

Regression Analysis

- It is a statistical tool to find the relationship between one dependent variable and one or more independent variables.
- Example: Consider a Sales Company Dataset and you are a Marketing Analyst of the company.
- Let the dataset has attributes like Adv. Cost (Rs.) and Sales Amount (Quantity)

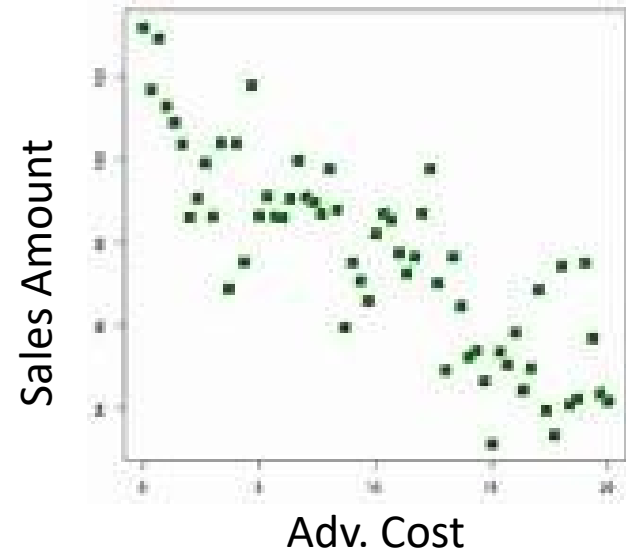
Adv. Cost (Rs.)	Sales Amount
1	1
2	1
3	2
4	2
5	4

Regression Analysis

- Here, we may decide how much money we can spent for advertising. So, amount of money spent is the controlled variable .
- But, we can't control the Sales amount, so it is not a controlled variable. This is dependent variable depends on Advertising Cost.
- This is not only the factor, Sales amount may also dependent on some other factors like no. of persons working, etc.
- Adv. Cost : $X \rightarrow$ Independent or Regressor Variable
- Sales Amount: $Y \rightarrow$ Dependent or Response or Random Variable

Regression Analysis

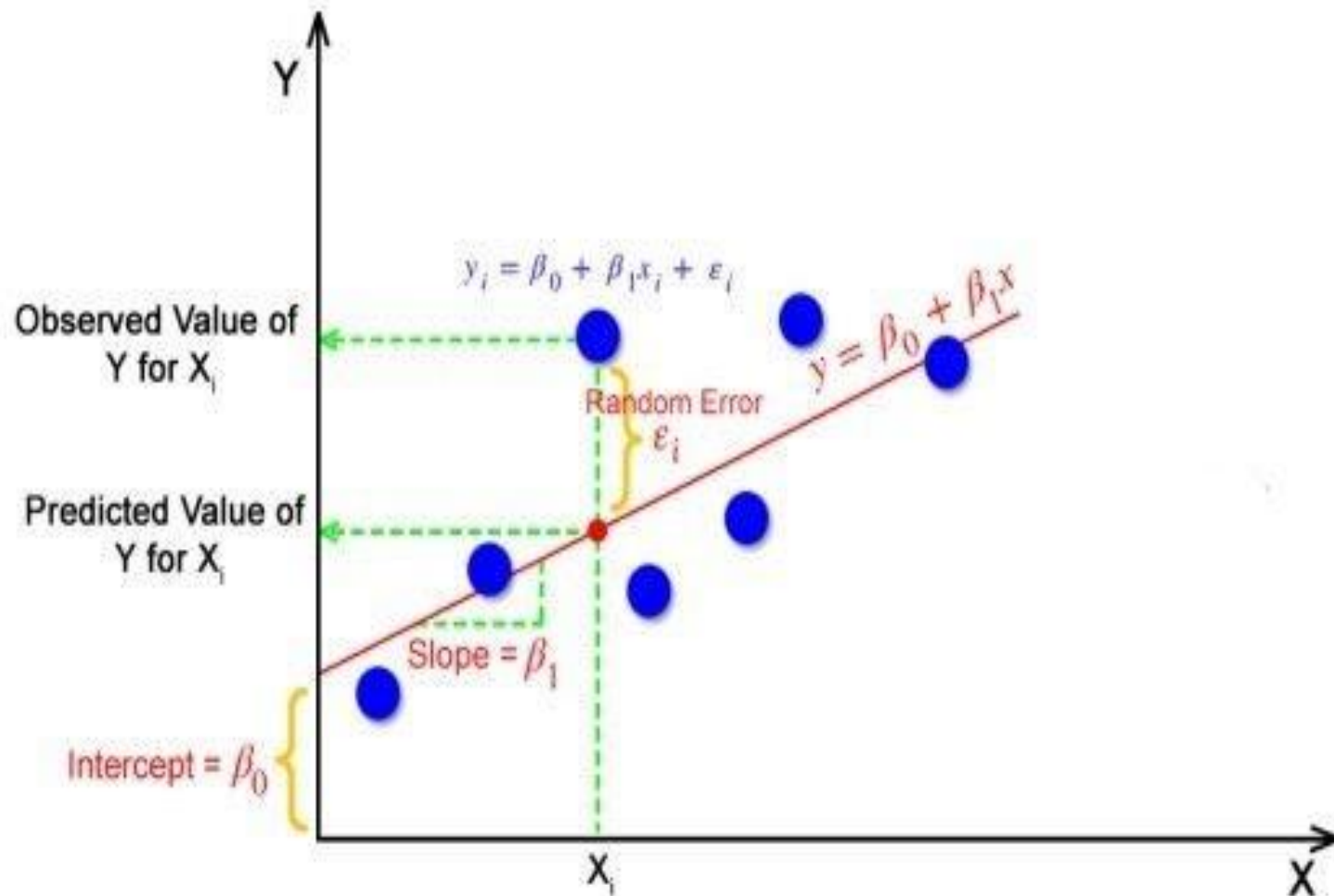
- Scatter plot is a mathematical diagram to display values of two variables for a set of data.
- It is used to investigate the possible relationship between the variables.
- If it indicates the linear relation then linear regression is considered; otherwise polynomial regression, and so on are considered.



