

L-9: Predictive Analytics & Revision





Agenda

- Review of last session
- Introduction to regression
- Method of least squares
- Simple linear regression

Covariance of X and Y

$$\text{Cov}(X, Y) =$$

$$E(X) = \sum x P(x) \\ = \int x f(x) dx$$

$$= [E(X - \mu_x)(Y - \mu_y)]$$

$$= \sum_x \sum_y (x - \mu_x)(y - \mu_y)$$

if discrete

$$= \iint (x - \mu_x)(y - \mu_y) f_{XY}(x, y)$$

if continuous

Joint
P.d.f

$P(x, y)$

Joint
prob
density
function

$$\checkmark \text{cov}(x, y)$$

$$= \frac{\sum_{i=1}^n (x - \mu_x)(y - \mu_y)}{n-1}$$

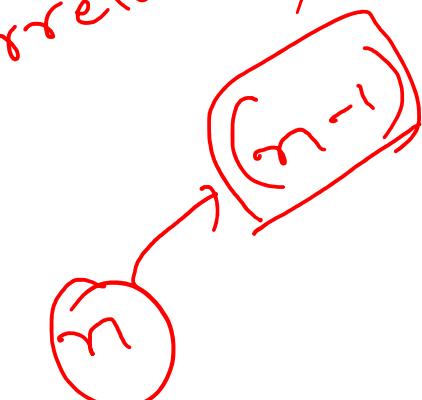
x
 y

income

expenditure

$n-1$

Correlation

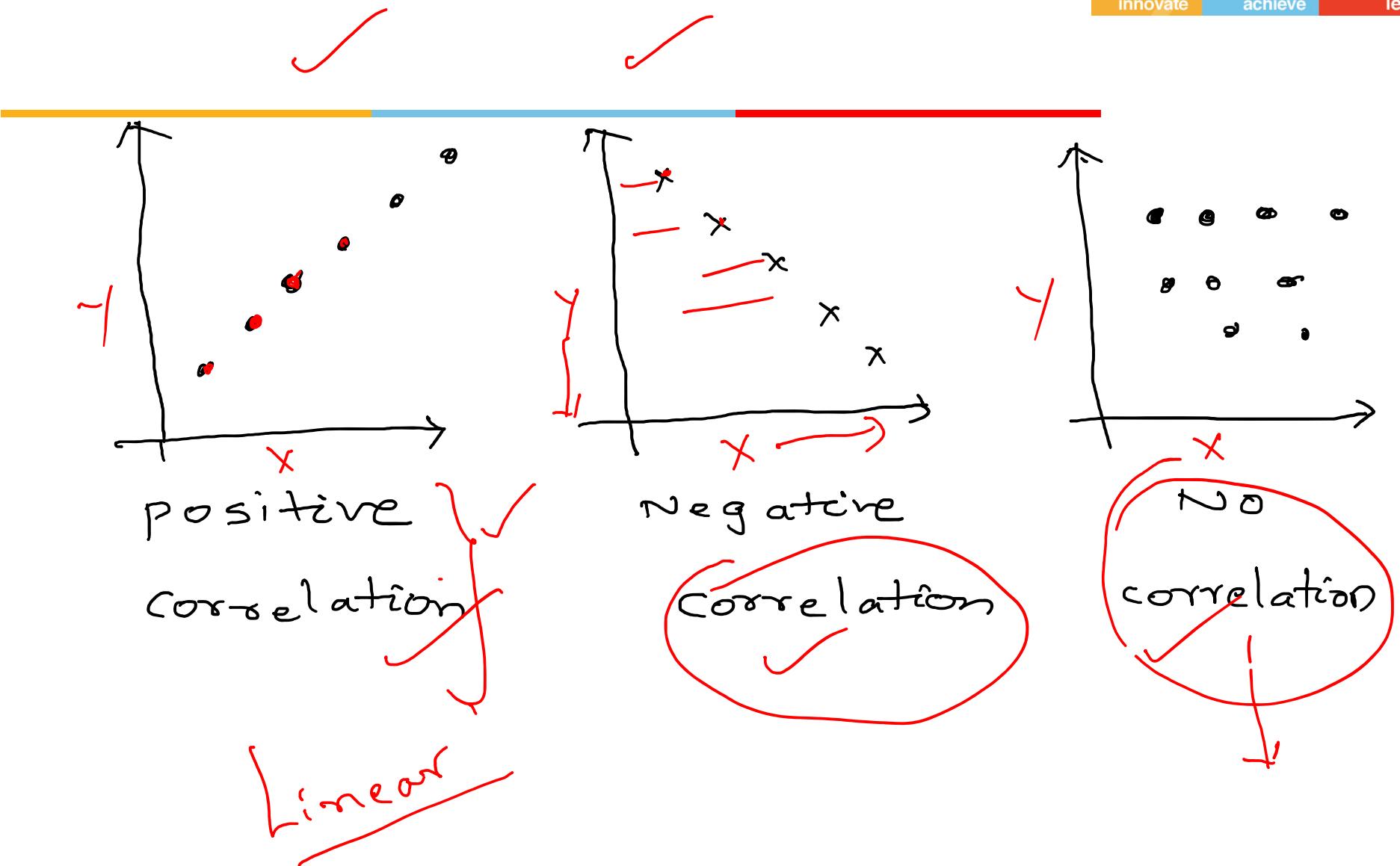


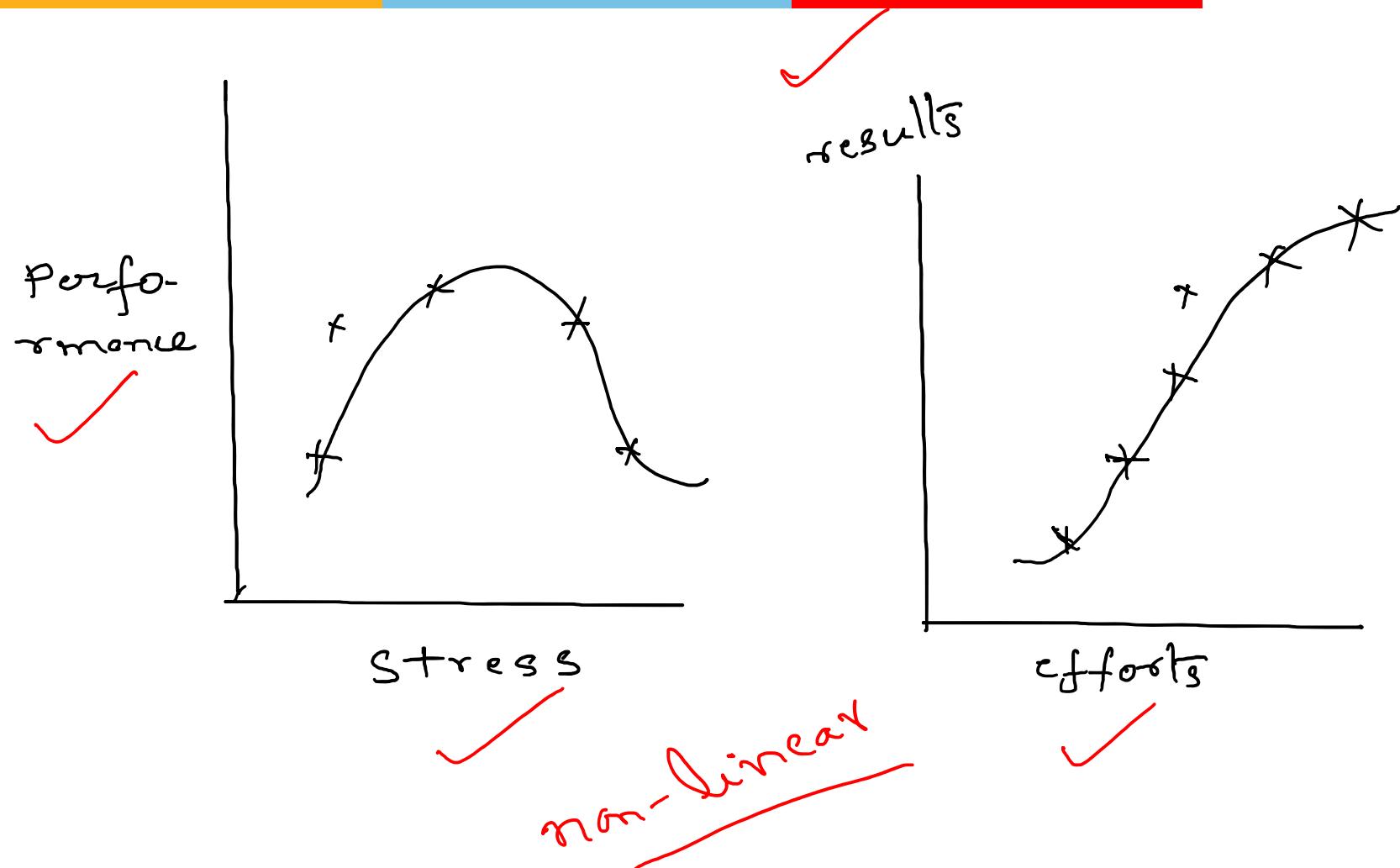
And also

⇒ Farmer has an impression that if he uses more fertilizers, then the crop yield increases.

we need to validate this?

How → ?





Coefficient of correlation:

$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{\sum xy}{\sqrt{\sum x^2 \cdot \sum y^2}}$$

where $x = x - \bar{x}$

$$y = y - \bar{y}$$

$$x^2 = (x - \bar{x})^2$$

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

Coefficient of Correlation

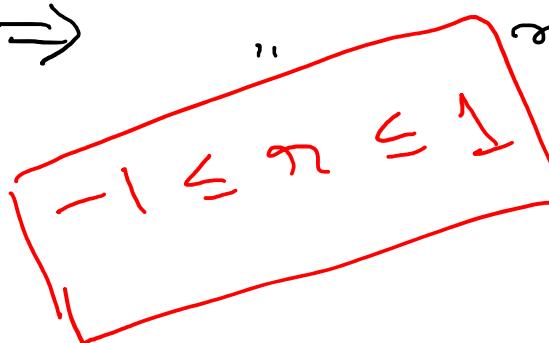
$r = 1 \Rightarrow$ Perfect and positive relation ✓

$r = -1 \Rightarrow$ " " negative relation

