

# Linear Regression

We have a function ( $f(x)$ ) for which given ( $x$ ) predict ( $y$ ). That is  $f(x) = y$ . And for regression,  $y$  is continuous.

Simplest type of function used in Linear Regression is Linear Function. That is:

$$f(x) = mx + b$$

And:

$$mx + b = y$$

*Note:  $mx + b$*

*is a linear function which represents the equation of a straight line. But  $f(x)$  can comprise of any function.*

There are two types of Regression:

1. Simple Regression.
2. Multiple Regression.

1. Simple Regression:

Where 'x' comprises of a single feature.

2. Multiple Regression

Where 'x' comprises of multiple feature.

## 1. Straight Line

The geometrical representation of an equation of the form:

$$Ax + By + C = 0$$

Where A, B, C are constants. When B is not zero, i.e. the equation contains:

When B is not zero, i.e. the equation contains a term in y, the equation of the straight line can be solved for y,

$$y = \left(-\frac{A}{B}\right)x + \left(-\frac{C}{B}\right)$$

which is the form of:

$$y = a + bx$$

The above form can also be written as:

$$y = mx + c$$

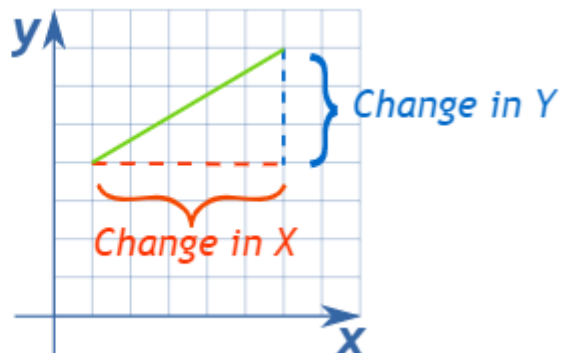
## 2. Slope of Straight Line

In the  $y = a + bx$ , the coefficient of 'x' on the right.

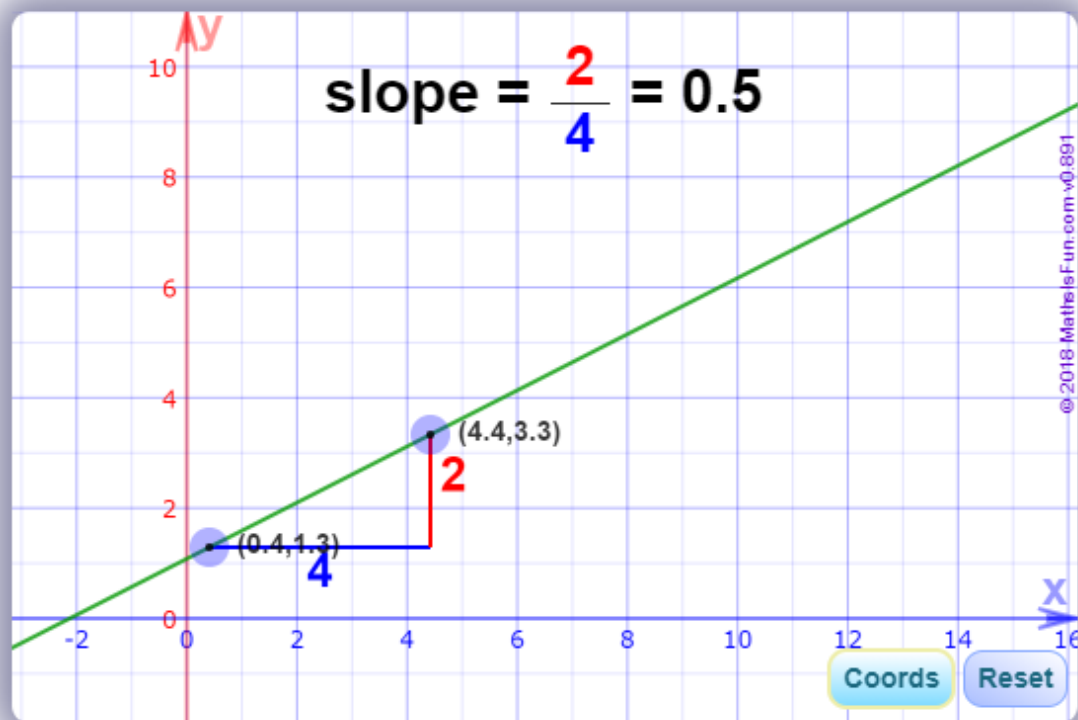
$$b = \text{'slope of the straight line'}$$

The slope may be positive, negative or zero.

$$\text{Slope} = \frac{\text{Change in Y}}{\text{Change in X}}$$

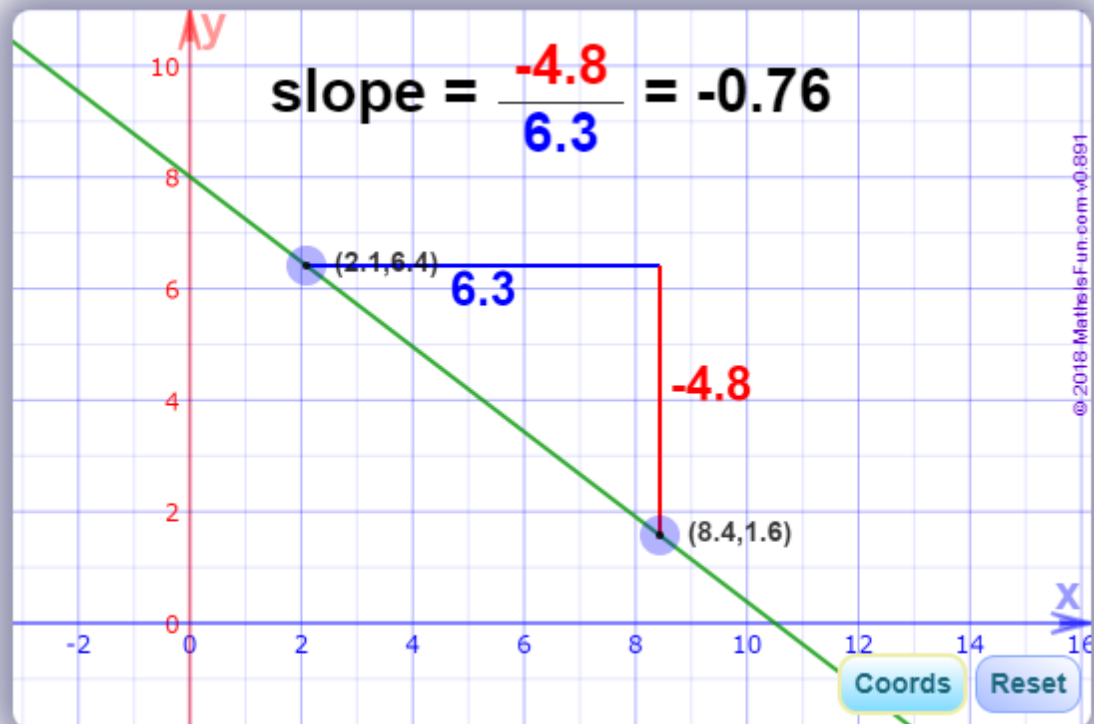


a) Slope is positive:

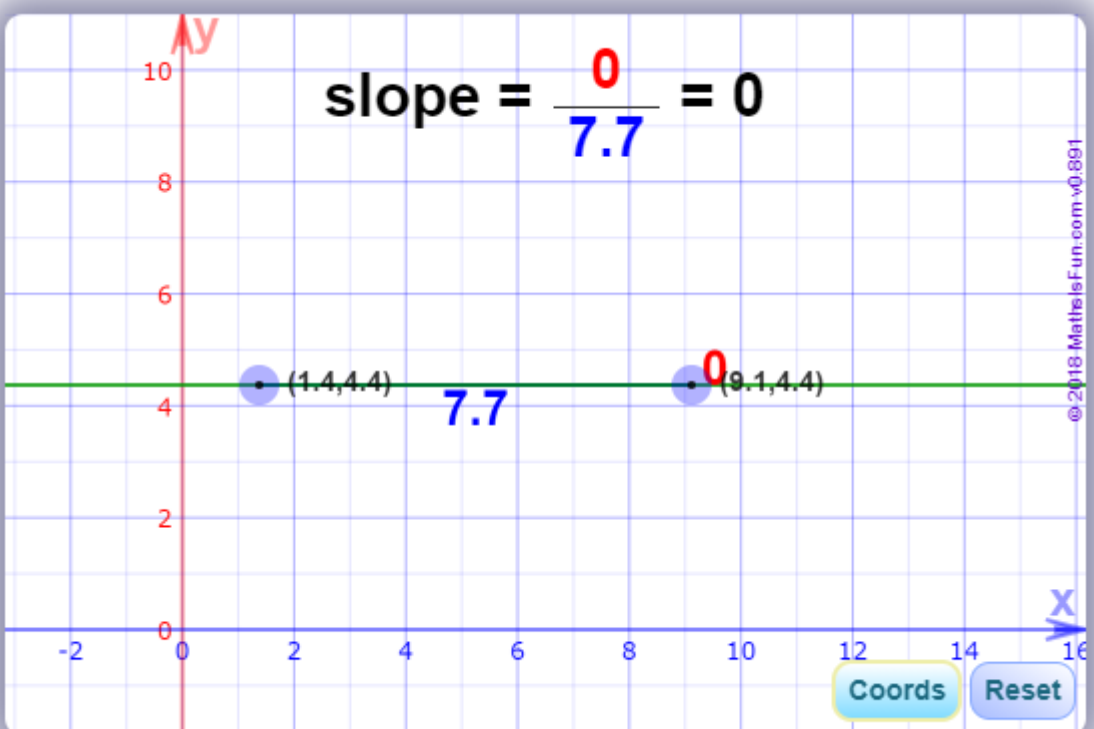


$$\begin{aligned}\text{slope} &= \frac{\text{Change in } Y}{\text{Change in } X} \\ &= \frac{3.3 - 1.3}{4.4 - 0.4} \\ &= \frac{2}{4} = 0.5\end{aligned}$$

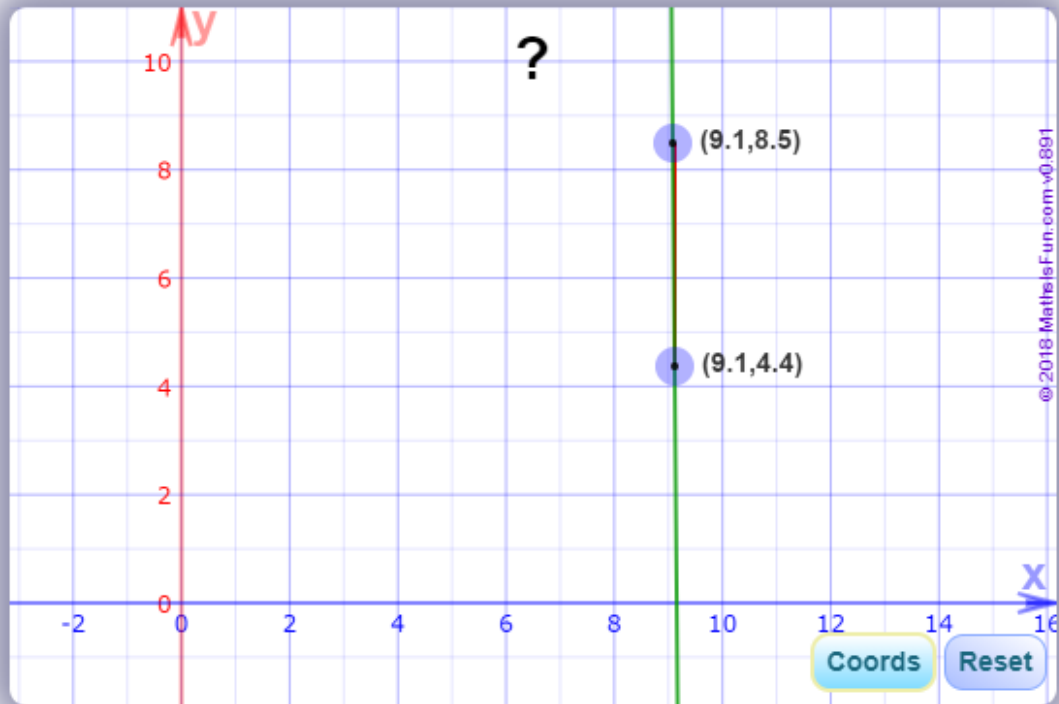
b) Slope is negative:



c) Slope is zero



d) Slope is Infinity:



Let there be an equation of a straight line:

$$5x + 2y - 6 = 0$$

$$\Rightarrow 5x + 2y = 6$$

$$\Rightarrow 2y = 6 - 5x$$

$$\Rightarrow y = 3 - \frac{5x}{2}$$

We know the equation:

$$y = mx + c$$

$$\text{therefore, } y = -\frac{5x}{2} + 3$$

$$\text{therefore, } mx = -\frac{5x}{2} \text{ and } c = 3$$

*Here,  $c$  is called the  $y$  – intercept.*

*$y$  – intercept is the point where the straight line*

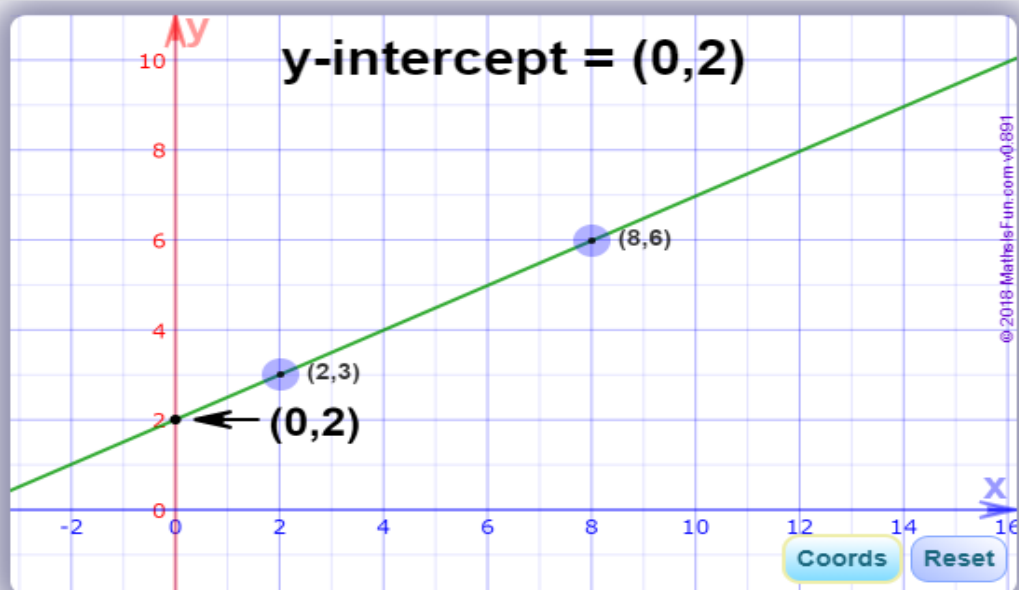
*touches(intercept) the y axis , where  $x = 0$ .*



In the above diagram the line crosses the y axis at  $y = 1$



Here the line crosses the y axis at  $y = -2$



Some useful information about straight lines:

a) If a straight line  $y = mx + c$  or  $y = a + bx$  passes through the point  $(x_0, y_0)$  the coordinate must satisfy the equation i. e.  $y_0 = a + bx_0$  or  $y_0 = mx_0 + c$

b) If a straight line passes through the points

$(x_1, y_1)$  and  $(x_2, y_2)$ , its slope is  $b = \frac{y_1 - y_2}{x_1 - x_2}$ . If a straight

line passes through the point  $(x_0, y_0)$  and has slope  $b$ ,

Its equation is:

$$y - y_0 = b(x - x_0)$$

## Parabola

Parabola is the geometrical representation of an equation of the form:

$$y = a + bx + cx^2$$

Where  $a, b, c$  are constants (the term in  $x^2$  must be present, i.e.  $c \neq 0$ ). The parabola is a special type of curve.

Following equations represent parabolas:

1.  $y = 3 - 2x + 7x^2$

2.  $y = 6 - 4x^2$

3.  $y = 0.5x^2$

4.  $y = a + bx + cx(x - 1)$

$$\Rightarrow a + bx + cx^2 - cx$$

$$\Rightarrow a + bx - cx + cx^2$$

$$\Rightarrow a + (b - c)x + cx^2$$

## Method of Least Square

Method of Least Squares is a device for finding the equation of a specified type of curve, which best fits, a given set of observations.

The method depends upon the *Principle of Least Squares*, which suggests that for the “best-fitting” curve, the sum of the squares of differences between the observed and the corresponding estimated values should be the minimum possible.

Suppose we are given ‘ $n$ ’ pairs of observation  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , and it is required to fit a straight line to these data. The general equation of a straight line  $y = a + bx$  is taken, where  $a$  and  $b$  are constants.

Any values for ‘ $a$ ’ and ‘ $b$ ’ would give a straight line, and once these values are obtained, an estimate of ‘ $y$ ’ can be had by substituting the values of ‘ $x$ ’.

That is to say, the estimated values of ‘ $y$ ’ when  $x = x_1, x_2, \dots, x_n$  would be :

$a + bx_1, a + bx_2, \dots, a + bx_n$  respectively.

In order that the equation  $y = a + bx$  gives a good representation of the relationship between  $x$  and  $y$ , it is desirable that the estimated values,

$a + bx_1, a + bx_2, \dots, a + bx_n$



are on whole, close enough to the corresponding observed values :  $y_1, y_2, \dots, y_n$ .

For the best fitting straight line ,therefore, our problem is only to choose such values of ' $a$ ' and ' $b$ ' for the equation :  $y = a + bx$  which will provide estimates of ' $y$ ' as close as possible to the observed values . This can be done in different ways.

