

Linear Regression      Intralit Slope

$$y = b + ax$$

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

Slope

$$b = \frac{1}{n} \left( \sum_{i=1}^n y_i - a \sum_{i=1}^n x_i \right)$$

Intralit

x	y
-2	-1
1	1
3	2

$$\begin{array}{cccc} x & y & xy & x^2 \\ \sum x_i = 2 & \sum y_i = 2 & \sum xy = 9 & \sum x^2 = 14 \end{array}$$

$$a = \frac{23}{38} \quad b = \frac{5}{19}$$

$$y = \frac{5}{19} + \frac{23}{38} x$$

$$x = 4$$

$$y = \frac{5}{19} + \frac{23}{38} \times 4 = \cancel{\dots} ?$$

$$= 0.26 + 0.82 \times 4$$

$$= 0.26 + 3.29$$

$$= 3.54$$

use the linear reg to find the sales of company

x	y	t = x - 2018	yt/t
2018	12	0	
2019	19	1	
2020	29	2	
2021	37	3	
2022	45	4	

$$y = 11.4 + 8.4t$$

$$y = 11.4 + 8.4(x - 2018)$$

## Multiple Linear Regression:-

multiple linear regression formula:-

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

$$b_0 = \bar{y} - b_1 \bar{x}_1 - b_2 \bar{x}_2$$

$$b_1 = \frac{[(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)]}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

if there are three variable

$$b_2 = \frac{[(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)]}{[(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2]}$$

$\hat{y}$	$x_1$	$x_2$
140	60	22
155	62	25
159	67	24
179	70	20
192	71	15
200	72	14
212	75	14
215	78	11
$\bar{y} = 181.5$	$69.38$	$18.13$
$S_y = 145.2$	$555$	$145$

if multiple linear  
 Regression Need to  
 calculate intercept only  
 no need to calculate  
 slope

→ more than one intercept  
 → How to calculate  $b$ ?

$$\text{Calculate } x_1^2 = 38767$$

$$x_2^2 = 2823$$

$$x_1 y = 101895$$

$$x_2 y = 25364$$

$$x_1 x_2 = 9859$$

$$\text{regression sum } (\sum x_1^2) \frac{\sum x_1^2 - (\sum x_1)^2}{n} = \frac{38767 - (555)^2}{8} = 263.875$$

$$(\sum x_2^2) \frac{\sum x_2^2 - (\sum x_2)^2}{n} = 194.88$$

$$(\sum x_1 y) = \frac{\sum x_1 y - \sum x_1 \sum y}{n} = 1162.5$$

$$\sum x_1 x_2 = \frac{\sum x_1 x_2 - \sum x_1 \sum x_2}{n} = -200.375$$

Calculate  $b_0, b_1, b_2$

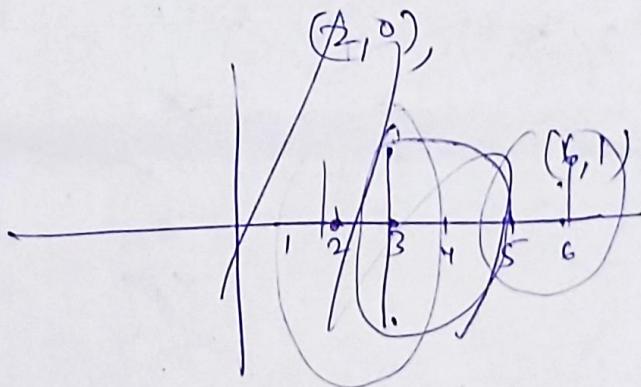
$$b_1 = 3.148$$

$$b_2 = -1.656$$

$$b_0 = 181.5 - 3.148(69.38) - (-1.656)(18.125)$$
$$= -6.867$$

$$b_1 = 3.148 \quad b_2 = \frac{-6.867}{-1.656}$$

$$\hat{y} = -6.867 + 3.148x_1 - 1.656x_2$$



a) Linear Regression - It establish a relationship b/w dependent var of one or more independent variables using best fit straight line

Least squares  
eqn. eqn.

$$y = b + ax \rightarrow \text{Independent}$$

dependent

$\therefore a$  = slope of line

$\therefore b$  = Intercept

$\therefore n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i$	(Intercept)
$\therefore a = \frac{1}{n} \left( \sum_{i=1}^n y_i - a \sum_{i=1}^n x_i \right)$	(slope) $n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2$

