{P+R} \rightarrow \text{Harmonic mean}$$

$$x \& y = \frac{2xy}{x+y}$$

$$\beta=0.5 \quad f_{0.5} = \frac{1.25 P \times R}{0.25 P + R}$$

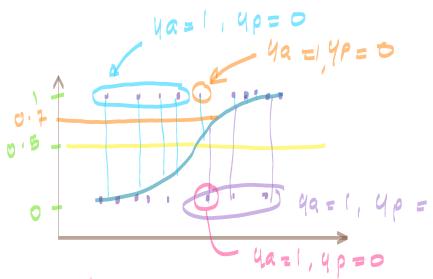
$$\beta=2 \quad f_2 = \frac{5PR}{4P+R}$$

		1	0
1	25	7	32
0	5	13	18

Precision =  $\frac{25}{32} \approx 0.78$   
 Recall =  $\frac{25}{30} = 0.83$   
 Precision =  $\frac{25}{25} = 1$   
 Recall =  $\frac{30}{30} = 1$

Loan → Approved / Declined  
 FP & FN

$$f_1\text{-score} \approx$$



$y_a$	$P$	$y_{p(0)}$	$y_{p(0.1)}$	$y_{p(0.2)}$	$y_{p(0.3)}$	$y_{p(0.5)}$	$y_{p(1)}$
1	0.8	1				0	0
✓	0.95	1	1			0	0
✓	0.42	1	1	1	0	0	0
✓	0.33	1	1	1	0	0	0
✓	0.72	1	1	1	0	0	0
✓	0.65	1	1	1	0	0	0
✓	0.45	1	1	1	0	0	0
✓	0.56	1	1	1	0	0	0
✓	0.15	1	1	1	0	0	0
✓	0.2	1	1	1	0	0	0

$$\begin{array}{c} \text{FP} \\ \text{FN} \end{array} \rightarrow \begin{array}{|c|c|} \hline 5 & 5 \\ \hline 0 & 0 \\ \hline \end{array} \quad \begin{array}{c} \text{FP} \\ \text{FN} \end{array} \rightarrow \begin{array}{|c|c|} \hline 5 & 4 \\ \hline 0 & 1 \\ \hline \end{array} \quad \begin{array}{c} \text{FP} \\ \text{FN} \end{array} \rightarrow \begin{array}{|c|c|} \hline 5 & 2 \\ \hline 0 & 3 \\ \hline \end{array}$$

$$\begin{array}{c} \text{FP} \\ \text{FN} \end{array} \rightarrow \begin{array}{|c|c|} \hline 4 & 1 \\ \hline 1 & 4 \\ \hline \end{array} \quad \begin{array}{c} \text{FP} \\ \text{FN} \end{array} \rightarrow \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 5 & 5 \\ \hline \end{array}$$

$$\begin{array}{c} y_{p(0)} \xrightarrow{\text{FP} \downarrow, \text{FN} \uparrow} y_{p(1)} \\ y_{p(0)} \xleftarrow{\text{FP} \downarrow, \text{FN} \uparrow} \text{Precision } f \\ y_{p(0)} \xleftarrow{\text{FP} \downarrow, \text{FN} \uparrow} \text{Recall } R \end{array}$$

## ROC & AOC

$$\begin{array}{c} \text{Sensi} \xleftarrow{\begin{array}{|c|c|} \hline \text{TP} & \text{FP} \\ \hline \text{FN} & \text{TN} \\ \hline \end{array}} \text{Specificity} \xrightarrow{\begin{array}{|c|c|} \hline 25 & 5 \\ \hline 5 & 15 \\ \hline \end{array}} \end{array}$$

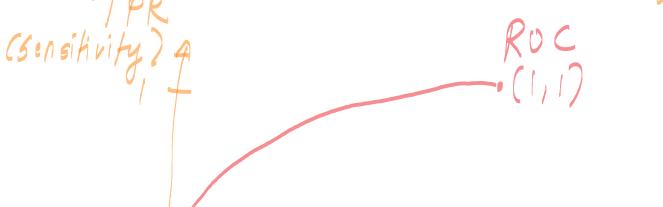
$$TNR = \frac{15}{20} = \frac{TN}{TN+FP} = 0.75$$

$$FPR = \frac{FP}{TN+FP} = \frac{5}{20} = 0.25$$

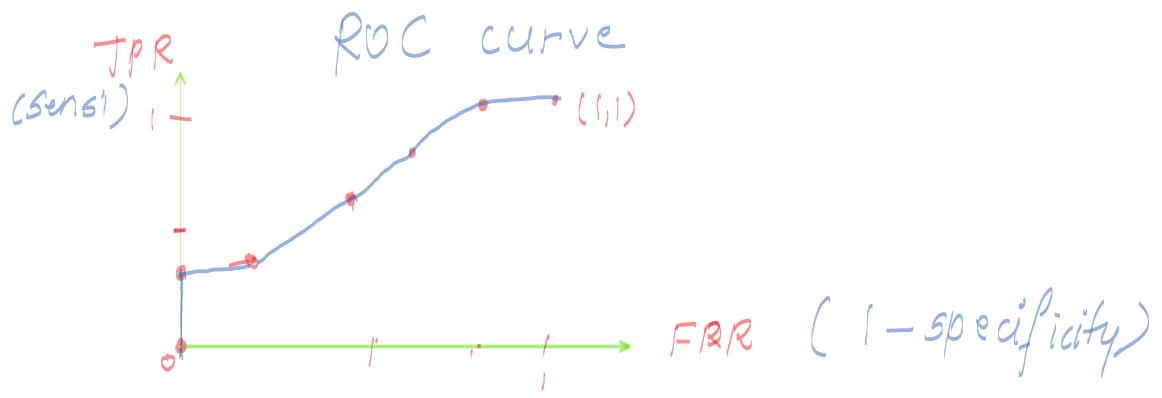
$$FPR = 1 - TNR = 1 - \text{specificity}$$

$\frac{\text{TPR}}{\text{Sensitivity}} \uparrow$

$\frac{\text{TPR}}{\text{ROC}} \uparrow$







①  $Th = 0$

	1	0
1	5	5
0	0	0

$$TPR = \frac{5}{5+0} = 1$$

$$FPR = \frac{5}{5+0} = 1$$

②  $Th = 0.2$

	1	0
1	5	4
0	0	1

$$TPR = 1$$

$$FPR = \frac{4}{4+1} = \frac{4}{5} = 0.8$$

③  $Th = 0.4$

	1	0
1	4	3
0	1	2

$$TPR = \frac{4}{5} = 0.8$$

$$FPR = \frac{3}{5} = 0.6$$

④  $Th = 0.5$

	1	0
1	3	2
0	2	3

$$TPR = \frac{3}{5} = 0.6$$

$$FPR = \frac{2}{5} = 0.4$$

⑤  $Th = 0.7$

	1	0
1	2	1
0	3	4

$$TPR = \frac{2}{5} = 0.4$$

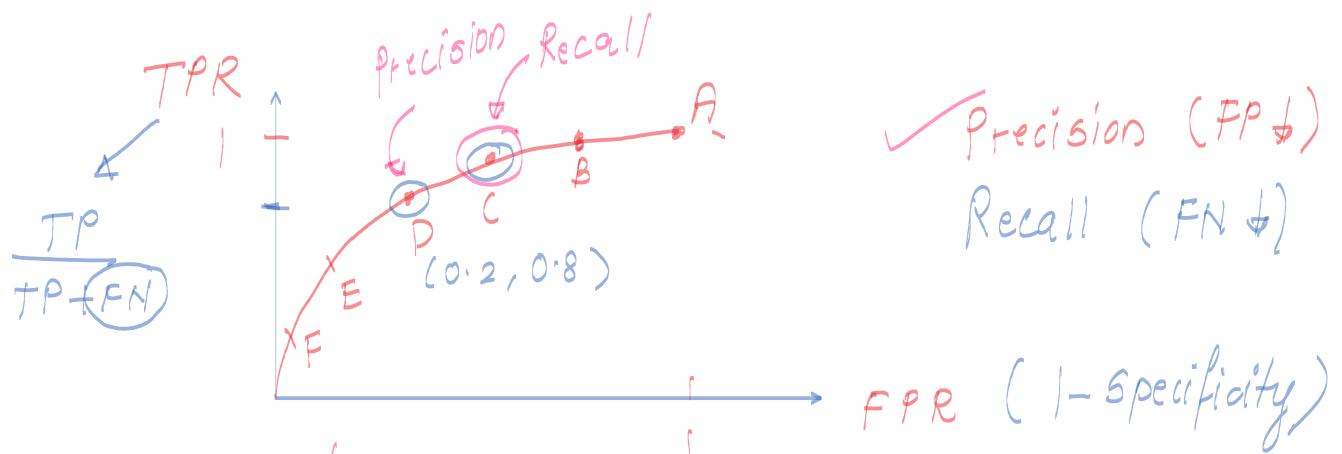
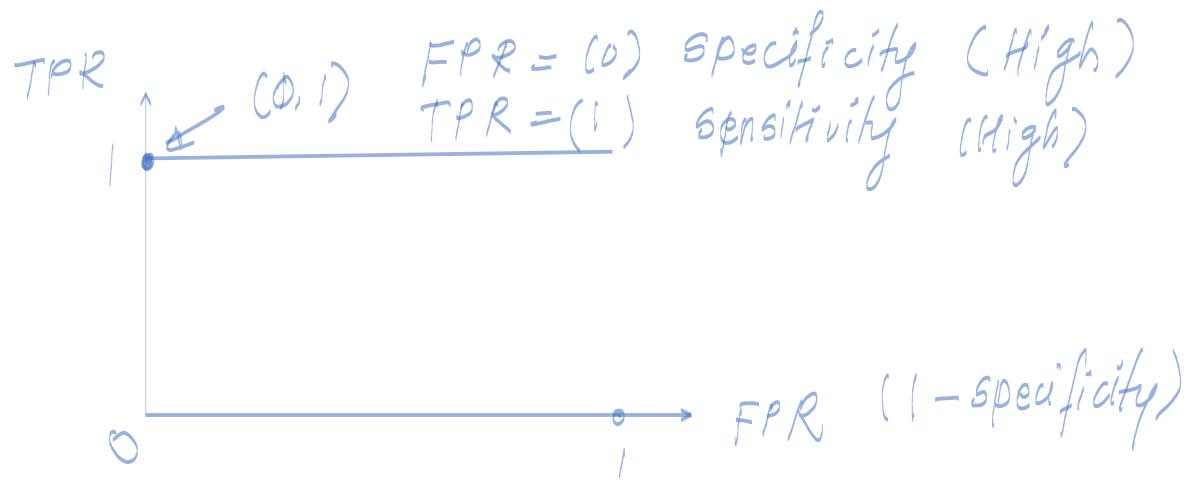
$$FPR = \frac{1}{5} = 0.2$$

⑥  $Th = 0.9$

	1	0
1	1	0
0	4	5

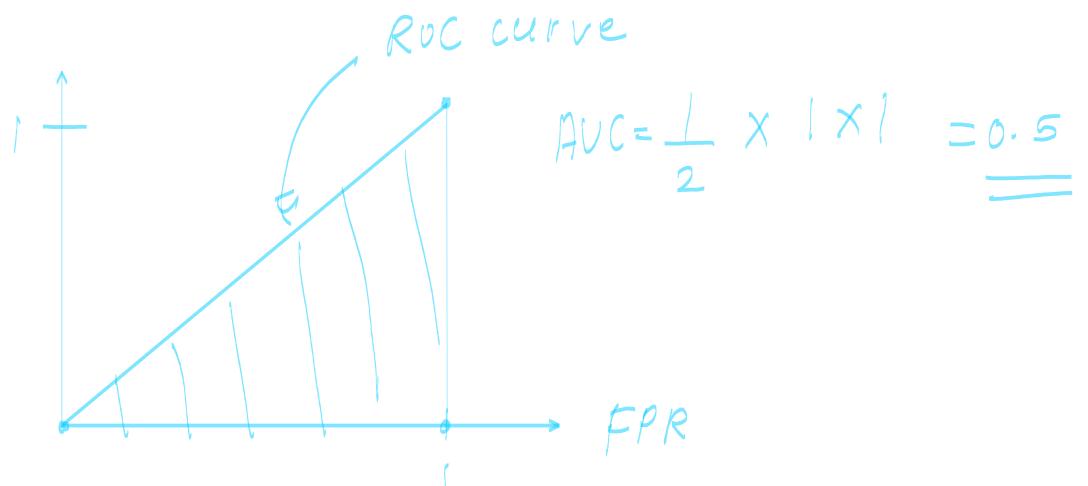
$$TPR = \frac{1}{5} = 0.2$$

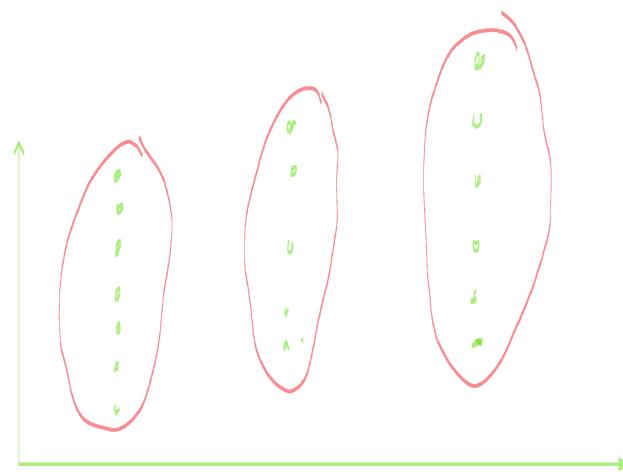
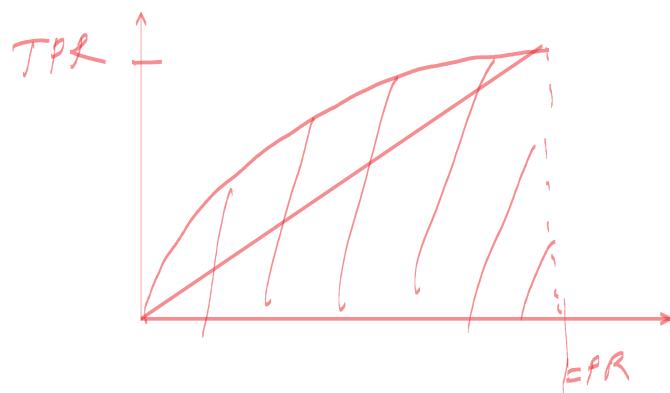
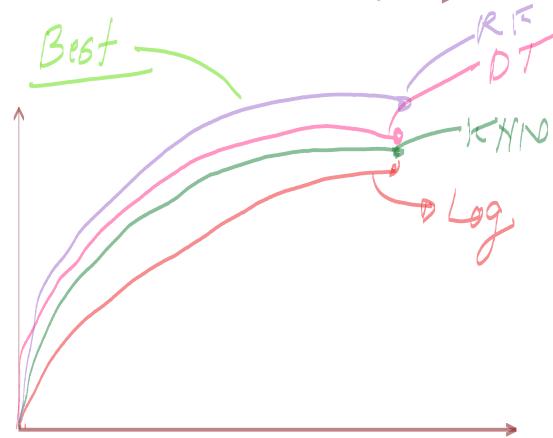
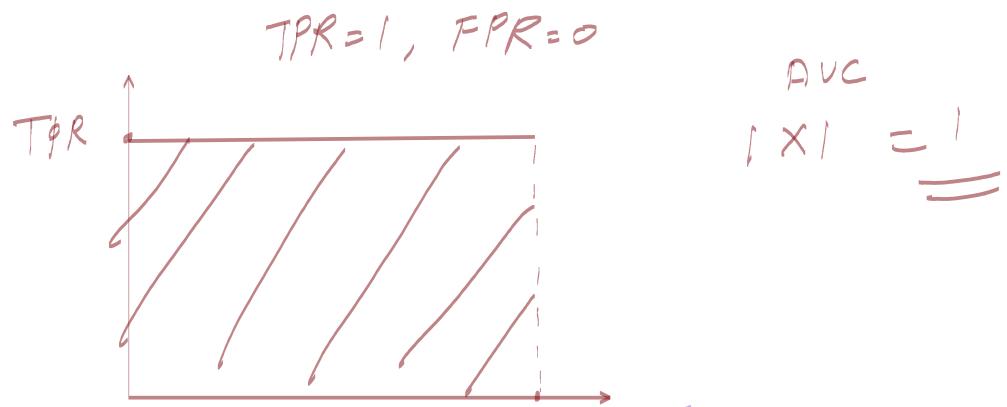
$$FPR = 0$$

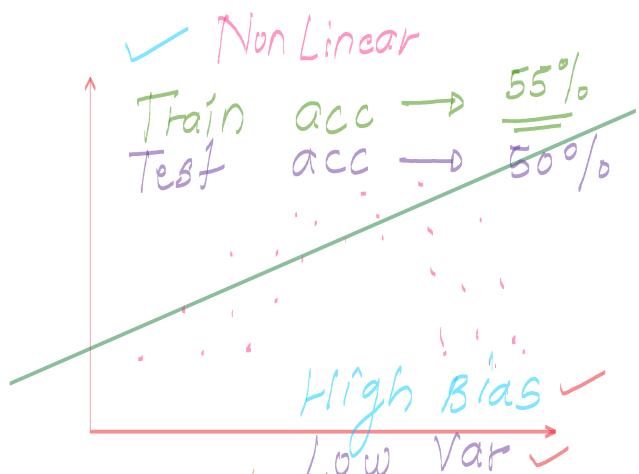
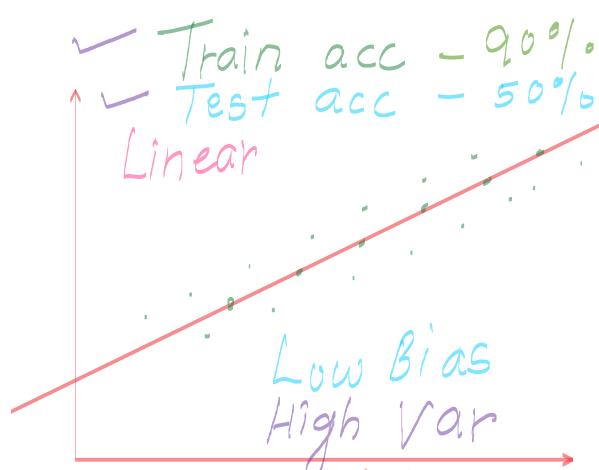
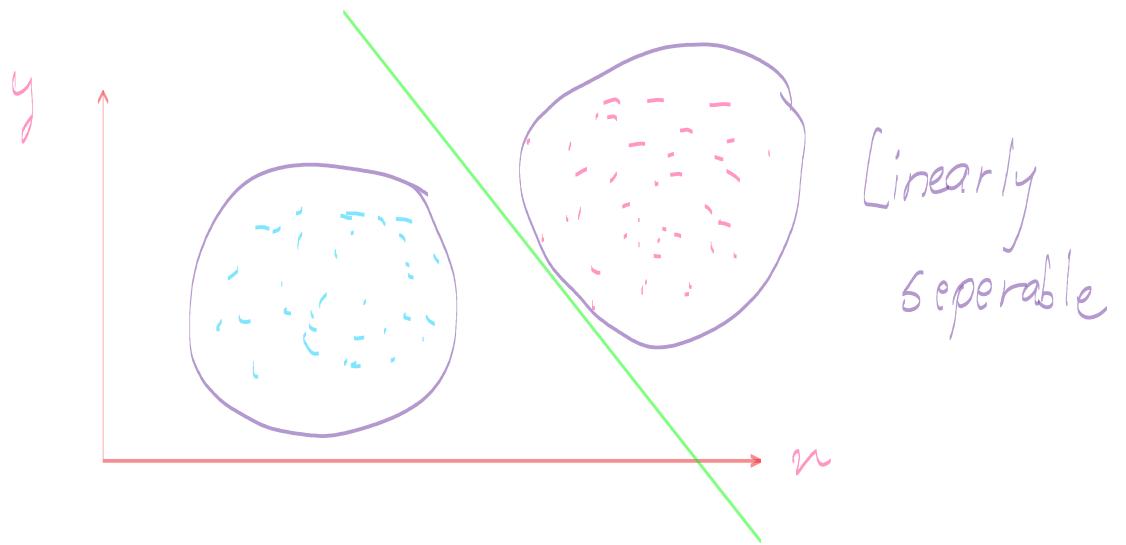