$r = 0 \Rightarrow$ no relation ✓

$0 < r < 1 \Rightarrow$ Partial positive relation

$-1 < r < 0 \Rightarrow$ " negative "



Example - 1

x	1	2	3	4	5	6	7	8	9
y	10	11	12	14	13	15	16	12	18

$$\bar{x} = \frac{\sum x}{n} = \frac{45}{9} = 5$$

$$\bar{y} = \frac{\sum y}{n} = \frac{126}{9} = 14$$

r

x	$x =$	x^2	y	$y = y - 14$	y^2	xy	\cdot
1	-4	16	10	-4	16	16	
2	-3	9	11	-3	9	9	
3	-2	4	12	-2	4	4	
4	-1	1	14	0	0	0	
5	0	0	13	-1	1	0	
6	1	1	15	1	1	1	
7	2	4	16	2	4	4	
8	3	9	17	3	9	9	
9	4	16	18	4	16	16	
		60		60	59		

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$= \frac{59}{\sqrt{60 \times 60}}$$

$$= 0.9833$$

$r = 0.9833$

→ Corrected
→ True corrected

x	$x =$	y	$y = y - 14$	y^2	xy
1	-4	16	10	-4	16
2	-3	9	11	-3	9
3	-2	4	12	-2	4
4	-1	14	0	0	0
5	0	0	13	-1	0
6	1	1	15	1	1
7	2	4	16	2	4
8	3	9	17	3	9
9	4	16	18	4	16

$\sum xy = 59$
 $n = 9$
 $\sum y^2 = 137$
 $\sum y = 105$

$$\text{cov}(x, y) = \frac{\sum xy}{n-1}$$

$$= \frac{59}{8}$$

$$= 7.375$$

Coefficient of Determination ✓

r is coeff. of correlation

r^2 is coeff of determination



Indicates the extent to which variation in one variable is explained by the variation in the other.

$$r = 0.9 \Rightarrow r^2 = 0.81$$

i.e. 81% of the variation in y

due to variation in x .

remaining 19% is due to some other factors.