Simple Linear Regression (SLR) Model

- It is a model with a single regressor variable X that has linear relationship with a response variable Y .
- The simple linear regression model is:
$$Y = a + cX + \epsilon$$
, where $a \rightarrow$ intercept, $c \rightarrow$ Slope, and ϵ is a random error component.
 \Rightarrow For a given X , the corresponding observation Y consists of the value $a + cX + \epsilon$.
- The same model may be written as:
$$y_i = a + cx_i + \epsilon_i$$
 for $i = 1, 2, \dots, n$; n is the no. of observations.

Simple Linear Regression (SLR) Model



Simple Linear Regression (SLR) Model

- Let the best fitted model is:

$$\hat{Y} = \hat{a} + \hat{c}X \quad \text{or}$$

$$\hat{y}_i = \hat{a} + \hat{c}x_i; i=1,2,\dots,n$$

- The line fitted by Least Square Method (LSM) which makes the sum of square of all vertical discrepancies as small as possible.
- The LSM estimates the parameters a and c using the dataset $\langle X, Y \rangle$ as \hat{a} and \hat{c} .

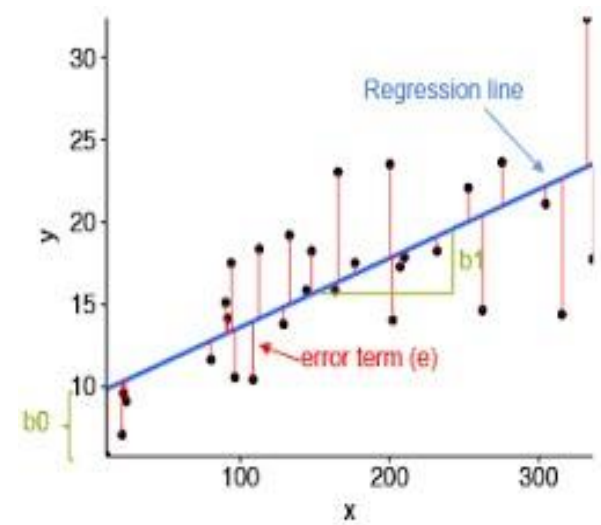
X Y

x_1 y_1

x_2 y_2

x_i y_i

x_n y_n



Simple Linear Regression (SLR) Model

- The cost function or error function of the LSM is:

$$S = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- We estimate a and c so that sum of square of all the differences between the observed y_i and the predicted \hat{y}_i is minimum, i.e., S is minimum.
- This S is called Sum of Square Residual, i.e., $SS_{\text{Res}} = S = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- Residual = Deviation between actual and predicted value
- Error = Deviation between actual value and mean of population.