	Act		Pred	
Pred	1	0	0	1
	$TP$	$FP$ (circled)	$TN$	$FP$
Act	1	$FN$	$FN$	$TP$
	0	$TN$	$TN$	$FP$

sklearn





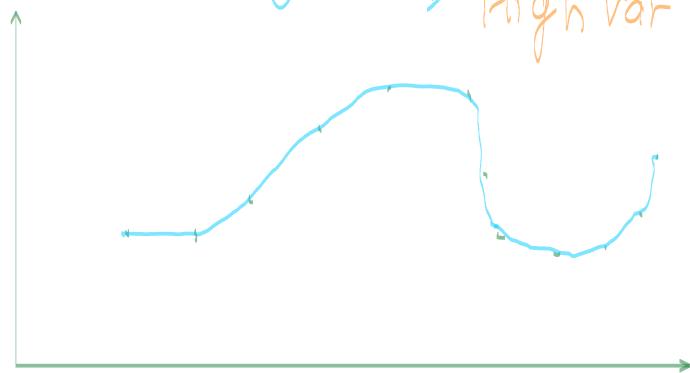


Bias → Low Bias (High Training)  
High Bias (Low Training)

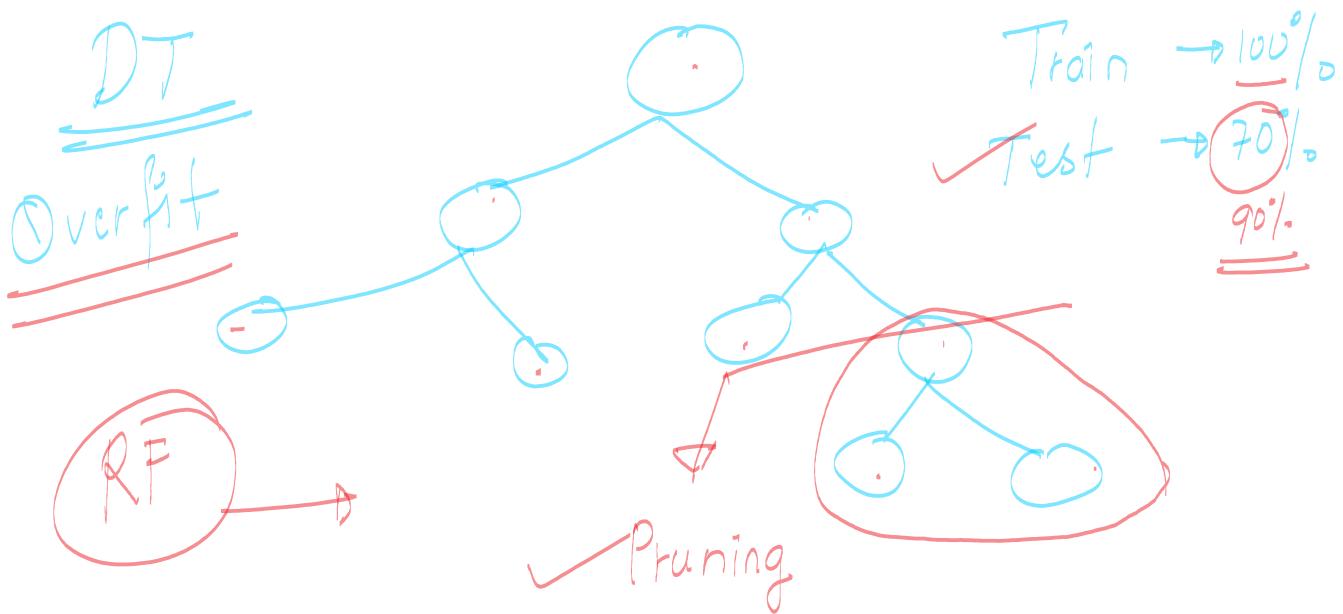
Variance → High Var → High Train, Low Test  
High Test, Low Train

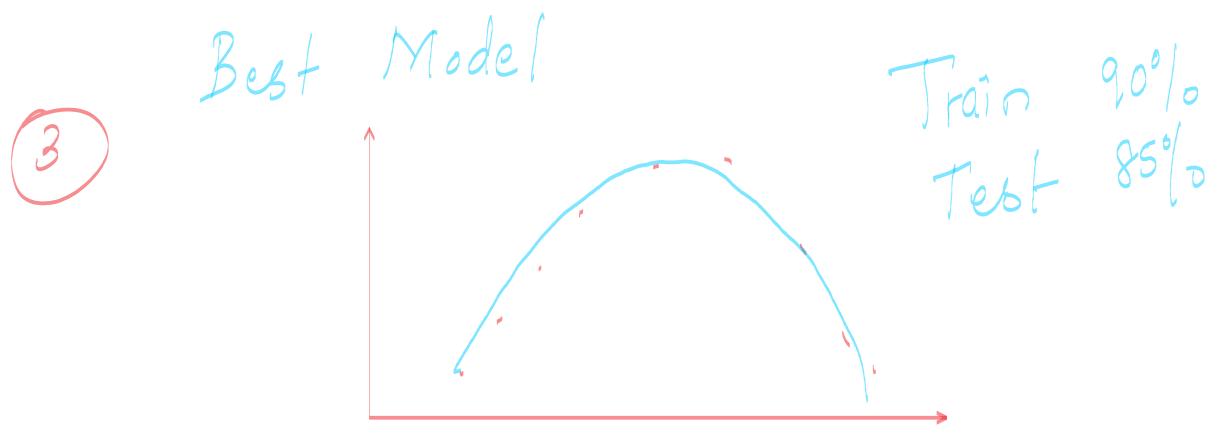
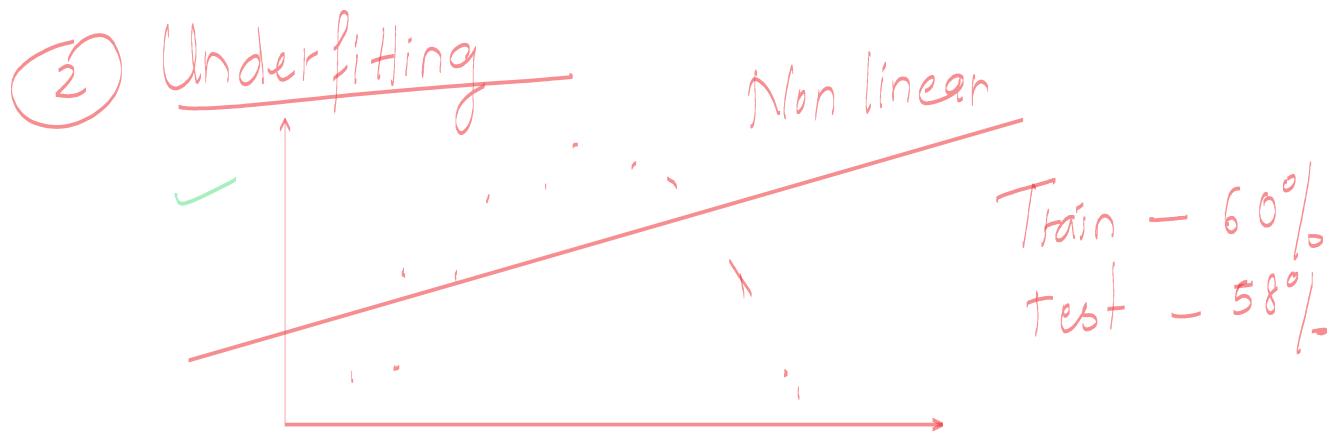
Variance →  $\sigma^2$  → High Test, Low Train  
 Low Var → Low Train, Low Test  
 High Var → High Train, High Test

① Over fitting → Low Bias  
 High Var

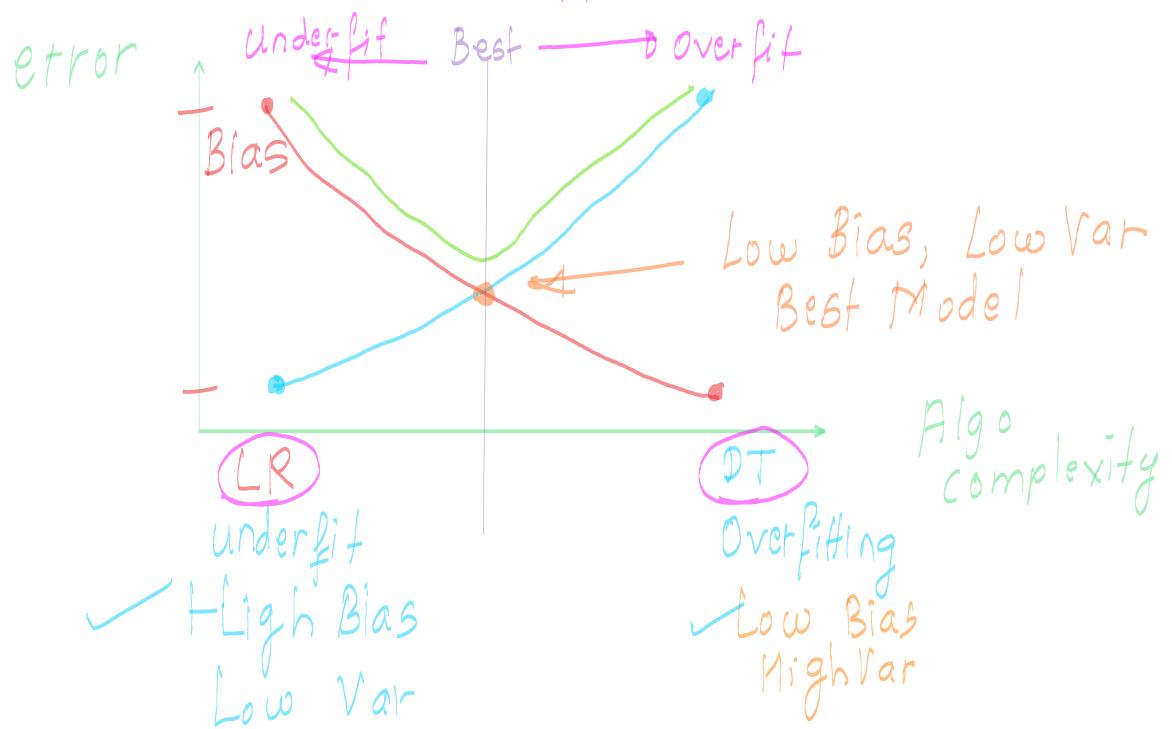


Train = 100%



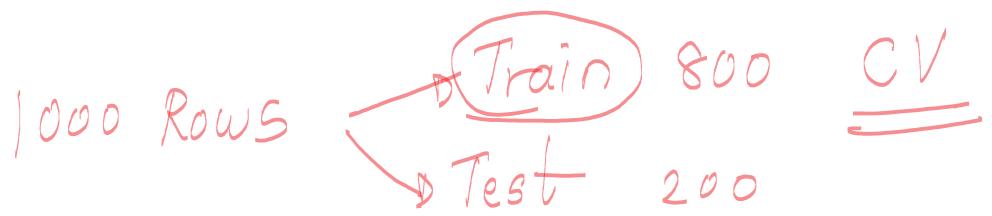


\* Bias-Variance - Trade off

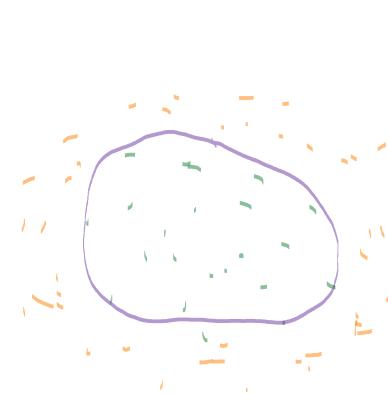
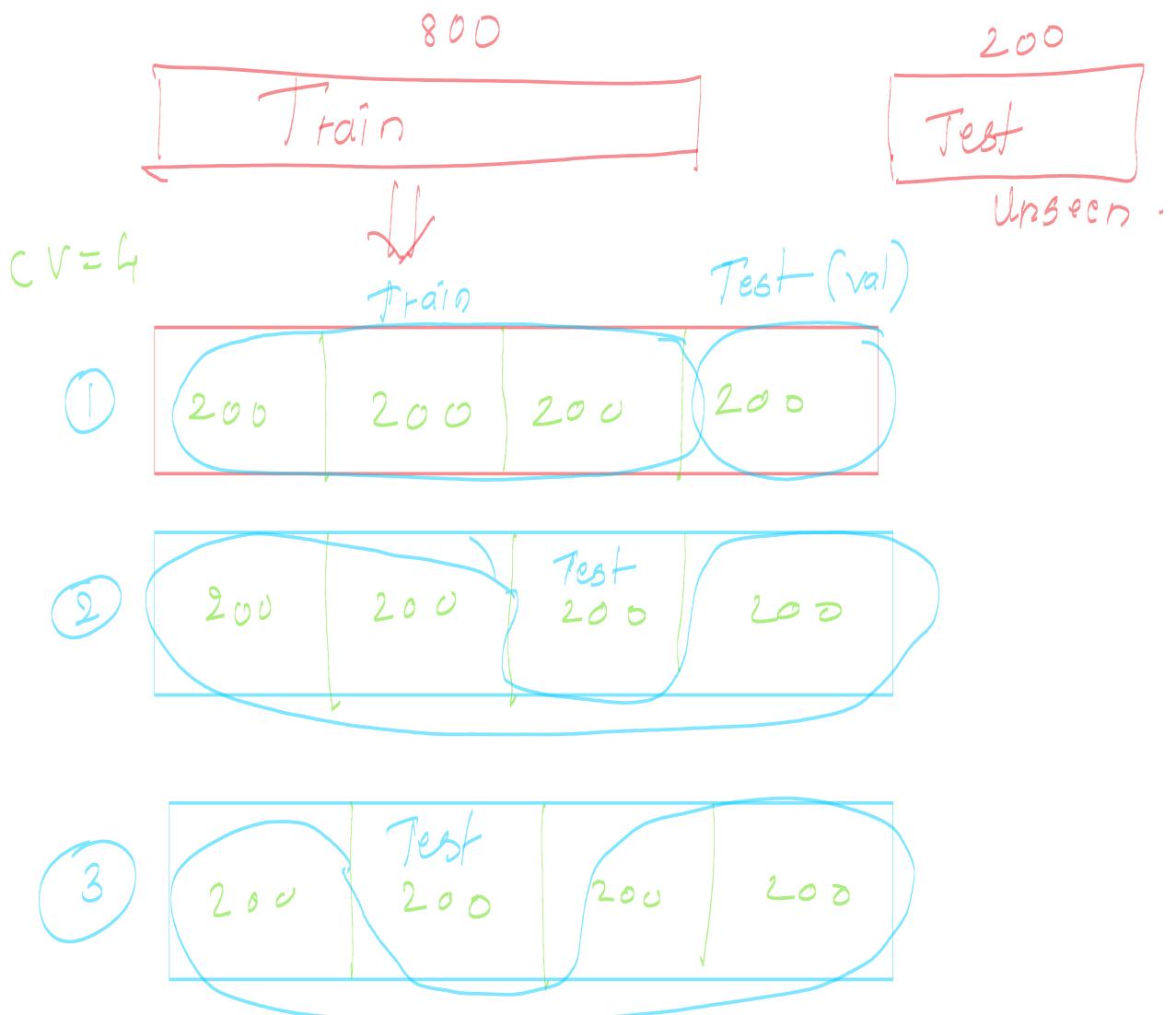


