# Simple Linear Regression

$\Rightarrow$ when we have a single input attribute (x) and we want to use linear regression, this is called Simple Linear Regression.

$\Rightarrow$ If we had multiple input attributes (eg. x1, x2, x3 etc.) This would be called multiple linear regression.

$\Rightarrow$ with Simple Linear Regression Model we want to model our data as follows :-

$$Y = B_0 + B_1 * x$$

&      O/P Variable      Coefficient      input Variable

Technically $B_0$ is called the intercept and $B_1$ is called the slope.

⇒ The goal is to find the best estimates for the coefficients to minimize the errors in predicting $y$ from $x$.

⇒ First we estimate the value of for $B1$ as:

$$B1 = \frac{sum((x_i - mean(x)) \times (y_i - mean(y)))}{sum((x_i - mean(x))^2)}$$

mean () → average value for the variable in our dataset.

⇒ we can calculate $B_0$ using $B1$ and some statistical from our dataset as follow:

$$B_0 = mean(y) - B1 \times mean(x).$$

Estimating the Slope (B1)

⇒ First we need to calculate the mean value of x and y.

$$mean = \frac{1}{n \times Sum(x)} \cdot \boxed{\frac{1}{n} \times Som(x)}$$

n ⇒ is the no. of values (5 in this case).
We can use the AVERAGE() function in your spreadshed.

mean (x) = 3
mean (y) = 2.8

$$\frac{1}{5} \times 15 = 3$$

$$\frac{1}{5} \times 14 = 2.8$$

| 1 | X | Y |  |
|---|---|---|---|
| 2 | 1 | 1 |  |
| 3 | 2 | 3 |  |
| 4 | 4 | 3 |  |
| 5 | 3 | 2 |  |
| 6 | 5 | 5 |  |

⇒ Now we calculate the error of each variable from the mean

| | 1 X | mean(x) | X-mean(x) |
|---|---|---|---|
| 2 | 1 | 3 | -2 |
| 3 | 2 | 3 | -1 |
| 4 | 3 | 3 | 0 |
| 5 | 4 | 3 | 1 |
| 6 | 5 | 3 | 2 |

| | Y | mean(y) | y-mean(y) |
|---|---|---|---|
| 1 | 1 | 2.8 | -1.8 |
| 2 | 3 | 2.8 | 0.2 |
| 3 | 3 | 2.8 | 0.2 |
| 4 | 2 | 2.8 | -0.8 |
| 5 | 5 | 2.8 | 2.2 |
| 6 | | | |

=> No $\omega$ we have all parts for calculating the numerator. All we need to do is multiple the error for each x with error for each y. and calculate the sum of multiplication

| | X-mean(x) | y-mean(y) | multiplication |
|---|---|---|---|
| 1 | -2 | -1.8 | 3.6 |
| 2 | 2 | 0.2 | -0.2 |
| 3 | -1 | 0.2 | 0.2 |
| 4 | | +0.8 | -0.8 |
| 5 | | 2.2 | 4.4 |
| 6 | | | ⑦ |

$$\begin{array}{r} 3.6 \\ 4.4 \\ \hline 8.0 \end{array}$$

8.0

=> Summing the final column we have the remesatch as 0.7

=> Now we calculate the bottom part of the eqⁿ for calculating B₁

| 1 | X-mean(x) | Squared |
|---|-----------|---------|
| 1 | -2 | 4 |
| 2 | -2 | 1 |
| 3 | -1 | 1 |
| 4 | 1 | 0 |
| 5 | 0 | 0 |
| 6 | 2 | 4 |

$$\frac{7}{10} = 0.7$$

$$B_1 = 8/10 = 0.8$$

Now estimating the intercept (B0)

