

Multiple linear Regression



→ Using Adjusted R^2 score is a good idea during Multiple Linear Regression.

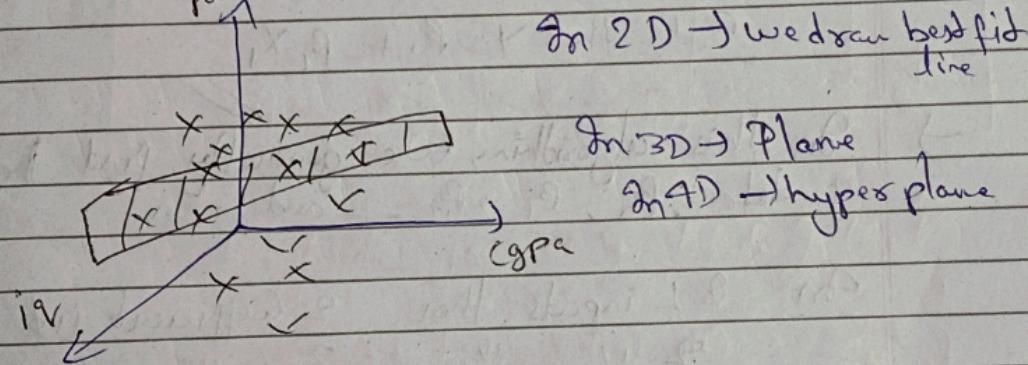
A Multiple Linear Regression

→ It is the extension form of Simple linear regression
→ we have more than one input in this.

Let's suppose $\text{GPA} | \text{IV} | \text{Package(Y)}$

(Now data is in 3D)

package(Y)



We take a plane or In 3D we fit a plane.

→ The plane is the surface that minimizes the distance (error) between actual and predicted value

Formulation → In Simple LR $\Rightarrow y = mx + c$

but for more inputs, noce the equation will look like :-

Suppose for two inputs

where y plane indepd at y.

$$y = m_1 X_1 + m_2 X_2 + b$$

m_2 is the slope of X_2 feature

m_1 is the slope of X_1 feature (deh)

Literature expression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$\downarrow \downarrow \downarrow$
 $m_1 X_1 \quad m_2$

We have to find $[\beta_0, \beta_1, \beta_2]$ m_1, m_2
 $\beta_0 = b$

→ If our data is in 4D

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

→ Same as for n inputs

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n$$

$$Y = \beta_0 + \sum_{i=1}^n \beta_i X_i \quad \text{General equation}$$

for 2D

$$\text{if } n=1 \text{ then } Y = \beta_0 + \beta_1 X_1$$

→ In this algorithm we generally find the coefficients of data $\beta_0, \beta_1, \beta_2, \dots$ and so on

For 2 inputs then 2 coefficients find

For " " " 3 Coeff - 1

If n inputs then n+1 coefficients, needs to find

• Taking example

$X_1 \rightarrow \text{GPA}$, $X_2 \rightarrow \text{IQ}$, $Y \rightarrow \text{LPA}$ (Package)

$$Y = \beta_0 + \beta_1 X_1 (\text{GPA}) + \beta_2 X_2 (\text{IQ})$$



These are the weights. Like this Y is dependent

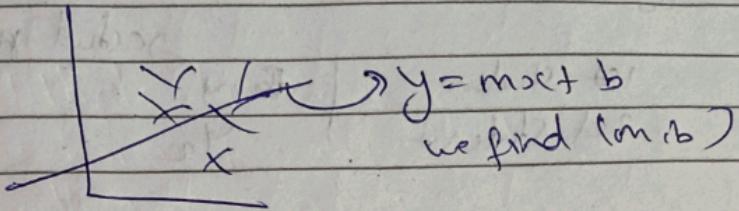
on β_2 of IQ

$\text{IQ} \uparrow \uparrow Y \uparrow$, $\text{IQ} \downarrow \downarrow B_2 \downarrow$ Just like this

