

# Multiple Linear Regression

It is a statistical method used to model the relationship between one continuous dependent variable and two or more independent variables, which may be continuous or categorical.

The goal is to predict the dependent variable based on the independent variables by fitting a linear equation to the data.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

$y$ : Dependent Variable  
 $\beta_0$ : Intercept (Offset)  
 $\beta_1, \beta_2, \dots, \beta_p$ : Coefficients of independent variable  
 $x_1, x_2, \dots, x_p$ : Independent Variables  
 $\epsilon$ : Error term

## Mathematical Foundation

Let's take sample data.

$x_1$ gpa	$x_2$ iq	$y$ salary
8 $x_{11}$	80 $x_{12}$	75
7 $x_{21}$	70 $x_{22}$	60
5 $x_{31}$	120 $x_{32}$	110

$$\begin{aligned}
 \hat{y}_1 &= \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} \dots + \beta_m x_{1m} \\
 \hat{y}_2 &= \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m} \\
 \hat{y}_3 &= \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} + \dots + \beta_m x_{3m} \\
 &\vdots \\
 \hat{y}_n &= \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm}
 \end{aligned}$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \hat{y}_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} \dots + \beta_m x_{1m} \\ \hat{y}_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m} \\ \hat{y}_3 = \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} + \dots + \beta_m x_{3m} \\ \vdots \\ \hat{y}_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_m x_{nm} \end{bmatrix}$$

$\hat{\mathbf{y}}$  is  $n \times 1$   
 $\mathbf{X}$  is  $n \times 1$

We can rewrite the matrix in different form by decomposing it.

$$\hat{y} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix} \rightarrow (m+1) \times 1$$

$\underbrace{\hspace{15em}}_X \quad \underbrace{\hspace{2em}}_{\beta}$   
 $\hspace{15em} \downarrow$   
 $\hspace{15em} n \times (m+1)$

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

$\therefore \boxed{\hat{y} = X\beta}$  - eq ① (This equation is true for any number of dimensions)

$$\underbrace{n \times (m+1) \leftrightarrow (m+1) \times 1}_{n \times 1}$$

$\hat{y}$  (shape of  $\hat{y}$ )

Error function for MLR

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} \quad \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$e = y - \hat{y} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

$$e^T e = \begin{bmatrix} y_1 - \hat{y}_1 & y_2 - \hat{y}_2 & \dots & y_n - \hat{y}_n \end{bmatrix} \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

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

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

$\beta_0$  = offset (reason)  
So, if we consider this data  
experience | grades | salary  
 $\beta_1$  |  $\beta_2$  |

The salary would be  
Salary = exp  $\times$   $\beta_1$  + grades  $\times$   $\beta_2$  +  $\beta_0$   
If the value of this are zero, then  
still there would be some base  
salary. So, that's what  $\beta_0$  gives

$\therefore E = e^T e$  (Loss Function) — eq (2)

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

(first we will prove that these both terms are the same)

(It's a tedious process, but it simply concludes that output is scalar)

$$E = y^T y - 2 y^T \hat{y} + \hat{y}^T \hat{y} \quad \text{— eq (3)}$$

Now, replace value of  $\hat{y}$  in eq (3) with eq (1)

$$E = y^T y - 2 y^T X \beta + (X \beta)^T (X \beta)$$

$$E = y^T y - 2 y^T X \beta + \beta^T X^T X \beta \quad \text{— eq (4)} \quad ((AB)^T = B^T A^T)$$

This is an important equation because it has  $\beta$  variable.

We have to find value of  $\beta$  Matrix for which  $E$  is minimum, So we will now have to differentiate eq (4) w.r.t  $\beta$

Differentiating the equation

$$\frac{dE}{d\beta} = \frac{d}{d\beta} (y^T y - 2 y^T X \beta + \beta^T X^T X \beta)$$

$$= 0 - 2 y^T X + 2 \beta^T X^T X$$

(Now we compare it with zero)

$$-2 y^T X + 2 \beta^T X^T X = 0$$

$$2 \beta^T X^T X = 2 y^T X$$

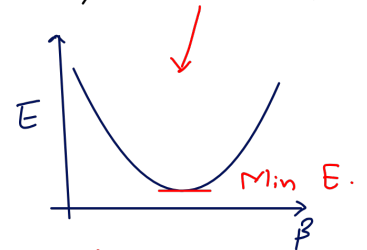
$$\beta^T X^T X = y^T X$$

(Now to find only for  $\beta^T$ , we multiply  $(X^T X)^{-1}$  on both sides to eliminate  $X^T X$ )

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

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

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



Because of  $\beta^T$ , the differentiation of  $\beta^T X^T X \beta$  is difficult. And it can be only solved by matrix differentiation)

(But the output for this is  $2 \beta^T X^T X$ )

(Note that the answer is only true when the matrix is symmetric, and here  $X^T X$  is symmetric).

Now Transposing both sides

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

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

$$\beta = \underbrace{[(x^T x)^{-1}]^T}_{\text{To prove: } (x^T x)^{-1} \text{ is symmetric}} (x^T y)$$

$$\boxed{\beta = (x^T x)^{-1} x^T y} \quad - \text{eq (5)}$$

$\therefore \beta$  values has shape of  $(m+1 \times 1)$

This entire method that we used, it's called as OLS (Ordinary Least square)

One of the reasons, Gradient descent is better than OLS is because in OLS, we are inverting a matrix, and in inverting a matrix, it normally takes time complexity of  $\boxed{O(n^3)}$  which requires a lot of Computation.

$$[(x^T x)^{-1}]^T = (x^T x)^{-1} \quad \leftarrow \begin{matrix} \text{To prove:} \\ (x^T x)^{-1} \text{ is symmetric} \end{matrix}$$

Assume that  $x^T x = A$

$$A A^{-1} = I$$

$$(A A^{-1})^T = I^T$$

$$(A^{-1})^T A^T = I$$

$$(A^{-1})^T A = I$$

$$(A^{-1})^T A A^{-1} = I A^{-1}$$

$$(A^{-1})^T I = A^{-1}$$

$$(A^{-1})^T = A^{-1}$$

$\therefore$  Proved that  $(A^{-1})^T = A^{-1}$

$$\left( (x^T x)^{-1} \right)^T = (x^T x)^{-1}$$