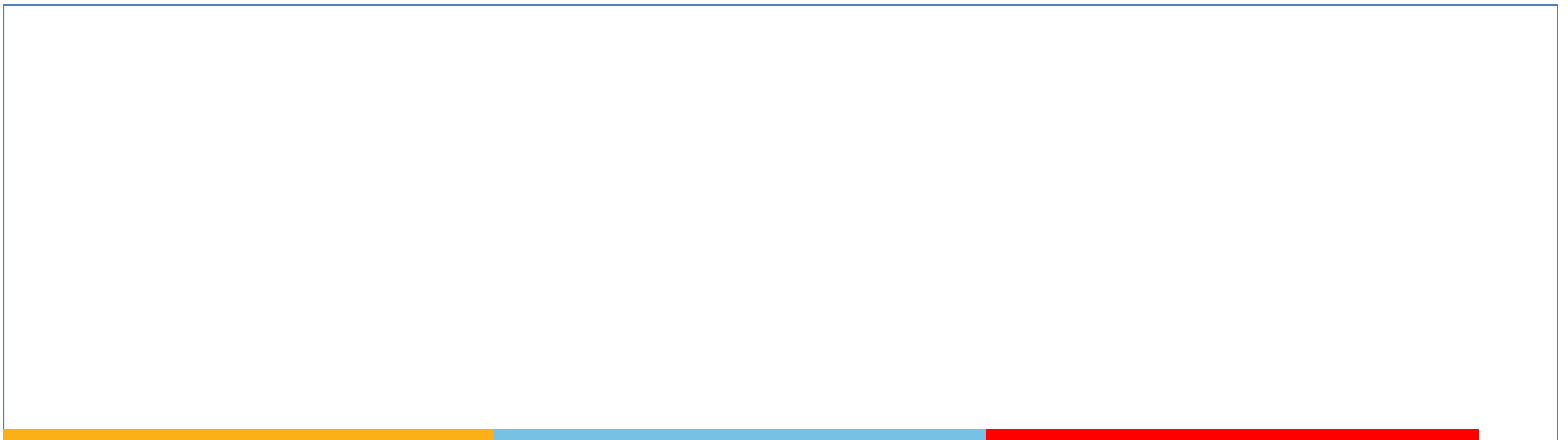# L-9: Predictive Analytics & Revision

# Agenda

➢ Review of last session

➢ Introduction to regression

➢ Method of least squares

➢ Simple linear regression

# Covariance of $x$ and $\gamma$

$$Cov(X, Y) =$$

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

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

$$= \sum_x \sum_y (x - \mu_x)(Y - \mu_y) \; \boxed{P(x, y)}$$

$\rightarrow$ joint p.d.f

if discrete

$\rightarrow$ joint prob. density fun

$$= \iint (x - \mu_x)(Y - \mu_y) f(x, y) \, dx \, dy$$

if continuous

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

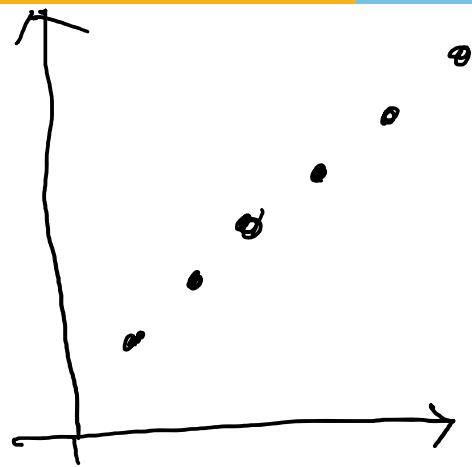$$= \frac{\sum_{i=1}^{n} (x - \mu_x)(y - \mu_y)}{n - 1}$$
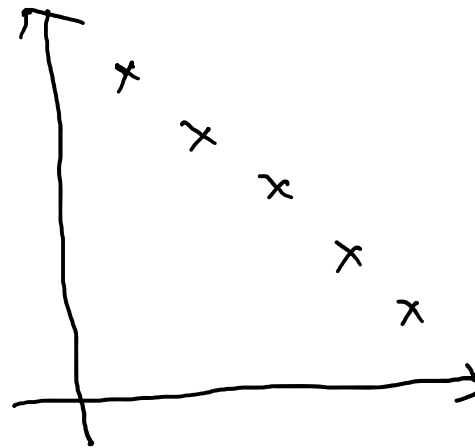
# And also

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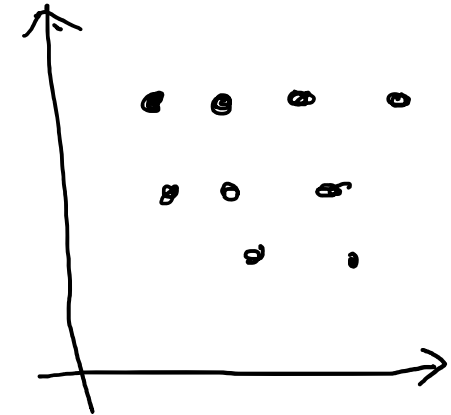crop yield increases.