Ques 1

find the Zineau Reg best fit line for the basis of given data  $x, y$

x	y	if draw best fit Reg line for the given eqn.
-2	-1	
1	1	
3	2	$\therefore n = 3$

x	y	xy	$x^2$
-2	-1	2	4
1	1	1	1
3	2	6	9

$\Sigma x_i$	$\Sigma y_i$	$\Sigma xy$	$\Sigma x^2$
2	2	9	14

$$a = \frac{3(9) - (2)(2)}{3((14) - (2)^2)} \rightarrow$$

$$\begin{array}{|c|c|} \hline & 23 \\ \hline 38 & : 30 \Rightarrow \frac{1}{2} \\ \hline \end{array}$$

$$a \rightarrow \frac{23}{38}$$

18

$$b = \frac{1}{3} \left( 2 - \frac{23}{38}(2) \right) \rightarrow \frac{1}{3} \left( \frac{76 - 56}{38} \right) \rightarrow \frac{20}{144}$$

$$\begin{array}{l} 1 \\ 4 \\ 2 \\ 8 \\ 3 \\ \hline + \frac{1}{x} \\ \hline 6 \end{array}$$

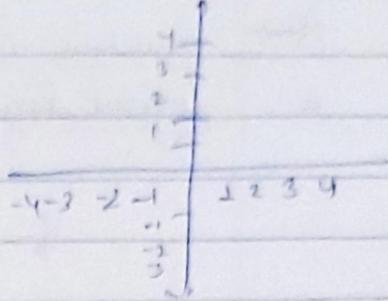
$$y = \frac{5}{19} + \frac{23}{38}x$$

Best fit line Reg eqn.

Y val of y at  $x=4$  prediction

$$y = \frac{5}{19} + \frac{23}{38}(4) \Rightarrow 51.19$$

$$y = \frac{5}{19} + \frac{92}{38} \rightarrow$$



- 1) LSRL (LRL)
- 2) Decrease?
- 3) prediction for Next

Q Use the Linear Reg to find the sales of the Company in 2012 &  
draw the trained line for the same  
(best fit)

$x$ (in years)	$y$ (sales) (mill dollar)	scby initial $t = x - 2005$	only data iter	$y = b + at \rightarrow$ $y = b + a(x - 2005)$
initial 2005	12	0	0	0
2006	19	1	19	1
2007	29	2	58	4
2008	37	3	111	9
2009	45	4	180	16
			$\sum y = 148$	$\sum t = 10$
			$\sum yt = 368$	$\sum t^2 = 30$

$$8.4 \quad a = \frac{5(368) - 10 \times 148}{5(30) - (10)^2} = \frac{1840 - 1450}{150 - 100} \Rightarrow \frac{390}{50} = 7.8$$

$$11.4 \quad b = \frac{1}{5} \left( 148 - 7.8 \times 10 \right) \Rightarrow 11.6 \quad 12.2 \\ 13.4$$

$$y = 11.4 + 8.4t \quad y = 11.4 + 8.4(x - 2005)$$

$$\text{for } 2012 \quad t = x - 2005 \Rightarrow 2012 - 2005 = 7$$

$$y = 11.6 + 8.4 \times 7$$

$$y = 11.4 + 8.4(2012 - 2005)$$

$$2012 \rightarrow \boxed{y = 70.4}$$

CLASSMATE  
Date \_\_\_\_\_  
Page \_\_\_\_\_

We can write multiple linear regression as in  
Linear Reg when  $n = \text{multiple}$  then  
 $[y = b + a_1x_1 + a_2x_2 + \dots + a_nx_n] \rightarrow y = b + a(x_1 + x_2 + \dots + x_n)$