However according to the 'Principle of Least Square', the "best-fitting" equation is interpreted as that which minimises the **sum of square of difference/ sum of square of errors (SSD/SSE)**:

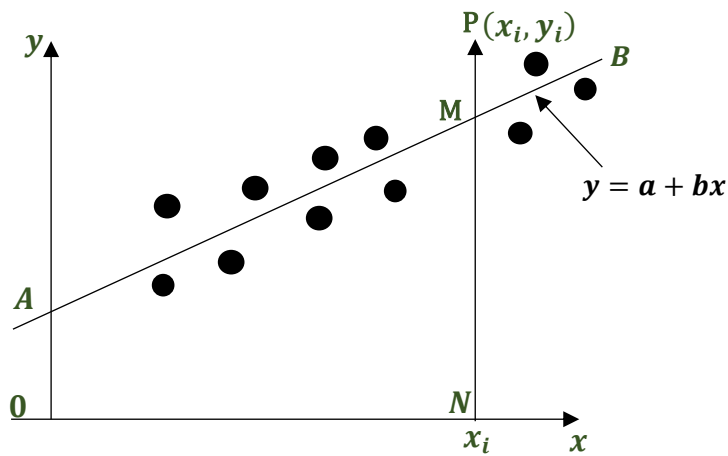
$$\sum_{i=1}^n (y_i - a - bx_i)^2$$

i.e.

$$(y_1 - a - bx_1)^2 + (y_2 - a - bx_2)^2 + \dots + (y_n - a - bx_n)^2$$

### Principle of Least Squares in Straight Line

$x$ -(1)	<i>observed y</i> -(2)	<i>Estimated</i> $y = a + bx$ -(3)	Difference (2)-(3)	(Difference) <sup>2</sup>
$x_1$	$y_1$	$a + bx_1$	$y_1 - a - bx_1$	$(y_1 - a - bx_1)^2$
$x_2$	$y_2$	$a + bx_2$	$y_2 - a - bx_2$	$(y_2 - a - bx_2)^2$
....	....	....	....	....
$x_n$	$y_n$	$a + bx_n$	$y_n - a - bx_n$	$(y_n - a - bx_n)^2$
$\sum_{i=1}^n x_i$	$\sum_{i=1}^n y_i$	$\sum_{i=1}^n a + bx_i$	$\sum_{i=1}^n y_i - a - bx_i$	$\sum_{i=1}^n (y_i - a - bx_i)^2$



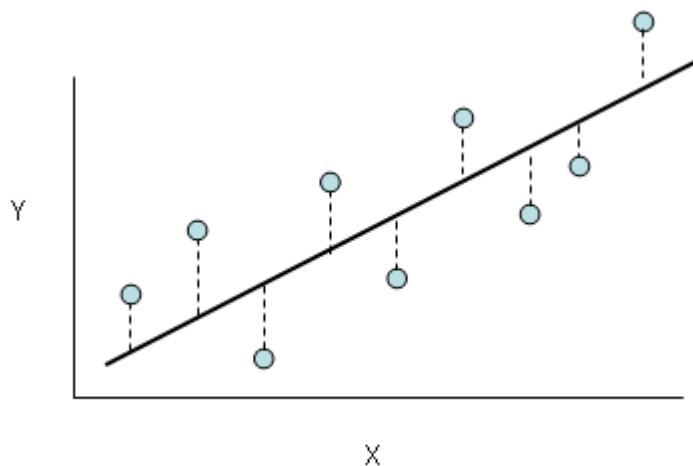
(Geometrical Interpretation)

Fig: Method of Least Squares

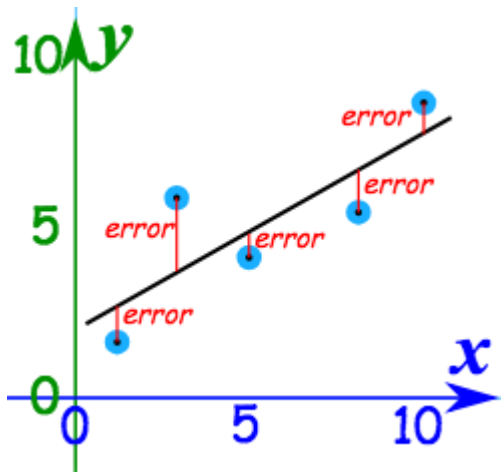
In the geometric sense, the problem of finding the best fitting straight line as follows:

If the pairs of observations:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

are plotted as points on a graph paper and all the possible straight lines are drawn on it, the straight line will be considered the “**best fitting**” for which the “*sum of the squares of vertical distances*” **PM** between the plotted points **P** and the line **AB** is the least.



Hence dotted points are **predicted y** and the straight line is the **actual y**, Therefore the distance between **actual y** and **predicted y** is the **difference / error**.



Such error or difference is known as **Sum of Squares of Difference (SSD)/Sum of Squares of Error (SSE)** i.e.

$$\sum_{i=1}^n (y_i - a - bx_i)^2$$

Obtained from **Least Square Method**.

Now we will do **partial derivative** of the above obtained difference:

$$SSE = \sum_{i=1}^n (y_i - a - bx_i)^2$$

Partial derivative of 'a' :

$$\begin{aligned} \frac{\partial SSE}{\partial a} &= \sum_{i=1}^n (y_i - a - bx_i)^2 \\ \Rightarrow \frac{\partial}{\partial a} \left( \sum_{i=1}^n (y_i - a - bx_i)^2 \right) \end{aligned}$$

By applying power rule  $\left(\frac{\partial}{\partial x}(x^n) = nx^{n-1}\right)$ :

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i)^{2-1} \times \frac{\partial}{\partial a}(y_i - a - bx_i) \right\}$$

We know while doing partial derivation except 'a' everything will be considered as constants:

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i) \times (0 - 1 - 0) \}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i) \times (-1) \}$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i) \times (-1) \cdots eqn(i)$$

Partial derivative of 'b':

$$\begin{aligned} \frac{\partial SSE}{\partial b} &= \sum_{i=1}^n (y_i - a - bx_i)^2 \\ \Rightarrow \frac{\partial}{\partial b} \left( \sum_{i=1}^n (y_i - a - bx_i)^2 \right) \end{aligned}$$

By applying power rule  $\left(\frac{\partial}{\partial x}(x^n) = nx^{n-1}\right)$ :

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i)^{2-1} \times \frac{\partial}{\partial b}(y_i - a - bx_i) \right\}$$

We know while doing partial derivation except 'b' everything will be considered as constants:

$$\Rightarrow \sum_{i=1}^n \{2(y_i - a - bx_i) \times \left(0 - 0 - (b \frac{\partial}{\partial b}(x_i) + (x_i) \frac{\partial}{\partial b}(b))\right)\}$$

$$\Rightarrow \sum_{i=1}^n \{2(y_i - a - bx_i) \times (0 - 0 - (0 + (x_i) \times 1))\}$$

$$\Rightarrow \sum_{i=1}^n \{2(y_i - a - bx_i) \times (-x_i)\}$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i) \times (-x_i) \cdots eqn(ii)$$