~~r_{xy}~~ ~~$r_{x,y}$~~

π 0.9833 \rightarrow *Coeff Correlation*

$\text{cov}(x, y)$ $=$ 7.375 \rightarrow *Covariance Interpretation*

$\pi^2 = 0.81 \rightarrow$ *Coeff of determination*

$-1 \leq \pi \leq 1$

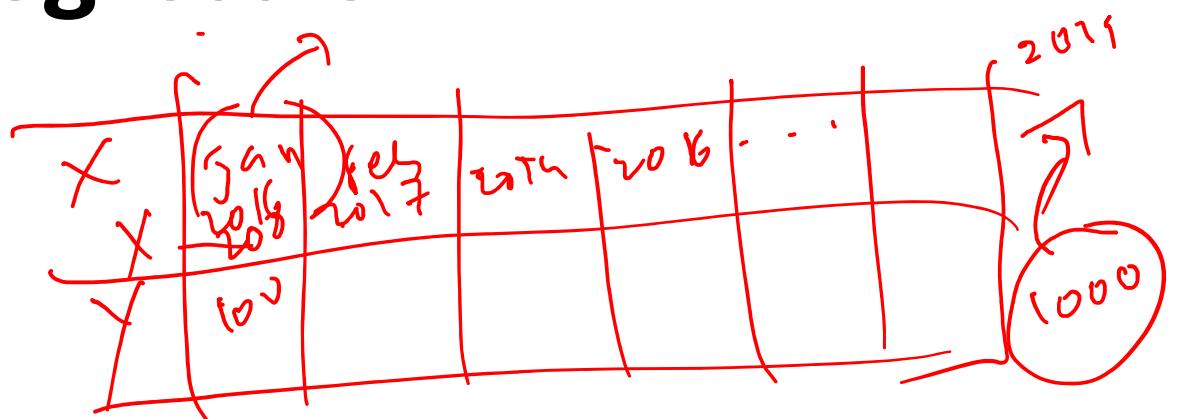
"9490233116"

farmer:

x : fertilizer
 y : crop yield

Correlation
 g

Regression



Regression :-

Correlation

r

x	1	2	3	4	5
y	1	4	9	16	25

$$y = x^2 \text{ when } x = ?$$

$y = ax$

when $a = 7$: $y = ?$

~~5.4 5.6 5.8
5.4 5.6 5.8~~

a	1	2	3	4	5
y	1	6	2	5	4

when $x = 7$, $y = ?$

$$y = ?$$

Correlation

- Measuring strength or degree of the relationship between two variables
- no estimation
- both variables are independent

Regression

- having an algebraic equation between two variables
- estimation
- one is dep't variable and other indept variables