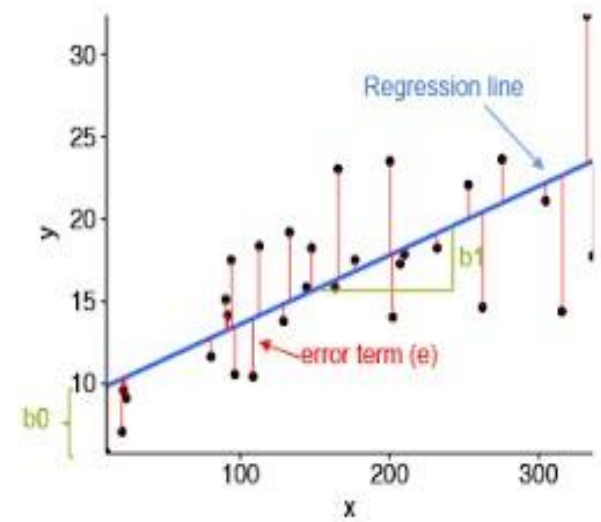
X Y

x_1 y_1

x_2 y_2

x_i y_i

x_n y_n



Parameter estimation using LSM

$$SS_{Res} = S = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{a} - \hat{c}x_i)^2$$

- By taking the partial derivatives = 0 and solving, we get the best fitted model is:

$\hat{Y} = \hat{a} + \hat{c}X$, where

$$\hat{a} = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{c} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Parameter estimation using LSM

The least square estimator of a & b
 (\hat{a} & \hat{b}) must satisfies: $\frac{\partial S}{\partial a} \Big|_{\hat{a}, \hat{b}} = 0$ and $\frac{\partial S}{\partial b} \Big|_{\hat{a}, \hat{b}} = 0$

$SS_{\text{res}} = S = \sum (y_i - \hat{y}_i)^2$
 $= \sum (y_i - \hat{a} - \hat{b}x_i)^2$

$\therefore \frac{\partial S}{\partial a} = 0$ and $\frac{\partial S}{\partial b} = 0$

$-2 \sum (y_i - \hat{a} - \hat{b}x_i) = 0$ and $-2 \sum x_i (y_i - \hat{a} - \hat{b}x_i) = 0$
 $\quad \quad \quad (1) \quad \quad \quad (2)$

$(1) \Rightarrow \sum (y_i - \hat{a} - \hat{b}x_i) = 0$
 $\Rightarrow \sum y_i = \sum \hat{a} + \sum \hat{b}x_i = 0$
 $\Rightarrow \sum y_i - n\hat{a} - \hat{b}\sum x_i = 0$
 $\Rightarrow \hat{a} = \frac{\sum y_i}{n} - \hat{b} \frac{\sum x_i}{n}$

$= \bar{y} - \hat{b}\bar{x}$ where $\bar{y} = \frac{\sum y_i}{n}$ and $\bar{x} = \frac{\sum x_i}{n}$
 $\quad \quad \quad (3)$

Parameter estimation using LSM

$$(2) \Rightarrow \sum x_i (y_i - \hat{a} - \hat{b} x_i) = 0$$

$$\Rightarrow \sum x_i (y_i - \bar{y} + \bar{b} - \hat{b} x_i) = 0 \quad [\text{using (3)}]$$

$$\Rightarrow \sum x_i (y_i - \bar{y}) + \bar{b} \sum x_i (\bar{x} - x_i) = 0$$

$$\Rightarrow \hat{b} \sum (x_i - \bar{x}) x_i = \sum (y_i - \bar{y}) x_i$$

$$\Rightarrow \hat{b} = \frac{\sum (y_i - \bar{y}) x_i}{\sum (x_i - \bar{x}) x_i}$$

DECEMBER 2018								JANUARY 2019							
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
48	30	31				1		1			1	2	3	4	5
49	2	3	4	5	6	7	8	2	6	7	8	9	10	11	12
50	9	10	11	12	13	14	15	3	13	14	15	16	17	18	19
51	16	17	18	19	20	21	22	4	20	21	22	23	24	25	26
52	23	24	25	26	27	28	29	5	27	28	29	30	31		

$$\hat{b} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})(x_i - \bar{x})} \quad (4)$$

$$\text{if } \sum (y_i - \bar{y}) \bar{x} = 0$$

$$\text{and } \sum (x_i - \bar{x}) \bar{x} = 0$$

Parameter estimation using LSM

$$\begin{aligned} \rightarrow \hat{b} &= \frac{\sum (x_i y_i) - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} \\ &= \frac{\sum x_i y_i - \frac{\sum y_i}{n} \sum x_i}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} \end{aligned}$$

Important :