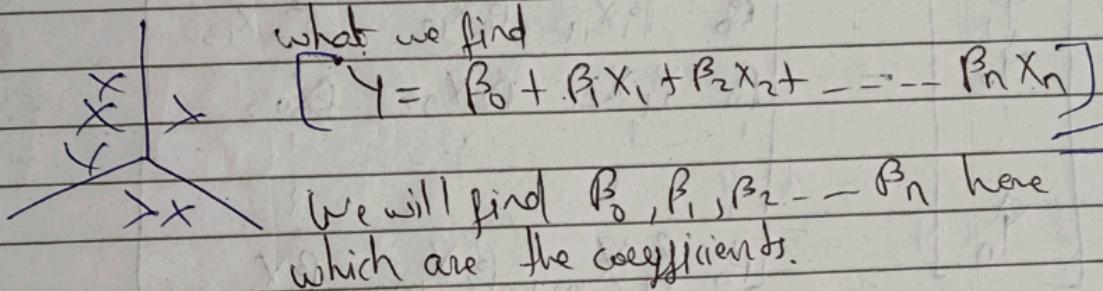
β_0 = offset (if all others are zero then)

A Mathematical formulation from Scratch

In Simple



here in 2D or 3D



Let's Suppose

GPA	IQ	Gender	LPA (Package) Target
X_1	X_2	X_3	Y

for only 1 input (In 2D) $\Rightarrow y = mx + b$

$$y = B_0 + B_1 X_1$$

Similarly for all

$$Y = B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3$$

↳ Predicted value

Suppose we have 100 students $(100, 4)$ 4 columns.

On next page

$\beta_0 \beta_1 \beta_2 \beta_3$ we have to find it
 like in original dataset each student has their CGPA
 package (y) give

1st student | $\beta_0 + \beta_1 x_1 + \beta_2 x_2$
 2nd student | $\beta_0 + \beta_1 x_1 + \beta_2 x_2$

let's make a matrix

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} \\ \vdots & \vdots & \vdots & \vdots \\ \beta_0 & \beta_1 x_{100,1} & \beta_2 x_{100,2} & \beta_3 x_{100,3} \end{bmatrix}$$

$y_1 \Rightarrow$ predicted package of 1st student
 $y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \beta_3 x_{23}$ 2nd student

$\beta_1 x_{11} =$ CGPA of 1st student

2nd student in the table 1st row 1st column

$\beta_3 x_{13} \Rightarrow$ CGPA of the 2nd student

$\beta_1 x_{100,1} \Rightarrow$ 100th student's CGPA

100th row 1st feature

$\beta_1 x_{100,2} \Rightarrow$ CGPA of 100th student

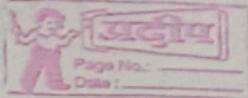
• Taking a small assumption for better

We don't have 100 rows we have n rows

" " " 3 input col + 1 output $\rightarrow m$ columns

100 rows $\rightarrow n$ rows

4 cols \rightarrow 3 col (input) \rightarrow m cols



$$\hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \beta_3 X_{13} + \dots + \beta_m X_{1m} \\ \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \beta_3 X_{23} + \dots + \beta_m X_{2m} \\ \vdots \\ \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \beta_3 X_{n3} + \dots + \beta_m X_{nm} \end{bmatrix}$$

wrong

$$\hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 X_{11} & \beta_2 X_{12} & \beta_3 X_{13} & \dots & \beta_m X_{1m} \\ \beta_0 & \beta_1 X_{21} & \beta_2 X_{22} & \beta_3 X_{23} & \dots & \beta_m X_{2m} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta_0 & \beta_1 X_{n1} & \beta_2 X_{n2} & \beta_3 X_{n3} & \dots & \beta_m X_{nm} \end{bmatrix}$$

$\beta_1 X_{11} \Rightarrow$ 1st student ka 1st value (feature 1 me)

$\beta_2 X_{22} \Rightarrow$ 2nd student ka 2nd value (feature 2 me)

We can also write the Matrix \hat{Y} as following
(Separating $\beta_0, \beta_1 - \beta_n$) Product ke form me

$$\hat{Y} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1m} \\ 1 & X_{21} & X_{22} & \dots & X_{2m} \\ 1 & X_{31} & X_{32} & \dots & X_{3m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix}$$

by finding its product we will also get \hat{Y}

Let's we will called this matrix as X and this β

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

\hat{Y} matrix = X matrix * β matrix

How we call $\hat{Y} = X\beta$ in English or professor

\hat{Y} is \rightarrow The prediction of all the input rows (of all student)
 \hat{Y} is the package of each student predicted