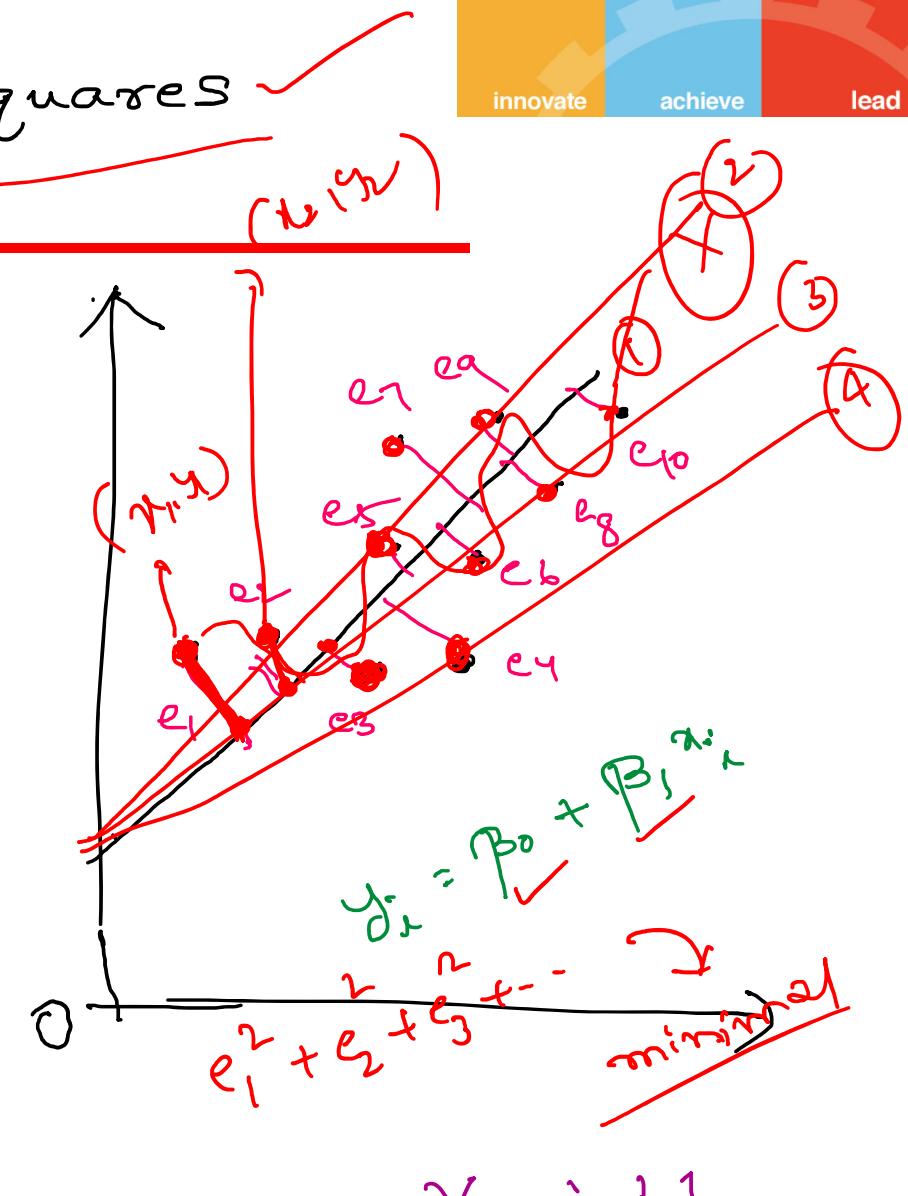
Method of Least squares

y : Dependent variable

x : Independent variable

predictor variable

response Variable

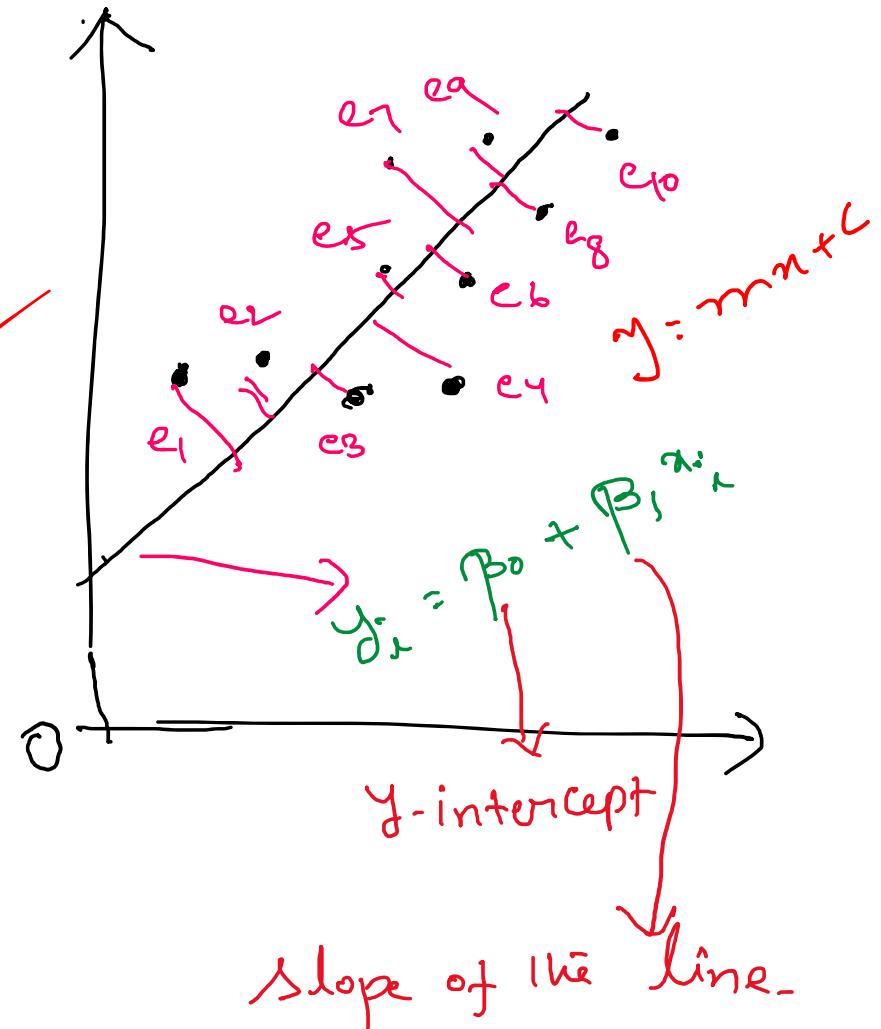


Method of Least squares

"minimizing the
error"

i minimize

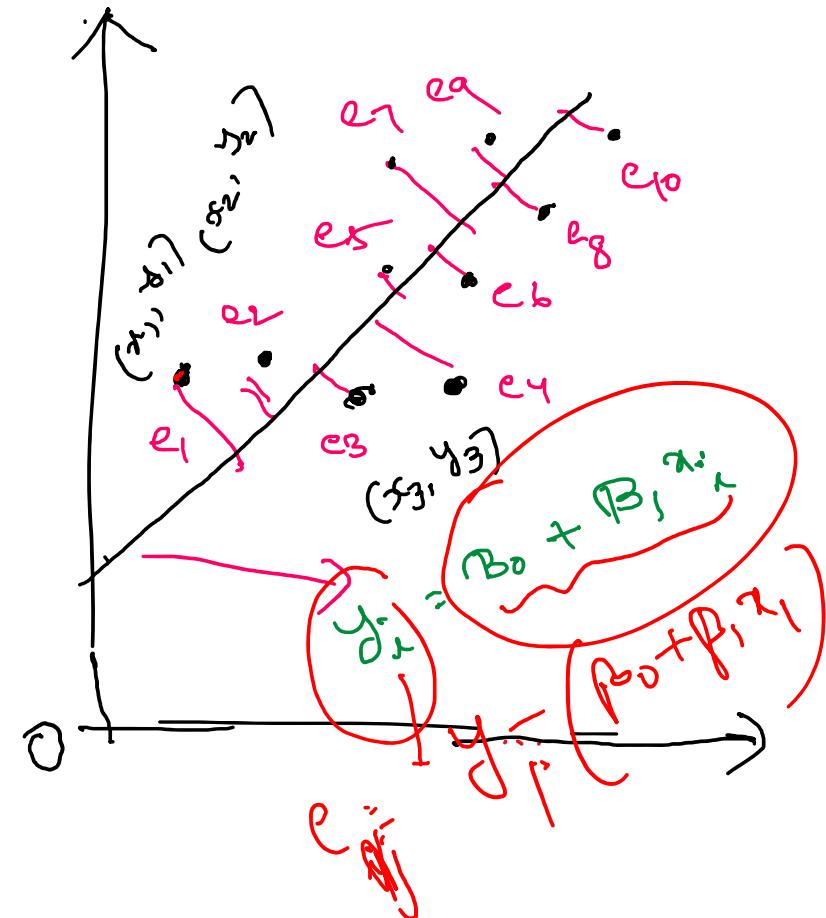
$$e_1^2 + e_2^2 + e_3^2 + \dots + e_{10}^2$$



Method of Least squares

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

we need to choose β_0 and β_1 , which minimizes the error.



Method of Least squares



$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S}{\partial \beta_0} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-1)$$

$$\Rightarrow \sum_{i=1}^n y_i = n \beta_0 + \beta_1 \sum_{i=1}^n x_i$$

$$\frac{\partial S}{\partial \beta_1} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (2)(-x_i)$$

$$\Rightarrow \sum_{i=1}^n x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2$$

on solving these, we get β_0 & β_1
which minimizes error.

