

## Explained Variance

Suppose  $X$  has a column of 1s, which we'll denote by  $\mathbf{1}_n$ . Since  $Y = X\hat{\beta} + e$  is the orthogonal decomposition of  $Y$  with respect to  $\text{col}(X)$ , we have  $\|Y\|^2 = \|X\hat{\beta}\|^2 + \|e\|^2$ . Now the sample mean of  $Y$  is  $\mathbf{1}'_n Y/n$  and the sample variance of  $Y$  is

$$\text{var}(Y) = \frac{\|Y\|^2}{n} - \left(\frac{\mathbf{1}'_n Y}{n}\right)^2.$$

We can express the sample sum as  $\mathbf{1}'_n Y = \mathbf{1}'_n (X\hat{\beta}) + \mathbf{1}'_n e = \mathbf{1}'_n X\hat{\beta}$ , so that the sample mean of  $Y$  is equal to the sample mean of  $X\hat{\beta}$ . (Another way of saying or seeing this is that the sample mean of  $e$  is 0.) Therefore,

$$\begin{aligned} \text{var}(Y) &= \frac{\|Y\|^2}{n} - \left(\frac{\mathbf{1}'_n Y}{n}\right)^2 = \frac{\|Y\|^2}{n} - \left(\frac{\mathbf{1}'_n X\hat{\beta}}{n}\right)^2 = \frac{\|X\hat{\beta}\|^2 + \|e\|^2}{n} - \left(\frac{\mathbf{1}'_n X\hat{\beta}}{n}\right)^2 \\ &= \text{var}(X\hat{\beta}) + \text{var}(e), \end{aligned}$$

where we have used the fact that the sample mean of  $e$  is 0. This is equation (4.22) in the book. It leads to the definition  $R^2 := \text{var}(X\hat{\beta})/\text{var}(Y)$ . The equation gives another way to write this:  $1 - R^2 = \text{var}(e)/\text{var}(Y)$ .

Now recall that with simple linear regression, the RMS error is  $\sqrt{1 - r^2} \cdot \text{sd}(Y)$ . The definition of the RMS error, in our new notation, is  $\sqrt{\|e\|^2/n} = \text{sd}(e)$ . In other words, we proved before for simple linear regression that  $\text{sd}(e) = \sqrt{1 - r^2} \cdot \text{sd}(Y)$ , i.e.,  $1 - r^2 = \text{var}(e)/\text{var}(Y)$ . Thus, we conclude that  $R^2 = r^2$ .

Exercise: Show that  $R$ , the non-negative square root of  $R^2$ , is the sample correlation between  $Y$  and  $\hat{Y}$ , where, as usual,  $\hat{Y} = X\hat{\beta}$ .