If we have 1000 students the \hat{Y} will contain 1000 predicted values in the matrix

$$\hat{Y} = X\beta \quad \text{for 2D}$$

Now from here what is $\beta \Rightarrow \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$ and so on

Q Now, what is X matrix

Suppose if we have CGPA | IQ | gender

$X_1 \quad X_2 \quad X_3$

This is X' matrix

for 2 extra col $X_1 \quad X_2 \quad X_3$

1	X_{11}	X_{12}	X_{13}
1	X_{21}	X_{22}	X_{23}
1	X_{31}	X_{32}	X_{33}

Q form a given df $\rightarrow XY$ if we extract Y

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \text{It is the actual value in the data}$$

Now we are creating a matrix $E = Y - \hat{Y} = \text{actual - Predicted}$

$$E = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} - \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix} = \begin{bmatrix} Y_1 - \hat{Y}_1 \\ Y_2 - \hat{Y}_2 \\ \vdots \\ Y_n - \hat{Y}_n \end{bmatrix}$$

Now in single LR

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

which is the error

And now we have $E = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$

if we will find $E^T E$ it will be equal to e

\downarrow
Proving this

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

$$E^T E$$

$$E^T$$

$$[y_1 - \hat{y}_1, (y_2 - \hat{y}_2), (y_3 - \hat{y}_3), \dots, (y_n - \hat{y}_n)] \begin{bmatrix} (y_1 - \hat{y}_1)^2 \\ (y_2 - \hat{y}_2)^2 \\ \vdots \\ (y_n - \hat{y}_n)^2 \end{bmatrix}$$

\uparrow shape $(1, n)$

Now after multiply we will get

shape $(n, 1)$

$$= [(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + \dots + (y_n - \hat{y}_n)^2]$$

$$\Rightarrow (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\text{Same as } e)$$

\downarrow Error function

So, we can write our Error function E (of matrix a)

$$E = E^T E = (y - \hat{y})^T (y - \hat{y})$$

And we know that (from 12th class notes)
 $(A+B)^T = A^T + B^T$
 $(A-B)^T = A^T - B^T$

$$E = (Y^T - \gamma^T)(Y - \gamma) \rightarrow$$

And from eq ① we know that $\gamma = XB$
 putting the val. of eq. ① in eq. ②

$$E = (Y^T - X^T(B^T))(Y - XB)$$

$$E = Y^T Y - Y^T X B - Y(XB)^T + (XB)^T * XB$$

Now we have to prove these both as eq.

A Detour for this

$$Y^T X B \quad (XB)^T Y \text{ is same}$$

$$\text{let's } Y = A, XB = B$$

$$A^T B = B^T A \rightarrow ①$$

$$② \text{ We know that } (AB)^T = B^T A^T \text{ same logic}$$

$$(A^T B)^T = B^T A^T \rightarrow ② \quad (A^T)^T = A$$

Mean, we have to prove $A^T B = (A^T B)^T$

lets say $A^T B = C$

And we want to prove

$$C = C^T$$

$$C = A^T \beta \quad (A = Y, \beta = X\beta)$$

$$C = Y^T \times \beta$$

We have to prove

$$Y^T \times \beta = (Y^T \times \beta)^T \quad \text{Proving this}$$

→ Let's say we have n studies then $Y = \begin{bmatrix} \vdots \\ \vdots \\ h \end{bmatrix}$ shape = $n \times 1$

Transpose of $Y \rightarrow$ shape $\Rightarrow (1 \times n)$

→ For $X \rightarrow$ for n stud → $\begin{bmatrix} 1 & x_1 & x_2 & \dots & x_n \end{bmatrix}$
 ↓ each will
 Shape $\Rightarrow n \times (m+1)$

→ For β for n st $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$
 ↓ n rows but
 with n inputs we have
 $m+1$ col

