

L- 8: Predictive Analytics



Agenda

- Covariance
- Correlation
- Introduction to regression
- Method of least squares
- Simple linear regression

Covariance of X and Y



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

$$= \left[E(X - \mu_X)(Y - \mu_Y) \right]$$

$$= \sum \sum (x - \mu_X)(y - \mu_Y) P(x, y)$$

if discrete

$$= \iint (x - \mu_X)(y - \mu_Y) f(x, y) dx dy$$

if continuous



consider the following

⇒ whether spending on advertising of a company is related to overall sales of the company.

→ If it is related, how it is related

⇒ Forecasting the sales, given the budget for advertising

And also



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

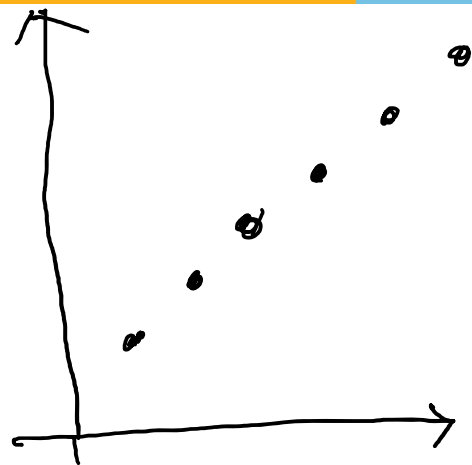
We need to validate this?

How → ?

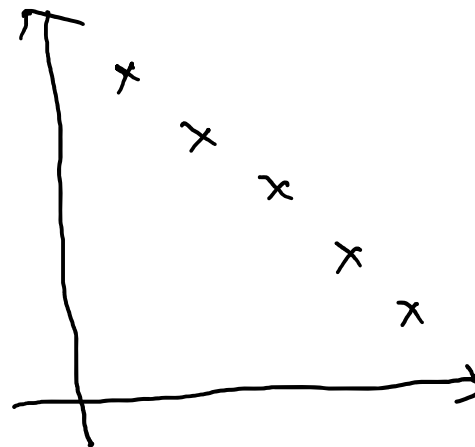
Correlation



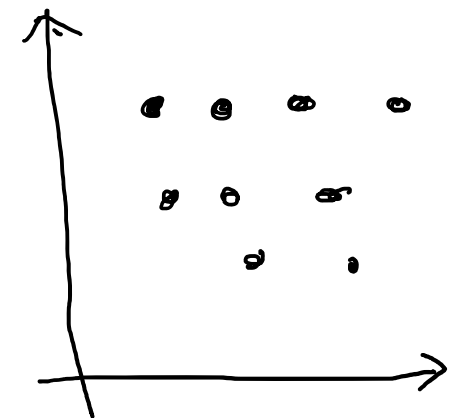
- Sales of a Company and Expenditure on advertisement
- Price and Demand of a product
- Inflation and Gold Price
- IQ and performance in Entrance.



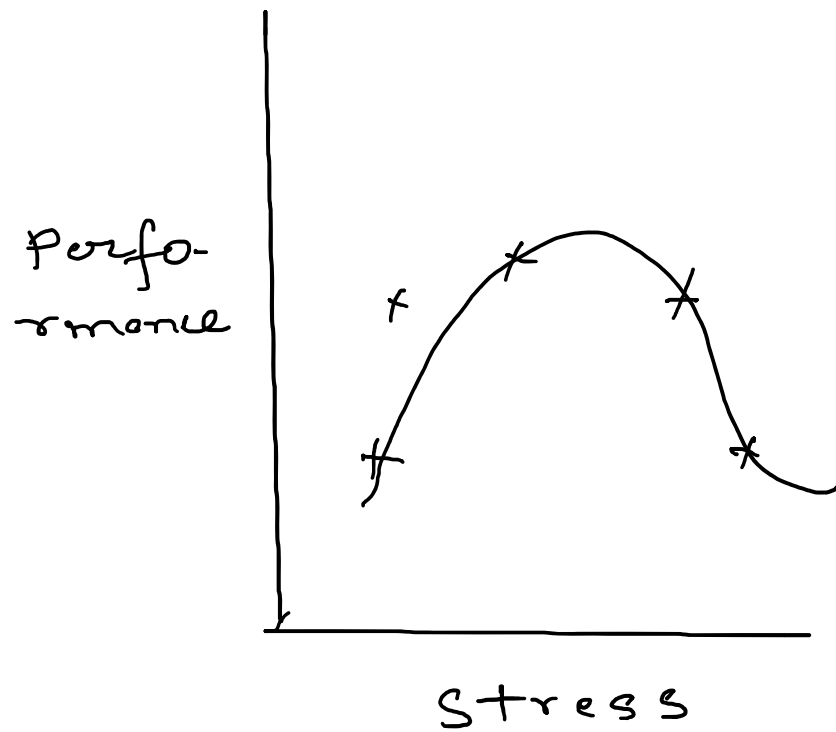
positive
correlation



Negative
correlation



No
correlation



Coefficient of correlation:



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

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

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

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

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



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

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

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

✓

x	$x - 5$	x^2	y	$y - 5$	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
		$\sum 60$		$\sum 60$		$\sum 59$

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

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

$$= 0.9833$$

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.

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

is 81% of the variation in y
due to variation in x

remaining 19% is due to some other factors.