We need to validate this?

How → ?

positive
correlation

Negative
correlation

No
correlation

# Coefficient of correlation:

$$r = \frac{cov(x, y)}{\sigma_x \, \sigma_y} = \frac{\Sigma XY}{\sqrt{\Sigma X^2 \cdot \Sigma Y^2}}$$

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 | 17 | 18 |

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

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

$r$

| $x$ | $X = x-5$ | $X^2$ | $y$ | $Y = y-14$ | $Y^2$ | $XY$ |
|---|---|---|---|---|---|---|
| 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 |

$$r = \frac{\Sigma XY}{\sqrt{\Sigma X^2 \, \Sigma Y^2}}$$

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

$$= 0.9833$$

innovate    achieve    lead

| x | X = x−5 | xy | Y | y = Y−14 | Y² | XY |
|---|---------|----|---|----------|----|----|
| 1 | −4 | 16 | 10 | −4 | 16 | 16 |
| 2 | −3 | 9  | 11 | −3 |    | 9  |
| 3 | −2 | 4  | 12 | −2 |    | 4  |
| 4 | −1 | 1  | 14 | 0  |    | 0  |
| 5 | 0  | 0  | 13 | −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 |    | 0  |    | 59 |

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

$$= \frac{\Sigma XY}{n-1}$$

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

$$= 7.375$$

# 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 81% of the variation in $y$ due to variation in '$x$'

remaining 19% is due to some other factors.

