

# How is the Ordinary Least Squares formula derived?

Statistics Linear Regression and Correlation Least Squares Regression Line (LSRL)

## 1 Answer



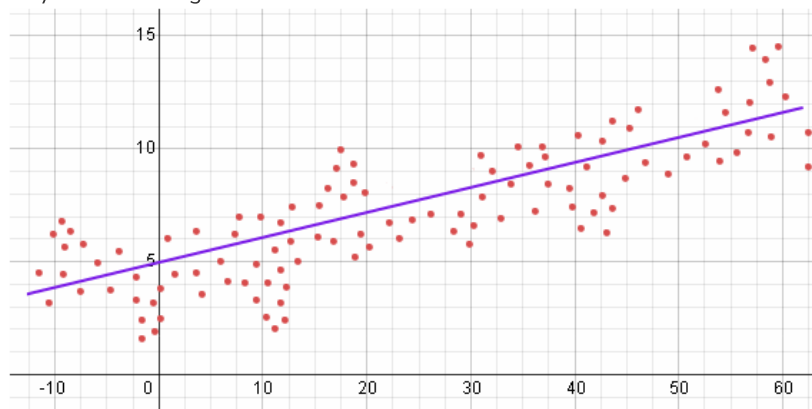
**Shwetank Mauria**  
Mar 19, 2018

Please see below.

### Explanation:

Let us say that we expect a linear relation, say  $y = ax + b$ , between two observable variables  $x$  and  $y$  and where  $a$  and  $b$  are not known. It is not necessary that the relation should be linear, but here we have assumed this only for the purpose of simplification.

However, actual observed data may be accompanied by errors or noise for various reasons and when we plot the observed data on a graph, it may not be a linear and may be some thing as shown below.



Observe that each point denotes observed value of  $x$  and  $y$  and  $n$  points provides us with  $n$  data points given by  $(x_i, y_i)$ , where  $i$  ranges from 1 to  $n$ .

We can draw many lines through these points, with varying slopes i.e.  $a$ 's and intercepts  $b$ 's, but how do we know which one is the best fit. This is done by the method of **Least squares**. In this method we define the relationship between observed data and expected data - **by minimizing the sum of squares of the deviation between observed and expected values**.

In other words, best line has minimum error between line and data points. Note that had we not squared, positive and negative errors would have almost cancelled out. We will talk more about this later.@

For a particular  $x_i$  observed value is  $y_i$ , expected value of  $y_i$  for is  $ax_i + b$  and difference between them  $d_i$  is given by  $y_i - ax_i - b$  and in **least square method** we seek to

minimise error  $E = \sum_{i=1}^n d_i^2$  i.e.  $\sum_{i=1}^n (y_i - ax_i - b)^2$ .

### Related questions

What is meant by the term "least squares" in linear regression?

What is the general formate for the equation of a least-squares regression line?

What is the primary use of linear regression?

What is regression analysis?

What is a regression analysis?

What does a regression analysis tell you?

Why must the R-Squared value of regression be less than 1?

What is the standard error?

What is the t-statistic?

In a regression analysis, if R-Squared = 1, then does SSE = SST

See all questions in Least Squares Regression Line (LSRL)

### Impact of this question

2733 views around the world



You can reuse this answer Creative Commons License

To get this minimum, we use calculus. For this, we must have first derivatives of  $a$

and  $b$  yielding zero. Differentiating  $\sum_{i=1}^n (y_i - ax_i - b)^2$  w.r.t.  $a$  and  $b$ , we get

$$\frac{\partial E}{\partial a} = -2 \sum_{i=1}^n x_i (y_i - ax_i - b) = 0$$

$$\text{and } \frac{\partial E}{\partial b} = -2 \sum_{i=1}^n (y_i - ax_i - b) = 0$$

To solve for  $a$  and  $b$ , we rewrite them as

$$a \sum x_i^2 + b \sum x_i = \sum x_i y_i \text{ and}$$

$$a \sum x_i + bn = \sum y_i$$

and solving them for  $a$  and  $b$  we get

$$a = \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\text{and } b = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

@ - Note that  $a \sum x_i + bn = \sum y_i$  can be expressed as  $a \sum \frac{x_i}{n} + b = \sum \frac{y_i}{n}$  is just a fit between averages of  $x_i$  and  $y_i$ . Hence, this alone may not give the best fit. The latter is arrived due to  $a \sum x_i^2 + b \sum x_i = \sum x_i y_i$ .

[Answer link](#)