$\star \rightarrow$  multiple Regression = multiple Linear Regression (multiple linear)  
 $\star$  Logistic Regression =  $a = \frac{a_1 + a_2 + \dots + a_n}{n}$

$$b = \frac{1}{n} \left( \sum_{i=1}^n y_i - a \sum x_i \right) , \quad b_1 = \frac{1}{n} \sum y_i - a \bar{x}_1$$

$$\boxed{y = b + ax}$$

$x_1$	$x_2$	$y$	$x_1y$	$x_2y$	$y = b + a_1x_1 + a_2x_2$	$x_1^2$
2	1	4	8	4		4
3	2	6	18	12		9
1	4	5	5	20		1
6	3	8	48	24		16
$\Sigma 12$	$\Sigma 10$	$\Sigma 23$	$\Sigma 79$	$\Sigma 60$		36
						9
						50
						30

$$\rightarrow a_1 = \frac{4(79) - (12 \times 23)}{4(50) - (10)^2} = \frac{316 - 276}{56} = \frac{40}{56} = 0.714$$

$$\rightarrow b_1 = \frac{1}{4} (23 - (0.714) 12) = \frac{14.432}{4} = 3.6$$

$$(x_2y) \rightarrow a_2 = \frac{4(60) - (10)(23)}{4(30) - (10)^2} = \frac{240 - 230}{10} = \frac{10}{20} = 0.5$$

$$\rightarrow b_2 = \frac{1}{4} (23 - (0.5)(10)) = 4.5$$

$$b = \frac{b_1 + b_2}{2} = \frac{3.6 + 4.5}{2} = 4.05$$

12.679

$$\rightarrow \boxed{y = 4.05 + (0.714)x_1 + 0.5x_2} \quad y = b + a_1x_1 + a_2x_2$$

$$y = 4.05 + 3.57 + 3.5 \Rightarrow 11.12$$

when  $x_1 = 5$

$x_2 = 7$

if intercept same line ||  
 $b$  if intercept diff then line ||

Classif

20-1

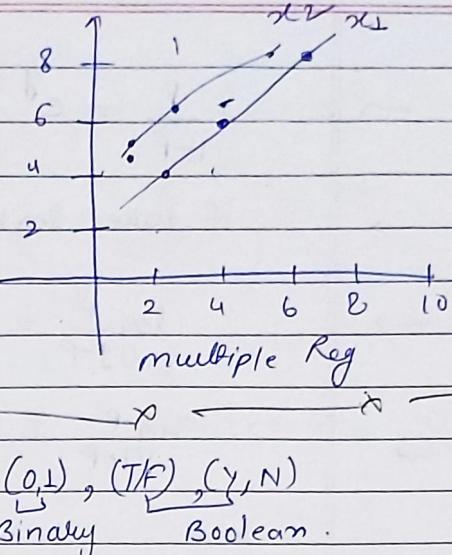
20-2

40 - final

20TA

\* Logistic Regression: when dependent ( $y$ ) var is in form of yes/No, 0/1 or T/F then we use Logistic Regression

→ It is used to find the Prob. of event success or event failure. we should use logistic Regression. when the dependent var is binary or boolean (0,1), (T/F), (Y, N)



∴ function → logit () → used for evaluating any Logistic Regression.

$$\rightarrow y = b_0 + b_1 x \quad (1) \quad \text{Linear regression.}$$

$$P(y) = b_0 + b_1 x \rightarrow P = b_0 + b_1 x \quad \text{so } 0 \leq P \leq 1$$

$$\rightarrow 0 \leq P \quad e^x \rightarrow +ve \quad e^{-x} \Rightarrow \frac{1}{e} \rightarrow +ve \quad \text{exponentiation.}$$

let  $P = e^{b_0 + b_1 x} \quad (2)$

$$P \rightarrow P \leq 1 \quad P = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}} \quad (3) \quad \text{so } b_0 + b_1 x = y$$

success

$$P = \frac{e^y}{1 + e^y} \rightarrow (3) \quad P + q = 1$$

$$q = 1 - P \rightarrow 1 - \frac{e^y}{1 + e^y} \rightarrow \frac{1 - q}{e^y + 1} \quad (4)$$

eqn (3) divide by (4)

$$\frac{P}{q} = \frac{P}{1-P} \Rightarrow \frac{\frac{e^y}{1+e^y}}{\frac{1-e^y}{1+e^y}} \Rightarrow \frac{e^y}{1-e^y} \Rightarrow e^y \quad \text{odd ratios}$$