$$TE = \text{Bias}^2 + \text{Var} + \text{Irreducible Error}$$

$\leftarrow$  Error

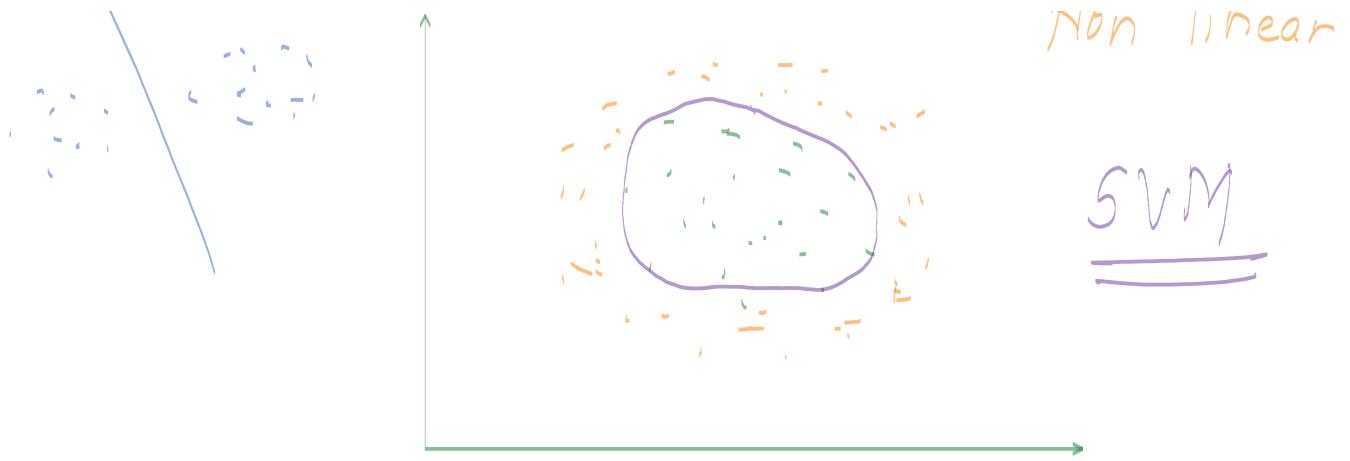


## Cross Validation

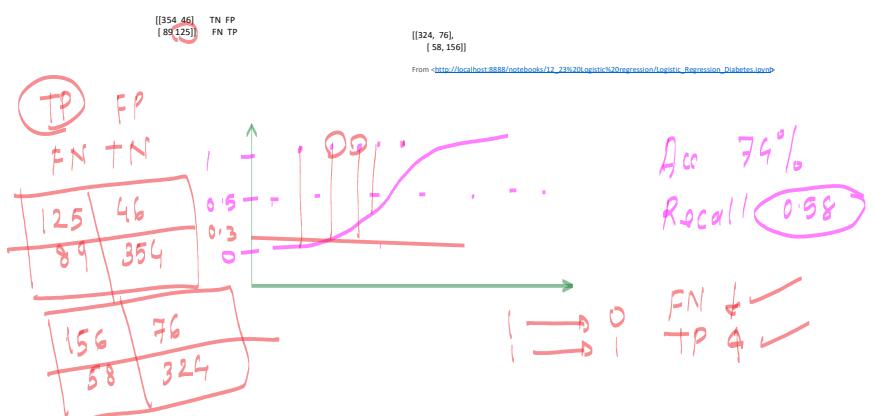


Non linear

SVM



SVM



$$\frac{125 + 354}{125 + 354 + 89 + 46} = 0.76$$

$$\frac{156 + 324}{156 + 324 + 58 + 76} =$$

$$\frac{6}{11} \approx 0.545$$

1	0
3	2
2	3

4	3
1	2

50 ✓ It is - Setosa → 0 | 0 | OVR

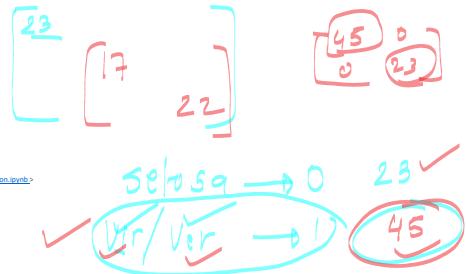
50 ✓ Iris - Versicolor → 1 | 0 |

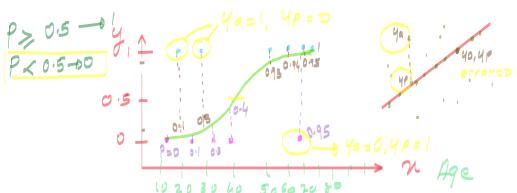
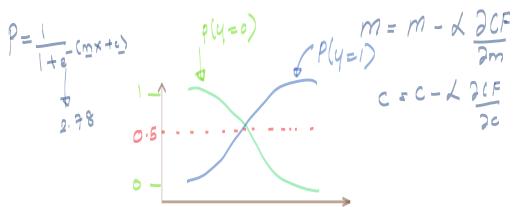
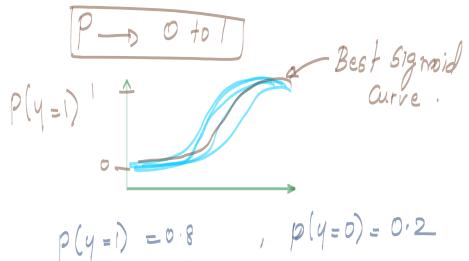
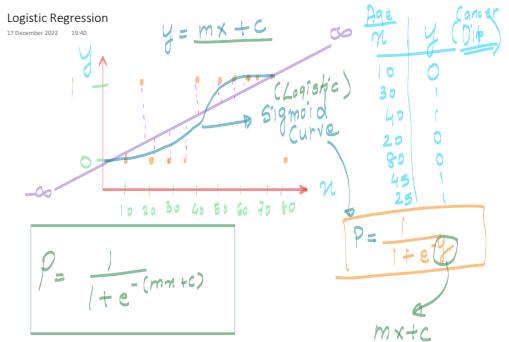
50 ✓ Iris - Virginica → 1 | 1 | 0



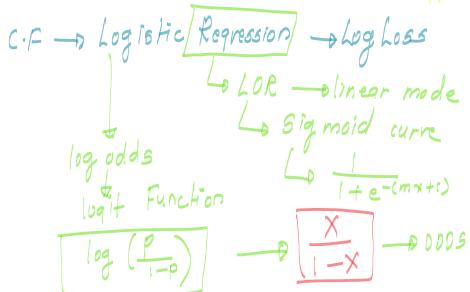
✓ Confusion Matrix :  
[[23 0 0]  
[0 17 5]  
[0 0 22]]  
Mutual Confusion Matrix :  
[[45 0]  
[0 23]]  
[[45 3]  
[3 15]]  
[[45 0]  
[0 23]]]

From <[http://localhost:8888/notebooks/12\\_73%20logistic%20regression/Multiclass\\_Logistic\\_Regression.ipynb](http://localhost:8888/notebooks/12_73%20logistic%20regression/Multiclass_Logistic_Regression.ipynb)>





$$CF \rightarrow \text{Linear Regression} = \text{MSE} \left( \frac{\sum (y_a - y_p)^2}{N} \right)$$

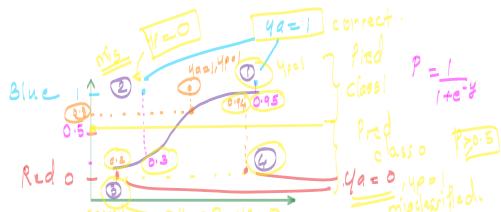


$$\text{Log Loss} = -\frac{1}{N} \left[ \sum \frac{y_a \log(p)}{o} + \frac{(1-y_a) \log(1-p)}{o} \right]$$

$$\frac{1}{1+e^{-y}} \rightarrow \text{Odds}$$

$$P \geq 0.5 \rightarrow 1$$

$$P < 0.5 \rightarrow 0$$



$$LL = -\frac{1}{N} \left[ \sum \frac{y_a \log(p)}{o} + \frac{(1-y_a) \log(1-p)}{o} \right]$$



$$L.L = -\frac{1}{N} \left[ \sum_{i=1}^N \left[ y_i \log(p) + (1-y_i) \log(1-p) \right] \right]$$

$$\boxed{y_i=1} \rightarrow L.L = -y_i \log(p) = -\log p$$

$$\boxed{y_i=0} \rightarrow L.L = -[(1-y_i) \log(1-p)]$$

~~6~~ ~~100~~ = ~~65%~~

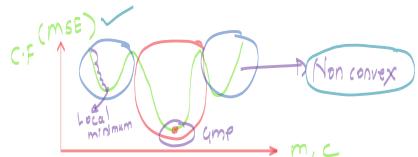
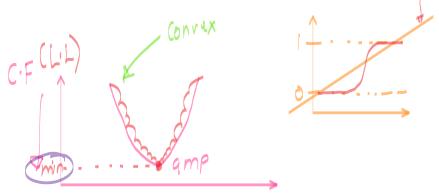
	$y_i$	$p$	$y_p$	$p \geq 0.5 \rightarrow 1$	$p < 0.5 \rightarrow 0$
Total = 6	1	0.2	0 ✓	✓	
Correct = 4	1	0.8	1 ✓		
	1	0.3	0 X	✗ (Mis)	
	0	0.9	1 X	✗ (Mis)	
	0	0.4	0 ✓		
	1	0.7	1 ✓		

$$P(y) = \frac{1}{1+e^{-\theta}} \rightarrow \text{eqn of linear}$$

(linear-model)

$$LR = |y| \rightarrow mx+c \quad \text{range} -\infty \text{ to } \infty$$

$$\text{Logistic (Sig. curve)} \rightarrow P = \frac{1}{1+e^{-\theta}}$$



\* Confusion Matrix

		Actual	
		1	0
Pred	1	TP ✓	FP
	0	FN	TN

Positive  $\rightarrow 1$   
Negative  $\rightarrow 0$

Actual 1

Pred	1	0
1	TP	FP
0	FN	TN

Actual 0

Pred	1	0
1	TP	FP
0	FN	TN

\* Confusion matrix

sklearn

Pred	1	0
1	TP 25	FP 5
0	FN 5	TH 16

$$\checkmark \text{class} 1 \rightarrow 30 \leftrightarrow \frac{TP \rightarrow 25 (cc)}{FN \rightarrow 5 (cc)}$$

$$\text{class} 0 \rightarrow 20 \leftrightarrow \frac{TN \rightarrow 16 (cc)}{FP \rightarrow 4 (cc)}$$

$$\textcircled{1} \quad TPR = \frac{TP}{TP+FN} = \frac{25}{25+5} = \frac{25}{30}$$

$$\textcircled{2} \quad FNR = \frac{FN \times}{FN+TP} = \frac{5}{5+25} = \frac{5}{30}$$

$$\textcircled{3} \quad TNR = \frac{TN}{TN+FP} = \frac{16}{16+4} = \frac{16}{20}$$



o TP | FN / TN

✓ Class 1 → 30       $\frac{TP \rightarrow 25 (C)}{FN \rightarrow 5 (D)}$  ✓  
✓ Class 0 → 20       $\frac{TN \rightarrow 16 (C)}{FP \rightarrow 4 (D)}$  X

① Sensitivity =  $\frac{TP}{TP+FN} = \frac{25}{25+5} = \frac{25}{30}$

② FNR =  $\frac{FN}{FN+TN} = \frac{5}{5+25} = \frac{5}{30}$

③ Specificity =  $\frac{TN}{TN+FP} = \frac{16}{16+4} = \frac{16}{20}$

④ FPR =  $\frac{FP}{TN+FP} = \frac{4}{16+4} = \frac{4}{20}$

		Act	
		1	0
Pred	1	(TP)	(FP)
	0	(FN)	(TN)

→ Precision =  $\frac{TP}{TP+FP}$

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

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

①  $\frac{900}{100} - \frac{100}{100}$

(100)	0
0	900

Accuracy =  $\frac{100+900}{100+900+0+0} = 1 \times 100 \Rightarrow 100\%$

②  $\frac{40}{100} + \frac{850}{900}$

$\frac{40+850}{1000} = \frac{890}{1000} = \underline{\underline{89\%}}$

③  $\frac{1}{0} \quad 0$   
 $\begin{array}{|c|c|c|} \hline & 0 & 0 \\ \hline 1 & 0 & 0 \\ \hline 0 & 100 & 900 \\ \hline \end{array}$

$\frac{900}{100+900} = \underline{\underline{90\%}}$



④

	1	0
1	900	100
0	0	0

$$\frac{900}{100+900} = \underline{\underline{90\%}}$$

50 - 50% Ideal Balanced

70 - 30

60 - 40

65 - 35

✓  $\begin{matrix} 90-10 \\ 80-20 \end{matrix}$  } Imbalanced  
 ✗  $\begin{matrix} 80-80 \\ 20-20 \end{matrix}$  }   
 ↳ SMOTE  
 ↳ Over Sampling  
 ↳ Under Sampling

### \* Classification Report

Precision   Recall   F1-score

Class 0

Class 1

① ✓ Recall(+) → Cancer / No Cancer

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

	1	0
1	40	10
0	20	30

FP

② ✓ Precision(+) → spam / Not spam

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

	1	0
1	40	15
0	10	35

FN

Offer letter  
 ↳ Spam → Not Spam

③ f1-score.

$$f_B\text{-score} = (1+\beta^2) \frac{P \times R}{(\beta^2 \times P) + R}$$

✓  $\beta=1$

$$f_1\text{-score} = \frac{2PR}{P+R}$$

$\frac{2XY}{X+Y}$   
 Harmonic Mean

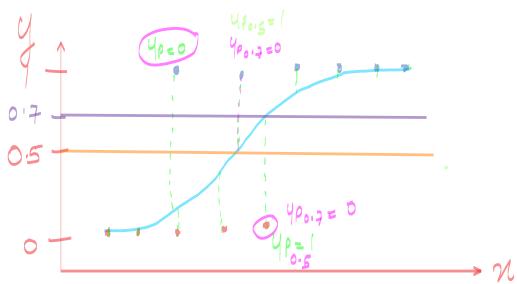
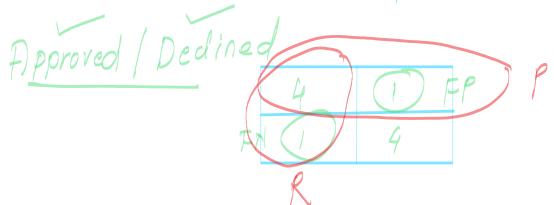


$$\checkmark f_1 - \text{score} = \frac{P+R}{P+R}$$

$\frac{x+y}{x+y}$   
Harmonic Mean

$$\beta = 0.5 \\ f_{0.5} - \text{score} = \frac{1.25 PR}{0.25 P + R}$$

$$\beta = 2, f_2 - \text{score} = \frac{5PR}{4P+R}$$



$y_0$	$P$	$y_{P(0)}$	$y_{P(0.1)}$	$y_{P(0.2)}$	$y_{P(0.4)}$	$y_{(0.5)}$	$y_{C(1)}$
-1	0.8	1	1	1	1	1	0
0	0.95	1	1	1	1	1	0
1	0.42	1	1	1	0	0	0
0	0.33	1	1	1	0	0	0
-1	0.72	1	1	1	0	0	0
-1	0.65	1	1	1	0	0	0
0	0.45	1	1	1	0	0	0
-1	0.56	1	1	1	0	0	0
0	0.15	1	1	1	0	0	0
0	0.26	1	1	1	0	0	0

$$P \geq 0.5$$

$$P < 0.5$$

$Th = 0$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 5 & 5 \\ \hline 0 & 0 \\ \hline \end{array}$	$Th = 0.1$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 5 & 5 \\ \hline 0 & 0 \\ \hline \end{array}$
$Th = 0.2$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 5 & 4 \\ \hline 0 & 1 \\ \hline \end{array}$	$Th = 0.4$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 5 & 2 \\ \hline 0 & 3 \\ \hline \end{array}$
$Th = 0.5$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 4 & 1 \\ \hline 1 & 4 \\ \hline \end{array}$	$Th = 1$	$\begin{array}{ c c }\hline 1 & 0 \\ \hline 0 & \text{FP} \\ \hline \end{array}$	$\begin{array}{ c c }\hline 0 & 5 \\ \hline 5 & 5 \\ \hline \end{array}$

Precision

$\checkmark Th \neq \text{FP} \neq FN \neq$

$\rightarrow Th \neq \text{FP}, FN \neq$

Recall

\* ROC - AUC Curve

TPR  
(Sensitivity)

Sensi  $\leftarrow$ 

TP	FP
FN	TN

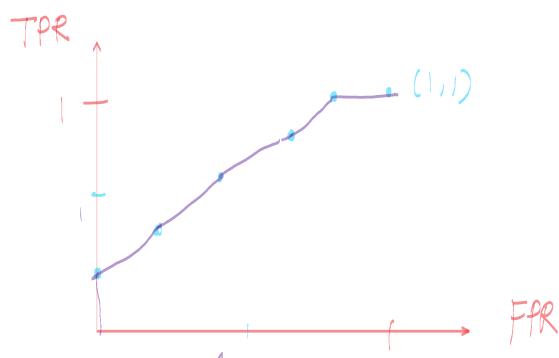
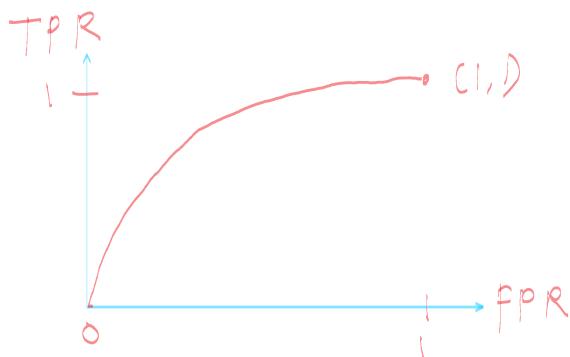
 $\rightarrow$  Specificity





TNR → Specificity

$$\boxed{\text{FPR} = 1 - \text{Specificity}}$$



$$Th = 0$$

5	5
0	0

$$\begin{aligned} \text{TPR} &= 1 \\ \text{FPR} &= 5 \end{aligned}$$

$$Th = 0.2$$

5	4
0	1

$$\begin{aligned} \text{TPR} &= 1 \\ \text{FPR} &= 0 \end{aligned}$$

$$Th = 0.5$$

3	2
2	3

$$\begin{aligned} \text{TPR} &\approx 0.6 \\ \text{FPR} &= 0.4 \end{aligned}$$

$$Th = 0.4$$

4	3
1	2

$$\begin{aligned} \text{TPR} &= 0.8 \\ \text{FPR} &= 0.6 \end{aligned}$$

$$Th = 0.7$$

2	1
3	4

$$\begin{aligned} \text{TPR} &\approx 0.4 \\ \text{FPR} &= 0.2 \end{aligned}$$

$$Th = 0.9$$

1	0
4	5

$$\begin{aligned} \text{TPR} &= \\ \text{FPR} &= \end{aligned}$$



✓ Recall

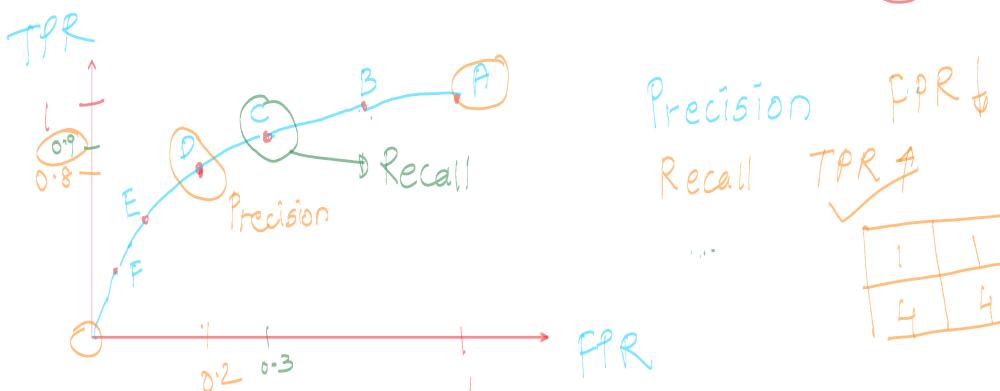
$$\frac{TP}{TP+FN}$$

$\checkmark FPR = 0 \rightarrow \text{High Specificity}$

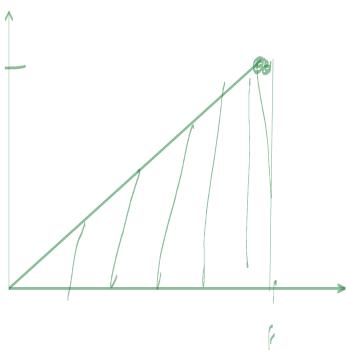
$TPR = 1 \rightarrow \text{High Sensitivity}$

$$FPR = (1 - \text{Specificity})$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