Shape $\Rightarrow (m+1) \times 1$

$\rightarrow (m+1) \times 1$

So $\rightarrow Y^T \times \beta$

$(1 \times n) \underset{x}{\underbrace{(n \times (m+1))}} ((m+1) \times 1)$

↓ after multiplying it because n is same

$(1 \times m+1) \underset{1}{\underbrace{((m+1) \times 1)}}$

$1 \times 1 = [1] \rightarrow$ This is a matrix

If we find its Transpose it will be same

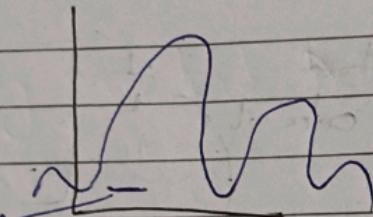
$[1]^T = [1]^T \cancel{\text{Same}}$

Mean, we have proved
 $\gamma^T X \beta = (\gamma^T X \beta)^T$ then $\gamma^T X \beta = (\beta^T \gamma)$
 Proved

Now from left case
 $E = \gamma^T \gamma - \gamma^T X \beta - (\beta^T \gamma) + (\beta^T X \gamma)$
 Now $E = \gamma^T \gamma - 2\gamma^T X \beta + \beta^T X^T X \beta$

The same concept we use in Simple linear Regress we will use
 If here also

The Concept



the m & b graph varying
 like that

and we have to find the minimum m & b value by
 partial derivative, same logic here and equating with 0

$$\frac{dE}{d\beta} = \frac{d}{d\beta} \left[\underbrace{\gamma^T \gamma}_{0} - \underbrace{2\gamma^T X \beta}_{2\gamma^T X} + \underbrace{\beta^T X^T X \beta}_{\beta^T \gamma} \right] = 0$$

$$= 0 - 2\gamma^T X + \frac{d}{d\beta} \left[\underbrace{\beta^T X^T X \beta}_{2\beta^T X^T X \beta} \right] = 0 \quad \text{here } \beta \neq \beta^T \\ \text{both occurs}$$

So we have to study Matrix differentiation but
 I will learn it later during revision, currently direct val

$Y = A^T X \beta$
$\frac{dY}{d\beta} = 2X^T A^T$

that's why it is called

$$= 0 - 2\gamma^T X + [2\beta^T X^T X]$$

same logic for e.g. $\frac{dG}{d\beta}$

$$\Rightarrow -2Y^T X + 2X X^T \beta^T = 0$$

$$\Leftrightarrow 2X X^T \beta^T - 2Y^T X = 0$$

$$X X^T \beta^T - Y^T X = 0$$

$$X X^T \beta^T = Y^T X$$

$$\beta^T = \frac{Y^T X}{X X^T}$$

Note

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

Transpose both side

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

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

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

↳ This is a square matrix, and after transposing we will also get the same matrix.

$$\beta = (X X^T)^{-1} (X^T Y) \rightarrow \text{find answer}$$

which we are finding

$$\text{here } X \rightarrow \begin{matrix} X_{\text{train}} \\ \vdots \\ \text{for one feature} \end{matrix} \Rightarrow \begin{matrix} \text{one} \\ \vdots \\ X_{\text{train}} \end{matrix}$$

$$Y \rightarrow Y_{\text{train}} \text{ (actual value)}$$

$\beta = n+1$, where n is no. of features, always β will be one more than β_0 .

Verifying

$$\hat{y}_\beta = \underbrace{(X^T X)^{-1} X^T Y}_{\text{J ustified}} \\ \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = \underbrace{(m+1) \times (m+1)}_{\text{J ustified}} \underbrace{[(m+1) \times n]}_{\text{J ustified}} [n \times 1] \\ (m+1) \times 1 = [(m+1) \times n] [n \times 1]$$

Shape

$$[(m+1) \times 1] = [(m+1) \times 1] \quad \text{Proved the shape is Same}$$

Q Why Gradient Descent?

Sklearn uses

- OLS for finding β_0 & β
- Gradient descent

OLS gives β by using the $\beta = (X^T X)^{-1} X^T Y$, but question is that why sklearn or most ML engine use Gradient descent
the answer is $(X^T X)^{-1}$ this inverse function

because OLS compute β with each, but the matrix inversion complexity is too high $O(n^3)$
So, for reducing it we will use Gradient descent

OLS for OLS → for Gradient Descent

Linear Regression SGD Regression

scratch code in Laptop