FEBRUARY 2019								MARCH 2019							
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
5						1	2	9	31					1	2
6	3	4	5	6	7	8	9	10	3	4	5	6	7	8	9

Parameter estimation using LSM

$$\begin{aligned}
 \text{from (3)} \quad \hat{a} &= \bar{y} - \hat{b} \bar{x} \\
 &= \bar{y} - \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \cdot \bar{x} \\
 &= \frac{\sum y_i}{n} - \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \cdot \frac{\sum x_i}{n} \\
 &= \frac{\sum y_i \{ n \sum x_i^2 - (\sum x_i)^2 \} - \sum x_i (n \sum x_i y_i - \sum x_i \sum y_i)}{n \{ n \sum x_i^2 - (\sum x_i)^2 \}} \\
 &= \frac{n \sum x_i^2 \sum y_i - (\sum x_i)^2 \sum y_i - n \sum x_i \sum x_i y_i + (\sum x_i)^2 \sum y_i}{n \{ n \sum x_i^2 - (\sum x_i)^2 \}} \\
 &= \frac{n \sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} \quad \text{--- (6).}
 \end{aligned}$$

Parameter estimation using LSM

2019. The simple regression model is

$Y = \hat{a} + \hat{b}X$ where

$\hat{a} = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$

$\hat{b} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$

$\epsilon \leftarrow$ random error variable

during computation some error may occur due to instrument.

JANUARY 17 THURSDAY 3rd Week • 017-348

two thousand Nineteen

09.00

10.00

11.00

Basic Assumptions on SLR Model

- In simple LRM, $y_i = a + cx_i + \epsilon_i$ for $i = 1, 2, \dots, n$:
 - i) ϵ_i is a random variable with zero mean and variance σ^2 (Unknown), i.e., $E(\epsilon_i)=0$ and $\text{Var}(\epsilon_i)=\sigma^2$
 - ii) ϵ_i and ϵ_j are uncorrelated, $i \neq j$, i.e., $\text{COV}(\epsilon_i, \epsilon_j) = 0$
 - iii) ϵ_i is a normally distributed random variable with zero mean and variance σ^2 . i.e., $\epsilon_i \sim N(0, \sigma^2)$
- So, ϵ_i 's are normally distributed and uncorrelated \Rightarrow ϵ_i 's are independent.

Consequences in terms of y_i

- $y_i = a + cx_i + \epsilon_i$ for $i = 1, 2, \dots, n$:

$$\begin{aligned} \text{i)} \quad E(y_i) &= E(a + cx_i + \epsilon_i) \\ &= E(a) + E(cx_i) + E(\epsilon_i) \\ &= a + x_i E(c) + 0 \text{ [as } x \text{ is controlled variable, and } E(\epsilon_i)=0] \\ &= a + cx_i \end{aligned}$$

$$\begin{aligned} \text{ii)} \quad \text{Var}(y_i) &= \text{Var}(a + cx_i + \epsilon_i) \\ &= \text{Var}(a) + \text{Var}(cx_i) + \text{Var}(\epsilon_i) \\ &= 0 + 0 + \sigma^2 = \sigma^2 \end{aligned}$$