$$B_0 = mean(y) - B_1 \times mean(x)$$

$$B_0 = 2.8 - 0.8 \times 3$$

$$B_0 = 0.4$$

making Predictions.
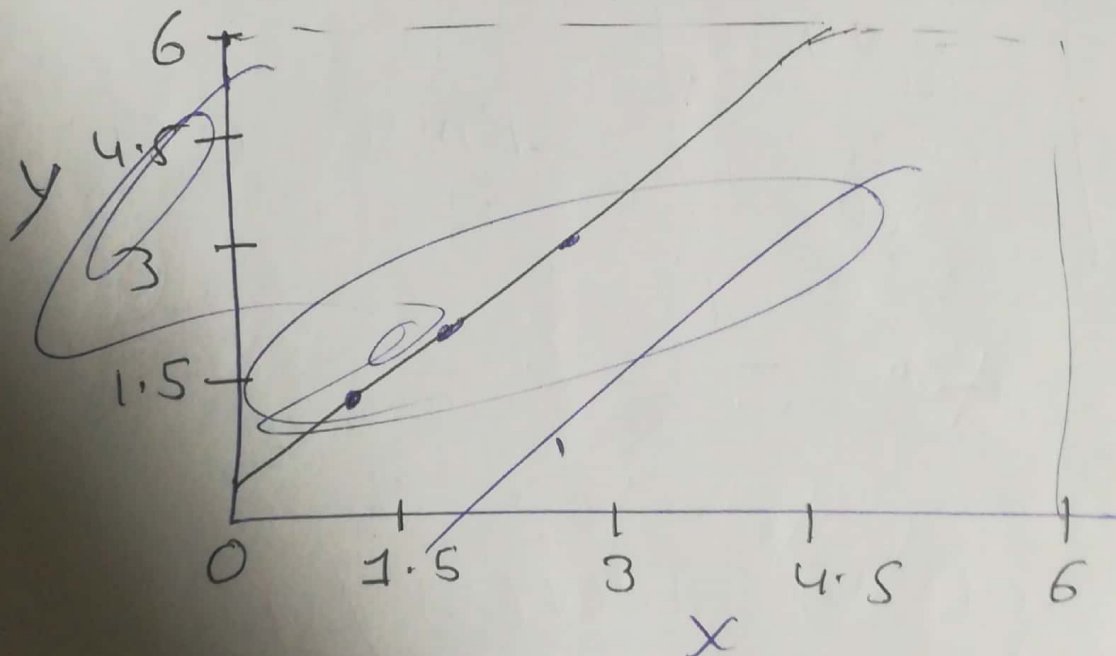
$$y = B_0 + B_1 * x$$

$$y = 0.7 + 0.8 * x$$

Lets try out the mode by making Prediction for our training data.

| | x | y | Predicted y |
|---|---|---|---|
| 1 | 1 | 1 | 1.2 |
| 2 | 2 | 3 | 2 |
| 3 | 3 | 3 | 3.6 |
| 4 | 4 | 2 | 2.8 |
| 5 | 5 | 5 | 4.4 |
| 6 | | | |

Estimating the Error.

Root mean Squared Error or RMSE

$$RMSE = sqrt\left(\frac{Sum((p_i - y_i)^2)}{m}\right)$$

$\downarrow$ Predicted value    $\searrow$ actual value.

| i | Pred | y | error |
|---|------|---|-------|
| 1 | | 1 | 0.2 |
| 2 | 1.2 | | |
| 3 | 2 | 3 | -1 |
| 4 | 3.6 | 3 | 0.6 |
| 5 | 2.8 | 2 | 0.8 |
| 6 | 4.4 | 5 | 0.6 |

| error | Squared error |
|-------|---------------|
| 0.2 | 0.04 |
| 1 | 1 |
| 0.6 | 0.36 |
| 0.8 | 0.64 |
| -0.6 | 0.36 |

RMSE = 0.692

Shortcut :-

$$B1 = \frac{Corr(x, y)}{\downarrow} * Stdev(y) / Stdev(x)$$

Pearson's correlation coefficient.