# SIMPLE LINEAR REGRESSION AND...

## SIEYE RYU

### 1. SIMPLE LINEAR REGRESSION

Suppose that we have $n$ pairs of data
$$\{(x_i, y_i) \in \mathbb{R}^2 : i = 1, \cdots, n\}.$$
We will construct a linear predictor function $y = \beta_1 x + \beta_0$ from the data, based on the method of least squares.

The *residual* $r_i$ $(i = 1, \cdots, n)$ is the difference between an observed value and the fitted value provided by a linear prediction:
$$r_i = y_i - (\beta_1 x_i + \beta_0) \qquad (i = 1, \cdots, n).$$

The method of least squares provides an approximation of the data by minimizing the sum of the squares of the residuals. In other words, if $S$ is a function on $\mathbb{R}^2$ defined by
$$S(b_0, b_1) = \sum_{i=1}^{n} (y_i - b_1 x_i - b_0)^2,$$
and if the point $(\beta_0, \beta_1)$ minimizes $S(b_0, b_1)$, then the linear predictor function is given by
$$y = \beta_1 x + \beta_0.$$

Intuitively, $S(b_0, b_1)$ has a global minimum and it has no global maximum or saddle points. Hence, it is enough to find a critical point of $S$. To do this we calculate the gradient of $S$. By chain rule, we have
$$\nabla S(b_0, b_1) = \left( -2 \sum_{i=1}^{n} (y_i - b_1 x_i - b_0), \ -2 \sum_{i=1}^{n} x_i(y_i - b_1 x_i - b_0) \right).$$
If $(\beta_0, \beta_1)$ is a critical point of $S$, then we have
$$\sum_{i=1}^{n} (y_i - \beta_1 x_i - \beta_0) = 0 \tag{1.1}$$

$$\sum_{i=1}^{n} x_i(y_i - \beta_1 x_i - \beta_0) = 0. \tag{1.2}$$

The first equation (1.1) can be written as follows:
$$\sum_{i=1}^{n} (y_i - \beta_1 x_i) - n\beta_0 = 0.$$
This leads us to
$$\beta_0 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \beta_1 x_i). \tag{1.3}$$

1

Substituting (1.3) into (1.2) yields

$$\sum_{i=1}^{n} x_i(y_i - \beta_1 x_i - \frac{1}{n} \sum_{j=1}^{n} (y_j - \beta_1 x_j)) = 0.$$

This is equivalent to

$$\beta_1 \left( -\sum_{i=1}^{n} x_i^2 + \frac{1}{n} \left( \sum_{i=1}^{n} x_i \right)^2 \right) + \sum_{i=1}^{n} x_i y_i - \frac{1}{n} \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right) = 0.$$

Thus, we have

$$\beta_1 = \frac{\sum x_i y_i - \frac{1}{n} \left( \sum x_i \right) \left( \sum y_i \right)}{\sum x_i^2 - \frac{1}{n} \left( \sum x_i \right)^2}. \tag{1.4}$$

We denote the means of $x_i$ and $y_i$ $(i = 1, \cdots, n)$ by $\bar{x}$ and $\bar{y}$, respectively. Then (1.3) is equivalent to

$$\beta_0 = \bar{y} - \beta_1 \bar{x}.$$

In the rest of the section, we prove that (1.4) is equivalent to

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}. \tag{1.5}$$

The following calculation shows that the numerator of (1.5) is the same as that of (1.4):

$$\begin{aligned}
\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^{n} \left( x_i - \frac{1}{n} \sum_{j=1}^{n} x_j \right) \left( y_i - \frac{1}{n} \sum_{j=1}^{n} y_j \right) \\
&= \sum_{i=1}^{n} x_i y_i - \frac{2}{n} \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right) + \frac{n}{n^2} \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right) \\
&= \sum_{i=1}^{n} x_i y_i - \frac{1}{n} \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right).
\end{aligned}$$

Replacing $y_i$ and $\bar{y}$ with $x_i$ and $\bar{x}$ shows that the denominators of (1.4) and (1.5) are the same.

## References

[1]     Linear regression on Wikipedia
[2]     Simple linear regression on Wikipedia
[3]     Least squares on Wikipedia