Linear regression

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

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

$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

Normal equations.

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

Regression Coefficients

$$y = a + b_1 x$$

↓

b_{yx} : Regression coeff
of y on x

regression line of y on x

$$x = c + d_1 y$$

↓

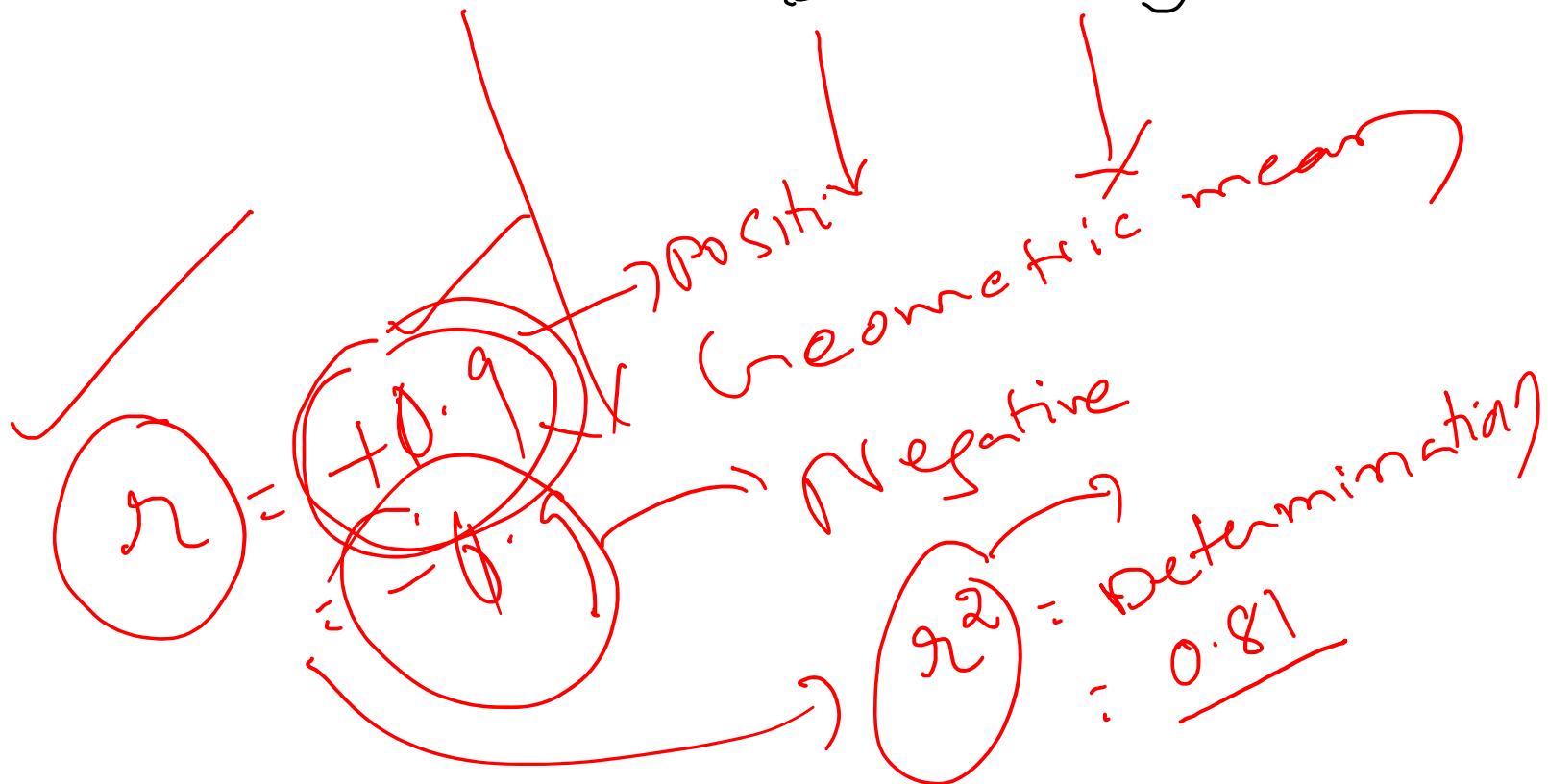
b_{xy} : regression coeff
of x on y

regression line of x on y

~~of y on x~~

Correlation coefficient

$$r = \sqrt{b_{yx} \times b_{xy}}$$



Example :-

company	Advt	Sales	Revenue
	Expt		
A	1	1	
B	3	2	
C	4	2	
D	6	4	
E	8	6	
F	9	8	
G	11	8	
H	14	9	

$$y = a + bx$$

$$\sum y = an + b \sum x$$

$$\sum xy = a \sum x + b \sum x^2$$

~~Example :-~~

Sales Revenue y	Advt expt. x	x^2	xy
1	1	1	1
2	3	9	6
2	4	16	8
4	6	36	24
6	8	64	48
8	9	81	72
8	11	121	88
9	14	196	126
$\sum 40$		$\sum 56$	$\sum 373$

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

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

$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

$$\Rightarrow 40 = 8\beta_0 + 56\beta_1$$

$$373 = 56\beta_0 + 524\beta_1$$

on solving

$$\beta_0 = 0.072$$

$$\beta_1 = 0.704$$

i.e. $y = (0.072) + (0.704)x$

$$i) y = (0.072) + (0.704)x$$

when $x = \underline{0.075}$, then

$$\begin{aligned} y &= (0.072) + (0.704)(0.075) \\ &= 0.1248 \quad \approx 12.48\% \end{aligned}$$

Example:

Consider the following data

x	1	2	4	0
y	0.5	1	2	0

Fit a linear regression line

Estimate y when $x = 5$.

x	y	xy	x^2
1	0.5	0.5	1
2	1	2	4
4	2	8	16
0	0	0	0
$\sum x = 7$		$\sum y = 3.5$	$\sum x^2 = 21$
$\sum xy = 10.5$			