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

$\rightarrow$  if take log both sides

$$\rightarrow \log \frac{P}{1-P} = \log(e^y)$$

$$\rightarrow \log \frac{P}{1-P} = \cancel{y} \rightarrow \boxed{\log \frac{P}{1-P} = \text{both, } x}$$

logit function

Logistic Regression eqn

Ques Consider the Data Set Given below of cricket players which contains the no of wickets taken in the last match played by them & on the basis of wicket taken in the previous match by each player, the player is either selected (1) or not selected (0). In the team for the next match to be played find the prob of selection of a player who has taken 2 wickets in the previous match that he played. Assume the value of  $b_0$  &  $b_1 = 1$

Last match (post data) no of wicket taken	Result(Selected 1 0 Not selected)
0	0
3	1
1	0
2	1
4	1

#

★

★

13/9/19

→ Q

same ques

Q

Q find the prob of the selection of the player who has taken 5 wickets in the previous match that he played. assume value of  $b_0$  &  $b_1 = 1$ .

$$\boxed{\log \frac{P}{1-P} = x}$$

when  $x=2$   $P=1$

# Likelihood Proof  $\rightarrow$  Binary  $\rightarrow$  Regression  
on features  $\rightarrow$  classification.

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

12

$$\frac{P}{1-P} = e^{b_0 + b_1 x_1} \rightarrow \frac{P}{1-P} = e^{\frac{x_1}{\beta}} \rightarrow \frac{P}{1-P} = e^{\frac{x_1}{2}}$$

$$\frac{P}{1-P} = 20.08$$

$$P = 20.08 - 20.08 P$$

$$P + 20.08 P = 20.08$$

$$\frac{21.08}{21.08} P = \frac{20.08}{21.08}$$

$$= 0.952 \approx 1$$

★  $\frac{P}{1-P} = e^{\frac{x_1 + x_2}{2}}$

$$P = (1-P) 403.428$$

$$P = \frac{403.428}{404.428}$$

$$P = 0.99 \approx 1 \text{ Selected}$$

$\rightarrow$  after Logistic we use classification like

0	1
0, 1	3, 2, 4, 5
	✓

# For multiple independent variable

$$\log \frac{P}{1-P} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

★ Data is Inconsistent,  $\rightarrow$  use Logistic regression

★ If Data is large  $\rightarrow$  use classification

13/9/19

multiple Logistic Regression /

last 2 match. selected (+) & selected (0)

find the prob of selection of a player who has taken 4 wkt in 1st match & 3 wkt in 2nd match. Assume value of  $b_0, b_1, b_2, \dots, b_n$

$$b_0, b_1, b_2, \dots, b_n = 1$$

$$x_1 = 0$$

$$x_2 = 4$$

$$\log \frac{P}{1-P} = b_0 + b_1 x_1 + b_2 x_2$$

$x_1$	$x_2$	f
0	1	0
1	0	0
3	2	1
4	3	1
2	4	1

130919

Logistic  
Regression

Universal or single dataset given thus  
we will find Training set & Test set init

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\log \frac{P}{1-P} = b_0 + b_1 x_1 + b_2 x_2$$

$$= 1 + 1 \cdot 0 + 1 \cdot 4$$

$$\left| \log \frac{P}{1-P} < 5 \right|$$

$$\frac{P}{1-P} = e^5 \rightarrow \frac{148.413}{149.413} \Rightarrow 1 = P$$

$$\frac{P}{1-P} = b_0 + b_1 x_1 + b_2 x_2$$

$$1 + 1 \cdot 0 + 1 \cdot 1$$

$$\frac{P}{1-P} = e^2$$

$$P = 7.38 = 0.98$$

~~7.38~~

$$P = 7.38 (1-P)$$

$$P = 7.38 - 7.38 \Rightarrow 0.87 \rightarrow \text{not decided}$$

∴ Logistic Regression → value if like 0.87 not necessarily true

\* Ridge & Lasso Regression - 2 Regularisation  $\xrightarrow{\frac{L_1}{L_2}}$

→ Ridge & Lasso Regression are powerful Technique's generally used for creating more accurate least square regression line or more accurate value for dependent variable.

$$y = b_0 + b_1 x_1 + \underbrace{A f}_{(\alpha) \text{Lasso Ridge w.r.t}} - \text{adjustment factor/error factor}$$

→ Ridge & Lasso regression depends on the value of  $\alpha$  (alpha) = level of significance - if the value of level of significance is given in the historical or quantitative data (dataset) Then we use the concept of Ridge & Lasso Regression.