For, 3D dataset,

$$\hat{y}_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12}$$

$$\hat{y}_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22}$$

datasets:-

| | $X_1$ | $X_2$ | $y$ |
|---|---|---|---|
| $P_1$ | $X_{11}$ | $X_{12}$ | $\hat{y}_1$ |
| $P_2$ | $X_{21}$ | $X_{22}$ | $\hat{y}_2$ |

For this,

$$b = \text{intercept} = \beta_0 = \beta_0 - \underset{\underset{\text{learning rate}}{\uparrow}}{\eta} (\beta_0\_\text{slope})$$

$$m_1 = \underset{\text{weight}}{\text{slope}} = \beta_1 = \beta_1 - \eta (\beta_1 - \text{slope})$$

$$m_2 = \text{weight} = \beta_2 = \beta_2 - \eta (\beta_2\_\text{slope})$$

Our work is to find $\beta_0\_\text{slope}$, $\beta_1\_\text{slope}$, $\beta_2\_\text{slope}$ from loss function.

$$L = \sum_{i=1}^{2} (y_i - \hat{y}_i)^2 \times \frac{1}{2} \quad (\because 2, \text{ no of person} = 2)$$

$$L = \left( (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 \right) \frac{1}{2}$$

$\Rightarrow$ we want to find $\dfrac{dL}{d\beta_0}$ to find ~~to~~ intercept,

$$L = \left( \left( y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12} \right)^2 + \left( y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22} \right)^2 \right) \cdot \frac{1}{2} \quad - \text{①}$$

$$\frac{dL}{d\beta_0} = \left( 2 \left( y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12} \right)(-1) + 2 \left( y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22} \right)(-1) \right) \cdot \frac{1}{2}$$

$$= -\frac{2}{2} \left( \left( y_1 - \hat{y}_1 \right) + \left( y_2 - \hat{y}_2 \right) \right) \qquad \times \text{ wrong.}$$

We cannot cut $2 \& 2$ as $2$ in denominator is the total no of person.

$$= -\frac{2}{n} \left( \left( y_1 - \hat{y}_1 \right) + \left( y_2 - \hat{y}_2 \right) \right)$$

In general for $n$-person data,

$$= -\frac{2}{n} \left( \left( y_1 - \hat{y}_1 \right) + \left( y_2 - \hat{y}_2 \right) + \cdots + \left( y_5 - \hat{y}_5 \right) + \cdots + \left( y_n - \hat{y}_n \right) \right)$$

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

Now for finding $\beta_1$ - slope, $\dfrac{\partial L}{\partial \beta_1}$

$$L = \left( (y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 \right) \tfrac{1}{2} \quad (\because \text{from } 1)$$

$$\frac{\partial L}{\partial \beta_1} = \left( 2(y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})(-x_{11}) + 2(y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})(-x_{21}) \right) \tfrac{1}{2}$$

$$= \frac{\cancel{2}}{2} \left( (y_1 - \hat{y}_1)\cancel{x_{11}} + (y_2 - \hat{y}_2)\cancel{x_{21}} \right)$$

$$= -\frac{2}{2} \left( (y_1 - \hat{y}_1) x_{11} + (y_2 - \hat{y}_2) x_{21} \right)$$

### In general,

$\beta_1$-slope

$$\boxed{ \frac{-2}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i) x_{i1} }$$

to find $\beta_2$-slope, $\dfrac{\partial L}{\partial \beta_2}$

$$L = \left( (y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 \right) \tfrac{1}{2}$$

$$\frac{\partial L}{\partial \beta_2} = \left( 2(y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})(-x_{12}) + 2(y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})(-x_{22}) \right) \tfrac{1}{2}$$

$$\frac{\partial L}{\partial \beta_2} = \frac{-2}{2}\left(\left(y_1 - \hat{y}_1\right)X_{12} + \left(y_2 - \hat{y}_2\right)X_{22}\right)$$

In general,

$$\frac{\partial L}{\partial \beta_2} = \frac{-2}{n}\left(\sum_{i=1}^{n}\left(y_i - \hat{y}_i\right)X_{i2}\right)$$

for $m^{th}$ column

$$\frac{\partial L}{\partial \beta_m} = \frac{-2}{n}\left(\sum_{i=1}^{n}\left(y_i - \hat{y}_i\right)X_{im}\right)$$

For $n^{th}$ column

$$\frac{\partial L}{\partial \beta_n} = \frac{-2}{n}\left(\sum_{i=1}^{n}\left(y_i - \hat{y}_i\right)X_{in}\right)$$

example, For $5^{th}$ column

$$\frac{\partial L}{\partial \beta_5} = \frac{-2}{5}\left(\sum_{i=1}^{n}\left(y_i - \hat{y}_i\right)X_{i5}\right)$$

Make it General,

According to Multiple linear Regression, we know that,

In $\sum_{i=1}^{n}(y_i - \hat{y}_i)$ If we take $y$ is the matrix then it would be,

$$\left(Y - \hat{Y}\right)_{n \times 1}$$

$n =$ No of person

$X$ would be the metrix then,

$$\left(X\right)_{n \times m}$$

$n =$ No of person / Row

$m =$ No of feature / column

~~one~~ = ~~for appending as column~~

as $\left(y - \hat{y}\right)$ is symmetric matrix then,

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

$\therefore$ In general it can be written as

$$\text{coff-slope} = \left( (Y-\hat{Y})^T_{1 \times n} \cdot (X)_{n \times m} \right) \left(\frac{-2}{n}\right)$$

In numpy it can be written as,

$$\text{coff-slope} = \frac{-2}{n}\left(np.dot\left((Y-\hat{y}), X\right)\right) \quad -$$

$\circ$ $\left((y-\hat{y})^T\right.$ taken case by numpy itself.

Disadvantages:- ① Large memory required, If our dataset is too large then It can through memory Exhausted Exception.

② More ~~accurate th~~ time required when Big dataset comes into the picture.

Advantages:- ① More accurate & focused compared any other Gradient descent.

:- To overcome This disadvantages, Stochastic Gradient Descent comes into the picture.