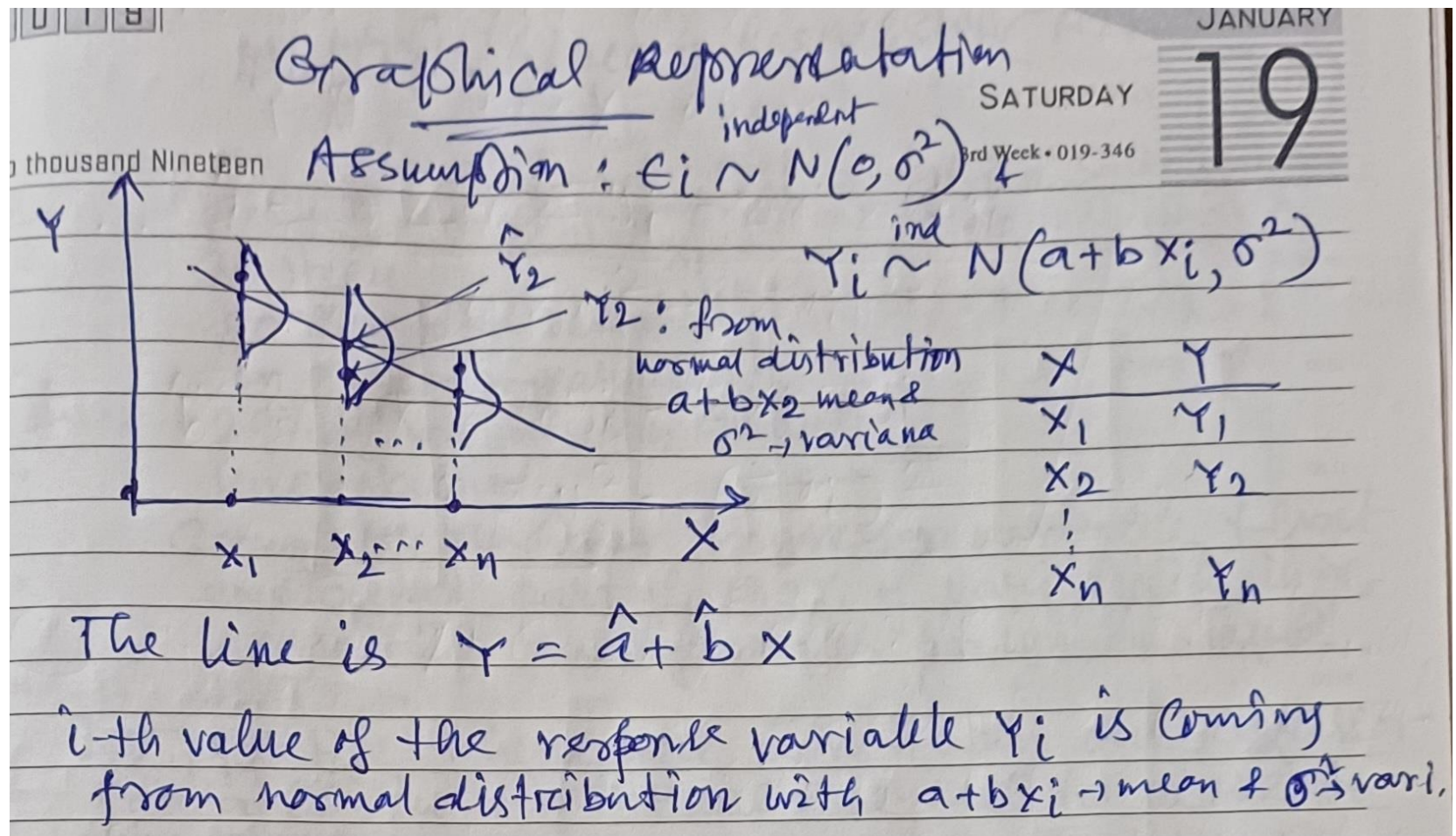
Since, ϵ_i follows normal distribution with zero mean and σ^2 variance, so the consequence of y_i is that y_i follows normal distribution with $a + cx_i$ mean and σ^2 variance,

i.e., $y_i \sim N(a + cx_i, \sigma^2)$

Consequences in terms of y_i

- So, we assume that, i-th observation, y_i is from normal distribution with mean = $a + cx_i$ and $\text{Var} = \sigma^2$ (Also, y_i 's are uncorrelated and independent)
- Therefore, given a set of data, the dataset must satisfy this assumptions.
- If the assumptions are not hold, then we should not apply this regression analysis. This is verify using topic modeling adequacy checking (study in your own, if interested).

Graphical interpretation of Assumptions in SLR Model



Multiple Linear Regression (MLR) Model

- Consider the same Company Sales Dataset.
- In case of SLR, we assume that the response variable “Sales Amount” is fully explained by the regressor variable “Adv. Cost”
- But in reality, it may be say, 80% explained by “Adv. Cost”. Remaining 20% may be explained by other factor, say “No. of sales person” employed.
- In practice, there are more than one regressor variables, in that case, we consider MLR.

MLR Model

MLR:

MLR Model:

more than one regressor variable,
say $k-1$ vari, x_1, x_2, \dots, x_{k-1}

General form of MLR is: $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-1} x_{k-1} + \epsilon$

or $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{k-1} x_{ik-1} + \epsilon_i$

for $i = 1, 2, \dots, n$

Here we have more
than 1 regressor variable
so called multiple regression
and the model is linear

in terms of ~~unknown~~ regressor variables.
parameters $\beta_0, \beta_1, \dots, \beta_{k-1}$

β_{k-1} (not in x_1, x_2, \dots, x_{k-1})
ith observation \rightarrow

	β_1	β_2	β_{k-1}	Y	ϵ
	x_1	x_2	\dots	x_{k-1}	Y
x_{11}	x_{12}	\dots	$x_{1,k-1}$	Y_1	ϵ_1
x_{21}	x_{22}	\dots	$x_{2,k-1}$	Y_2	ϵ_2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{i1}	x_{i2}	\dots	$x_{i,k-1}$	Y_i	ϵ_i
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{n1}	x_{n2}	\dots	$x_{n,k-1}$	Y_n	ϵ_n

important:

The model is called
MLR

FEBRUARY 2019								MARCH 2019							
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
5						1	2	9	31					1	2
6	3	4	5	6	7	8	9	10	3	4	5	6	7	8	9
7	10	11	12	13	14	15	16	11	10	11	12	13	14	15	16
8	17	18	19	20	21	22	23	12	17	18	19	20	21	22	23

MLR Model

JANUARY 23 WEDNESDAY 4th Week • 023-342

Assumption: error $\epsilon_i \sim N(0, \sigma^2)$ and ϵ_i are independent

Define matrices:

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1}$$

vector of observations on Y

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{K-1} \end{bmatrix}_{K \times 1}$$

vector of parameters

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

vector of errors

and

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

Dataset

x_1	x_2	\dots	x_{K-1}	Y
x_{11}	x_{12}	\dots	$x_{1,K-1}$	Y_1
\vdots	\vdots	\ddots	\vdots	\vdots
x_{n1}	x_{n2}	\dots	$x_{n,K-1}$	Y_n

MLR Model

This model can be expressed as:

$$Y = X\beta + \epsilon \text{ in matrix form.}$$

$$X: n \times k, \beta: k \times 1$$

$$\Rightarrow X\beta: n \times 1$$

$$\epsilon: n \times 1$$

$$\therefore Y: n \times 1$$

we have to fit this model means
we have to estimate the parameters.

Estimation of Model parameters (MLR)

Like Simple Linear Regression (SLR), we will estimate

using Least Square ~~estimation~~ method, where
the parameters are ~~estimated~~ determined by

minimizing the
SS Residual (sum of square
residual)

DECEMBER								2018	JANUARY								2019
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S		
48	30	31					1	1			1	2	3	4	5		
49	2	3	4	5	6	7	8	2	6	7	8	9	10	11	12		
50	9	10	11	12	13	14	15	3	13	14	15	16	17	18	19		
51	16	17	18	19	20	21	22	4	20	21	22	23	24	25	26		
52	23	24	25	26	27	28	29	5	27	28	29	30	31				

Important:

MLR Model

Least square method (LSM) THURSDAY 24
 4th Week • 024-341
 determines the parameters by minimizing SS_{res} where

$$SS_{res} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Let the fitted model is: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_{k-1} X_{k-1}$

$$\therefore SS_{res} = \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \hat{\beta}_2 X_{i2} - \dots - \hat{\beta}_{k-1} X_{i(k-1)})^2$$

We now represent SS_{res} in matrix form for that we consider residual vector e as

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad e_i = i\text{-th residual}$$

$$e_i = Y_i - \hat{Y}_i \quad \forall i = 1, 2, \dots, n$$

$$e' = (e_1, e_2, \dots, e_n)$$

$\therefore e = Y - \hat{Y}$

vector of observation \rightarrow vector of observations for the fitted value

MLR Model

05.00 $\therefore SS_{Res} = \sum_{i=1}^n e_i^2 = e'e$ $e'e = [e_1 e_2 \dots e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$

06.00 $= (Y - \hat{Y})'(Y - \hat{Y})$

$\therefore Y = XB + e \Rightarrow e = (Y - X\hat{\beta})'(Y - X\hat{\beta})$ $= e_1^2 + e_2^2 + \dots + e_n^2 = \sum_{i=1}^n e_i^2$

and $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_{k-1} X_{k-1}$

$= (Y' - \hat{\beta}' X') (Y - X\hat{\beta})$

$= X\hat{\beta}$ (in matrix form) $= Y'Y - Y'X\hat{\beta} - \hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$

$= 1 \times n, n \times 1 \Rightarrow 1 \times 1 \Rightarrow$ scalar quantity $1 \times n, n \times k, k \times 1 \Rightarrow 1 \times 1 \Rightarrow$ scalar quantity $1 \times k, k \times n, n \times 1 \Rightarrow 1 \times 1 \Rightarrow$ scalar quantity $1 \times k, k \times n, n \times k, k \times 1 \Rightarrow 1 \times 1 \Rightarrow$ scalar quantity

Important
Transpose of
the scalar
is itself.

