

Linear Regression Demystified

Linear regression is an important subject in statistics. In elementary statistics courses, formulae related to linear regression are often stated without derivation. This note intends to derive these formulae for students with more advanced math background.

1 Simple Regression: One Variable

1.1 Least Square Prescription

Suppose we have a set of data points (x_i, y_i) ($i = 1, 2, \dots, n$). The goal of linear regression is to find a straight line that best fit the data. In other words, we want to build a model

$$\hat{y}_i = \beta_0 + \beta_1 x_i \quad (1)$$

and tune the parameters β_0 and β_1 so that \hat{y}_i is as close to y_i as possible for all $i = 1, 2, \dots, n$.

Before we do the math, we need to clarify the problem. How do we judge the “closeness” of \hat{y}_i and y_i for all i ? If the data points (x_i, y_i) do not fall exactly on a straight line, y_i and \hat{y}_i is not going to be the same for all i . The deviation of y_i from \hat{y}_i is called the *residual* and is denoted by ϵ_i here. In other words,

$$y_i = \hat{y}_i + \epsilon_i = \beta_0 + \beta_1 x_i + \epsilon_i. \quad (2)$$

The least square prescription is to find β_0 and β_1 that minimize the sum of the square of the residuals SSE :

$$SSE = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2. \quad (3)$$

If you are familiar with calculus, you will know that the minimization can be done by setting the derivatives of SSE with respect to β_0 and β_1 to 0. The resulting equations are¹

$$\sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n \epsilon_i = 0 \quad \text{and} \quad \sum_{i=1}^n x_i (y_i - \hat{y}_i) = \sum_{i=1}^n x_i \epsilon_i = 0. \quad (4)$$

If you are not familiar with calculus, you will have to take my words. We will try to give you an intuition of the meaning of the above equations below. Before we do that, let's remind you some concepts in statistics.

The mean and standard deviation of a set of points $\{u_i\}$ are defined as

$$\bar{u} = \frac{1}{n} \sum_{i=1}^n u_i, \quad SD_u = \sqrt{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^2}. \quad (5)$$

The Z score of u_i is defined as

$$Z_{ui} = \frac{u_i - \bar{u}}{SD_u}. \quad (6)$$

The *correlation* of two set of points (with equal number of elements) $\{x_i\}$ and $\{y_i\}$ is defined as

$$r_{xy} = r_{yx} = \frac{1}{n} \sum_{i=1}^n Z_{xi} Z_{yi}. \quad (7)$$

¹Careful students may realize that this only guarantees that SSE is stationary but not necessarily minimum. There is also a concern whether the resulting solution is a global minimum or a local minimum. Detail analysis reveals that the solution in our case is a global minimum.

The Cauchy-Schwarz inequality states that for two set of real numbers $\{u_i\}$ and $\{v_i\}$,

$$\left(\sum_{i=1}^n u_i v_i\right)^2 \leq \left(\sum_{i=1}^n u_i^2\right) \left(\sum_{i=1}^n v_i^2\right)$$

and the equality holds if and only if $v_i = k u_i$ for all i , where k is a constant. You can find the proof of the inequality in, e.g., Wikipedia.

It follows from the Cauchy-Schwarz inequality that the correlation $|r_{xy}| \leq 1$, with $|r_{xy}| = 1$ if and only if x_i and y_i fall exactly on a straight line, i.e. $\epsilon_i = 0$ for all i . So the correlation r_{xy} measures how tightly the points (x_i, y_i) are clustered around a line.

