

Multiple Linear Regression

Multiple Linear Regression is a Supervised machine Learning algorithm used to predict a continuous target variable using two or more independent variables.

- It models a linear relationship b/w features & target.
- It finds a straight-line relationship between dependent variable (Target) → what we want to predict.
- Independent Variable (Features) → factors impacting prediction.

Ex:- Suppose we want to predict house price using:

<u>Feature</u>	<u>meaning</u>
i) Size	Bigger house → higher price
ii) Bedrooms	more rooms → higher price
iii) Distance to city	far → lower price.

So the model learns:-

$$\text{House Price} = \text{Function}(\text{Size}, \text{Bedrooms}, \text{Distance})$$

Mathematical Equation :-

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

y → Predicted values

β_0 → intercept (base prediction) Hyperplane in

x_1, x_2, \dots → Features (base pr) n-dim co-ordinate

β_1, β_2, \dots → Co-efficients (importance of Each feature)

coefficients tell how much the target changes when one feature increases by 1 unit.

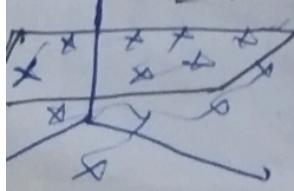
→ The model finds the best coefficients, so the line fits the data points.

→ It reduces the error between Actual value vs Predicted

cgpa	iq	package
-	-	-

to find the values of $\beta_0, \beta_1, \beta_2$

(coefficients)



$$\begin{aligned}
 E &= e^T e = (y - \hat{y})^T (y - \hat{y}) \\
 &= (y^T - \hat{y}^T) (y - \hat{y}) \quad \text{symmetric means same} \\
 &= y^T y - [\hat{y}^T y - \hat{y}^T \hat{y}] + \hat{y}^T \hat{y} \\
 E &= y^T y - 2 \hat{y}^T y + \hat{y}^T \hat{y} \rightarrow \text{Eq(3)}
 \end{aligned}$$

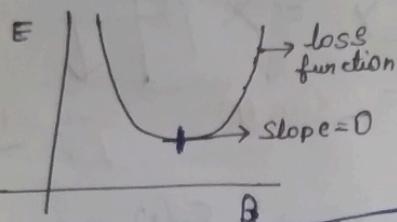
now replace $\hat{y} \rightarrow x\beta$.

$$\begin{aligned}
 E &= y^T y - 2 y^T x\beta + (x\beta)^T (x\beta) \\
 E &= y^T y - 2 y^T x\beta + \beta^T x^T x\beta \rightarrow \text{Eq(4)}
 \end{aligned}$$

here E is a function of β

$y \rightarrow$ data output
 $x \rightarrow$ data input } data will not change, so here
 E is will change, when β change

$E(\beta) \Rightarrow$ find such value of β matrix
for which E is min



$$\frac{dE}{d\beta} = 0$$

$$E = 0 - 2 y^T x + 2 \beta^T x^T x = 0$$

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

$$\beta^T x^T x = y^T x$$

multiply by $(x^T x)^{-1}$ on b.s

$$\beta^T x^T x \underbrace{(x^T x)^{-1}}_{I} = y^T x (x^T x)^{-1}$$

$$\underbrace{\beta^T}_{\beta^T} I = y^T x (x^T x)^{-1}$$

apply (T) transpose on b.s

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

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

$$\beta = [(x^T x)^{-1} x^T y] \rightarrow \text{Eq(5)}$$

$$\begin{aligned}
 \frac{d}{d\beta} (\beta^T x^T x \beta) &= 2 \beta^T x^T x \\
 \therefore x^T x \beta &= 2 x^T \beta \\
 \text{only if } A \text{ symmetric:} \\
 A &= (x^T x)^T = x^T x \\
 \underbrace{A^T}_{\text{Symmetric}} &= A
 \end{aligned}$$

$$\begin{aligned}
 (x^T x)^{-1} &\xrightarrow{T} \\
 [A^{-1}]^T &= A^{-1} \\
 (x^T x)^{-1} &\xrightarrow{\text{Symmetric}}
 \end{aligned}$$