\Rightarrow every
term is
scalar
quantity

FEBRUARY							2019	MARCH							2019
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
5						1	2	9	31					1	2
6	3	4	5	6	7	8	9	10	3	4	5	6	7	8	9
7	10	11	12	13	14	15	16	11	10	11	12	13	14	15	16
8	17	18	19	20	21	22	23	12	17	18	19	20	21	22	23
9	24	25	26	27	28			13	24	25	26	27	28	29	30

MLR Model

JANUARY 25 FRIDAY 4th Week • 025-340

2019 two thousand Nineteen

$$Y'X\hat{\beta} = (X'Y)'\hat{\beta} = (\hat{\beta}'(X'Y))'$$

1x1 matrix

For 1x1 matrix A , $A=A'$

$$= (\hat{\beta}'X'Y)' = \hat{\beta}'X'Y$$

1x1 matrix

$$SS_{res} = Y'Y - Y'X\hat{\beta} - \hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$$

$$= Y'Y - \hat{\beta}'X'Y - \hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$$

$$= Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$$

∴ We have two SS_{res} in different forms.

For estimation of K -unknowns $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}$ we get,

$$\frac{\partial SS_{res}}{\partial \beta_0} \bigg|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}} = 0$$

$$\frac{\partial SS_{res}}{\partial \beta_1} \bigg|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}} = 0$$

$$\vdots$$

$$\frac{\partial SS_{res}}{\partial \beta_{K-1}} \bigg|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}} = 0$$

So we have K normal equations with K -unknowns \Rightarrow solution gives the estimated values of $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}$

MLR Model

Process LSM for both forms;

1. Matrix form

$$SS_{Res} = Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$$

$$\frac{\partial SS_{Res}}{\partial \hat{\beta}} = 0 \Rightarrow$$

$$\frac{\partial}{\partial \hat{\beta}} (Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}) = 0$$

$$0 - 2X'Y + 2X'X\hat{\beta} = 0$$

$$\Rightarrow -2X'Y + 2X'X\hat{\beta} = 0$$

2nd Form

$$SS_{Res} = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_{k-1} X_{i,k-1})^2$$

normal eqs.

$$\frac{\partial SS_{Res}}{\partial \hat{\beta}_0} = 0 \Rightarrow$$

$$\sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_{k-1} X_{i,k-1}) = 0$$

$$\Rightarrow \sum_{i=1}^n e_i = 0$$

DECEMBER 2018								JANUARY 2019							
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
48	30	31						1	1	2	3	4	5	6	7
49	2	3	4	5	6	7	8	2	6	7	8	9	10	11	12
50	9	10	11	12	13	14	15	3	13	14	15	16	17	18	19
51	16	17	18	19	20	21	22	4	20	21	22	23	24	25	26

Important:

MLR Model

2019

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

two thousand Nineteen

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

09.00

10.00

11.00

12.00

01.00

02.00

03.00

04.00

05.00

this will give all the unknown parameters.

$$Y = X\beta$$

$$(Y'X)(X'X)^{-1}(X'X)\beta + (Y'X)(X'X)^{-1}(X'Y)$$

$$\frac{d}{dx}(X'B) = B$$

$$\frac{d}{dx}(X'b) = b$$

$$\frac{d}{dx}(X'X) = 2X$$

$$\frac{d}{dx}(X'BX) = 2BX$$

SATURDAY

4th Week • 026-339

JANUARY

26

$$\frac{\partial SS_{Res}}{\partial \beta_1} = 0 \Rightarrow$$

$$\sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_{k-1} X_{i,k-1}) \cdot X_{i1} = 0$$

$$\Rightarrow \sum (Y_i - \hat{Y}_i) X_{i1} = 0$$

$$\Rightarrow \sum_{i=1}^n e_i X_{i1} = 0$$

similarly we get

$$\sum_{i=1}^n e_i X_{i,k-1} = 0$$

so we have k-normal equations

$$\sum e_i = 0, \sum e_i X_{i1} = 0, \dots, \sum e_i X_{i,k-1} = 0$$

All are independent which gives the soln of k-unknowns $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{k-1}$

Degree of Freedom (DF) of SS_{Res}

2019 FEBRUARY 09 SATURDAY 6th Week • 040-325

Linear regression

two thousand Nineteen

We may represent SS_{residual} i.e., SS_{res} in many different forms, like

i) $SS_{\text{res}} = \sum_{i=1}^n e_i^2$

ii) $SS_{\text{res}} = Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}$
where $\hat{\beta} = (X'X)^{-1}X'Y$

$\therefore SS_{\text{res}} = Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X(X'X)^{-1}X'Y$
 $= Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'Y$
 $= Y'Y - \hat{\beta}'X'Y$