We can now rewrite the first equation of (4) as

$$\bar{\epsilon} = 0. \quad (8)$$

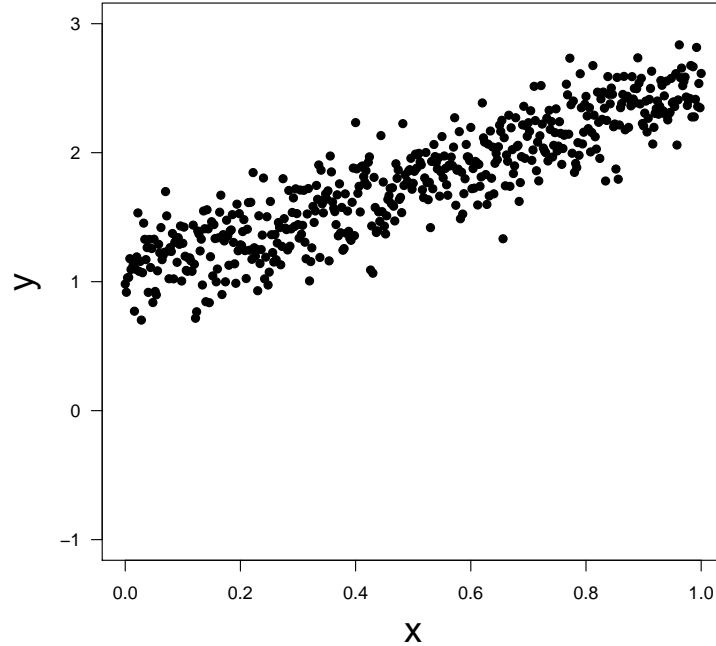
That is to say that the mean of ϵ_i vanishes. We can also combine the two equations of (4) as

$$\sum_{i=1}^n x_i \epsilon_i - \bar{x} \sum_{i=1}^n \epsilon_i = 0 \Rightarrow \sum_{i=1}^n (x_i - \bar{x})(\epsilon_i - \bar{\epsilon}) = 0 \Rightarrow nSD_x SD_\epsilon r_{x\epsilon} = 0 \Rightarrow r_{x\epsilon} = 0.$$

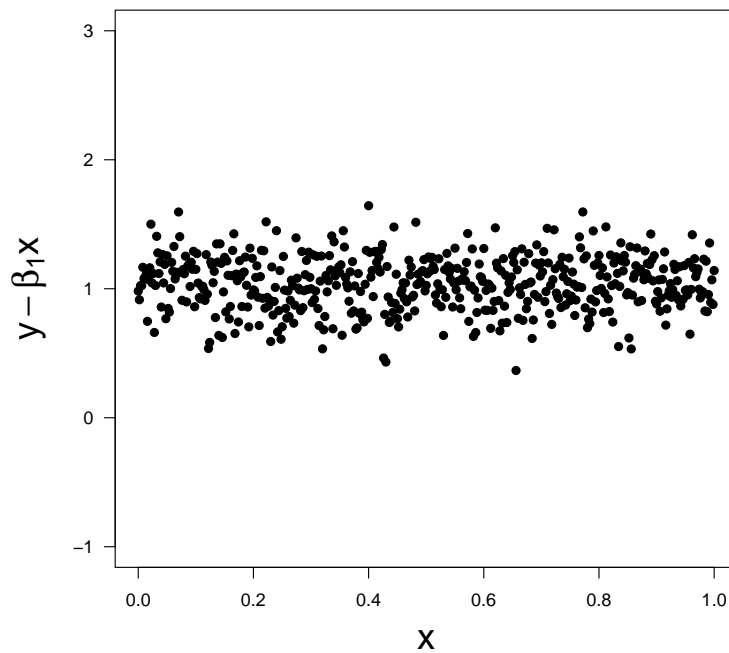
That is to say that the correlation coefficient between x_i and ϵ_i is 0.

Thus the least square prescription is to find β_0 and β_1 to make the residuals having zero mean and zero correlation with x_i . Let's illustrate how this can be done graphically.