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

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

$$\sum xy = \beta_0 \sum x_1 + \beta_1 \sum x^2$$

$$3.5 = 4\beta_0 + \beta_1 \quad (?)$$

$$10.5 = 7\beta_0 + \beta_1 \quad (?)$$

on solving these

$$\beta_0 = 0$$

$$\beta_1 = 0.5$$

$$\text{i.e. } y = 0 + (0.5)x$$

$$\text{when } x = 5, \quad y = (0.5)5 \\ = 0.25$$

Linear regression (multiple regression)

Example:-

x_0	size	No of rooms	No of floors	Age of home	Price Lakh
1	2000	5	2	45	4000
1	1400	3	1	40	2000
1	1600	3	2	30	3000
1	800	2	1	35	2000

House → Price

$y \rightarrow n$

increment per

$y =$

20 years
1 floor
2 rooms
1200 sqft



Multiple Linear Regression

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

$$\sum y = \beta_0 n + \beta_1 \sum x$$
$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

$$\sum y = \beta_0 n + \beta_1 \sum x_1 + \beta_2 \sum x_2 + \beta_3 \sum x_3 + \beta_4 \sum x_4$$

$$\sum xy = \beta_0 \sum x_1 + \beta_1 \sum x_1^2 + \beta_2 \sum x_1 x_2 + \beta_3 \sum x_1 x_3 + \beta_4 \sum x_1 x_4$$

$$\sum x_2 y = \beta_0 \sum x_2 + \beta_1 \sum x_1 x_2 + \beta_2 \sum x_2^2 + \beta_3 \sum x_2 x_3 + \beta_4 \sum x_2 x_4$$

$$\sum x_3 y = \beta_0 \sum x_3 + \beta_1 \sum x_1 x_3 + \beta_2 \sum x_2 x_3 + \beta_3 \sum x_3^2 + \beta_4 \sum x_3 x_4$$

$$\sum x_4 y = \beta_0 \sum x_4 + \beta_1 \sum x_1 x_4 + \beta_2 \sum x_2 x_4 + \beta_3 \sum x_3 x_4 + \beta_4 \sum x_4^2$$

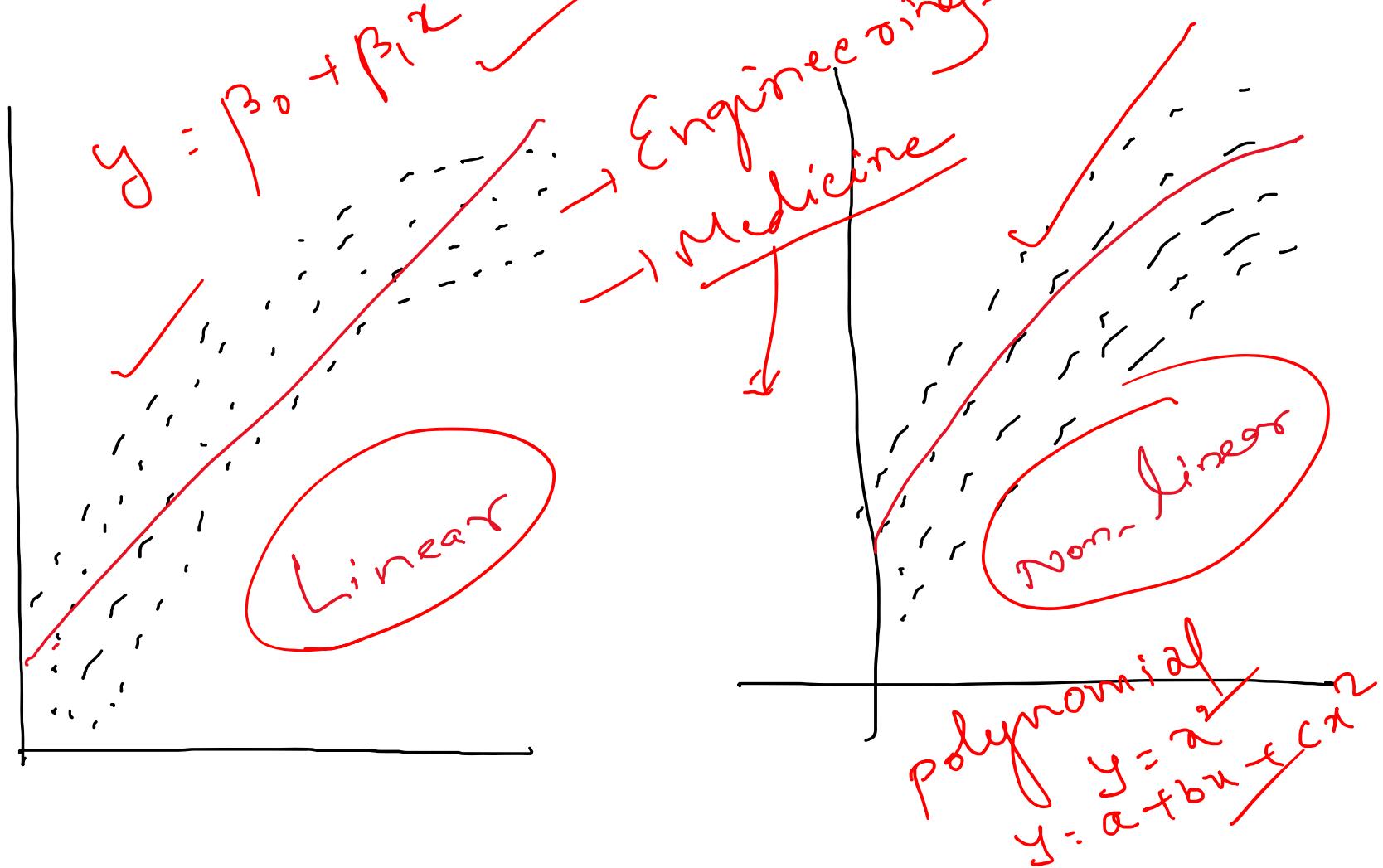
solve for $\beta_0, \beta_1, \beta_2, \beta_3, \beta_4$

Other regressions



just a look

Tech
Revenue



Suppose

$$y = a e^{bx}$$

exponential curve

$$\log y = \log a + b \log x$$

Y X

i.e.