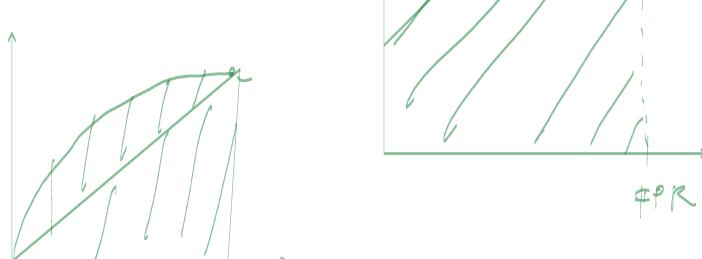


$$\frac{1}{2} \times 1 \times 1 = \underline{\underline{0.5}}$$

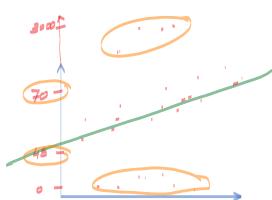


$$0.5 < \text{AUCC} /$$

$$\begin{matrix} 0.8 \\ 0.7 \\ 0.6 \end{matrix}$$

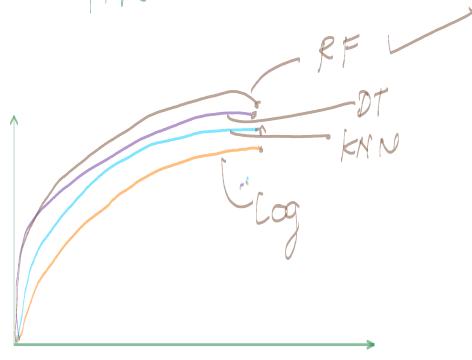


$$1 \times 1 = \underline{\underline{1}}$$



\* Classification → Target  
Imbalanced  $\begin{cases} 0 \rightarrow 90, 80 \\ 1 \rightarrow 10, 20 \end{cases}$

--- Ideal balanced





\* Classification → Target

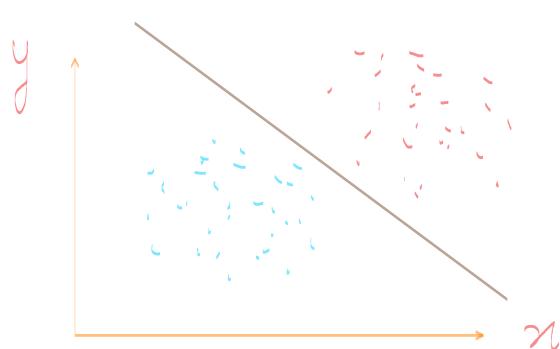
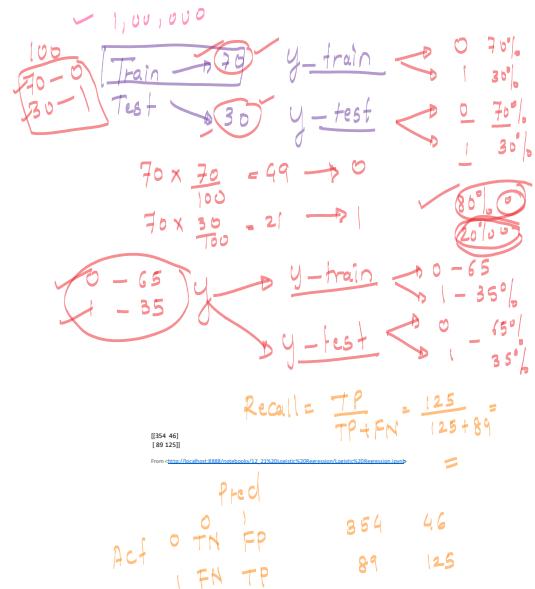
Imbalanced ↗ 0 → 90, 80 ✓  
                   ↘ 1 → 10, 20 ✓



Balanced { 50 - 50      Ideal      Balanced  
           60 - 40  
           65 - 35

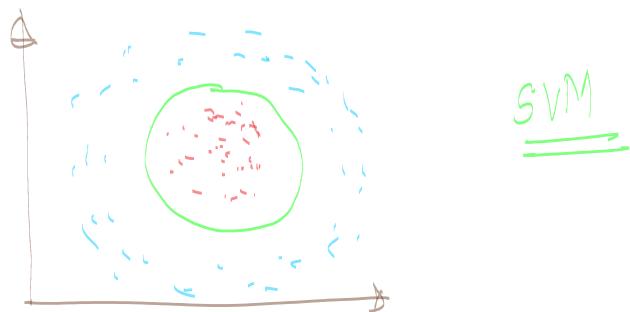
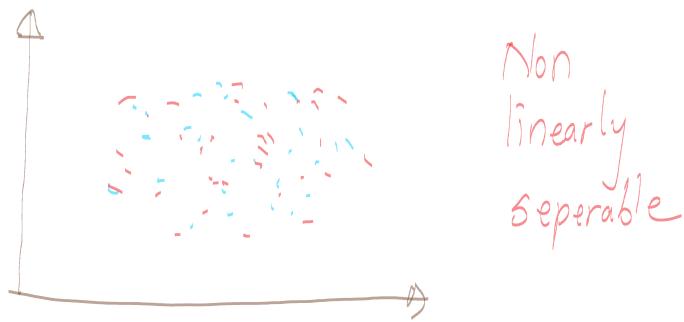
SMOTE

OverSampling (4)	80 - 80	Balanced
UnderSampling (4)	20 - 20	Balancing



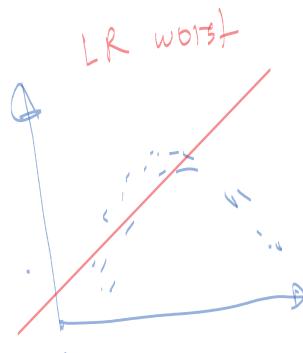
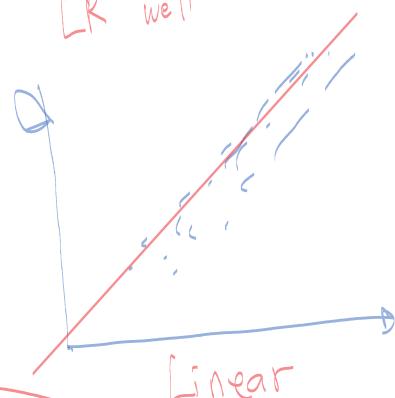
linearly  
separable





LR well

LR worst



### \* Assumption

Train  $\rightarrow$  10 Features  $\leftrightarrow$  Train  $\checkmark 0.78 \checkmark$   
 test  $\times 0.75 \times$

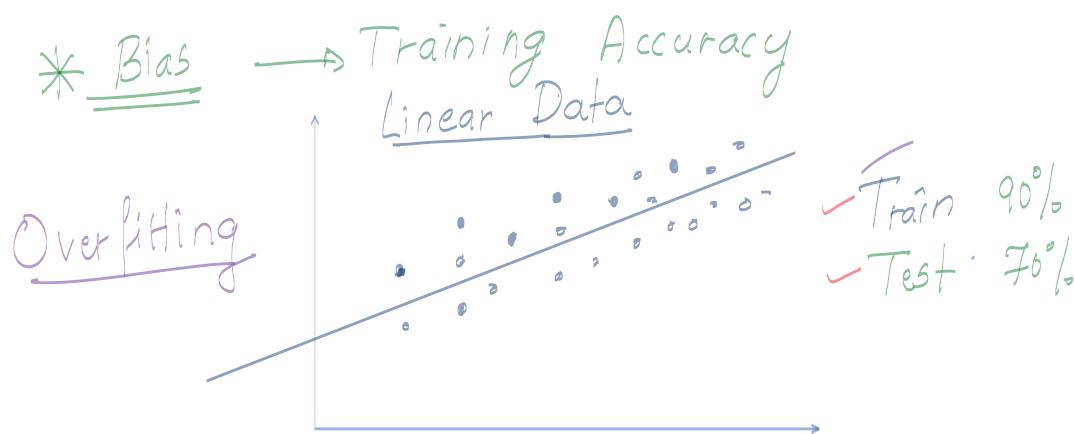
Linearly 5 Corr., 5 Non Corr  $\rightarrow$  [drop] (if) score  
 $\begin{cases} 1\text{drop} \\ 2\text{drop} \\ 3\text{drop} \end{cases}$

### \* Variance $\rightarrow$ Difference bet<sup>n</sup> Accuracies

High Var  $\rightarrow$  High Train and Low Test  
 Low Train and High Test

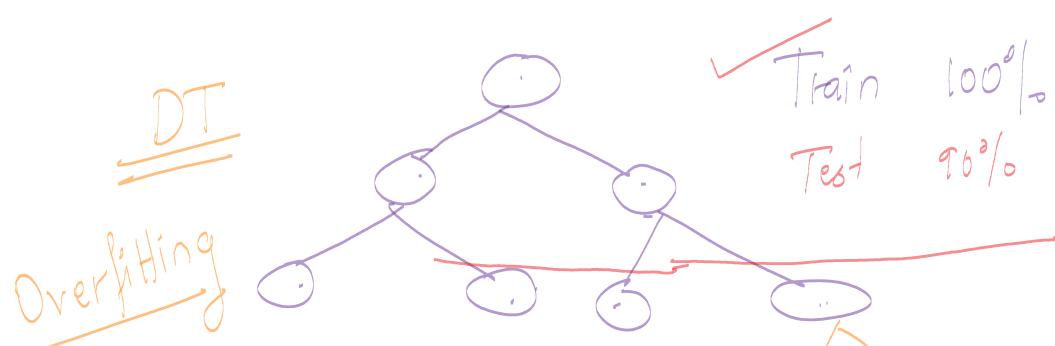
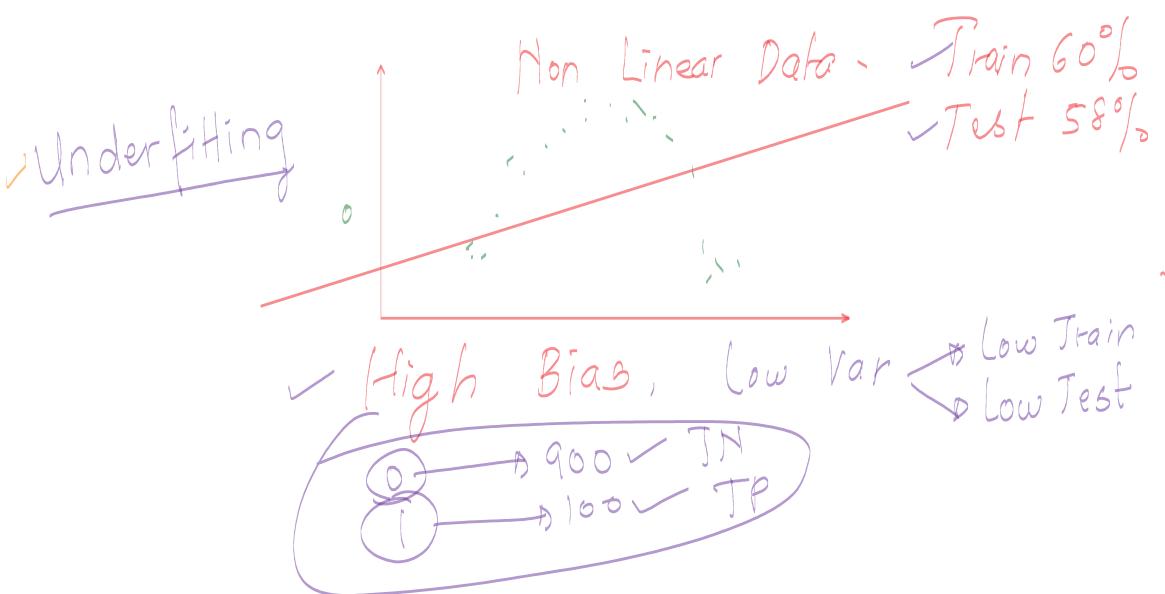
Low Var  $\rightarrow$  Low Train and Low Test  
 High Train and High Test





High Training Accuracy → Low Bias

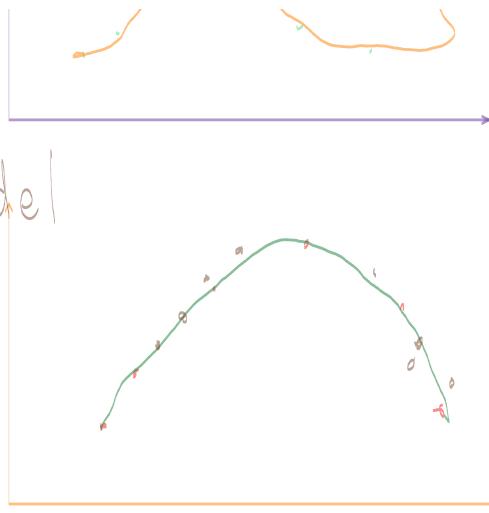
High Variance → High Train  
Low Test



Pruning (Cut)  
Hyper Parameter Tuning



DT)

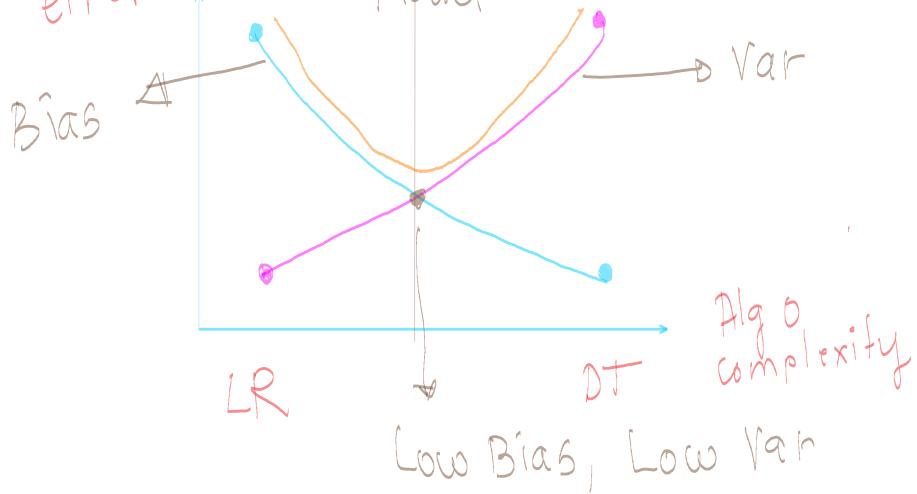


## \* Best Model

Train 98%  
Test 95%

High Training → Low Bias ✓  
→ Low Variance ✓

## \* Bias - Variance Trade-off



Non-linear  
LR → Und  
High Bias, L  
DT → Ove  
Low Bias, H

$$TE = \text{Bias}^2 + \text{Var} + \text{Irreducible Error}$$

Noise

## \* Cross Validation

1000 Rows

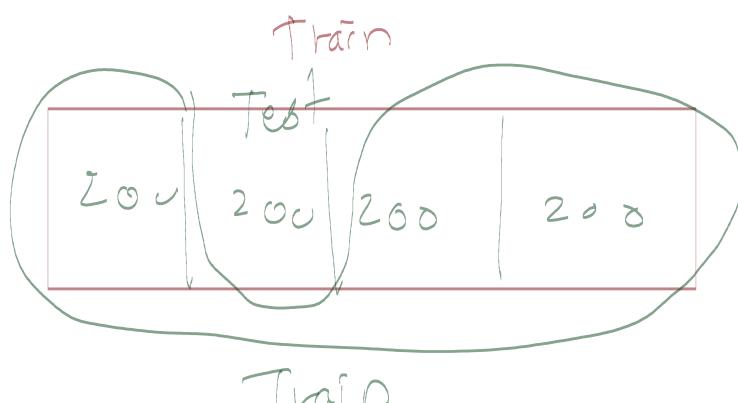
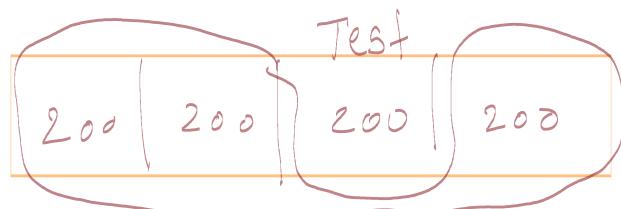
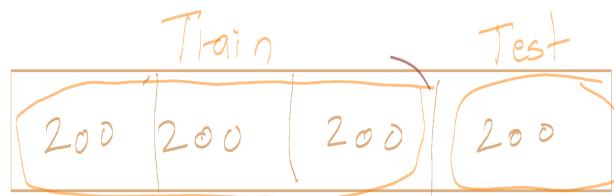
er fit  
low var

er fit  
high var

$CV = 4$

800 Train

200 Test



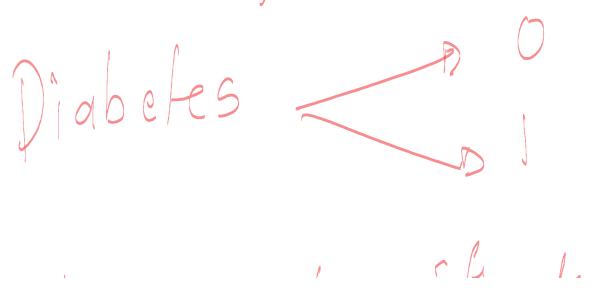
Decision\_Tree (min-sample-leaf = 2)

$$= \text{min-sample-leaf} = [2, 3, 4]$$

Logistic - Train



\* Binary Classification



5,6...10

Diabetes

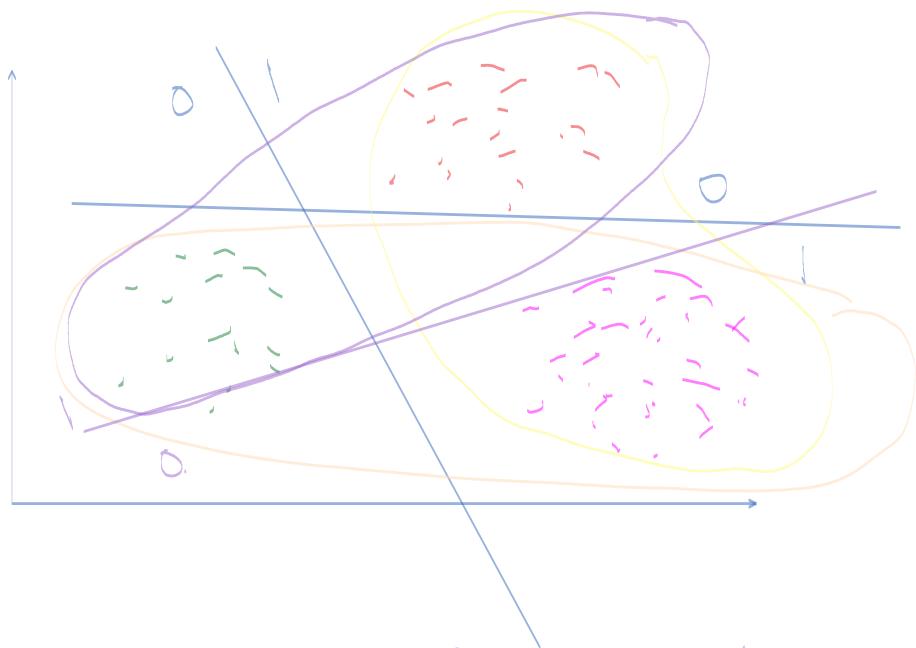
## \* Multiclass Classification One Vs Rest

