

THE MATH BEHIND LINEAR REGRESSION

By *Clutch Data Science*

The simple linear regression model:

$$Y = \beta_0 + \beta_1 X + \epsilon \quad (1)$$

We use $\hat{\beta}_0, \hat{\beta}_1$ as the estimator of β_0, β_1 . To find the best line that approximate the linear relationship, we use the Least Squared method, finding the $\hat{\beta}_0, \hat{\beta}_1$ that minimize the sum of squared residuals (SSR).

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

1. THE CLOSED FORM SOLUTION FOR $\hat{\beta}_0$ AND $\hat{\beta}_1$

To derive the solution, we need to take the partial derivative w.r.t $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\frac{\partial SSR}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad (3)$$

$$\frac{\partial SSR}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad (4)$$

To get $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize the SSR , we Set (3), (4) to 0. Then we can derive:

$$\bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x} = 0 \quad (5)$$

$$\sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0 \quad (6)$$

From (5) we have:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (7)$$

Substitute (7) into (6):

$$\begin{aligned}
\sum x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum x_i - \hat{\beta}_1 \sum x_i^2 &= 0 \\
\hat{\beta}_1 (\bar{x} \sum x_i - \sum x_i^2) &= - \sum x_i y_i + \bar{y} \sum x_i \\
\hat{\beta}_1 \sum x_i (\bar{x} - x_i) &= \sum x_i (\bar{y} - y_i) \\
\hat{\beta}_1 &= \frac{\sum x_i (\bar{y} - y_i)}{\sum x_i (\bar{x} - x_i)}
\end{aligned}$$

Furthermore, we normally write $\hat{\beta}_1$ in the following format:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})(x_i - \bar{x})} \quad (8)$$

Because $\sum \bar{x}(y_i - \bar{y}) = \bar{x} \sum (y_i - \bar{y}) = 0$ and $\sum \bar{x}(x_i - \bar{x}) = \bar{x} \sum (x_i - \bar{x}) = 0$

Therefore, (7) and (8) gives the mathematical solution for $\hat{\beta}_0$ and $\hat{\beta}_1$.

2. THE VARIANCE OF $\hat{\beta}_0$ AND $\hat{\beta}_1$

From (7), (8) we have:

$$\begin{aligned}
\text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) \\
&= \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2\bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1)
\end{aligned}$$

$$\text{Var}(\bar{y}) = \text{Var}\left(\frac{1}{n} \sum y\right) = \frac{1}{n^2} \sum \text{Var}(y) = \frac{\sigma^2}{n}$$

$$\text{Var}(\hat{\beta}_1) = \frac{1}{[\sum (x_i - \bar{x})^2]^2} \sum (x_i - \bar{x})^2 \text{Var}(y_i) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

The covariance term is:

$$\begin{aligned}
\text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}\left(\frac{1}{n} \sum y_i, \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2}\right) \\
&= \frac{1}{n} \frac{1}{\sum (x_i - \bar{x})^2} \text{Cov}\left(\sum_i y_i, \sum_j (x_j - \bar{x}) y_j\right) \\
&= \frac{1}{n \sum (x_i - \bar{x})^2} \sum_j (x_j - \bar{x}) \text{Cov}\left(\sum_i y_i, y_j\right) \\
&= \frac{1}{n \sum (x_i - \bar{x})^2} \sum_j (x_j - \bar{x}) \sum_i \text{Cov}(y_i, y_j) \\
&= \frac{1}{n \sum (x_i - \bar{x})^2} \sum_j (x_j - \bar{x}) n \sigma^2 \\
&= \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \sum_j (x_j - \bar{x}) \\
&= 0
\end{aligned}$$

Therefore, we have:

$$\text{Var}(\beta_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right) \quad (9)$$

$$\text{Var}(\beta_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \quad (10)$$