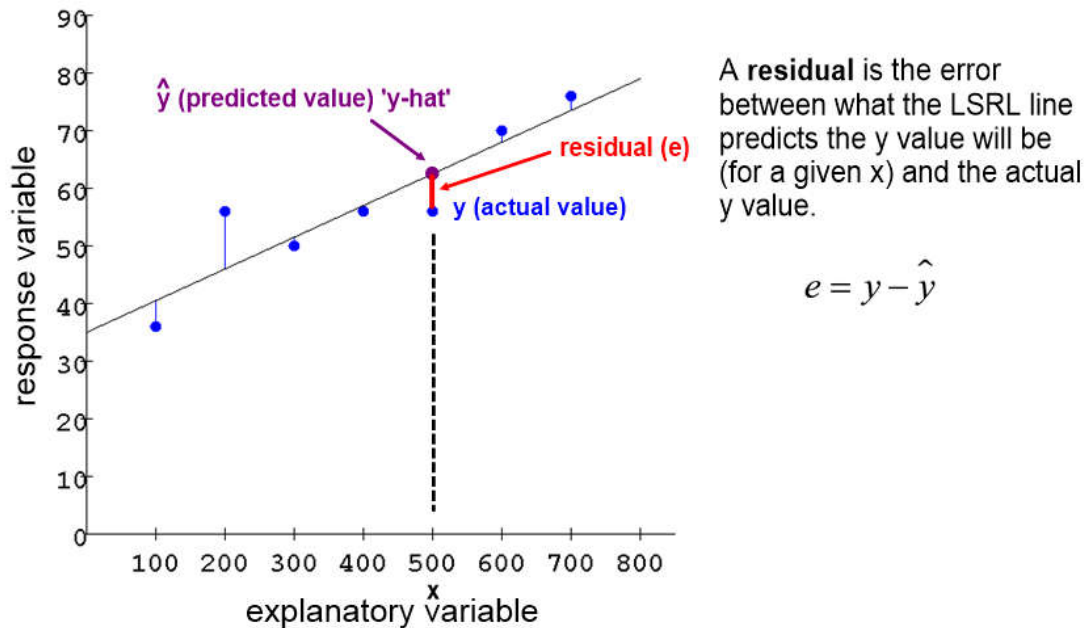


## Derivation of equations used for calculating the Least Squares Regression Line (LSRL) for a data set

In statistics, we often need to 'fit a line' to a set of data in such a way that the errors between the predicted  $y$  value for each  $x$  and the actual  $y$  value for that  $x$  are as small as possible. We can use multivariable optimization techniques to derive the needed equations for computation.

We start by defining the residual,  $e$ , as the difference between actual  $y$  and predicted  $y$  (actual – predicted) at any given  $x$ :



We then define a Least Squares Regression Line (LSRL) which will be the resulting equation:  $\hat{y} = a + bx$   
Where  $y$ -hat represents the predicted  $y$  output from the LSRL equation at input value  $x$ .

The error is then given by:  $e = y - \hat{y} = y - (a + bx)$

What we are really solving for are the two coefficients,  $a$  and  $b$  which result in a minimum total error across all the points. We can define a function for error in terms of  $a$  and  $b$ :

$$\text{total error (squared)} = f(a, b) = \sum_{i=1}^n (y_i - (a + bx_i))^2$$

The reason we square is to force all errors both above and below the line to be positive (we don't want the negative errors cancelling out the positive errors).

If we want to optimize this function (find its minimum value) we will take the two partial derivatives, set them equal to zero and solve for critical points, then use the Hessian determinant to verify that this is a minimum.

First, expanding the function:

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

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

$$f(a,b) = \sum_{i=1}^n (y_i^2 - 2y_i a - 2y_i b x_i + a^2 + 2abx_i + b^2 x_i^2)$$

We take the first derivatives and set them equal to zero:

$$\begin{cases} f_a = \sum_{i=1}^n (-2y_i + 2a + 2bx_i) = 0 \\ f_b = \sum_{i=1}^n (-2y_i x_i + 2ax_i + 2bx_i^2) = 0 \end{cases}$$

These form a system. Working further with the  $f_a$  equation:

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

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

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

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

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

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

$$\left( \sum_{i=1}^n (1) \text{ is just } n \right)$$

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

Working further with the  $f_b$  equation:

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

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

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

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

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

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

So now we can write the system like this:

$$\left\{ \begin{array}{l} f_a = \sum_{i=1}^n (-2y_i + 2a + 2bx_i) = 0 \\ f_b = \sum_{i=1}^n (-2y_i x_i + 2ax_i + 2bx_i^2) = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} \sum_{i=1}^n (y_i) = na + b \sum_{i=1}^n (x_i) \\ \sum_{i=1}^n (y_i x_i) = a \sum_{i=1}^n (x_i) + b \sum_{i=1}^n (x_i^2) \end{array} \right.$$

$$\left\{ \begin{array}{l} [n]a + \left[ \sum_{i=1}^n (x_i) \right] b = \left[ \sum_{i=1}^n (y_i) \right] \\ \left[ \sum_{i=1}^n (x_i) \right] a + \left[ \sum_{i=1}^n (x_i^2) \right] b = \left[ \sum_{i=1}^n (y_i x_i) \right] \end{array} \right.$$

The expressions in the brackets are numerical values computed from the data set. To make this easier to see, we can replace each with a constant (which would be calculated from the data set):

*Define:*

$$Q = \left[ \sum_{i=1}^n (x_i) \right]$$
$$R = \left[ \sum_{i=1}^n (y_i) \right]$$
$$S = \left[ \sum_{i=1}^n (x_i^2) \right]$$
$$T = \left[ \sum_{i=1}^n (y_i x_i) \right]$$

Then the system would become:

$$\begin{cases} [n]a + \left[ \sum_{i=1}^n (x_i) \right] b = \left[ \sum_{i=1}^n (y_i) \right] \\ \left[ \sum_{i=1}^n (x_i) \right] a + \left[ \sum_{i=1}^n (x_i^2) \right] b = \left[ \sum_{i=1}^n (y_i x_i) \right] \end{cases}$$
$$\begin{cases} na + Qb = R \\ Qa + Sb = T \end{cases}$$

Which can then be solved using RREF for the LSRL coefficients  $a$  and  $b$ .

How do we know this is where total squared error is minimized? (As opposed to maximized or a saddle point)

To verify this is a minimum, we would need to determine the value of the Hessian Determinant, so first we need expressions for the 2<sup>nd</sup> derivatives:

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

$$f_{aa} = \sum_{i=1}^n (2) = 2n$$

$$f_{ab} = \sum_{i=1}^n (2x_i)$$

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

$$f_{bb} = \sum_{i=1}^n (2x_i^2)$$

The Hessian determinant would then be:

$$\begin{aligned} D &= f_{aa}f_{bb} - (f_{ab})^2 \\ &= (2n) \left( \sum_{i=1}^n (2x_i^2) \right) - \left( \sum_{i=1}^n (2x_i) \right)^2 \\ &= 4n \sum_{i=1}^n (x_i^2) - 4 \sum_{i=1}^n (x_i^2) \\ &= 4 \sum_{i=1}^n (x_i^2) (n-1) \end{aligned}$$

...which is definitely positive, so a maximum or a minimum. Looking at concavity  $f_{aa}$  to check:  $f_{aa} = 2n$  is also definitely positive, which means concave up, which verifies that these values of  $a$  and  $b$  produce the minimum total square error.