50 Iris - Setosa - 0 (50)

50 Iris - Versicolor

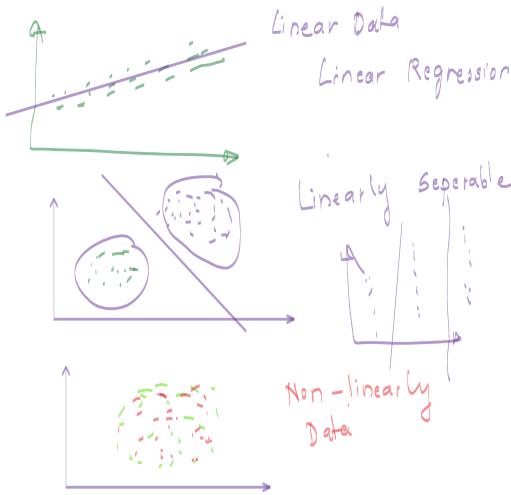
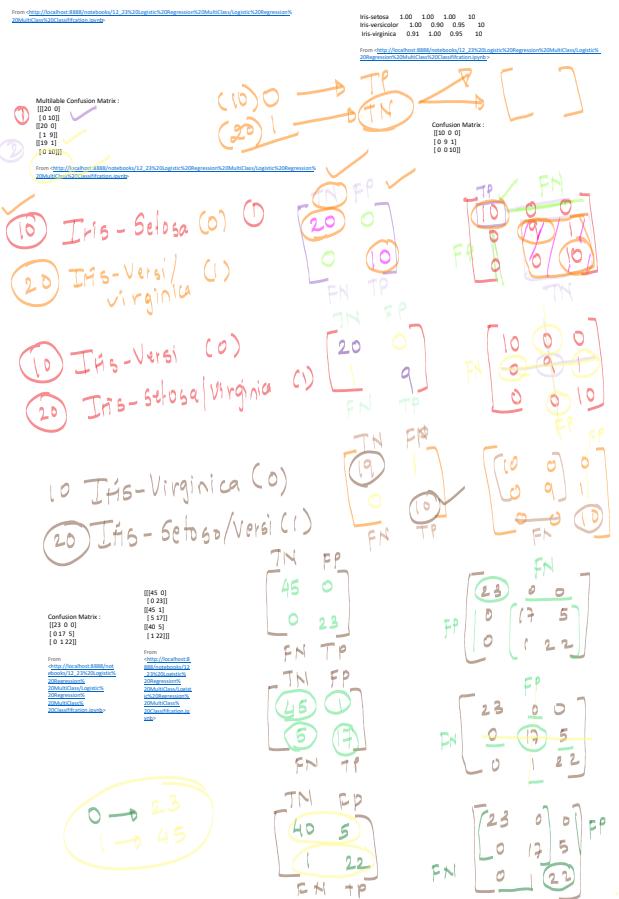
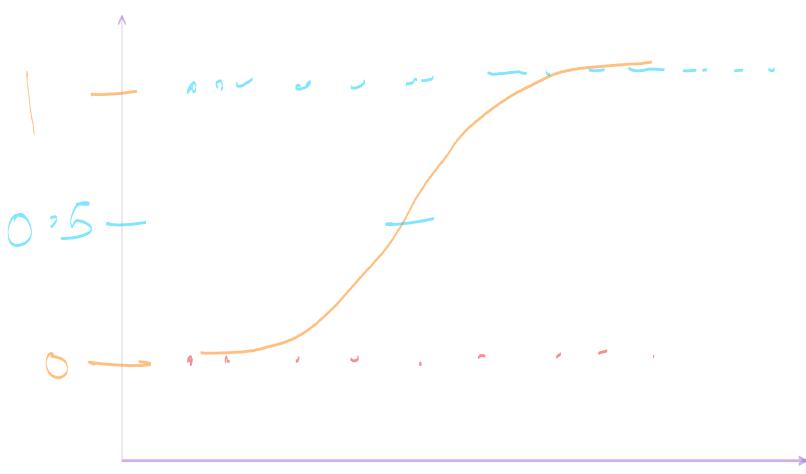
50 Iris - Virginica

1 (100) |  
0 (50) |



Multiclass confusion matrix







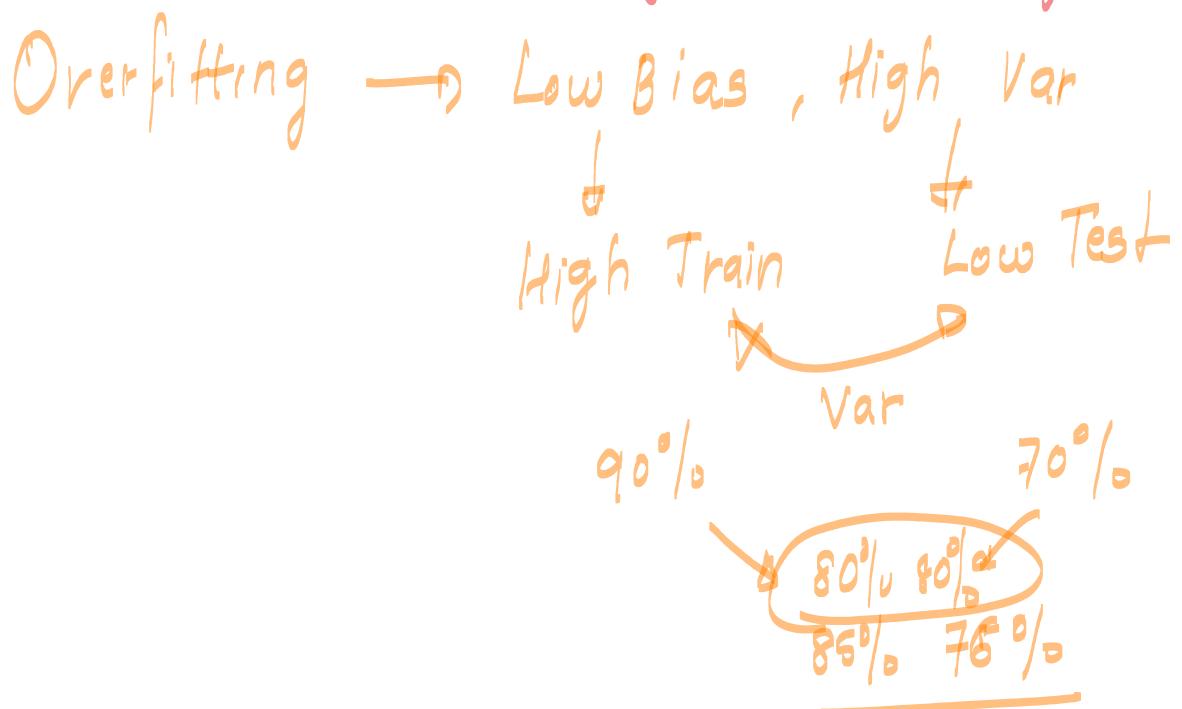
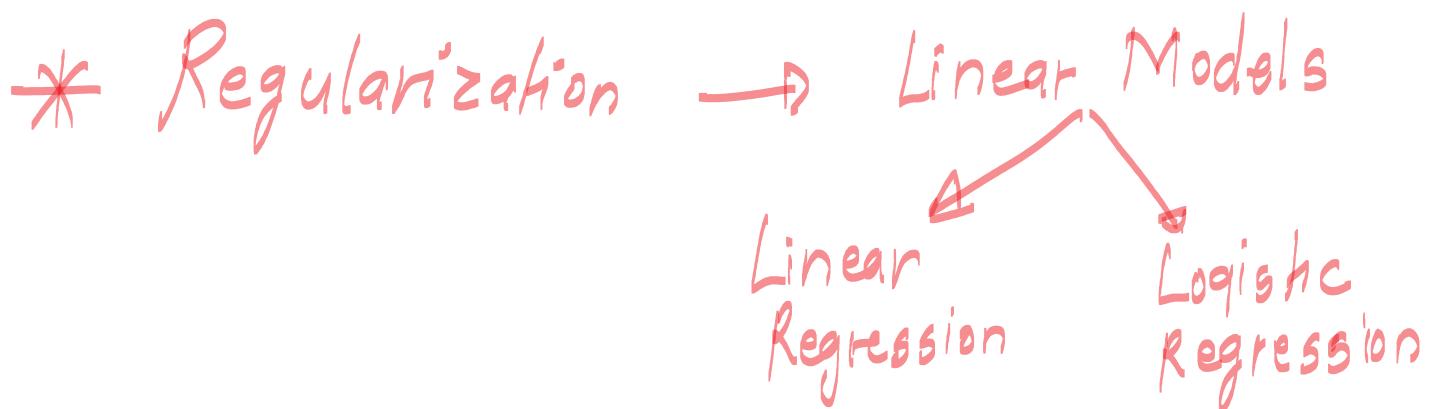
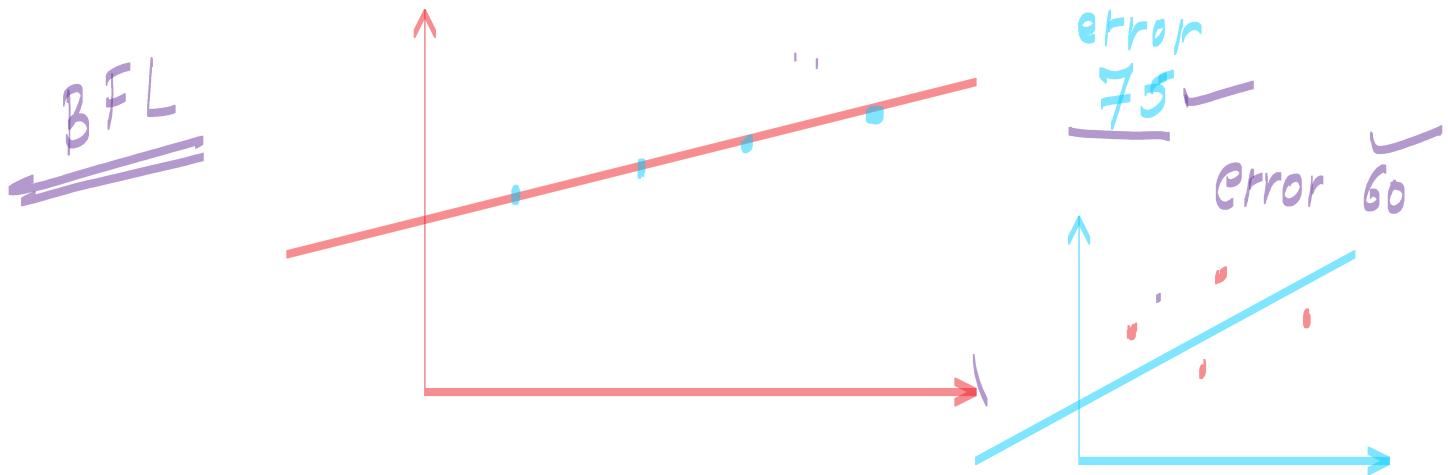
(50) Iris-Setsosa → 0  
(100) (Iris-Versicolor / virginica) → 1

(100): 1  
(50): 0



## Regularization

26 December 2022 07:30



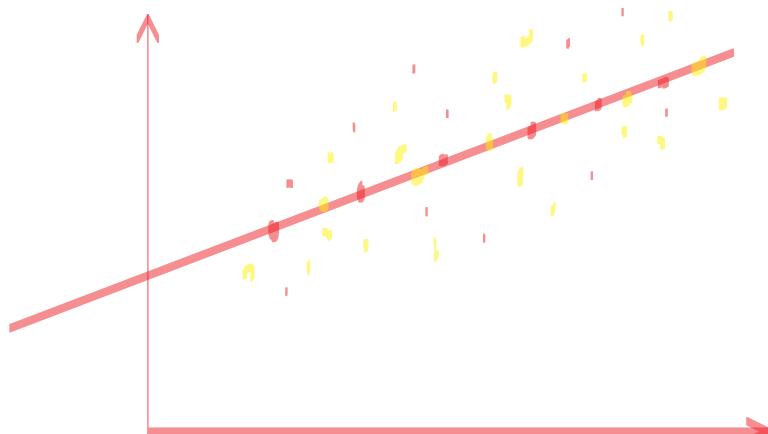
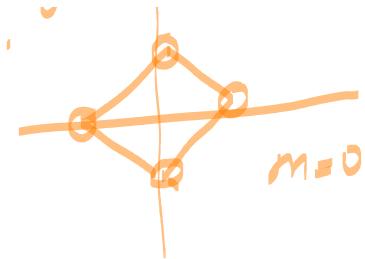
- ① Lasso Regression
- ② Ridge Regression

(L<sub>1</sub> Regularization)  
(L<sub>2</sub> Regularization)

m ≠ 0

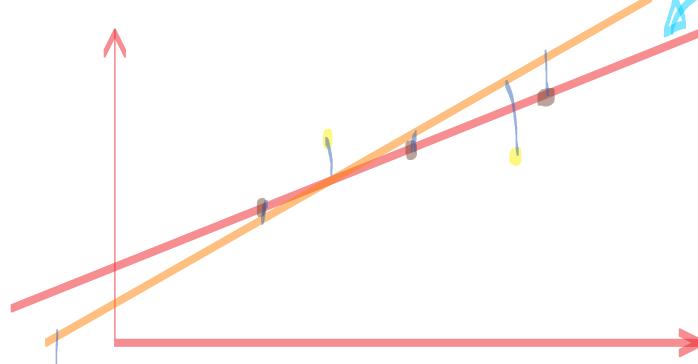
# \* Lasso Regression

$m \neq 0$



0.8 Train  
0.75 Test

LOR  $\rightarrow$  Linear.

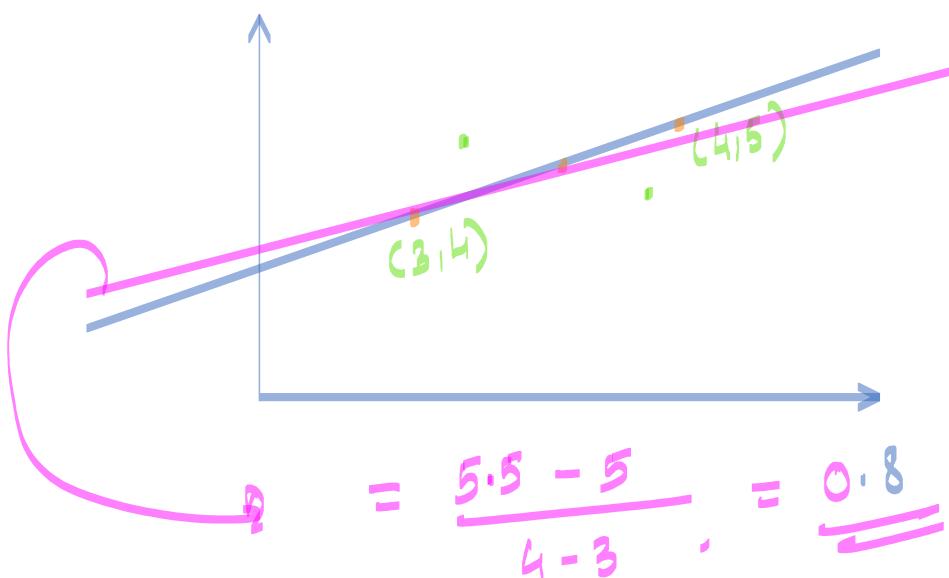


Train error = 0  
 $\hookrightarrow 100\%$

Test error = 10  
 $\hookrightarrow 80\%$

Train  $\rightarrow$  8 = error  
Test  $\rightarrow$  6 = error

$91\% \quad 85\% \quad \} \text{var} \downarrow$



LOR  
 $m = \frac{5-4}{4-3}$

$$= \frac{1}{1} = 1$$

$$= \frac{5.5 - 5}{4 - 3} = 0.8$$

$CF = MSE \rightarrow$  Linear Reg

$$CF = SSE + \lambda |m|$$

slope

①  $CF = SSE + \lambda |m| \xrightarrow{\text{LR}} m^2 \text{ Ridge} \rightarrow \text{Linear model}$

