

Regression Math Explained

Assume dataset below:

X	Y
0	4
1	12
2	28
3	52
4	80

To find the line of the best fit, we can choose two points and calculate $y = mx + c$ easily:

$X_1, Y_1 = 0, 4$

$X_2, Y_2 = 4, 80$

Firstly, calculate the slope using two values:

$$m = \frac{y_1 - y_2}{x_1 - x_2}$$

$$m = \frac{4 - 80}{0 - 4}$$

$$m = -\frac{76}{-4} = \mathbf{19}$$

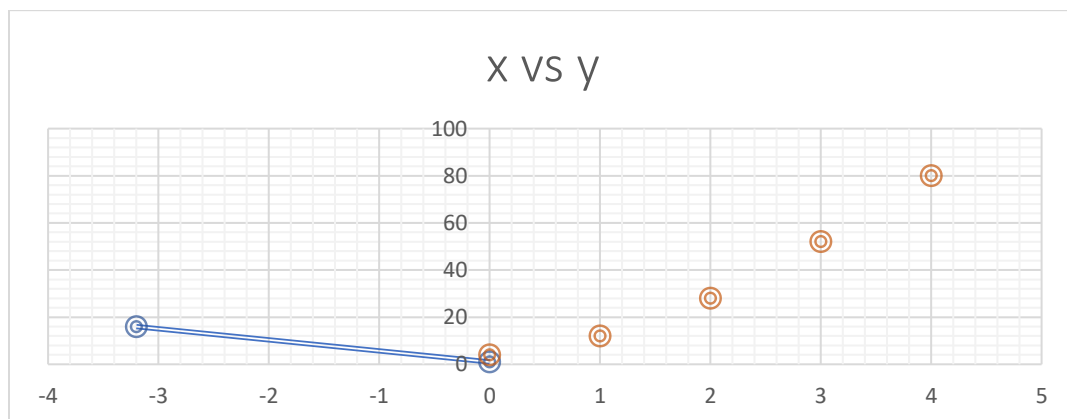
Y-intercept or coefficient bias is what the value of y is when x is 0. Given $m = 19$ and $x, y = 1, 4$:

$$y - y_1 = m * (x - x_1)$$

$$y - 4 = 19x - 19$$

$$\mathbf{y = 19x - 15}$$

This formula above is fairly accurate however, we need to minimize the errors. Notice graph below. It perfectly describes points used, however ones in between are not described as well.



To minimize the error, we use Least Square Regression method which minimizes the error in such a way that all square errors are minimized. First, the formula for calculating slope is:

$$m = \frac{\text{sum}((x - x_{\text{mean}}) * (y - y_{\text{mean}}))}{\text{sum}(x - x_{\text{mean}})^2}$$

Firstly, calculate x and y means:

$$x_{\text{mean}} = \frac{0 + 1 + 2 + 3 + 4}{5} = 2$$

Additionally, calculate $x - x_{\text{mean}}$ for each value of x :

$$x_1 - x_{\text{mean}} = 0 - 2 = -2$$

$$x_2 - x_{\text{mean}} = 1 - 2 = -1$$

$$x_3 - x_{\text{mean}} = 2 - 2 = 0$$

$$x_4 - x_{\text{mean}} = 3 - 2 = 1$$

$$x_5 - x_{\text{mean}} = 4 - 2 = 2$$

Do the same for y values:

$$y_{\text{mean}} = \frac{4 + 12 + 28 + 52 + 80}{5} = 35.2$$

Additionally, calculate $y - y_{\text{mean}}$ for each value of y :

$$y_1 - y_{\text{mean}} = 4 - 35.2 = -31.2$$

...

$$y_5 - y_{\text{mean}} = 80 - 35.2 = 44.8$$

Sum of all these multiplied together will give us numerator which equals to **192**. Sum of $(x - x_{\text{mean}})^2$ values will give us **10**. Therefore, slope results in **19.2**.

To calculate the y -intercept we use:

$$\begin{aligned} c &= y_{\text{mean}} - m * x_{\text{mean}} \\ &= 35.2 - 19.2 \\ &\quad * 2 = -22.4 \end{aligned}$$

Overall formula equals to:

$$y = 19.2x - 22.4$$

Notice the difference in the plot on the right.

