① { multiple linear regression }

MLR is used when there are multiple input features and a single output feature. (more than 1 predictors)

Suppose we have $n$ input features. Let them be $X_1, \ldots, X_n$ and let the output feature be $y$. Then MLR is:

$$\hat{y} = \beta_0 + \beta_1 X_1 + \ldots + \beta_n X_n$$

The objective is to find $(\beta_0, \ldots, \beta_n) \ni$ we have minimum loss

② { Mathematical formulation }

Let us suppose we have $m$ input features $(X_1, \ldots, X_m)$ and $n$ such data points. Let $y$ be the output feature. So our dataset looks like this.

m predictors

|  | $X_1$ | $X_2$ | | $X_m$ | $y$ |
|---|---|---|---|---|---|
| | $X_{11}$ | $X_{12}$ | . . . . | $X_{1m}$ | $y_1$ |
| | $X_{21}$ | $X_{22}$ | . . . . | $X_{2m}$ | $y_2$ |
| n data points | ⋮ | | | | ⋮ |
| | $X_{n2}$ | | | $X_{nm}$ | $y_n$ |

output feature

Then

$$\hat{y}_1 = \beta_0 + \beta_1 X_{11} + \ldots + \beta_m X_{1m} \qquad - ①$$

⋮ ⋮

$$\hat{y}_n = \beta_0 + \beta_1 x_{n1} + \cdots + \beta_m x_{nm} \qquad \text{---}\textcircled{n}$$

These $n$ equations can then be written as:

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \cdots + \beta_m x_{1m} \\ \vdots \\ \beta_0 + \beta_1 x_{n1} + \cdots + \beta_m x_{nm} \end{bmatrix} \Rightarrow (n \times 1)_{shape}$$

$(n \times 1)$ shape

$$= \begin{bmatrix} 1 & x_{11} & \cdots & x_{1m} \\ & \vdots & & \\ 1 & x_{n1} & \cdots & x_{nm} \end{bmatrix}_{n \times (m+1)} \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_m \end{bmatrix}_{(m+1) \times 1}$$

Thus, $\boxed{\hat{y} = X\beta} \Rightarrow$ super important equation.

③ { Mathematical formulation of error function }

Let $\quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad$ and $\quad \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$

We define $\quad E = d_1^2 + d_2^2 + \cdots + d_n^2$

$$= (y_1 - \hat{y}_1)^2 + \cdots + (y_n - \hat{y}_n)^2$$

$$= \begin{bmatrix} y_1 - \hat{y}_1 & \cdots & y_n - \hat{y}_n \end{bmatrix}_{1 \times n} \begin{bmatrix} y_1 - \hat{y}_1 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}_{n \times 1}$$

$$= e^T e \; , \quad e = \begin{bmatrix} y_1 - \hat{y}_1 \\ \\ y_n - \hat{y}_n \end{bmatrix}_{n \times 1}$$

$$= \|e\|^2 \qquad \longrightarrow \text{(norm of } e\text{)}$$

Thus, error function $E$ can also be represented as a product of two matrices, namely $e^T$ and $e$.

We may also write $e$ as $y - \hat{y}$

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

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

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

These two are equal.

Claim :- $y^T \hat{y} = (\hat{y})^T y$ i,e $y^T \hat{y}$ is a symmetric matrix

Proof :- $y^T$ has shape $1 \times n$ and $\hat{y}$ has shape $n \times 1$.

Thus, $y^T y$ has shape $1 \times 1$. Hence, we have shown that it is a scalar matrix.

We know that every scalar matrix is symmetric. Thus

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

$$= (\hat{y})^T y \qquad \boxed{QED}$$

Thus, $\quad E = y^T y - 2 y^T \hat{y} + (\hat{y})^T \hat{y}$

④ { final derivation }

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

putting this in the equation of $E$, we get

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

$$\boxed{E = y^T y - 2 y^T X \beta + \beta^T X^T X \beta} \qquad \Rightarrow \; \text{main equation}$$

Thus, $E$ is a function of $\beta$. as $X$ and $y$ are already fixed.

$$E(\beta) = y^Ty - 2y^TX\beta + \beta^TX^TX\beta$$

We need to find $\beta \ni E(\beta)$ has minimum value.

Here
$$y = (y_1 \ldots y_n)^T$$

$$\beta = (\beta_0 \ldots \beta_m)^T$$

$$X = \begin{bmatrix} 1 & X_{11} & \cdots & X_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & \cdots & X_{nm} \end{bmatrix}_{n \times (m+1)}$$

find such value of $\beta$ matrix for which $E$ is min. for minimum, we should have.

$$\frac{\partial E}{\partial \beta} = 0$$

Now

$$\frac{\partial E}{\partial \beta} = \frac{\partial}{\partial \beta} \left( y^T y - 2y^T x \beta + \beta^T x^T x \beta \right)$$

$$= 0 - 2y^T x + \frac{\partial}{\partial \beta} \left( \beta^T x^T x \beta \right)$$

For this we need to study matrix calculus which is not feasible as of now. However, we have a special result in matrix calculus which says that if A is a symme-- matrix and $\alpha = x^T A x$, then

$$\frac{\partial \alpha}{\partial x} = 2x^T A \ , \text{ where } x \text{ is}$$

a $(n \times 1)$ vector.

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

Thus, $\frac{\partial E}{\partial \beta} = 2\beta^T x^T x - 2y^T x$

Equating this to 0, we get

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

Assuming that $x^T x$ is invertible, we have

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

$$\Rightarrow \quad \beta^T = y^T x x^{-1} (x^T)^{-1}$$

$$\Rightarrow \quad (\beta^T)^T = \{ y^T x x^{-1} (x^T)^{-1} \}^T$$

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

$$\Rightarrow \quad \beta = x^{-1} (x^{-1})^T x^T y$$

$$\Rightarrow \quad \beta = x^{-1} (x^T)^{-1} x^T y$$

$$= [x^T x]^{-1} x^T y$$

Thus, $\boxed{\beta = (x^T x)^{-1} x^T y}$

$x \longrightarrow n \times (m+1)$

$x^T \longrightarrow (m+1) \times n$

$[x^T x]^{-1} \longrightarrow (m+1) \times (m+1)$

Now $[x^T x]^{-1} x^T$ has shape $(m+1) \times n$. Also, since $y$ has shape $n \times 1$. Thus, $\beta = [x^T x]^{-1} x^T y$ has shape $(m+1) \times 1$. ∎

⑤ { Problem with OLS solution }

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

Matrix multiplication is a computationaly expansive operation. In higher dimension, this is a slow procedure. Thus in higher dimension, we prefer to use non-closed approximation techniques like gradient descent which is computationaly less expansive.