

Computing Variance of Predictions in Linear Regression

Suppose we have independent variables

$$\{x_1, \dots, x_n\}$$

and dependent variables

$$\{y_1, \dots, y_n\}$$

Recall that linear regression uses predictions

$$\hat{y} = \bar{y} + r \frac{SD(y)}{SD(x)} (x - \bar{x})$$

Here

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$SD(y) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}, \quad SD(x) = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{SD(x) \cdot SD(y)}$$

Remember that we can use any expression of the form

$$\hat{y} = a + bx$$

However we choose parameters a for the intercept and b for the slope with the average loss

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

We want to minimize the average loss by choosing \hat{a} and \hat{b} . We can use the derivative to determine when the average loss attains the minimum.

Rule of change with respect to a

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

Rule of change with respect to b

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

Remember that we call the errors

$$y_i - \hat{y}_i = y_i - (a + bx_i)$$

between observed values and predicted values the residuals. We can interpret these equations

Rule of change with respect to a

The average of the residuals = 0

Rule of change with respect to b

The covariance of the residuals and independent variables = 0

We want both expressions to be 0 because the rate of change should be 0 at the minimum. By setting both expressions equal to 0, we have two equations and two unknowns.

These equations allow us to solve for \hat{a} and \hat{b} . Additionally these equations help us to relate the variance of $\{y_1, \dots, y_n\}$ to the variance of $\{\hat{y}_1, \dots, \hat{y}_n\}$. Note

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n \hat{y}_i \\ &= \frac{1}{n} \sum_{i=1}^n \bar{y} + r \frac{SD(y)}{SD(x)} (x_i - \bar{x}) \\ &= \frac{1}{n} \sum_{i=1}^n \bar{y} + \frac{1}{n} \sum_{i=1}^n r \frac{SD(y)}{SD(x)} (x_i - \bar{x}) \\ &= \bar{y} + 0 \end{aligned}$$

So we can calculate

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{y})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ \text{use } \bar{y} = \bar{\hat{y}} &\rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{\hat{y}})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{\hat{y}}) \end{aligned}$$

We can simplify

$$\begin{aligned} \sum_{i=1}^n 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n 2(y_i - \hat{y}_i) \left[\frac{r SD(y)}{SD(x)} (x_i - \bar{x}) \right] \\ &= \frac{2r SD(y)}{SD(x)} \sum_{i=1}^n (y_i - \hat{y}_i)(x_i - \bar{x}) \end{aligned}$$

With

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{y}_i)(x_i - \bar{x}) &= \sum_{i=1}^n (y_i - \hat{y}_i) x_i - \bar{x} \sum_{i=1}^n y_i - \hat{y}_i \\ &= 0 + 0 \\ &\quad \text{use sum of residuals is 0} \quad \text{use covariance of residuals and independent variables is 0} \end{aligned}$$

Therefore

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{\hat{y}})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{\hat{y}})^2 \end{aligned}$$

We conclude that

$$\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2$$

\uparrow Variance of observed values \uparrow mean square error \uparrow Variance of predicted values

Note that the variance of predicted values is less than the variance of observed values.

This is called regression to the mean.