②  $CF = 0.2 + (\times 0.8) \xrightarrow{\text{reducing slope}}$

③  $CF = 0.3 + (\times 0.6) \xleftarrow{\text{reducing slope}}$

④  $CF = 0.4 + (\times 0.4) \xrightarrow{\text{optimum}}$

⑤  $CF = 0.6 + (\times 0.3) \xleftarrow{\text{reducing slope}}$

$m_2 \rightarrow 1 \xrightarrow{\text{reduce}} 0$

$$\check{y} = m_1 n_1 + [m_2 n_2] + m_3 n_3 + c$$

Lasso Regression

Feature Selection Techniques

# Feature Selection Technique

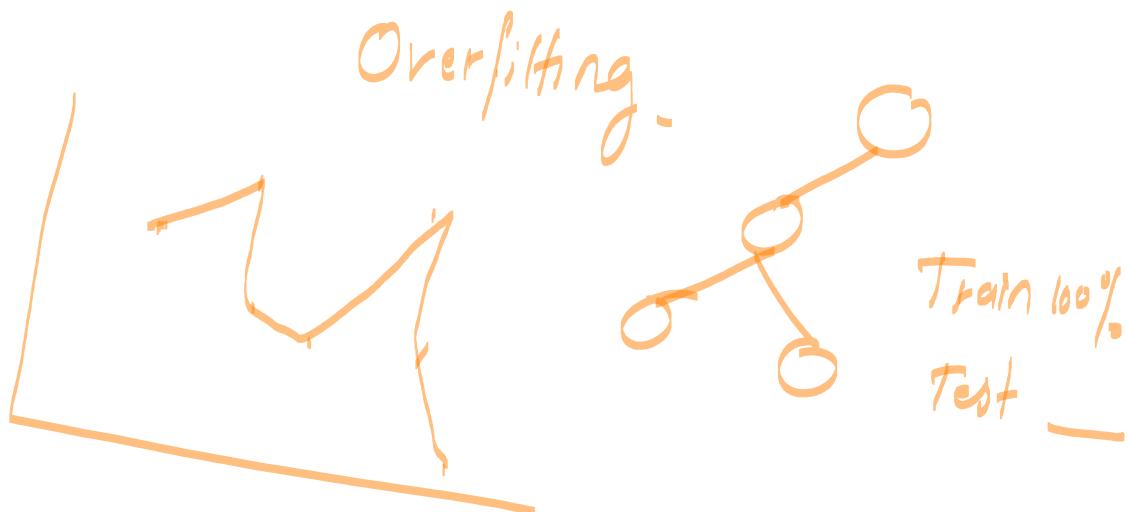
## Embedded

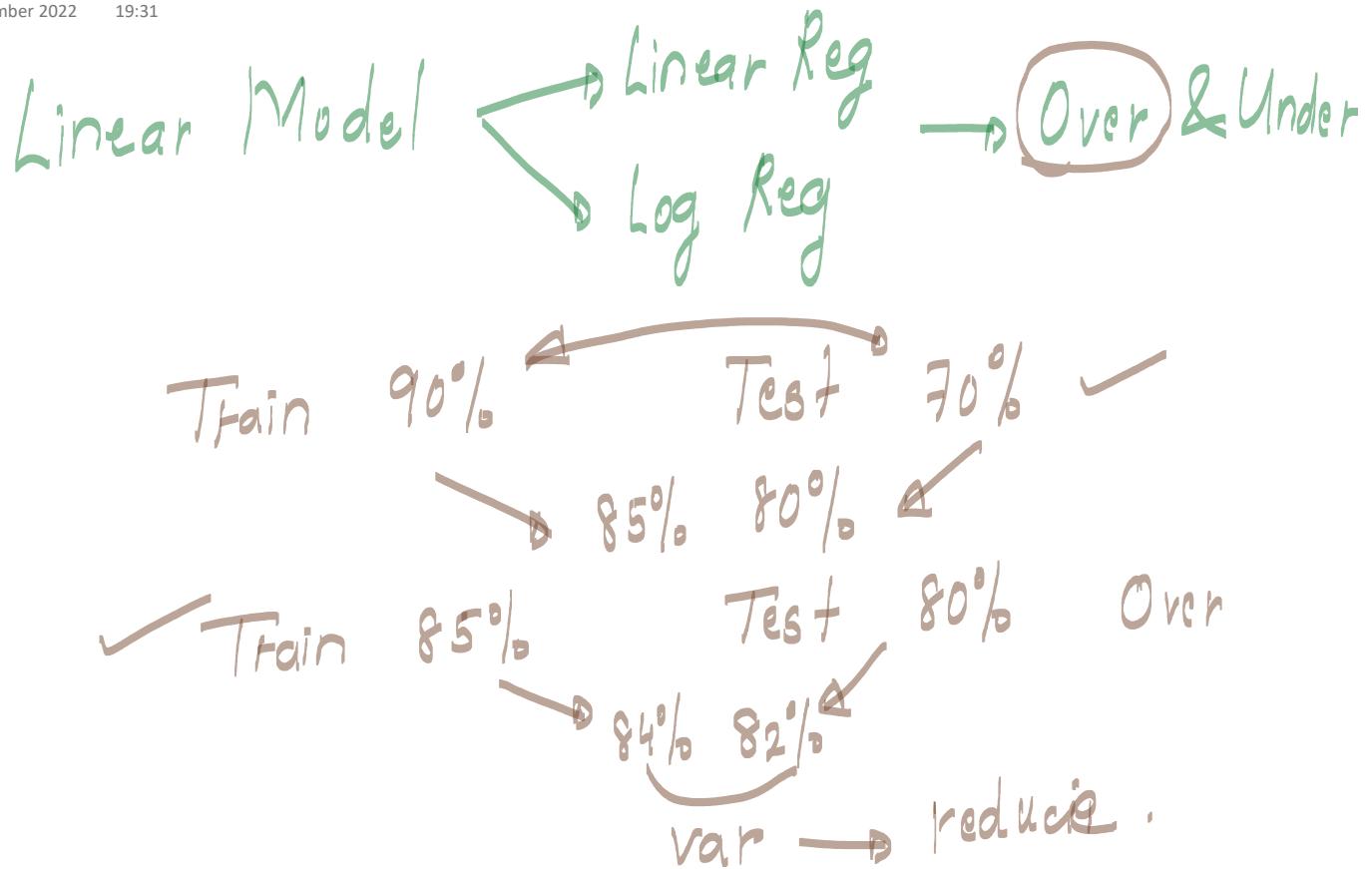
Correlation → Before Training

### \* Ridge Regression

$$\left\{ \begin{array}{l} CF = SSE + \lambda \underline{m^2} \rightarrow L_2 \text{ Reg.} \\ CF = SSE + \lambda \underline{|m|} \rightarrow \text{Ridge} \\ CF = SSE + \lambda \underline{|m|} \rightarrow \text{Lasso} \\ CF = SSE + \lambda \underline{m^2} \rightarrow L_1 \text{ Reg.} \end{array} \right.$$

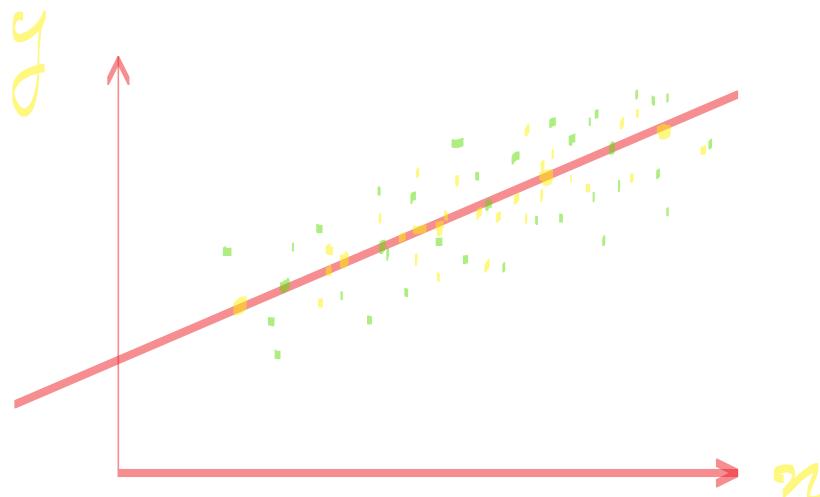
$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 \quad [m \neq 0]$$



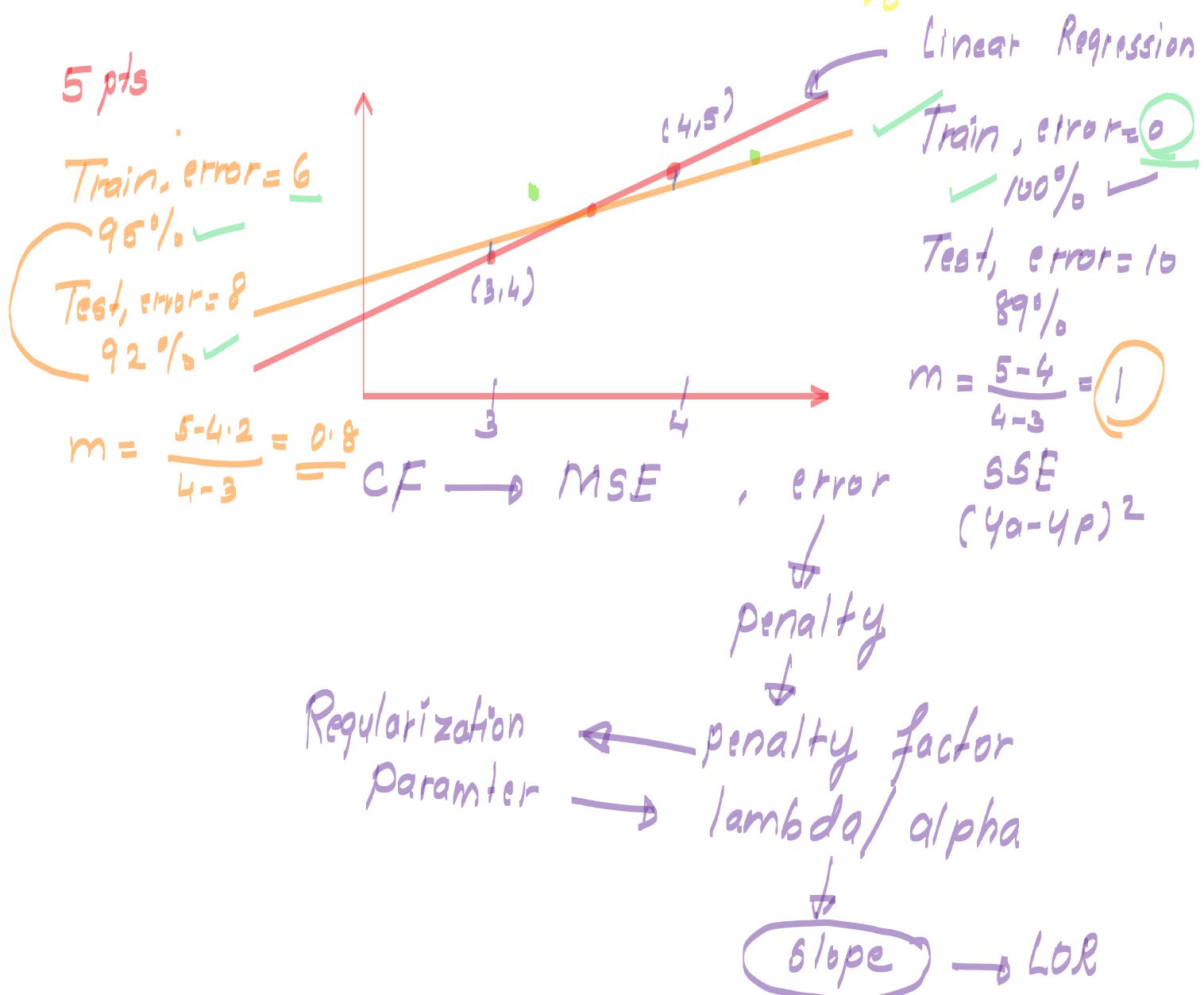


\* Lasso Regression (L<sub>1</sub> Regularization)

↳ Least Absolute Shrinkage and Selection Operator.



Train → 90%  
Test → 88%



## \* Lasso Regression

$$CF = SSE + \lambda |m|$$

$$C.F = SSE + \lambda |m|$$

- ①  $C.F = 10 + 1 \times 1.4$   $m = 1.4$
- $$= \underline{\underline{1.4}} \rightarrow \text{error (LOR)}$$
- ↓
- ②  $C.F = 0.2 + 1 \times 0.8$   $m = 0.8$
- $$= \underline{\underline{1}}$$
- ↓
- ③  $C.F = 0.3 + 1 \times 0.6$   $m = 0.6$
- $$= \underline{\underline{0.9}}$$
- ↓
- ④  $E.F = 0.4 + 1 \times 0.3$   $m = 0.4$
- $$= \underline{\underline{0.8}}$$
- ↓
- ⑤  $C.F = 0.6 + 1 \times 0.3$   $m = 0.3$
- $$= \underline{\underline{0.9}}$$

$$y = m_1 n_1 + \underline{m_2 n_2} + m_3 n_3 + C$$

↳ Lasso Reg → Feature Selection  
 Ensembled → Random Forest

## \* Ridge Regression

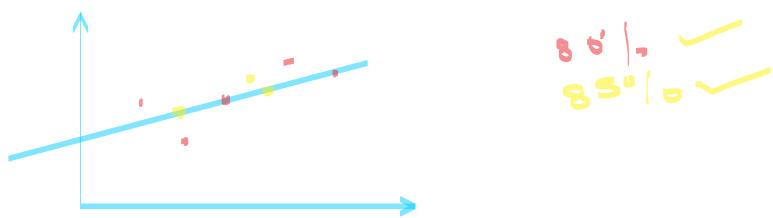
$$C.F = SSE + \lambda m^2$$

$\lambda$   $L_2$   
 Regularization

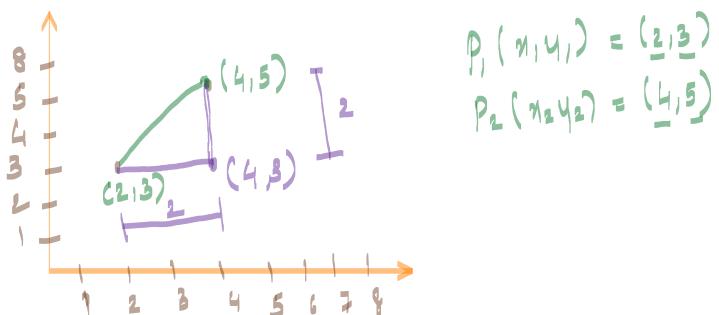
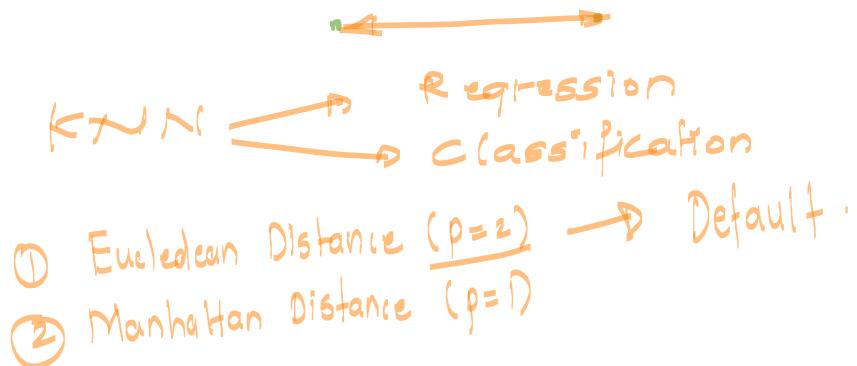
↓

non-zero

$$y = m_1 n_1 + m_2 n_2 + m_3 n_3$$



## \* Distance Based Algorithm



① Euclidean distance

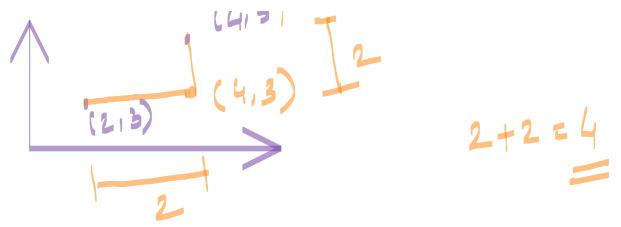
$$\begin{aligned} ED &= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \\ &= \sqrt{(4-2)^2 + (5-3)^2} = \sqrt{4+4} = \sqrt{8} = 2\sqrt{2} \end{aligned}$$

$$p_1(5,1), p_1(4,5)$$

$$ED = \sqrt{(4-5)^2 + (5-1)^2} = \sqrt{1+16} = \sqrt{17} = 4.1$$

## ② Manhattan distance





MD > ED

$$P_1(2, 4, 3) \quad , \quad P_2(6, 7, 8)$$

$$ED = \sqrt{16 + 9 + 25} = \sqrt{50} = 5\sqrt{2}$$

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

$$\begin{aligned} MD &= |x_1 - x_2| + |y_1 - y_2| + |z_1 - z_2| \\ &= |6 - 2| + |7 - 4| + |8 - 3| \\ &= 4 + 3 + 5 = 12 \end{aligned}$$

MD > ED

### \* Minkowski Distance

$$D_1 \quad D_2 \\ (\sum (D_1 - D_2)^p)^{1/p}$$

$$ED \Rightarrow p=2 \\ \sum ((D_1 - D_2)^2)^{1/2}$$

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

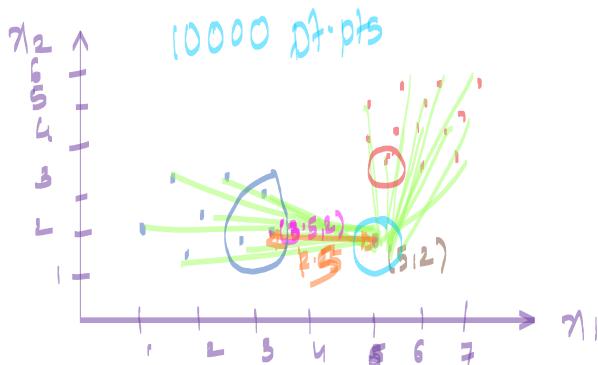
$$MD \Rightarrow p=1 \\ \sum |(D_1 - D_2)|^{1/1}$$

$$\frac{|MD|}{MD} = \frac{|D_1 - D_2|}{|x_1 - x_2| + |y_1 - y_2|}$$

$$INDEX = \frac{1}{W_1 - W_2}$$

$$INDEX = |W_1 - W_2| + |Y_1 - Y_2|$$

KNN



$$\begin{aligned} ED &= \sqrt{(5-3.5)^2 + (2-2)^2} \\ &= \sqrt{1.25} = \underline{\underline{1.25}} \\ &= \underline{\underline{1.5}} \end{aligned}$$



	\$n_1\$	\$n_2\$	\$y\$
1	1	1	B
2	2	3	D
3	3	6	R
4	4	2	B
5	5	2	B
6	6	5	R
7	7	2	B
8	8	8	R