Now,  $eqn(i)$  and  $eqn(ii)$  is a linear equation i.e. equals to zero('0')

$$2 \sum_{i=1}^n (y_i - a - bx_i) \times (-1) = 0 \dots eqn(iii)$$

$$2 \sum_{i=1}^n (y_i - a - bx_i) \times (-x_i) = 0 \cdots eqn(iv)$$

$eqn(iv)$  can be written as:

$$2 \sum_{i=1}^n (y_i - a - bx_i) \times ((-1) \times (x_i)) = 0$$

Now we will make all the constants i.e. -1 and 2 to right side and divide with zero in eqn(iii) and eqn(iv):

$$\sum_{i=1}^n (y_i - a - bx_i) = 0 \dots \text{eqn}(v)$$

$$\sum_{i=1}^n (y_i - a - bx_i) \times (x_i) = 0 \dots \text{eqn}(vi)$$

Now we multiply summation( $\Sigma$ ) in eqn(v) and  $x_i$  and summation( $\Sigma$ ):

$$\sum_{i=1}^n (y_i - a - bx_i) = 0 \dots \text{eqn}(v)$$

$$\Rightarrow \sum_{i=1}^n y_i - a \sum_{i=1}^n 1_i - b \sum_{i=1}^n x_i = 0$$

We know,

$$\sum_{i=1}^n \mathbf{1}_i = \mathbf{1}_1 + \mathbf{1}_2 + \cdots + \mathbf{1}_n = n$$

Therefore,

$$\Rightarrow \sum_{i=1}^n y_i - an - b \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \sum_{i=1}^n y_i = an + b \sum_{i=1}^n x_i \dots \text{eqn(vii)}$$

Similarly eqn(vi):

$$\sum_{i=1}^n (y_i - a - bx_i) \times (x_i) = 0 \dots \text{eqn(vi)}$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i - ax_i - bx_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \dots \text{eqn(viii)}$$

Now eqn(vii) and eqn(viii)

$$\sum_{i=1}^n y_i = an + b \sum_{i=1}^n x_i \dots \text{eqn(vii)}$$

$$\sum_{i=1}^n y_i x_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \dots \text{eqn(viii)}$$

is considered the best-fit equations where least squared values can be fitted with minimum no. of errors hence these two equations is known as **best fit of a straight line**.

Hence equations for **best fit of a straight line** are:

$$\sum_{i=1}^n y_i = an + b \sum_{i=1}^n x_i \dots \text{eqn(i)}$$

$$\sum_{i=1}^n y_i x_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \dots \text{eqn(ii)}$$

### **Best fit on Curve/Curve fit though Method of Least Square:**

We know equation of parabola:

$$y = a + bx + cx^2$$



## Principle of Least Squares in Curve

$x$ -(1)	<i>observed y</i> -(2)	<i>Estimated</i> $y = a + bx + cx^2$ -(3)	Difference (2)-(3)	(Difference) <sup>2</sup>
$x_1$	$y_1$	$a + bx_1 + cx_1^2$	$y_1 - a - bx_1 - cx_1^2$	$(y_1 - a - bx_1 - cx_1^2)^2$
$x_2$	$y_2$	$a + bx_2 + cx_2^2$	$y_2 - a - bx_2 - cx_2^2$	$(y_2 - a - bx_2 - cx_2^2)^2$
....	....	....	....	....
$x_n$	$y_n$	$a + bx_n + cx_n^2$	$y_n - a - bx_n - cx_n^2$	$(y_n - a - bx_n - cx_n^2)^2$
$\sum_{i=1}^n x_i$	$\sum_{i=1}^n y_i$	$\sum_{i=1}^n a + bx_i + cx_i^2$	$\sum_{i=1}^n y_i - a - bx_i - cx_i^2$	$\sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$

Therefore the **sum of square of difference/ sum of square of errors (SSD/SSE)** obtained from **Least Square Method**:

$$\sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

$$= (y_1 - a - bx_1 - cx_1^2)^2 + (y_2 - a - bx_2 - cx_2^2)^2 + \dots + (y_n - a - bx_n - cx_n^2)^2$$

Now we will do **partial derivative** of the above-obtained difference:

$$SSE = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

1)

$$\frac{\partial SSE}{\partial a} = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

$$\Rightarrow \frac{\partial SSE}{\partial a} = \frac{\partial}{\partial a} \sum_{i=1}^n \{(y_i - a - bx_i - cx_i^2)^2\}$$

By applying power rule  $\left( \frac{\partial}{\partial x} (x^n) = nx^{n-1} \right)$ :

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i - cx_i^2)^{2-1} \times \frac{\partial}{\partial a} (y_i - a - bx_i - cx_i^2) \right\}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (0 - 1 - 0 - 0) \}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (-1) \}$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) \dots eqn(i)$$

2)

$$\frac{\partial SSE}{\partial b} = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

$$\Rightarrow \frac{\partial SSE}{\partial b} = \frac{\partial}{\partial b} \sum_{i=1}^n \{(y_i - a - bx_i - cx_i^2)^2\}$$

By applying power rule  $\left( \frac{\partial}{\partial x} (x^n) = nx^{n-1} \right)$ :

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i - cx_i^2)^{2-1} \times \frac{\partial}{\partial b} (y_i - a - bx_i - cx_i^2) \right\}$$

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i - cx_i^2) \times \left( 0 - 0 - \left( b \frac{\partial}{\partial b} x_i + x_i \frac{\partial}{\partial b} b \right) - 0 \right) \right\}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (0 - 0 - (0 + x_i \times 1) - 0) \}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (-x_i) \}$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i) \dots \dots \text{eqn(ii)}$$

iii)

$$\frac{\partial SSE}{\partial c} = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

$$\Rightarrow \frac{\partial SSE}{\partial c} = \frac{\partial}{\partial c} \sum_{i=1}^n \{(y_i - a - bx_i - cx_i^2)^2\}$$

By applying power rule  $\left( \frac{\partial}{\partial x} (x^n) = nx^{n-1} \right)$ :

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i - cx_i^2) \frac{\partial}{\partial c} (y_i - a - bx_i - cx_i^2) \right\}$$

$$\Rightarrow \sum_{i=1}^n \left\{ 2(y_i - a - bx_i - cx_i^2) \times \left( 0 - 0 - 0 - \left( c \frac{\partial}{\partial c} x_i^2 + x_i^2 \frac{\partial}{\partial c} c \right) \right) \right\}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (0 - 0 - 0 - (0 + x_i^2 \times 1)) \}$$

$$\Rightarrow \sum_{i=1}^n \{ 2(y_i - a - bx_i - cx_i^2) \times (-x_i^2) \}$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i^2) \dots \text{eqn(iii)}$$