$$r = 0.9833$$

$$cov(x, y) = 7.375$$

$$r^2 = 0.81$$

Interpretation

# Regression

# Regression :-

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

when $x = 7 : y = ?$

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

when $x = 7, y = ?$

# 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

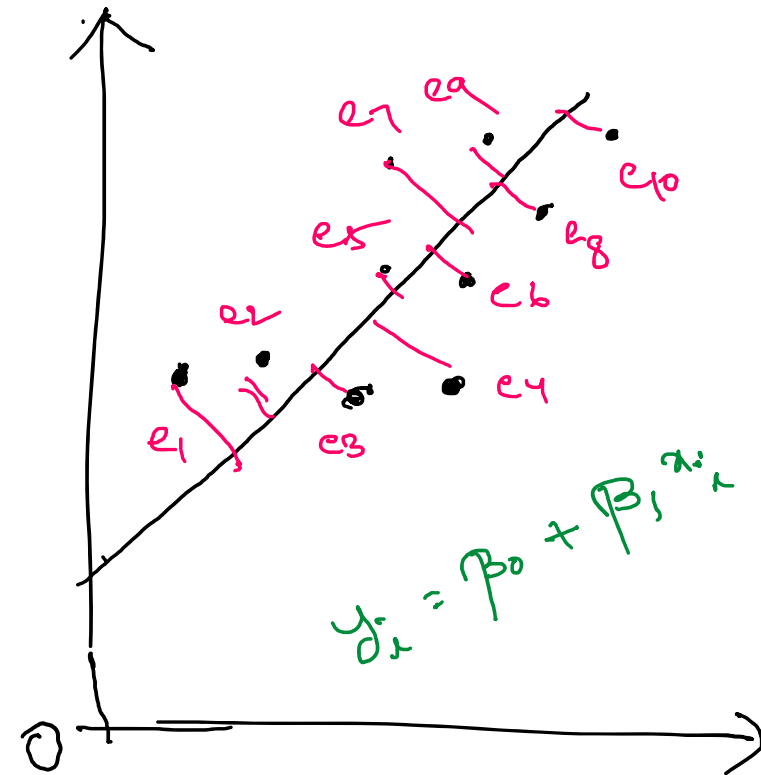→ estimation

→ one is dept variable and other indept variables

# Method of Least squares

$y$ : Dependent
        Variable

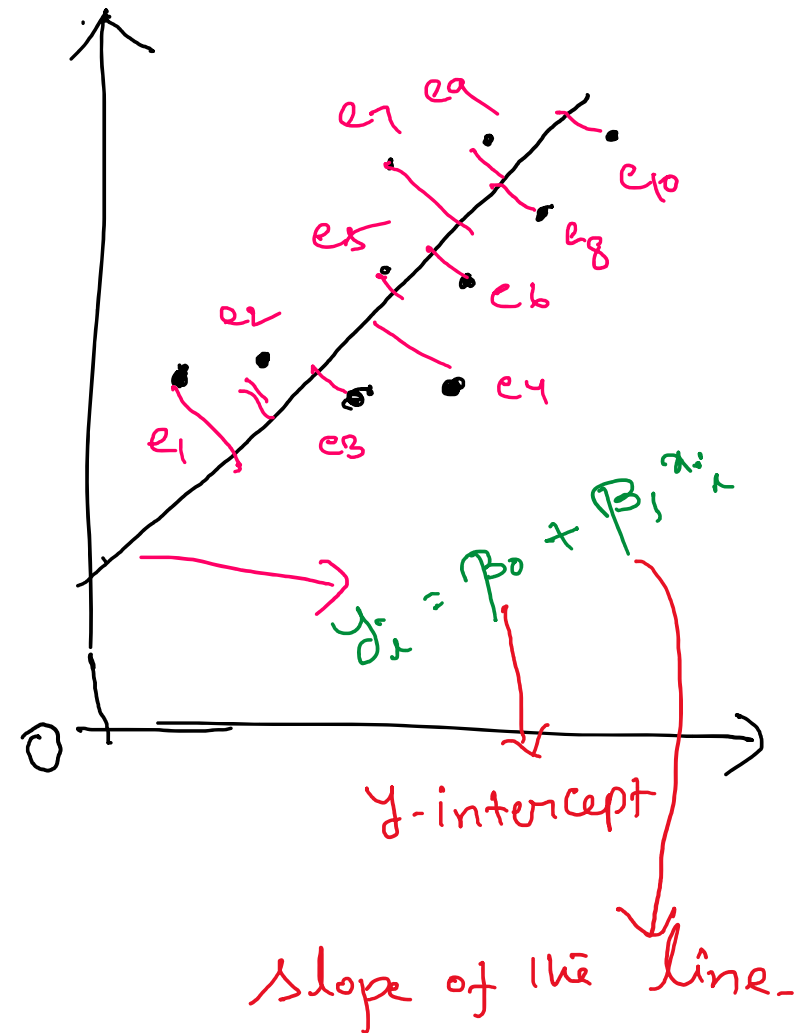$x$ : Independent
        Variable

predictor
        variable

response Variable

$y_i = \beta_0 + \beta_1 x_i$

$e_1$  $e_2$  $e_3$  $e_4$  $e_5$  $e_6$  $e_7$  $e_8$  $e_9$  $e_{10}$

O

# Method of Least squares

" minimizing the error "

$$\text{minimize}$$

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

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

y-intercept

Slope of the line.

$e_7$  $e_9$  $e_{10}$  $e_5$  $e_8$  $e_b$  $e_2$  $e_4$  $e_1$  $e_3$
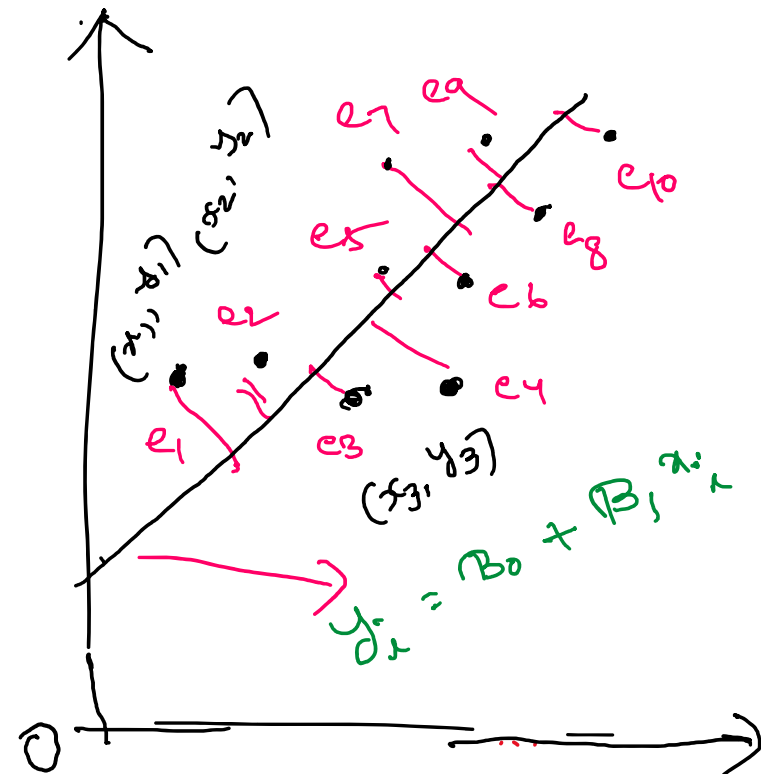
$O$