$\boxed{K=5}$  ✓ Odd.  
 No. of neighbours  
 $B \rightarrow 3$  ✓ 4p  
 $R \rightarrow 2$   
 $B \rightarrow 4$  ✓ 4p  
 $R \rightarrow 2$

$\boxed{K=6}$   
 $B = 3 \rightarrow$  No. of Dt.pt = 5 ✓  
 $R = 3 \rightarrow$  No. of Dt.pt = 3

Tens	Units	Thousands
10	2	15000
20	4	20000
30	6	30000

$P_1 (10, 2, 15000)$   
 $P_2 (20, 4, 20000)$   
 $P_3 (30, 6, 30000)$

$$\begin{array}{ccc} 20 & 4 & 20000 \\ 30 & 6 & 30000 \\ \downarrow 25 & \downarrow 5 & \downarrow 25000 \end{array} \quad \rightarrow \begin{array}{c} P_3(30, 6, 30000) \\ \downarrow 4 \quad \downarrow 5 \quad \downarrow 25000 \\ P_L(25, 5, 25000) \end{array}$$

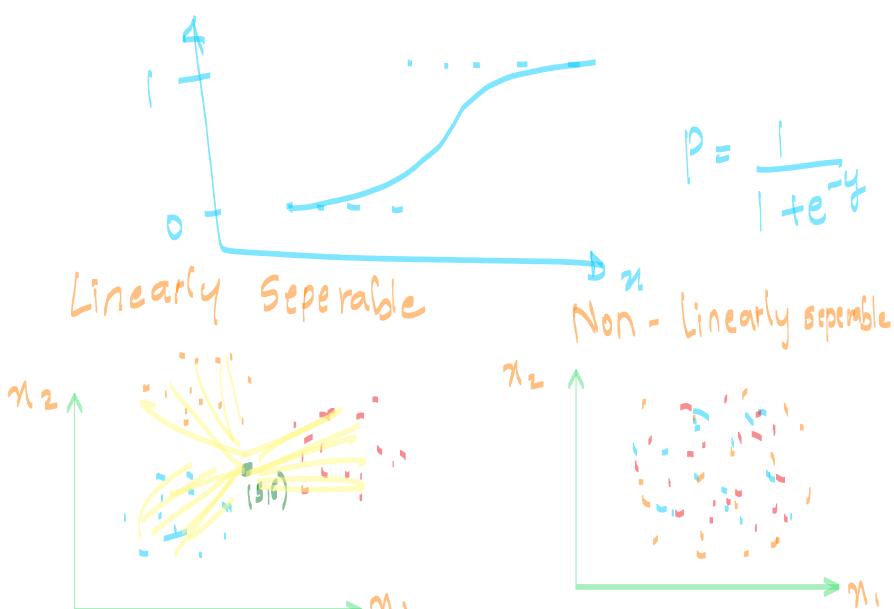
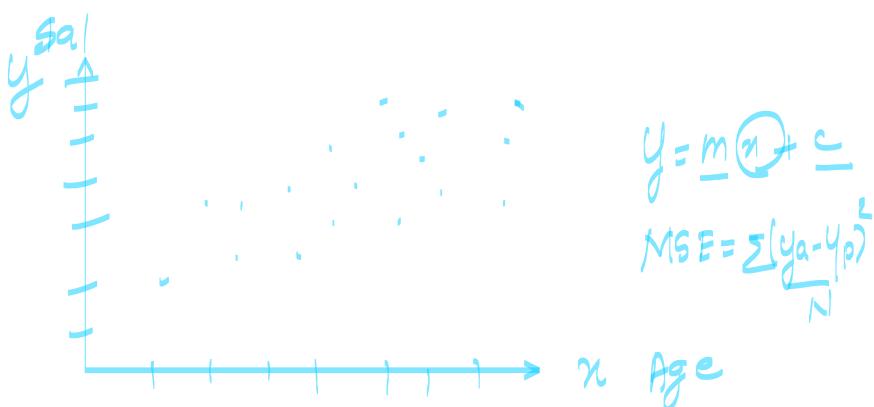
$$ED = \sqrt{(26 - 10)^2 + (5 - 2)^2 + (25000 - 15000)^2}$$

$$= \sqrt{225 + 9 + 1000000}$$

$$= \sqrt{1000234} =$$

## ✓ Feature Scaling.

[ 41    61    101    20    15 ]



① It will save the data (training)

(1) It will do the following:

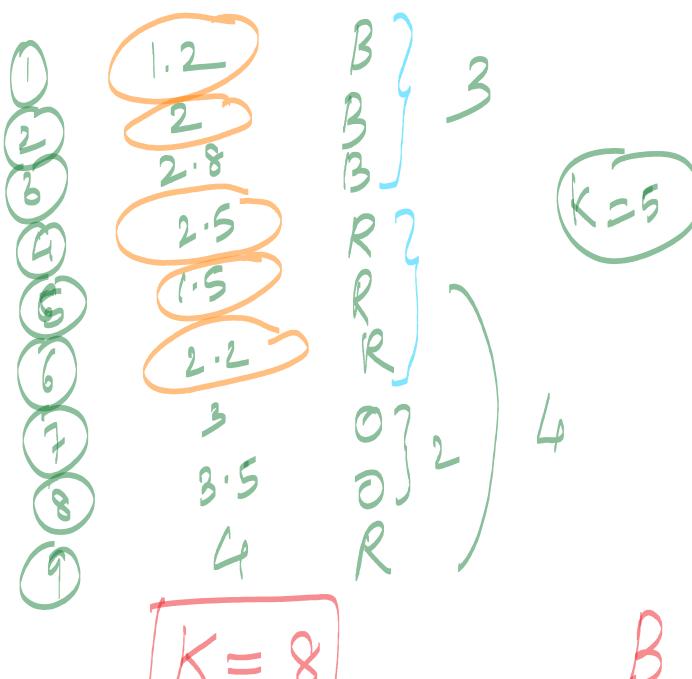
$n_1$	$n_2$	$B$	$R$	Training
2	10			
3	5	O		
4	8	O		
5	6	?		Test

(2) It will find distance between Test pt. & all the training dt pts.

(3) Arrange distances in ascending

(4) nearest neighbours,  $\boxed{K=5}$

(5) Majority classification  $\rightarrow y_p$



$B \rightarrow 2$   
 $R \rightarrow 3$   
 $O \rightarrow 0$

$$y_p = \begin{cases} B = 3 \\ R = 3 \end{cases} \quad \frac{\text{No. of Dt pts}}{\text{No. of Dt pts}} = \frac{3}{4}$$

$$O = 2$$

$$\boxed{K=8}$$

$$- R = 3 \quad \text{No. of pts} = 4$$

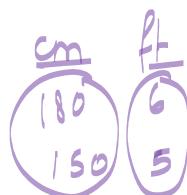
$$K=8$$

$y_p$

$B = 3$     No. of pts = 4  
 $R = 3$     No. of pts = 4

$\{ 1.2 + 1.8 + 2 = 5.0 \}$   
 $2 + 1.2 + 2.5 = 5.7$

km	m
2	2000
3	3000
4	4000



## \* KNN Regression



Distance  $n_1(\text{Age})$   $n_2(\text{Exp})$   $y$  ( $\text{Sal} / k$ )

1.2	30	4	30
2	40	5	40
3	50	6	50
1.8	20	2	25
4	25	2.5	30
3.2	45	5	55
4.1	35	4	50

$K=5$

Training

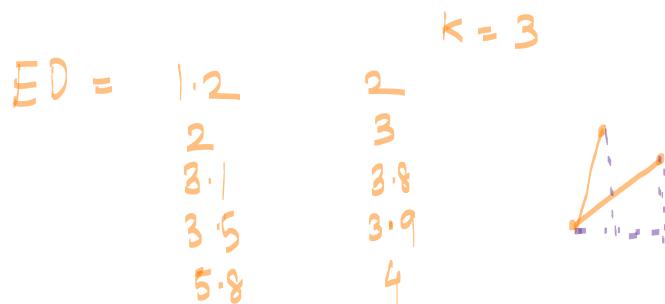
$$\begin{array}{c} 35 \\ 3 \\ ? \\ y_p = 40 \end{array}$$

$$y_p \text{ mean} = \frac{30 + 40 + 50 + 25 + 55}{5} = 40$$

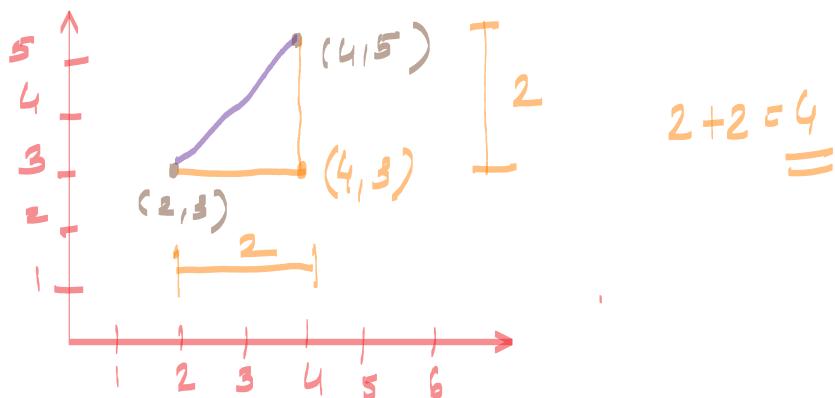
$$\begin{array}{ll}
 X_i & X_{\text{new}} \\
 \hline
 1 & 148 \quad 0.743 \\
 2 & 85 \quad 0.427 \\
 3 & 183 \quad 0.91 \\
 4 & 150 \quad 0.753 \\
 \hline
 \end{array}$$

Norm  $\rightarrow \frac{X_i - X_{\min}}{X_{\max} - X_{\min}}$   
 $= \frac{148 - 85}{183 - 85}$   
 $= \underline{\underline{0.056}}$

Sensitive  
to Outliers



\* KNN  Continuous  
Classification Categorical



\* Euclidean Distance  $P_1(n_1, y_1) = (2, 3)$ ,  $P_2(n_2, y_2) = (4, 5)$

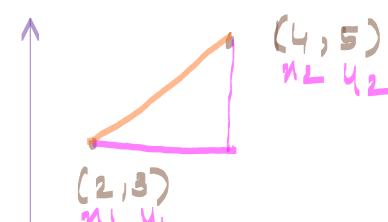
$$\begin{aligned} ED &= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \\ &= \sqrt{(4-2)^2 + (5-3)^2} = \sqrt{4+4} \\ &= \sqrt{8} = \underline{\underline{2\cdot 8}} \end{aligned}$$

$P_1(5, 1)$ ,  $P_2(4, 5)$

$$ED = \sqrt{(4-5)^2 + (5-1)^2} = \sqrt{1+16} = \sqrt{17}$$

\* Manhattan Distance

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



$$\begin{aligned} MD &= |2-4| + |3-5| \\ &= 2 + 2 \\ &= \underline{\underline{4}} \end{aligned}$$

$P_1(n_1, y_1, z_1) = (2, 4, 3)$ ,  $P_2(n_2, y_2, z_2) = (6, 7, 8)$

**MD > ED**

$$\rho_1(2, 4, 3) \quad , \quad \rho_2(6, 7, 8)$$

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

$$= \sqrt{(6-2)^2 + (7-4)^2 + (8-3)^2}$$

$$= \sqrt{16 + 9 + 25} = \sqrt{50} = \underline{\underline{5\sqrt{2}}}$$

$$MD = |6-2| + |7-4| + |8-3|$$

$$= 4 + 3 + 5 = \underline{\underline{12}}$$

$$\boxed{MD > ED}$$

## ① Minkowski Distance.

$$(\sum |D_1 - D_2|^P)^{1/P}$$

$P=2 \rightarrow$  Euclidean Distance

$$ED = (\sum |D_1 - D_2|^2)^{1/2}$$

$$= \sqrt{\sum_{\substack{|D_1 - D_2|^2 \\ n_1 u_1 \\ n_2 u_2}}}$$

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

$P=1 \rightarrow$  Manhattan Distance

$$MD = (\sum_{\substack{|D_1 - D_2| \\ n_1 u_1 \\ n_2 u_2}})^{1/1}$$

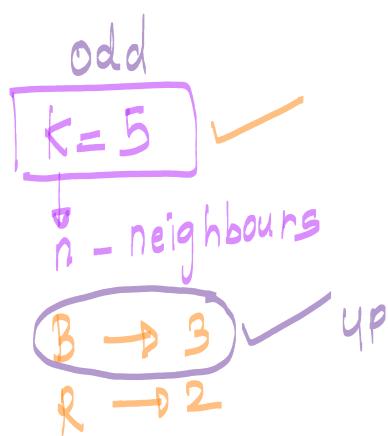
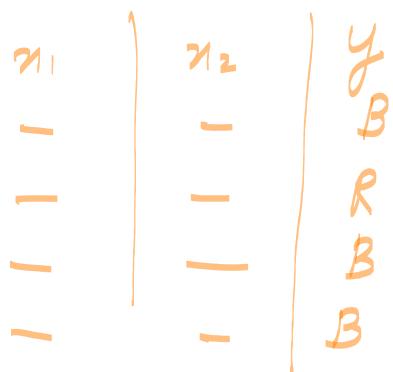
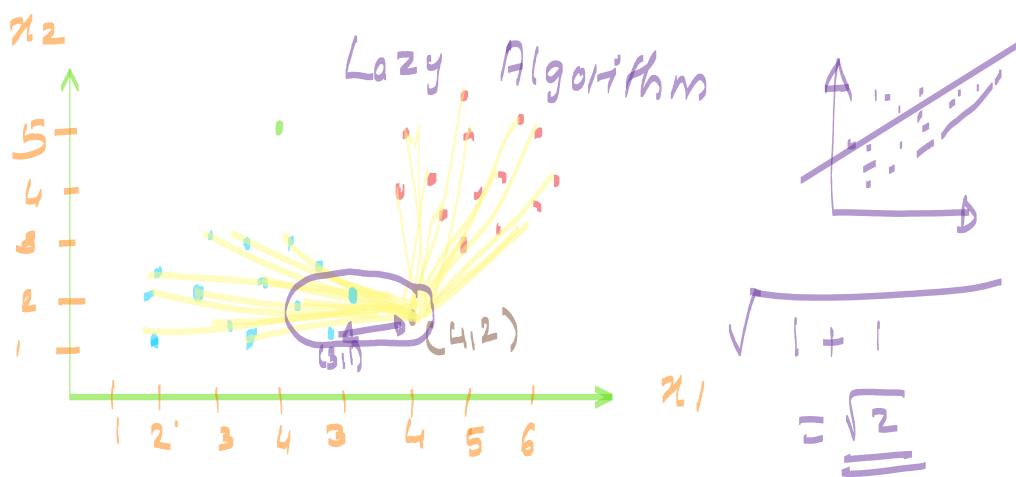
$$\boxed{MD = |x_1 - x_2| + |y_1 - y_2|}$$

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

$$|MD \geq ED|$$

$P=2 \rightarrow \text{Default}$

## kNN



$$\begin{matrix} \text{Get} \\ \text{Data} \end{matrix} \quad \begin{matrix} \text{8.5} \\ 1.2 \\ 1.8 \end{matrix} \quad \begin{matrix} \text{R} \\ \text{P} \end{matrix} \quad \begin{matrix} \text{y} \\ 3 \end{matrix} \quad R \rightarrow 2$$

\* Distance Based Algorithm.

$$K=6$$

$$\begin{matrix} B \\ R \end{matrix} \rightarrow \begin{matrix} L \\ 2 \end{matrix} \quad \checkmark \quad y_p$$

$$\begin{matrix} 3 \\ 3 \end{matrix} \rightarrow \begin{matrix} B \\ R \end{matrix} \rightarrow \begin{matrix} \text{No. of Data points} \\ \text{No. of Data points} \end{matrix} \rightarrow \begin{matrix} 5 \\ 3 \end{matrix}$$

$$\begin{matrix} 10 \\ 6 \\ 3 \\ 2 \\ 3 \end{matrix} \quad K=5$$

$$\begin{matrix} R \\ B \end{matrix} \rightarrow \begin{matrix} 2 \\ 2 \end{matrix}$$



$$\begin{matrix} 3 \rightarrow R \\ 1 \rightarrow P \\ 1 \rightarrow B \end{matrix}$$

$$\begin{matrix} K=5 \\ K=6 \end{matrix} \quad \begin{matrix} 3 \\ 2 \\ 3 \\ 3 \end{matrix}$$

\* Feature Scaling.

$$\begin{matrix} n_1 & n_2 & n_3 \\ 10 & 1 & 15000 \\ 20 & 2 & 20000 \\ 30 & 3 & 30000 \end{matrix} \quad \begin{matrix} p_1 & (10, 1, 15000) \\ p_2 & (20, 2, 20000) \\ p_3 & (30, 3, 30000) \end{matrix}$$

$\begin{matrix} 10 \\ 25 \end{matrix}$       3      30000       $P_3$   
 $\begin{matrix} 25 \\ 4 \end{matrix}$       25000       $T(25, 4, 25000)$

$$\begin{aligned}
 & \sqrt{(25-10)^2 + (4-1)^2 + (25000-15000)^2} \\
 & = \sqrt{15^2 + 3^2 + 10000^2} \\
 & = \sqrt{225 + 9 + 10000000} \\
 & = \sqrt{1000234}
 \end{aligned}$$

np.log

$$\begin{bmatrix} 1, 2, 3, \dots, 4, 9, 15, 20 \\ 0, 0.1, 0.3, \dots, 0.5, 0.8, 0.9, 0.1 \end{bmatrix}$$

## \* Feature Scaling

- ① Normalization  $\rightarrow$  MinMax Scalar
- ② Standardization  $\rightarrow$  Standard Scaler.