Now we get three equations:

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) \dots \text{eqn}(i)$$

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i) \dots \text{eqn}(ii)$$

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i^2) \dots \text{eqn}(iii)$$

Now , we putting eqns = 0 we get,

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) = 0 \dots \text{eqn}(iv)$$

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i) = 0 \dots \text{eqn}(v)$$

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i^2) = 0 \dots \text{eqn}(vi)$$

from eqn (iv) we get:

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) = 0 \dots eqn(iv)$$

$$\Rightarrow \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i) = \sum_{i=1}^n a + bx_i + cx_i^2 \dots eqn(vii)$$

Similarly from eqn(v),

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) \times (x_i) = 0 \dots eqn(v)$$

$$\Rightarrow \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i - ax_i - bx_i^2 - cx_i^3) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i) = \sum_{i=1}^n ax_i + bx_i^2 + cx_i^3 \dots eqn(viii)$$

Similarly from eqn(vi),

$$2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-x_i^2) = 0 \dots eqn(vi)$$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (-1) \times (x_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i - a - bx_i - cx_i^2) \times (x_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i^2 - ax_i^2 - bx_i^3 - cx_i^4) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i x_i^2) = \sum_{i=1}^n ax_i^2 + bx_i^3 + cx_i^4 \dots eqn(ix)$$

Therefore we get three equations through which we get **best fit in a curve**:

$$\sum_{i=1}^n (y_i) = \sum_{i=1}^n a + bx_i + cx_i^2 \dots eqn(i)$$

$$\sum_{i=1}^n (y_i x_i) = \sum_{i=1}^n ax_i + bx_i^2 + cx_i^3 \dots eqn(ii)$$

$$\sum_{i=1}^n (y_i x_i^2) = \sum_{i=1}^n ax_i^2 + bx_i^3 + cx_i^4 \dots eqn(iii)$$

Through the above three equations we achieve best fit in a curve  
i.e. **less error in experimental values very nearest to the curve**  
also known as **curve fitting or curve fit** obtained from **least  
square method**.

**Linear regression:**



There are two types of Linear Regression:

1. Simple Linear Regression.
2. Multiple Linear Regression.

### Simple Linear Regression:

Simple linear regression is used to estimate the relationship between two quantitative variables.

$$y = a + bx$$

Or,

$$y = mx + c$$

Or,

$$y = \beta_0 + \beta_1 x + \epsilon$$

' $\epsilon$ ' Stands for 'epsilon' is the error while constructing a linear regression model, where:

*$\beta_0$  or  $a$  or  $c$  is the 'y – intercept'*

*$\beta_1$  or  $m$  or  $b$  is the 'slope'*

What we got in best fit while calculation for a straight line from partial derivation:

$$\sum_{i=1}^n y_i = an + b \sum_{i=1}^n x_i \dots eqn(i)$$

$$\sum_{i=1}^n y_i x_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \dots eqn(ii)$$

can be written as :

$$\sum_{i=1}^n y_i = \beta_0 n + \beta_1 \sum_{i=1}^n x_i \dots eqn(i)$$

$$\sum_{i=1}^n y_i x_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \dots eqn(ii)$$

Therefore, multiplying 'n' in equation (ii) and  $\sum_{i=1}^n x_i$  in eqn(i) we get:

$$\sum_{i=1}^n x_i y_i = n \times \beta_0 \sum_{i=1}^n x_i + \beta_1 \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) \dots eqn(iii)$$

$$n \sum_{i=1}^n x_i y_i = n \times \beta_0 \sum_{i=1}^n x_i + n \times \beta_1 \sum_{i=1}^n x_i^2 \dots \text{eqn(iv)}$$

**Subtracting from eqn(iv) from eqn(iii) we get:**

$$\sum_{i=1}^n x_i y_i = n \times \beta_0 \sum_{i=1}^n x_i + \beta_1 \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) \dots \text{eqn(iii)}$$

-

$$n \sum_{i=1}^n x_i y_i = n \times \beta_0 \sum_{i=1}^n x_i + n \times \beta_1 \sum_{i=1}^n x_i^2 \dots \text{eqn(iv)}$$

---


$$\begin{aligned} \sum_{i=1}^n x_i y_i - n \sum_{i=1}^n x_i y_i \\ = 0 + \beta_1 \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) - \left( n \times \beta_1 \sum_{i=1}^n x_i^2 \right) \end{aligned}$$


---

$$\begin{aligned} \sum_{i=1}^n x_i y_i - n \sum_{i=1}^n x_i y_i \\ = \beta_1 \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) - \left( n \times \beta_1 \sum_{i=1}^n x_i^2 \right) \end{aligned}$$

Interchanging sides we get:

$$\begin{aligned} \Rightarrow \left( n \times \beta_1 \sum_{i=1}^n x_i^2 \right) - \beta_1 \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) \\ = n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i y_i \end{aligned}$$

$$\begin{aligned} \Rightarrow \beta_1 \left( n \sum_{i=1}^n x_i^2 \right) - \left( \sum_{i=1}^n x_i \times \sum_{i=1}^n x_i \right) \\ = n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i y_i \end{aligned}$$

$$\Rightarrow \beta_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i y_i}{(n \sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i \times \sum_{i=1}^n x_i)}$$

Hence, we got slope as:

$$\beta_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i y_i}{(n \sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i \times \sum_{i=1}^n x_i)}$$