# Method of Least squares

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

We need to choose $\beta_0$ and $\beta_1$ which minimizes the error.

# Method of Least Squares

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

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

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

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

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

on solving these, we get $\beta_0$ & $\beta_1$
which minimizes error.

# Linear regression

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

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

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

Normal equations.

# Regression Coefficients

$$y = a + bx$$

regression line of $y$ on $x$

$b_{yx}$ : Regression coeff of $y$ on $x$

$$x = c + dy$$

regression line of $x$ on $y$

$b_{xy}$ : regression coeff of $y$ on $x$

Correlation coefficient

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

## Example :-

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

$$y = a + bx$$

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

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

# Example :-

| Sales Revenue $y$ | Advt expr. $x$ | $x^2$ | $xy$ |
|---|---|---|---|
| 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 |
| $\Sigma\,40$ | $\Sigma\,56$ | $\Sigma\,524$ | $\Sigma\,373$ |

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

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

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

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

on solving

$$\beta_0 = 0.072$$
$$\beta_1 = 0.704$$

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

$\text{i}\dot{e} \quad y = (0.072) + (0.704)x$

$\text{when} \quad x = 0.075, \text{ then}$

$$y = (0.072) + (0.704)(0.075)$$

$$= 0.1248 \approx 12.48\%$$

# 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 |
| $\Sigma = 7$ | $\Sigma 3.5$ | $\Sigma 10.5$ | $\Sigma 21$ |

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

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

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

$$3.5 = 4\beta_0 + \beta_1 \ (7)$$

$$10.5 = 7\beta_0 + \beta_1 \ (21)$$

on solving these

$$\beta_0 = 0$$

$$\beta_1 = 0.5$$