## \* Gradient Descent Algorithm

$$\begin{array}{c}
 \text{y} = \Theta x + \Theta \\
 \frac{\partial C}{\partial \Theta} \\
 \text{MSE} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2
 \end{array}$$

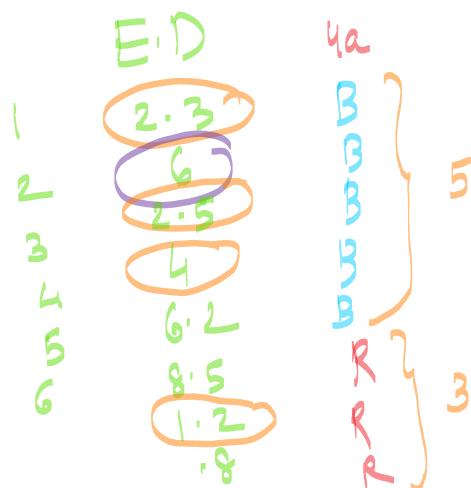


$$P = \frac{1}{1 + e^{-y}}$$

$$\circ \quad \begin{array}{c} t \\ \downarrow \end{array} \rightarrow \boxed{y = mx + c}$$

$$\checkmark LL = -\log P$$

$$\checkmark LL = -\log(1-P)$$



Ascending

1.2	R
1.8	R
2.3	R
2.5	R
4	R
6	R
6.2	R
8.5	R
8.8	R

$k=5$

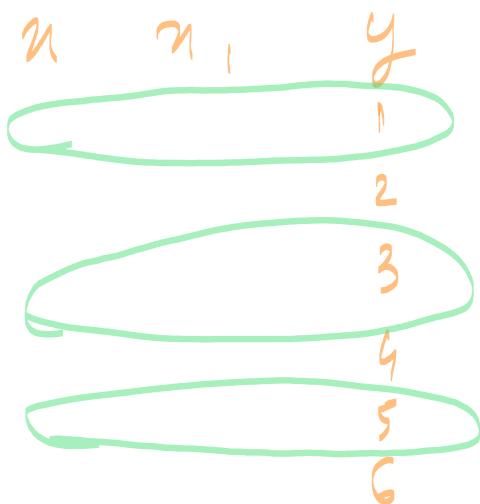
$y_p$

$B \rightarrow 3$

$R \rightarrow 2$

No. of neighbours

$n$  - neighbours

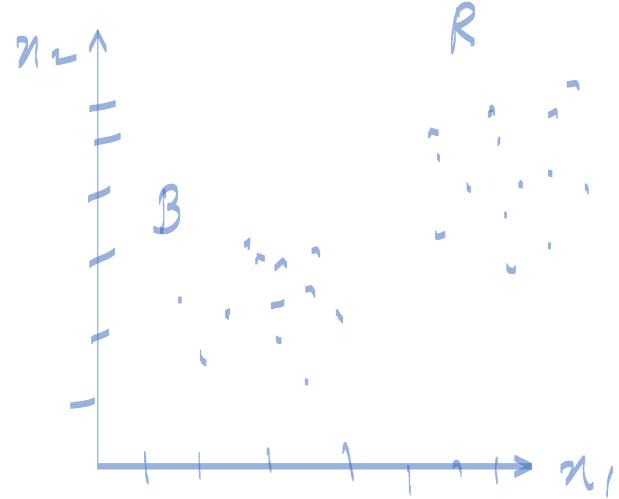


$$\boxed{k=3}$$

$$\frac{1+3+5}{3} = \frac{9}{3} = 3$$

① It will scatter plot and save training data.

$n_1$	$n_2$	$y$
1	1	B
1	1	R
1	2	B
1	2	R
1	3	B
1	3	R
1	4	B
1	4	R
1	5	B
1	5	R
1	6	B
1	6	R
1	7	B
1	7	R
1	8	B
1	8	R
1	9	B
1	9	R
1	10	B
1	10	R
1	11	B
1	11	R
1	12	B
1	12	R
1	13	B
1	13	R
1	14	B
1	14	R
1	15	B
1	15	R
1	16	B
1	16	R
1	17	B
1	17	R
1	18	B
1	18	R
1	19	B
1	19	R
1	20	B
1	20	R
1	21	B
1	21	R
1	22	B
1	22	R
1	23	B
1	23	R
1	24	B
1	24	R
1	25	B
1	25	R
1	26	B
1	26	R
1	27	B
1	27	R
1	28	B
1	28	R
1	29	B
1	29	R
1	30	B
1	30	R
1	31	B
1	31	R
1	32	B
1	32	R
1	33	B
1	33	R
1	34	B
1	34	R
1	35	B
1	35	R
1	36	B
1	36	R
1	37	B
1	37	R
1	38	B
1	38	R
1	39	B
1	39	R
1	40	B
1	40	R
1	41	B
1	41	R
1	42	B
1	42	R
1	43	B
1	43	R
1	44	B
1	44	R
1	45	B
1	45	R
1	46	B
1	46	R
1	47	B
1	47	R
1	48	B
1	48	R
1	49	B
1	49	R
1	50	B
1	50	R

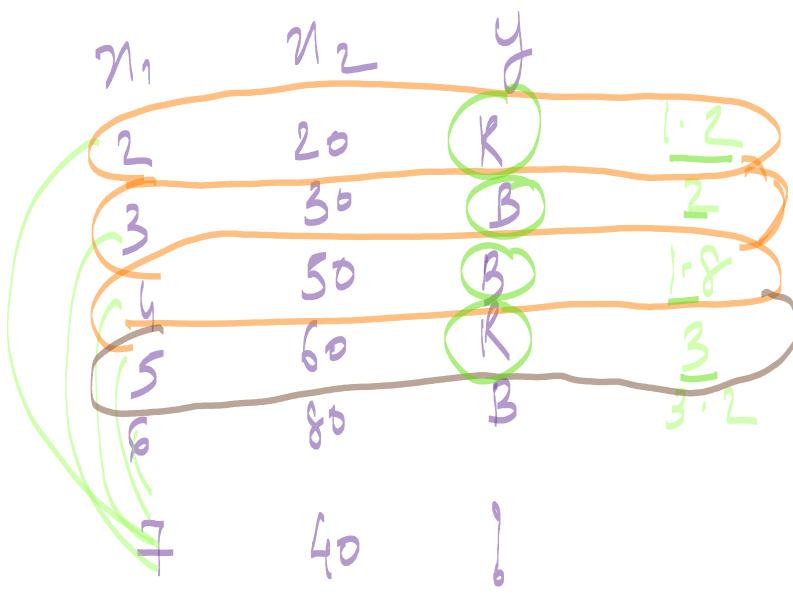


② Testing pt. for prediction

Training data



$n_1$	$n_2$	$y$
2	20	R
3	30	B
5	50	R
6	60	B
8	80	R
9	90	B
10	100	R
11	110	B
12	120	R
13	130	B
14	140	R
15	150	B
16	160	R
17	170	B
18	180	R
19	190	B
20	200	R
21	210	B
22	220	R
23	230	B
24	240	R
25	250	B
26	260	R
27	270	B
28	280	R
29	290	B
30	300	R



$K=3$

$n$ -neighbors

$B \rightarrow 2 \checkmark$   
 $R \rightarrow 1$

$K=4$   
Training data

$$4p \quad \checkmark \quad R = 2 \\ B = 2$$

١٢ - ٦

## Training data

? No · v· d· p· ts

No of off. pts 60

# Setosa

# Versicolor

ၫၦ

$$\begin{array}{ccc} 1 & 0 & \} \\ & 0 & \end{array} \quad \begin{array}{c} R \rightarrow 2 \\ B \rightarrow 2 \end{array}$$

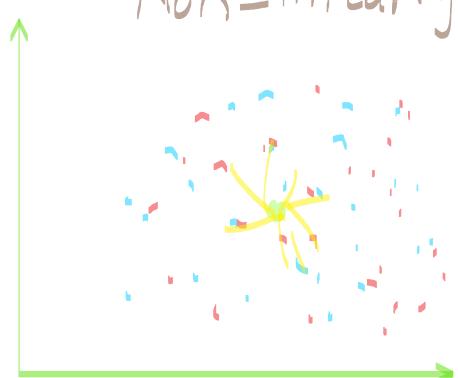
$$\Rightarrow 1 \cdot 2 + 3 = 4 \cdot 2$$

$$\rightarrow 2+8=3\cdot 8$$

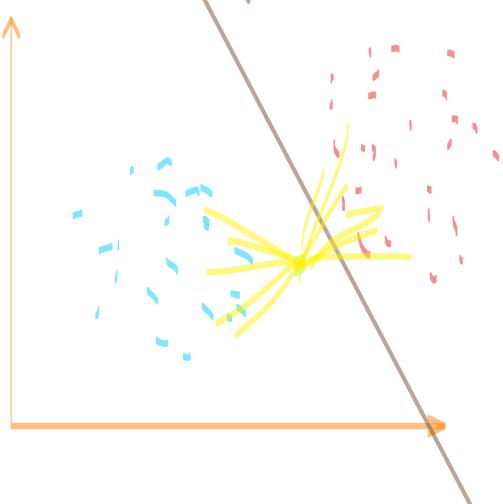
48

# \* Regression

## Non-linearity



## Linearly separable



Age	Exp	Sal	
Y1	Y2	Y3	
20	2	20k	2
30	3	30k	2.3
40	4	40k	3
50	5	45k	5
60	6	50k	5.4
70	7	60k	4
80	8	80k	4.5
35	3.5	9	11

$y_p$  → mean of  $y$  value  
of  $n$ -neighbours

$$y_p = 20 + 30 + 40 + 60 + 80$$

$$y_p = \frac{20 + 30 + 40 + 60 + 80}{5} \\ = \underline{\underline{45 \text{ K}}}$$

## \* Feature Scaling

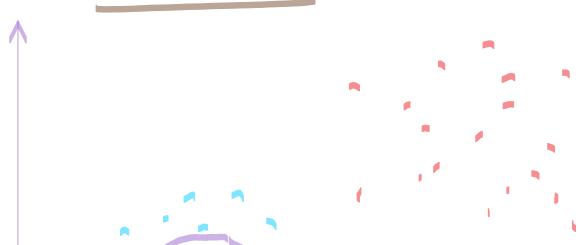
	$x_1$	$x_2$	$x_3$
1	100	15000	
2	200	20000	
3	300	30000	

$$\sqrt{(2-1)^2 + (200-100)^2 + (20000-15000)^2}$$

	<u>km</u>	<u>m</u>	<u>cm</u>	<u>ft</u>	<u>impact of higher scale values</u>
1	2000	15000	180	6	
2	3000	20000	150	5	
3	4000	30000	120	4	

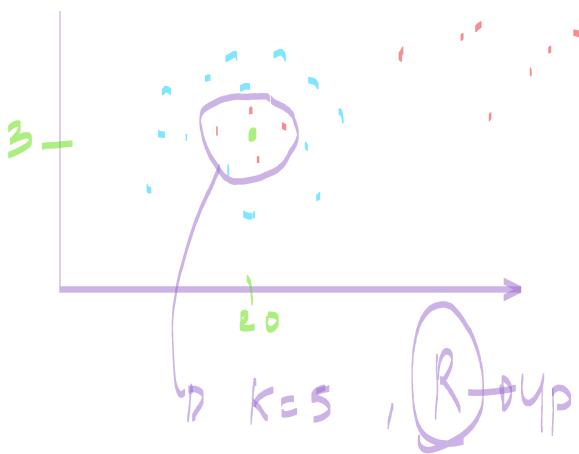
## \* Outliers

Sensitive to outliers



Not Sensitive





	$X_i$	$X_{\text{new}}$
0	50	0.409
1	660	0.540
2	64	0.524
3	66	0.540
4	40	0.327

0.016

$$x_{\min} = 40$$

$$x_{\max} = 660$$

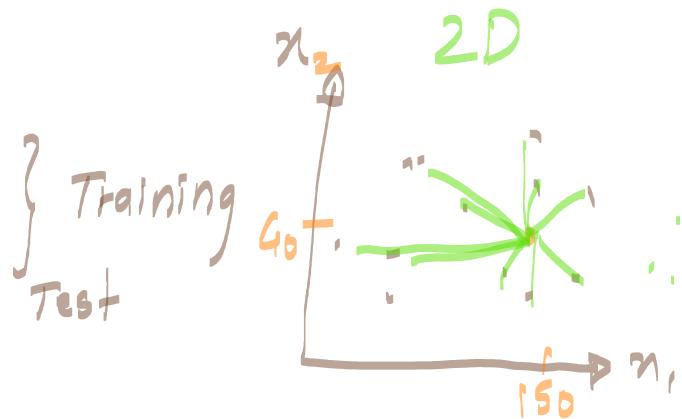
$$\frac{50 - 40}{660 - 40} = 0.016$$

From  
[http://localhost:8888/notebooks/12\\_29%20KNN/K%20-%20Nearest%20Neighbour.ipynb](http://localhost:8888/notebooks/12_29%20KNN/K%20-%20Nearest%20Neighbour.ipynb)

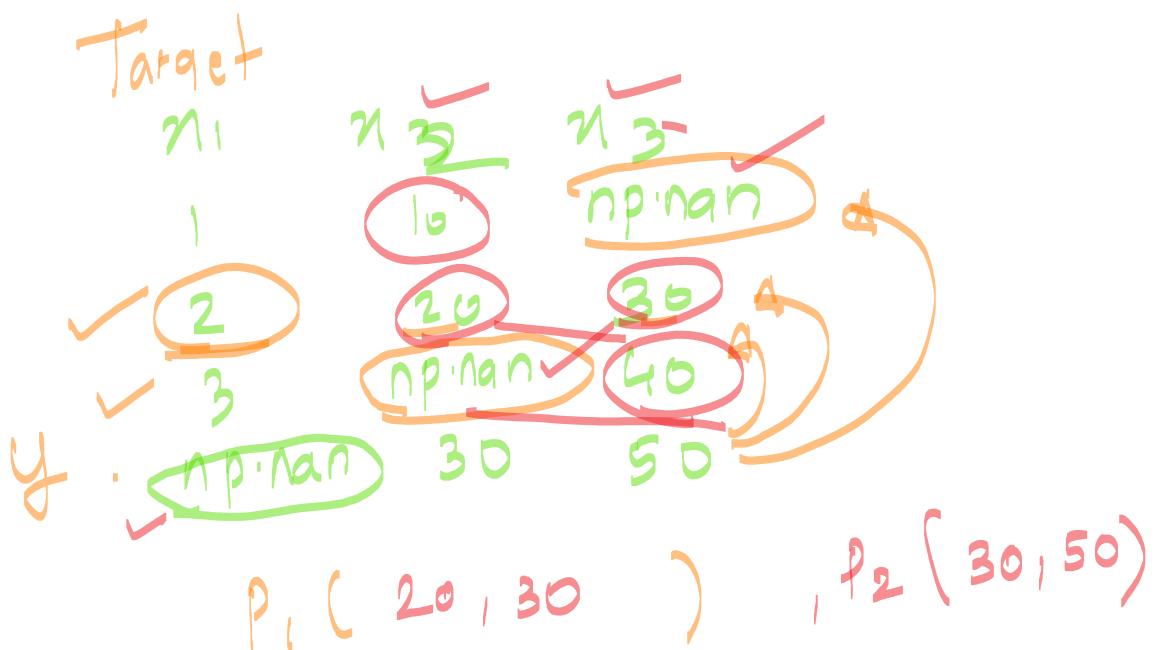
## KNN IMPUTER

29 December 2022 20:40

	$n_1$	$n_2$	$y$						
1	148	50	50	35	0	33.6	0.627	50	1
2	85	66	66	29	0	26.6	0.351	31	0
3	183	64	64	0	0	23.3	0.672	52	1
4	150	66	66	23	94	28.1	0.167	21	0
5	150	40	?	35	168	43.1	2.288	33	1



$$\frac{50+66}{2} \rightarrow \text{missing value} = 58$$



$$ED = \sqrt{10^2 + 20^2}$$

$$= \sqrt{500}$$

$$* \text{nan\_ED} = \sqrt{\omega_1 \times ED}$$

$$= \sqrt{\omega_1 \times \sqrt{(n_2 - n_1)^2 + (y_2 - y_1)^2}}$$

$$= \sqrt{\omega_1 \times 500}$$

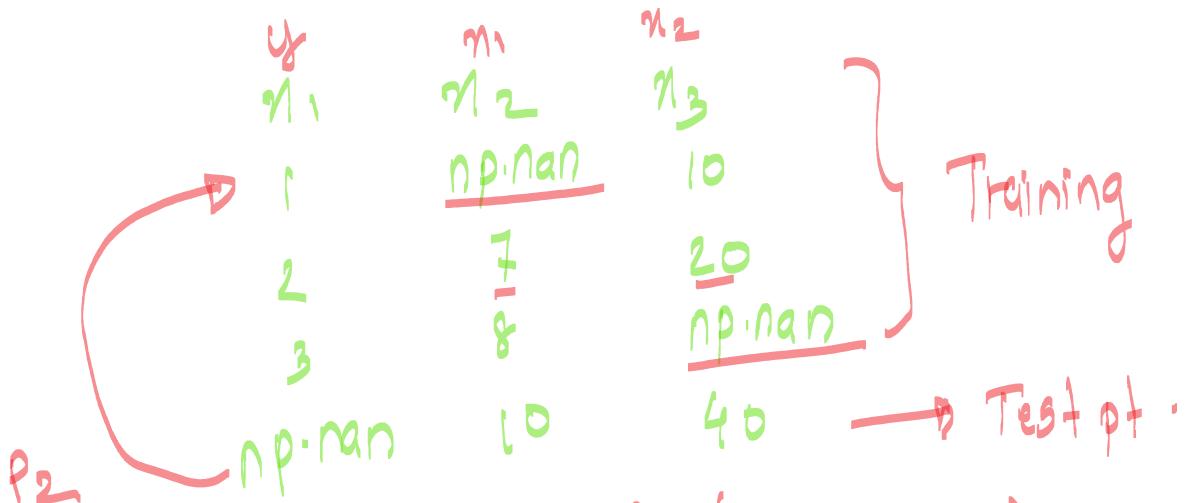
$$= \sqrt{\frac{\text{Total of coordinates}}{\text{Present cordin}} \times ED}$$

$K=3$  $P_1(148, 50), P_2(150, 74)$ 

	$n_1$	$n_2$	$y$
Glucose	BloodPressure	Target	
148	50	50	
85	66	66	
183	64	64	
150	66	66	
150	40	40	
150	74	?	

Training  
PL, Testing

$$\frac{50 + 66 + 64}{3} = \underline{\underline{66}}$$

 $P_2(10, 40), P_1(7, 8)$  $P_2(10, 40), P_1(20, 10)$ 

$$ED = \sqrt{(10 - 7)^2 + (40 - 20)^2}$$

$$= \sqrt{\underline{\underline{9}} + \underline{\underline{400}}} = \underline{\underline{22}}$$

 $P_2(10, 40)$   
 $P_1(8, np\cdot nan)$  $\text{Non-Euclidean} \rightarrow \sqrt{wt.} \times ED$ 

$$= \sqrt{wt.} \times \sqrt{(n_2 - n_1)^2 + (y_2 - y_1)^2}$$

$$= \sqrt{\omega_t \times ((x_2 - x_1)^2 + (y_2 - y_1)^2)}$$
$$= \sqrt{\frac{\text{Total coordinates}}{\text{Present coordinates}}} \times$$

## Outliers

02 January 2023 08:06

$Z\text{-score}$   
 $+3 >$   
 $-3 <$