$$Y = A + BX \quad \rightarrow \text{linear eqn}$$

$$\sum Y = A n + B \sum X \rightarrow 0$$

$$\sum XY = A \sum X + B \sum X^2 \rightarrow 2$$

$$\begin{aligned} A &= ? \\ B &= ? \end{aligned}$$

$\Rightarrow A$

Hence, we get
 $y = ae^{bx}$

Suppose $y = a x^b$ non Linear Power Curve

$$\log y$$

$$(\log y) = (\log a) + b \log x$$

$$\text{ie } Y = A + bX$$

$$\sum Y = A n + b \sum X$$

$$\sum XY = A \sum X + b \sum X^2$$

$$y = a x^b$$

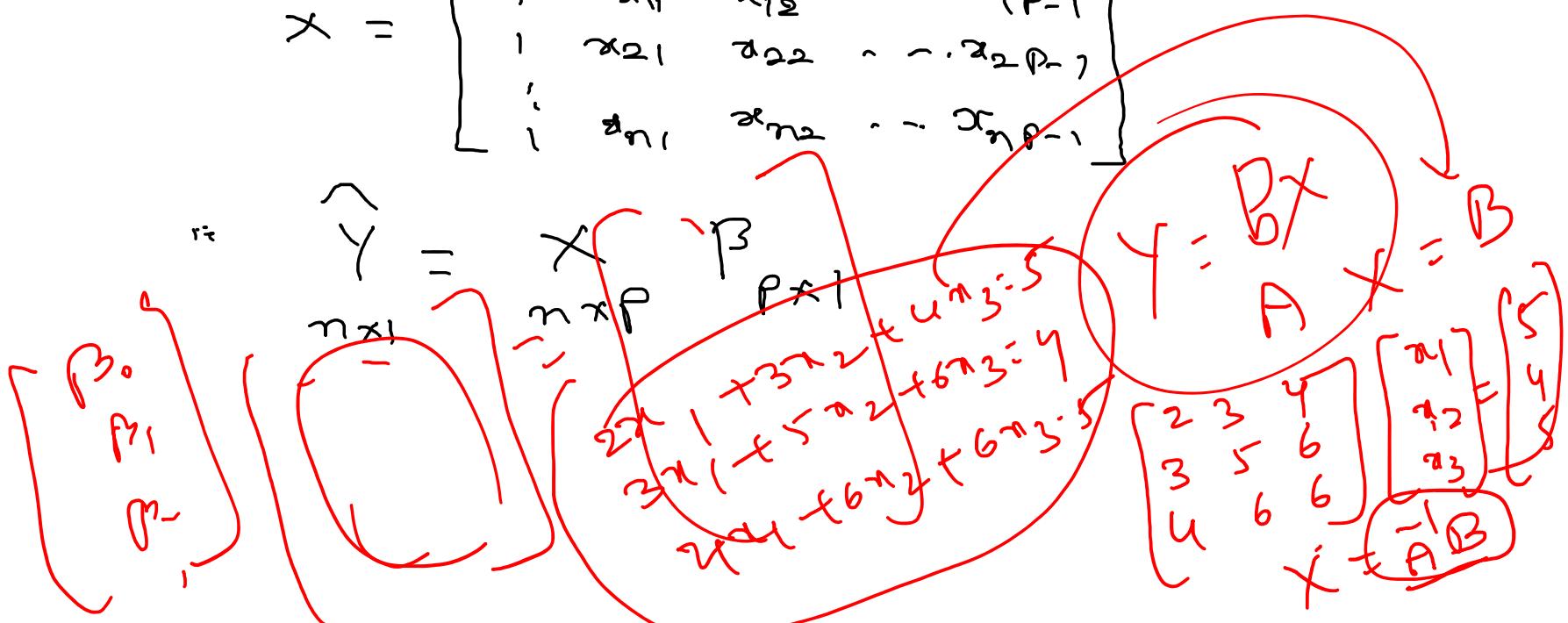
Matrix Approach:

Let $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1}$

observations $y_i = 1, 2, \dots, n \rightarrow$ by a vector γ

Unknowns $\beta_0, \beta_1, \dots, \beta_{p-1} \rightarrow \dots \beta$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1(p-1)} \\ 1 & x_{21} & x_{22} & \dots & x_{2(p-1)} \\ \vdots & x_{n1} & x_{n2} & \dots & x_{n(p-1)} \end{bmatrix}$$



Find β to minimize

$$S(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1 - \beta_2 x_2 - \dots)^2$$

$$= \|y - x\beta\|^2 = \|y - \hat{y}\|^2$$

Diff S wrt to each β we get linear eqns

$$x^T x \hat{\beta} = x^T y \rightarrow \text{normal eqns}$$

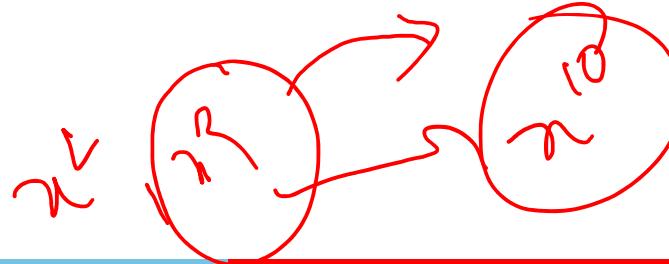
If $x^T x$ is non-singular, the soln is

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

$$x = f$$

$$A x = B$$

$$x = A^{-1} B$$



Computationally, it is sometimes unwise even to form the normal equations because forming $\mathbf{x}^T \mathbf{x}$ can introduce undesirable round-off error.