* Now, find degree of freedom (DF) of SS_{res}

$\therefore SS_{\text{res}} = \sum_{i=1}^n e_i^2$ where $e_i \sim N(0, \sigma^2)$ and e_i satisfies k -Constraints:

$$\left. \begin{aligned} \sum e_i &= 0 \\ \sum e_i x_{i1} &= 0 \\ &\vdots \\ \sum e_i x_{i,k-1} &= 0 \end{aligned} \right\}$$

Degree of Freedom (DF) of SS_{Res}

6.00
 ∴ To compute SS_{res} , we don't have freedom to choose all e_i for $i = 1, 2, \dots, n$ independently. We can choose $n-k$ of e_1, e_2, \dots, e_n independently if, we have the degree of freedom to choose $n-k$ of n e_i 's and the remaining k e_i 's have to be chosen in such a way that they satisfy above k -constraints.

Important:

∴ We are losing k degree of freedom for k -constraints on the residuals.

$$\therefore DF(SS_{res}) = n - k.$$

MARCH								APRIL							
2019								2019							
WK	S	M	T	W	T	F	S	WK	S	M	T	W	T	F	S
9	3					1	2	14		1	2	3	4	5	6
10		4	5	6	7	8	9	15	7	8	9	10	11	12	13
11	10	11	12	13	14	15	16	16	14	15	16	17	18	19	20
12	17	18	19	20	21	22	23	17	21	22	23	24	25	26	27
13	24	25	26	27	28	29	30	18	28	29	30				

Degree of Freedom (DF) of SS_T

FEBRUARY 11 MONDAY 7th Week • 042-323 two thousand nineteen

SS_{Total} is SS_T is the variation i.e., variability in the response variable (Y).

$\therefore SS_T = \sum_{i=1}^n (Y_i - \bar{Y})^2$ in matrix form.

$$\begin{aligned} &= \sum (Y_i^2 - 2Y_i\bar{Y} + \bar{Y}^2) \\ &= \sum Y_i^2 - 2\bar{Y} \sum Y_i + \bar{Y}^2 \sum 1 \\ &= \sum Y_i^2 - 2n\bar{Y} + n\bar{Y}^2 \\ &= \sum Y_i^2 - n\bar{Y}^2 \end{aligned}$$

partial derivative $\frac{\partial}{\partial \bar{Y}} \sum (Y_i - \bar{Y})^2 = 0$
choose $n-1$ free

Here, $SS_T =$ sum of n -terms i.e., $\sum (Y_i - \bar{Y})^2$

To compute SS_T we must satisfy the constraint $\sum_{i=1}^n (Y_i - \bar{Y}) = 0 \leftarrow \text{obvious}$

Degree of Freedom (DF) of SS_T

\therefore To compute SS_T we don't have freedom to choose all the n -terms independently. out of all the n -terms $(Y_1 - \bar{Y})$, $(Y_2 - \bar{Y})$, \dots , $(Y_n - \bar{Y})$ we may choose at most $(n-1)$ terms independently and 1-term should be chosen in such a way that the constraint $\sum (Y_i - \bar{Y}) = 0$

$\therefore DF(SS_T) = n-1$

Degree of Freedom (DF) of SS_{Reg}

$SS_{\text{regressor}}$ is, SS_{reg} defines how much of the variability in response variable Y is explained by the regressor model.

$\therefore SS_T = SS_{\text{reg}} + SS_{\text{res}}$

\downarrow Total variability in response variable \downarrow variability in response variable explained by model \downarrow variability in Y not explained by model (i.e., error)

\therefore Matrix form of SS_{reg} is

$$\begin{aligned} SS_{\text{reg}} &= SS_T - SS_{\text{res}} \\ &= \sum Y_i^2 - n\bar{Y}^2 - (Y'Y - \hat{\beta}'X'Y) \quad \leftarrow \text{From previous page} \\ &= \cancel{Y'Y} - n\bar{Y}^2 - \cancel{Y'Y} + \hat{\beta}'X'Y \\ &= \hat{\beta}'X'Y - n\bar{Y}^2 \end{aligned}$$

$\therefore SS_T = SS_{\text{reg}} + SS_{\text{res}}$

$\therefore DF(SS_T) = DF(SS_{\text{reg}}) + DF(SS_{\text{res}})$

$\Rightarrow n-1 = DF(SS_{\text{reg}}) + n-K$

$\Rightarrow DF(SS_{\text{reg}}) = K-1$

Thank you