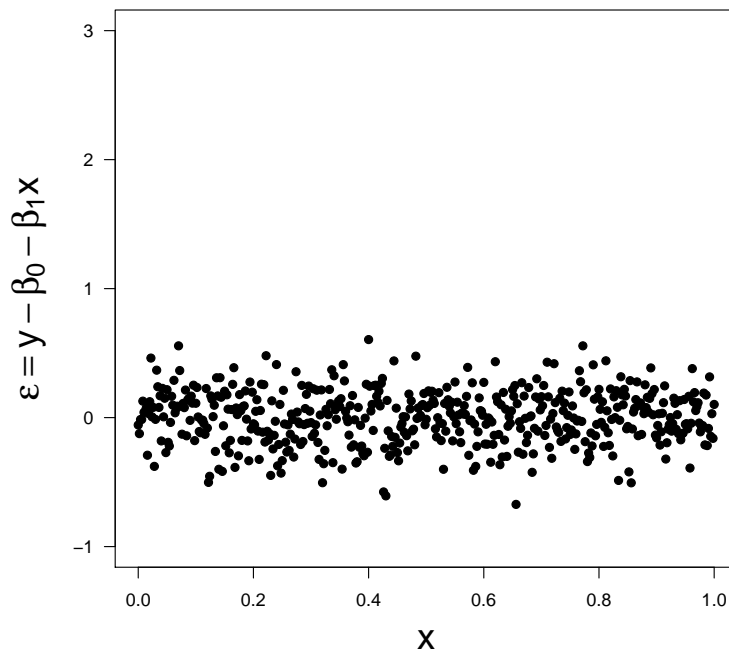
Suppose we have the following data set of x_i and y_i :



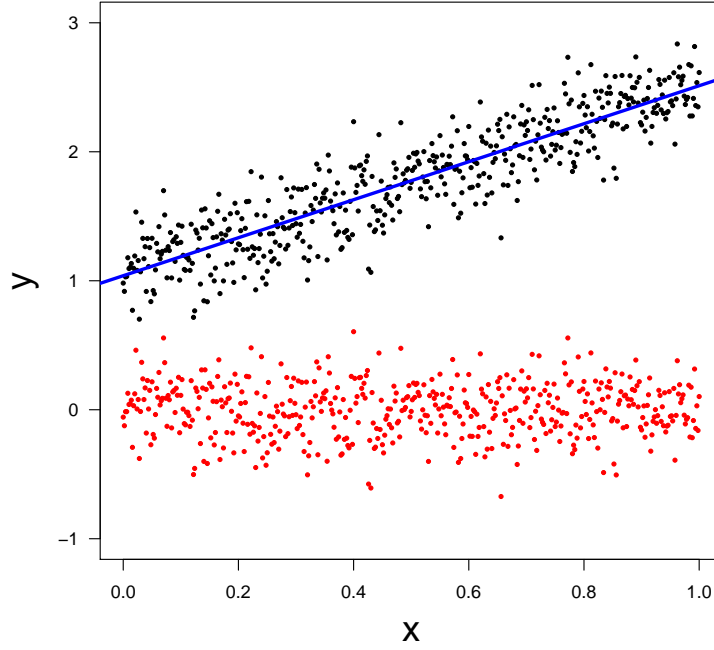
We want to fit a straight line $y = \beta_0 + \beta_1 x$ to the data points. We know the line has to make the residuals having zero mean and zero correlation with x_i . First, let's imagine pick various values of β_1 and see what the plot $y - \beta_1 x$ versus x look like. In general, we will have the plot similar to the one above but with different overall slope. If we pick the right value of β_1 , we will see that the residual $y - \beta_1 x$ vs x is flat:



We know have achieved the first goal: by choosing the right value of β_1 , the residual $y - \beta_1 x$ has zero correlation with x . However, the residuals do not have zero mean. So next we want to choose β_0 so that $y - \beta_0 - \beta_1 x$ has zero mean. When the right value of β_0 is chosen, the residuals become this:



We have accomplished our goal. The residuals now have zero mean and zero correlation with x_i . The resulting regression line $y = \beta_0 + \beta_1 x$ is the straight line that is best fit to the data points. The graph below shows the regression line (blue), data points (black) and the residuals.