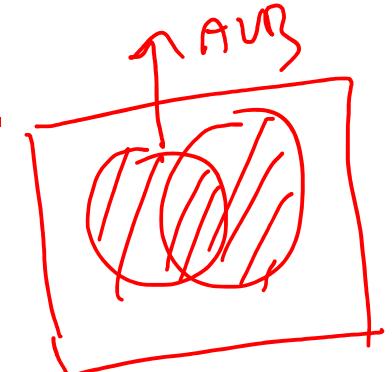
- If $\mathbf{x}^T \mathbf{x}$ is non-invertible ... ?
- ✓ Redundant features
 - ✓ too many features
- ↳ scaling

Revision

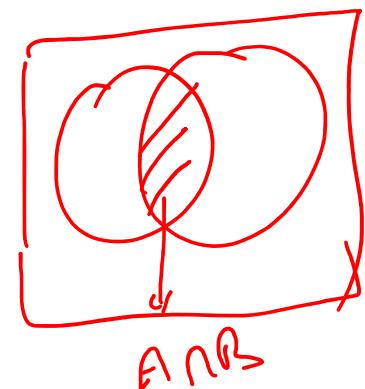
→ probability

$$\rightarrow P(A \cup B)$$

$$P(A \cap B)$$



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



→ Conditional probability:

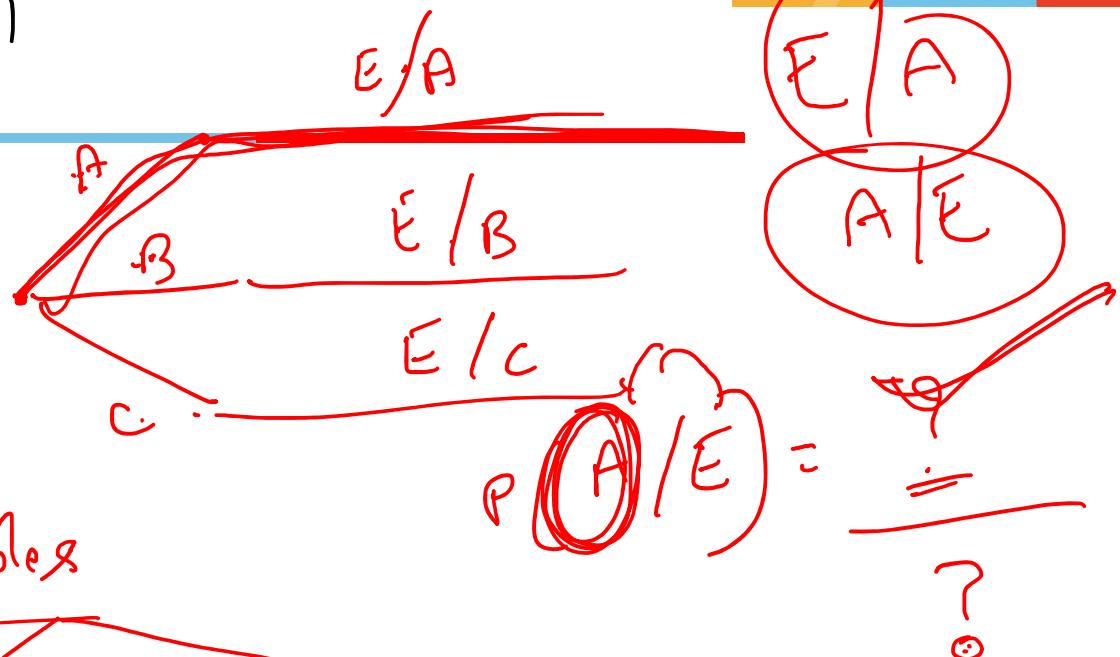
$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B | A) = \frac{P(A \cap B)}{P(A)}$$

$$\text{i.e. } P(A \cap B) = P(A | B)P(B) - P(B | A)P(A)$$

Revision

→ Bayes Theorem :-



→ Random variables

Discrete

$$P(x)$$

$$\text{i, } 0 \leq P(x) \leq 1$$

$$\text{ii, } \sum P(x) = 1$$

$$\text{Mean} = E(x) = \sum x P(x) \\ = \int x f(x) dx$$

Continuous

$$\text{i, } f(x) \\ \text{ii, } 0 \leq f(x) \leq 1$$

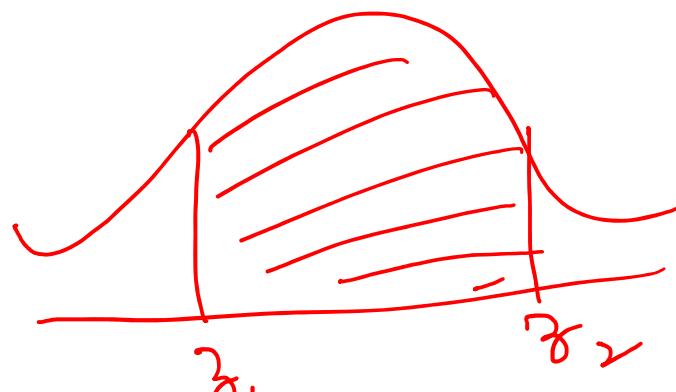
$$\text{iii, } \int f(x) dx = 1$$

$$\text{variance: } E(x-\mu)^2 \\ = E(x^2) - \mu^2 \\ = E(x^2) - [E(x)]^2$$

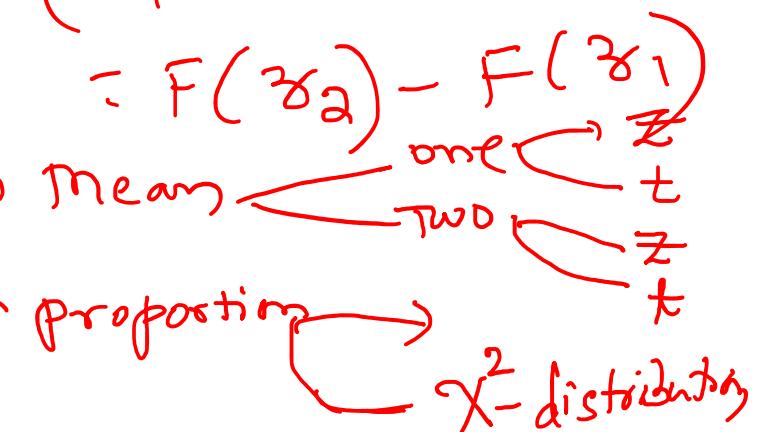
Revision

→ Binomial dist $P(x) = {}^n C_x P^x Q^{n-x}$, $x=0,1,2 \dots n$
 poisson dist $P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$, $x=0,1,2, \dots \infty$

→ Normal distribution :- $P(30 \leq x \leq 50)$



$$z = \frac{x - \mu}{\sigma}$$

$$P(z_1 \leq z \leq z_2) = F(z_2) - F(z_1)$$




Thanks