Regression

Regression :-

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

when $x = 7$: $y = ?$

x	1	2	3	4	5
y	1	6	2	5	4

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

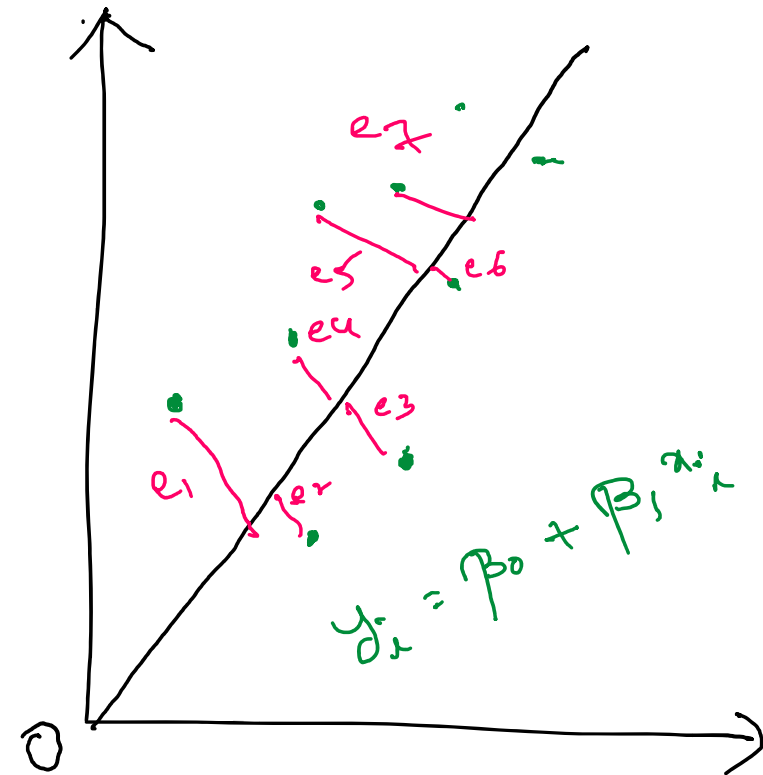
Method of Least Squares



y : Dependent Variable

x : Independent Variable

predictor variable



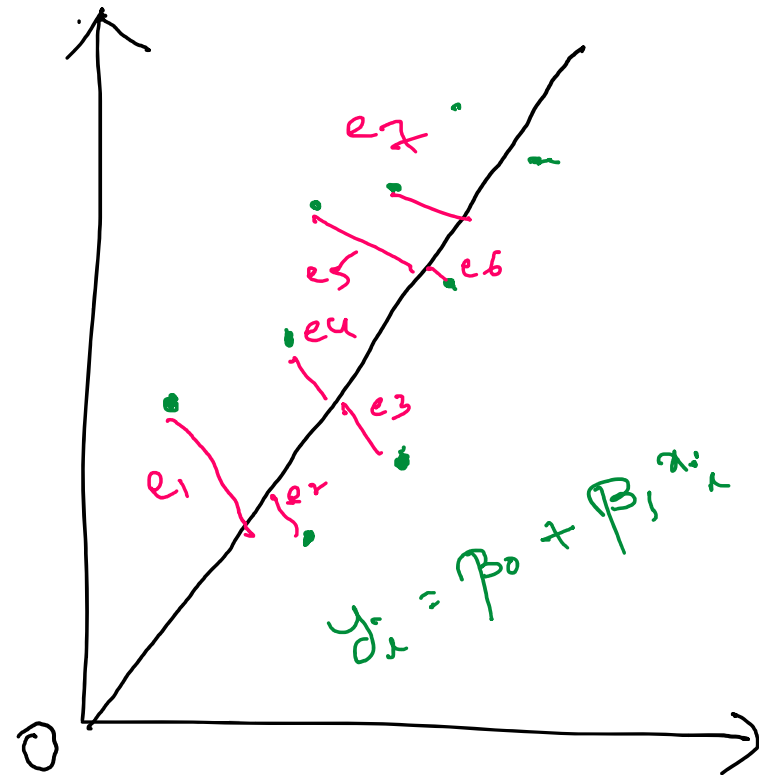
response variable

Method of Least Squares



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

↓
 β_0 and β_1 are chosen to minimize the sum of squared errors.



Method of Least squares



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

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

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

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

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

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

Linear regression

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

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

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

normal equations.

Matrix Approach:



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

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

unknowns $\beta_0, \beta_1, \dots, \beta_{p-1} \rightarrow$ " " β

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

$$\hat{\gamma}_{n \times 1} = X_{n \times p} \beta_{p \times 1}$$

Find β to minimize

$$S(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots)^2$$

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

Diff S w.r.t 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

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

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.

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
	↓ x_1	↓ x_2	↓ x_3	↓ x_4	↓ y

Linear regression (multiple regression)



example:-

$X =$

y

1	2000	5	2	45	4000
1	1400	3	1	40	2000
1	1600	3	2	30	3000
1	800	2	1	35	2000

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

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

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

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

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

$$3.5 = 4\beta_0 + \beta_1 \quad (1)$$

$$10.5 = 7\beta_0 + \beta_1 \quad (2)$$

on solving these

$$\beta_0 = 0$$

$$\beta_1 = 0.5$$

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

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



Thanks