Similarly,

Again taking the best fit of a straight line:

$$\sum_{i=1}^n y_i = \beta_0 n + \beta_1 \sum_{i=1}^n x_i \dots \text{eqn}(i)$$

$$\Rightarrow \sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i = \beta_0 n$$

$$\Rightarrow \beta_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

*We got y – intercept as: –*

$$\beta_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

Hence for 2-d problem:

$$y = \beta_0 + \beta_1 x$$

We got two learn parameter  $\beta_0$  and  $\beta_1$  where:

$$\beta_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{(n \sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i \times \sum_{i=1}^n x_i)}$$

$$\beta_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

## How Simple Linear Regression works:

Suppose we have 'x' data and 'y' data:

x	1	2	3	4	5
y	3	4	2	4	5

we know :

$$y = \beta_0 + \beta_1 x$$

and

$$\beta_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{(n \sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i \times \sum_{i=1}^n x_i)}$$

$$\beta_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

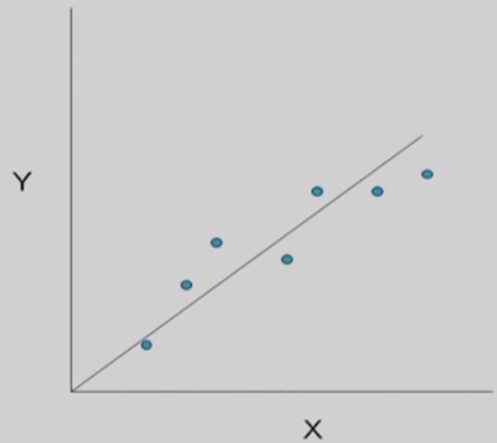
The linear regression function will take each value of 'x', calculate 'y', and plot them in the graph with a **random error** ' $\epsilon$ ' making the equation:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

But in practical sense , we do plotting in huge number of data size such as :

## Linear regression

- Given an input  $x$  compute an output  $y$
- For example:
  - Predict height from age
  - Predict house price from house area
  - Predict distance from wall from sensors



House Number	Y: Actual Selling Price	X: House Size (100s ft <sup>2</sup> )
1	89.5	20.0
2	79.9	14.8
3	83.1	20.5
4	56.9	12.5
5	66.6	18.0
6	82.5	14.3
7	126.3	27.5
8	79.3	16.5
9	119.9	24.3
10	87.6	20.2
11	112.6	22.0
12	120.8	.019
13	78.5	12.3
14	74.3	14.0
15	74.8	16.7
Averages	88.84	18.17

Sample 15 houses from the region.

Therefore 'Y' becomes:

$$Y = \beta_0 x^0 + \beta_1 x^1 + \beta_2 x^2 + \dots + \beta_p x^p + \varepsilon$$

When we have 'p' predictor variables.

The above equation is the **equation of the population line** i.e. the equation from which the examples are actually drawn.

Therefore **expected value of 'y' given 'x'** represented as:

$$(E(Y|X)) =$$



$$\hat{\beta}_0 x^0 + \hat{\beta}_1 x^1 + \hat{\beta}_2 x^2 + \dots + \hat{\beta}_p x^p$$

$\hat{\beta}$  is known as hat.

We need to find out  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_n$  so that **sum of squared difference or sum of squared error** is minimised through least squared line.

Now about error,

**Error( $\epsilon$ )**

We have to make assumption about the error:

$$d_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

$$d_2 = \beta_0 + \beta_1 x_2 + \epsilon_2$$

.

.

$$d_n = \beta_0 + \beta_1 x_n + \epsilon_n$$

we assume that expected value of epsilon 'i':

$$E(\epsilon_i) = 0$$

Residual error will have a Mean of Zero (Zero Conditional Mean),

And standard deviation of the error =

$$\sigma(\varepsilon_i) = \sigma(\varepsilon)$$

Errors  $(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$  are independent. And we can also assume that these errors are normally distributed.

This type of noise is known as Gaussian Noise or white noise.

We know sum of the squared errors:

$$\sum_{i=1}^n (y_i - a - bx_i)^2$$

Also written as:

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Or,

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

Sum of the squared errors for estimated value of 'y' on 'x':

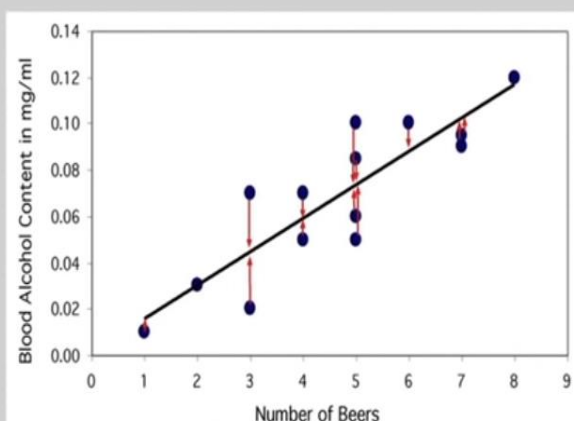
$$\sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

Now , if  $d_i = (x_i, y_i)$ , the above equation can also be written as:

$$\sum_{d \in D} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

## The regression line

The least-squares regression line is the unique line such that the sum of the squared vertical (y) distances between the data points and the line is the smallest possible.



## How do we "learn" parameters

- For the 2- $d$  problem

$$Y = \beta_0 + \beta_1 X$$

- To find the values for the coefficients which minimize the objective function we take the partial derivatives of the objective function (SSE) with respect to the coefficients. Set these to 0, and solve.

$$\beta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\beta_0 = \frac{\sum y - \beta_1 \sum x}{n}$$

17

i.e.:

$$\beta_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{(n \sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i \times \sum_{i=1}^n x_i)}$$

$$\beta_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

Now coming to **Multiple Linear Regression**:

## Multiple Linear Regression

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

$$h(x) = \sum_{i=0}^n \beta_i x_i$$

- There is a closed form which requires matrix inversion, etc.
- There are iterative techniques to find weights
  - delta rule (also called LMS method) which will update towards the objective of minimizing the SSE.

Alternative to matrix inversion etc. methods, the most famous rule called **delta rule** also called **LMS (least minimum slope)** method.

## Linear Regression

$$h(x) = \sum_{i=0}^n \beta_i x_i$$

To learn the parameters  $\theta$  ( $\beta_i$ ) ?

- Make  $h(\mathbf{x})$  close to  $y$ , for the available *training examples*.
- Define a cost function  $J(\theta)$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x)^{(i)} - (y)^{(i)})^2$$

- Find  $\theta$  that minimizes  $J(\theta)$ .