We have outlined the idea of linear regression. Let's now figure out how to do this in math.

Taking the average of equation (2) and using $\bar{\epsilon} = 0$, we obtain

$$\bar{y} = \beta_0 + \beta_1 \bar{x} \quad (9)$$

This means that the regression line passes through the point of average (\bar{x}, \bar{y}) . We can express β_0 in terms of β_1 as

$$\beta_0 = \bar{y} - \beta_1 \bar{x}. \quad (10)$$

What we need to do next is to calculate β_1 . There are many ways to do it. We will use a method that will prove useful for the multiple regression. This method use the concept of a vector, which we introduce next.

1.2 Vectors

An n -dimensional vector \mathbf{V} contains n real numbers v_1, v_2, \dots, v_n , often written in the form

$$\mathbf{V} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix}$$

The numbers v_1, v_2, \dots, v_n are called the *components* of \mathbf{V} . Here we adopt the convention to write vector variables in boldface to distinguish them from numbers. It is useful to define a special vector \mathbf{X}_0 whose components are all equal to 1. That is,

$$\mathbf{X}_0 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad (11)$$

The scalar product of two vectors \mathbf{U} and \mathbf{V} is defined as

$$\mathbf{U} \cdot \mathbf{V} = \mathbf{V} \cdot \mathbf{U} = \sum_{i=1}^n u_i v_i,$$

where v_i and u_i are components of \mathbf{V} and \mathbf{U} . Note that the result of the scalar product of two vectors is a number, not a vector. Two vectors \mathbf{U} and \mathbf{V} are said to be *orthogonal* if and only if $\mathbf{U} \cdot \mathbf{V} = 0$.

It follows from the definition that

$$\mathbf{V} \cdot \mathbf{V} = \sum_{i=1}^n v_i^2 \geq 0,$$

with $\mathbf{V} \cdot \mathbf{V} = 0$ if and only if all components of \mathbf{V} are 0. It follows from the Cauchy-Schwarz inequality that

$$(\mathbf{U} \cdot \mathbf{V})^2 \leq (\mathbf{U} \cdot \mathbf{U})(\mathbf{V} \cdot \mathbf{V})$$

with the equality holds if and only if $\mathbf{U} = k\mathbf{V}$ for some constant k .

Now that we have introduced the basic concept of vectors, we are now ready to express the regression equations in vector form.

1.3 Regression Equations in Vector Form

Equation (2) can be written as

$$\mathbf{Y} = \hat{\mathbf{Y}} + \boldsymbol{\epsilon} = \beta_0 \mathbf{X}_0 + \beta_1 \mathbf{X} + \boldsymbol{\epsilon}. \quad (12)$$

The sum of square of residuals SSE in (3) can be written as

$$SSE = \boldsymbol{\epsilon} \cdot \boldsymbol{\epsilon} = (\mathbf{Y} - \hat{\mathbf{Y}}) \cdot (\mathbf{Y} - \hat{\mathbf{Y}}) = (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}) \cdot (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}) \quad (13)$$

To minimize SSE , we set the derivatives of SSE with respect to β_0 and β_1 to 0. The resulting equations are

$$\mathbf{X}_0 \cdot (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}) = 0 \quad , \quad \mathbf{X} \cdot (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}) = 0$$

or

$$\mathbf{X}_0 \cdot \boldsymbol{\epsilon} = \mathbf{X} \cdot \boldsymbol{\epsilon} = 0. \quad (14)$$

This is equation (4) written in vector form. We see that the least square prescription can be interpreted as finding β_0 and β_1 to make the residual vector $\boldsymbol{\epsilon}$ orthogonal to both \mathbf{X}_0 and \mathbf{X} .