ie $\quad y = 0 + (0.5)x$

when $x = 5$, $\quad y = (0.5)5$

$$= 0.25$$

# Linear regression
## (Multiple regression)

Example:-

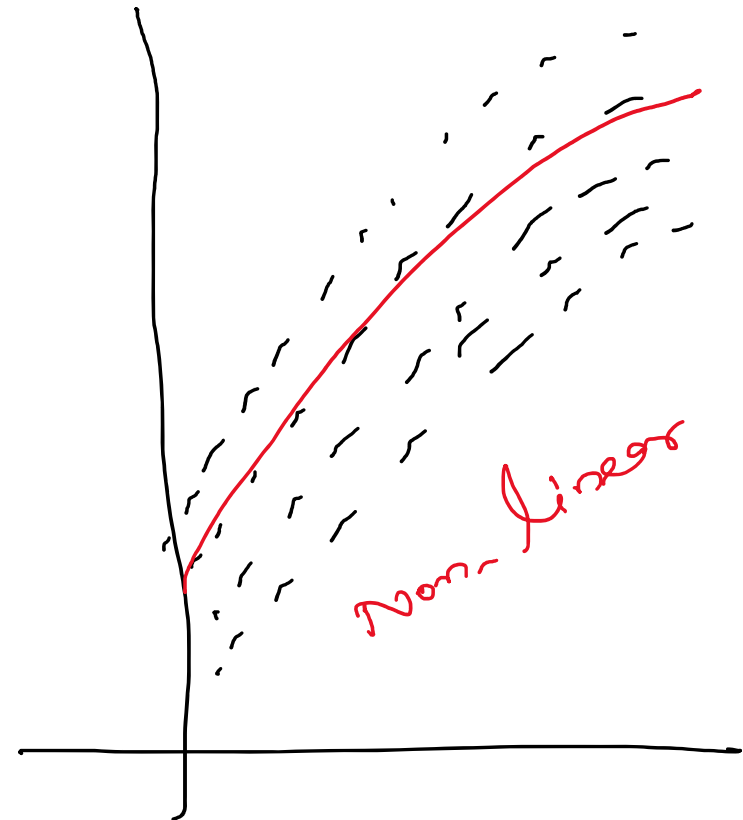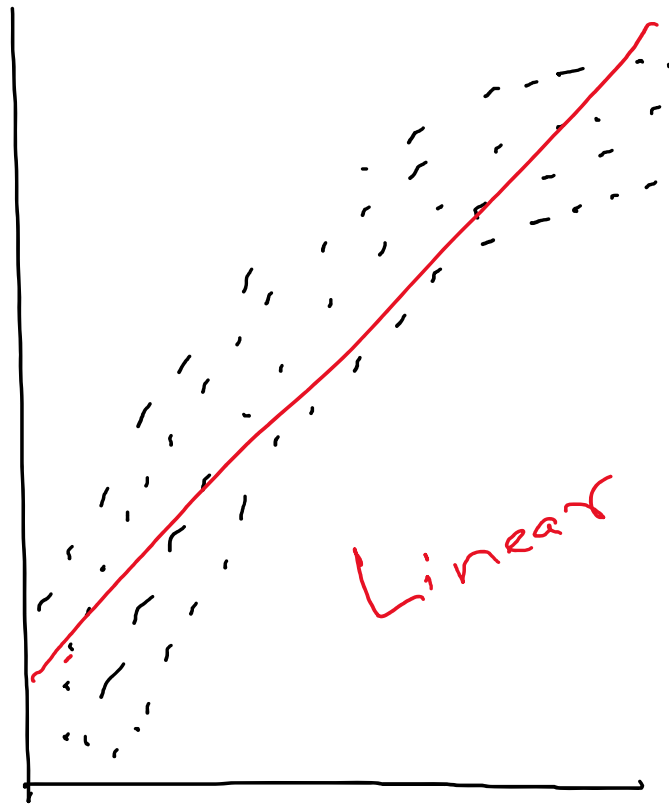| $x_0$ | Size | No of rooms | No of floors | Age of home | Price Lakhs |
|---|---|---|---|---|---|
| 1 | 2000 | 5 | 2 | 45 | 4000 |
| 1 | 1400 | 3 | 1 | 40 | 2000 |
| 1 | 1600 | 3 | 2 | 30 | 3000 |
| 1 | 800 | 2 | 1 | 35 | 2000 |

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

# Multiple Linear regression

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

# Other regressions
## just a look



Linear

Non-linear

Suppose $y = ae^{bx}$ $\rightarrow$ exponential curve

$\boxed{\log y} = \boxed{\log a} + b\boxed{\log x}$

$\underset{Y}{} \quad \underset{A}{} \quad \underset{X}{}$

i.e. $Y = A + bX$

$\Sigma Y = An + b\Sigma X \longrightarrow \boxed{1}$

$\Sigma XY = A\Sigma X + b\Sigma X^2 \longrightarrow \boxed{2}$

$A = ? \implies \boxed{a}$ ✓

$\boxed{b}$ ✓

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

Suppose $y = ax^b$ $\longrightarrow$ Power Curve

# Matrix Approach:

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

observations $y_i = 1, 2, \cdots n \rightarrow$ by a vector $Y$

unknowns $\beta_0, \beta_1 \cdots \beta_{p-1} \rightarrow$  "  .. $\beta$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p-1} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p-1} \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np-1} \end{bmatrix}$$

$\hat{Y}$ $= X$ $\beta$

$\underset{n \times 1}{Y} = \underset{n \times p}{X} \ \underset{p \times 1}{\beta}$

Find $\beta$ to minimize

$$S(\beta) = \sum_{i=1}^{m} \left( y_i - \beta_0 - \beta_1 x_1 - \beta_2 x_2 \cdots \right)^2$$

$$= \| Y - X\beta \|^2 = \| Y - \hat{Y} \|^2$$

Diff $S$ wrt to each $\beta$ we get Linear eqns

$$X^T X \hat{\beta} = X^T Y \longrightarrow \text{normal eqns}$$

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

$$\hat{\beta} = \left( X^T X \right)^{-1} X^T Y$$

Computationally, it is sometimes unwise even to form the normal equations because the multiplications involved in forming $X^T X$ can introduce undesirable round-off error.

→ If $X^T X$ is non-invertible ... ?

✓ Redundant features

✓ too many features

# Revision

# Revision

# Revision

# Thanks