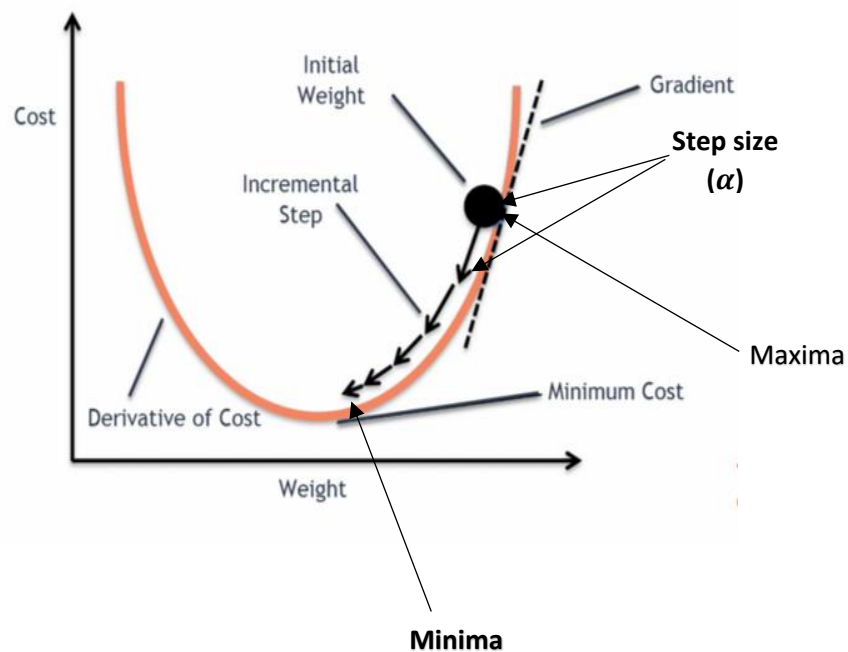
## LMS Algorithm

- Start a search algorithm (e.g. gradient descent algorithm,) with initial guess of  $\theta$ .
- Repeatedly update  $\theta$  to make  $J(\theta)$  smaller, until it converges to minima.

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\theta)$$

- $J$  is a convex quadratic function, so has a single global minima. gradient descent eventually converges at the global minima.
- At each iteration this algorithm takes a step in the direction of steepest descent(-ve direction of gradient).

*$\alpha$  is the step size, if  $\alpha$  is small we take smaller steps and if  $\alpha$  is large we take larger steps.*



$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\theta)$$

Where

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x)^i - y^i)^2$$

Therefore putting value of  $(\theta)$  , if we have only one training example  $(x,y)$ :

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \theta_j} \times \frac{1}{2} \sum_{i=1}^m (h(x)^i - y^i)^2$$

$$\Rightarrow \beta_j - \alpha \frac{\partial}{\partial \theta_j} \times \frac{1}{2} \sum_{i=1}^m (h(x)^i - y^i)^2$$

By applying power rule  $\left( \frac{\partial}{\partial x} (x^n) = nx^{n-1} \right)$ :

$$\Rightarrow \beta_j - \alpha \times \frac{1}{2} \sum_{i=1}^m \{2(h(x)^i - y^i)^{2-1} \times \frac{\partial}{\partial \theta_j} (h(x_i)^i - y^i)\}$$

$$\Rightarrow \beta_j - \alpha \times \frac{1}{2} \sum_{i=1}^m \{2(h(x)^i - y^i) \times \frac{\partial}{\partial \theta_j} (h(x_i)^i - y^i)\}$$

we know  $\theta$  is the change in the gradient for  $i_{th}$  value , hence

$$h = \sum_{i=1}^m \theta_i$$

*similarly we know,*

$$h(x) = \sum_{i=0}^m \beta_i x_i$$



*hence we also can write ,* 
$$h(x) = \sum_{i=0}^m \theta_i x_i$$

Similarly,

$$h(x)^i = \sum_{i=0}^m \theta_i (x_i)^i$$

$$\Rightarrow \beta_j - \alpha \times \left[ \frac{1}{2} \sum_{i=1}^m \{ 2 ( h(x)^i - y^i ) \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^m \theta_i (x_i)^i - y^i \right) \} \right]$$

Then for  $\theta_j$

$$\Rightarrow \beta_j - \alpha \times \left[ \frac{1}{2} \sum_{i=1}^m \{ 2 ( h(x)^i - y^i ) \frac{\partial}{\partial \theta_j} \left( \sum_{j=0}^m \theta_j (x_j)^i - y^i \right) \} \right]$$

$$\begin{aligned} \Rightarrow \beta_j - \alpha & \times \left[ \frac{1}{2} \sum_{i=1}^m \{ 2 ( h(x)^i \right. \\ & \left. - y^i ) \left( \sum_{j=0}^m ( \theta_j \frac{\partial}{\partial \theta_j} (x_j)^i + (x_j)^i \frac{\partial}{\partial \theta_j} \theta_j - 0 \right) \} \right] \end{aligned}$$

$$\Rightarrow \beta_j - \alpha \times \left[ \frac{1}{2} \sum_{i=1}^m \{ 2 ( h(x)^i - y^i ) \left( \sum_{j=0}^m ( 0 + (x_j)^i \times 1 - 0 ) \right) \} \right]$$

$$\Rightarrow \beta_j - \alpha$$

$$\times [\frac{1}{2} \sum_{i=1}^m \{2(\mathbf{h}(\mathbf{x})^i - \mathbf{y}^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) + (\mathbf{h}(\mathbf{x})^i - \mathbf{y}^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) \}]$$

$$\Rightarrow \beta_j - \alpha \times [\frac{1}{2} \sum_{i=1}^m \{2(\mathbf{h}(\mathbf{x})^i - \mathbf{y}^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) \}]$$

$$\Rightarrow \beta_j - \alpha \times [\frac{2}{2} \sum_{i=1}^m \{(\mathbf{h}(\mathbf{x})^i - \mathbf{y}^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) \}]$$

$$\Rightarrow \beta_j - \alpha \times [\sum_{i=1}^m \{(\mathbf{h}(\mathbf{x})^i - \mathbf{y}^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) \}]$$

$$\Rightarrow \beta_j + \alpha \times [\sum_{i=1}^m \{(\mathbf{y}^i - \mathbf{h}(\mathbf{x})^i) \left( \sum_{j=0}^m (\mathbf{x}_j)^i \right) \}]$$

If we remove  $\Sigma$  from the equations, we get:

$$\beta_j + \alpha \times (y^i - h(x)^i)(x_j)^i$$

### LMS Update Rule

- If you have only one training example  $(x, y)$

$$\begin{aligned}\frac{\partial}{\partial \theta} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h(x) - y) \frac{\partial}{\partial \theta_j} (h(x) - y) \\ &= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h(x) - y) x_j\end{aligned}$$

- For a single training example, this gives the update rule:

$$\beta_j = \beta_j + \alpha (y^{(i)} - h(x^i)) x_j^{(i)}$$