OLS \Rightarrow ordinary least square methodology

Mathematical Formulation:-

Ex:-	x_1^1		x_2^1		Placement
	CGPA	x_{11}	x_{12}	8	
	8	x_{21}	x_{22}	7	
	7	x_{31}	x_{32}	15	
	5				

now for m cols $\rightarrow n$ students/group

$$\begin{aligned} y_1 &= 8 \quad ; \quad y_2 = 7 \quad ; \quad y_3 = 15 \\ \hat{y}_1 &= ? \quad , \quad \hat{y}_2 = ? \quad , \quad \hat{y}_3 = ? \\ \hat{y}_1 &= \beta_0 + \beta_1 8 + \beta_2 80 \\ \hat{y}_2 &= \beta_0 + \beta_1 7 + \beta_2 70 \\ \hat{y}_3 &= \beta_0 + \beta_1 5 + \beta_2 120 \end{aligned}$$

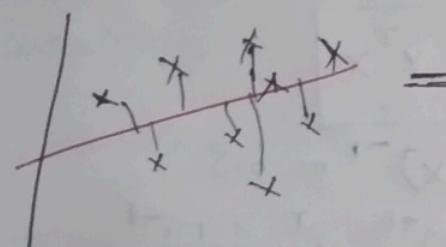
$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \beta_3 x_{13} + \dots + \beta_m x_{1m} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_m x_{2m} \\ \beta_0 + \beta_1 x_{31} + \beta_2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

$[n \times 1]$

$$\hat{Y} = \begin{bmatrix} x_{11} + x_{12} + \dots + x_{1m} \\ x_{21} + x_{22} + \dots + x_{2m} \\ x_{n1} + x_{n2} + \dots + x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}^{m+1} \Rightarrow \hat{Y} = X \beta$$

$(n \times 1) = n \times (m+1)$ $(m+1)$

Loss function for SLR



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

for mLR

need to write matrix form.

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{[n \times 1]} \quad \hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}_{[n \times 1]}$$

$$e = y - \hat{y} \Rightarrow \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}_{n \times 1}$$

$$e^T \cdot e = [y_1 - \hat{y}_1 \quad y_2 - \hat{y}_2 \dots y_n - \hat{y}_n]_{n \times 1} \cdot \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}_{n \times 1} \Rightarrow \frac{\partial E}{\partial \beta}$$

$$e^T \cdot e = [(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2] = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$E = e^T \cdot e$$

\Rightarrow error function
MLR Eq(2)

Problems

linear regression \rightarrow OLS method
 \rightarrow gradient descent

for 10 i/p:-

$$\text{OLS} \rightarrow [(X^T X)^{-1} X^T Y]$$

huge inverse function will consume more time
 Time complexity = $O(n^3)$

$$\frac{(X^T X)^{-1}}{(m+1) \times n} \xrightarrow{\text{shape}} \underbrace{(m+1) \times (m+1)}_{\substack{n \\ \times \\ n}} \quad \leftarrow \begin{array}{l} \text{for 10 i/p} \\ \text{Time complexity} = O(n^3) \end{array}$$

$$\frac{(m+1)(m+1)}{(100)(100)} \Rightarrow (1000000)$$

for high dimensional data like 100 i/p OLS
 $(100 \times 100)^3 \Rightarrow (10000)^3$ → will give result
 → but too slow

i) Sensitive to outliers :-

- OLS Squares errors → outliers get huge weight
- Even one extreme point can change the regression line completely.

ii) Multicollinearity :-

- when independent variables are highly correlated
- Coefficients become unstable and interpretation becomes misleading
- Hard to identify which feature actually impacts the target.

iii) Heteroscedasticity (Non-constant variance of errors)

- residual variance changes with predictors
- standard errors become biased → hypothesis tests become unreliable.

iv) Linearity Assumption :- works only if the relationship b/w variables is linear. If the relation is non-linear → OLS gives poor predictions.