The mean and standard deviation in (5) can be expressed as

$$\bar{u} = \frac{1}{n} \mathbf{X}_0 \cdot \mathbf{U} \quad , \quad SD_u = \sqrt{\frac{1}{n} (\mathbf{U} - \bar{u} \mathbf{X}_0) \cdot (\mathbf{U} - \bar{u} \mathbf{X}_0)} \quad (15)$$

and the Z score in (6) becomes

$$\mathbf{Z}_u = \frac{\mathbf{U} - \bar{u} \mathbf{X}_0}{SD_u} \quad (16)$$

It follows that

$$\mathbf{Z}_u \cdot \mathbf{Z}_u = n \quad (17)$$

The correlation r_{xy} in (7) is proportional to the scalar product of \mathbf{Z}_x and \mathbf{Z}_y :

$$r_{xy} = \frac{1}{n} \mathbf{Z}_x \cdot \mathbf{Z}_y \quad (18)$$

The equation $r_{x\epsilon} = 0$ is equivalent to

$$\mathbf{Z}_x \cdot \boldsymbol{\epsilon} = 0 \quad (19)$$

We see that the equations in vector form are simpler and more elegant. We are now ready to solve the regression equations (12) and (14) for β_0 and β_1 .

1.4 Regression Coefficients

Start with equation (12). Taking the scalar product of (12) with \mathbf{X}_0 , using equations (15), (14) and $\mathbf{X}_0 \cdot \mathbf{X}_0 = n$, we obtain

$$n\bar{y} = n\beta_0 + n\beta_1\bar{x} \Rightarrow \bar{y} = \beta_0 + \beta_1\bar{x}, \quad (20)$$

which is equation (9) we derived earlier. Multiply both sides of the above equation by \mathbf{X}_0 gives

$$\bar{y}\mathbf{X}_0 = \beta_0\mathbf{X}_0 + \beta_1\bar{x}\mathbf{X}_0$$

Subtracting the above equation from (12) gives

$$\mathbf{Y} - \bar{y}\mathbf{X}_0 = \beta_1(\mathbf{X} - \bar{x}\mathbf{X}_0) + \boldsymbol{\epsilon}$$

It follows from equation (16) that

$$\mathbf{Y} - \bar{y}\mathbf{X}_0 = SD_y\mathbf{Z}_y \quad \text{and} \quad \mathbf{X} - \bar{x}\mathbf{X}_0 = SD_x\mathbf{Z}_x.$$

Hence we have

$$SD_y\mathbf{Z}_y = \beta_1 SD_x\mathbf{Z}_x + \boldsymbol{\epsilon}$$

Taking the scalar product with \mathbf{Z}_x and using $\mathbf{Z}_x \cdot \boldsymbol{\epsilon} = 0$, we obtain

$$SD_y\mathbf{Z}_x \cdot \mathbf{Z}_y = \beta_1 SD_x\mathbf{Z}_x \cdot \mathbf{Z}_x$$

It follows from equation (18) and (17) that $\mathbf{Z}_x \cdot \mathbf{Z}_y = nr_{xy}$ and $\mathbf{Z}_x \cdot \mathbf{Z}_x = n$ and so the above equation reduces to

$$SD_y r_{xy} = \beta_1 SD_x \Rightarrow \beta_1 = r_{xy} \frac{SD_y}{SD_x}.$$

Combining the above equation with (20), we finally solve β_0 and β_1 :

$$\boxed{\beta_1 = r_{xy} \frac{SD_y}{SD_x} \quad , \quad \beta_0 = \bar{y} - \beta_1\bar{x}} \quad (21)$$

The regression line is the straight line with slope $r_{xy}SD_y/SD_x$ passing through the point (\bar{x}, \bar{y}) .

1.5 Root Mean Square Error

Having found a straight line that best fit the data points, we next want to know how good the fit is. One way to characterize the “goodness” of the fit is to calculate the mean square error RMSE, defined as

$$RMSE = \sqrt{\frac{SSE}{n}} = \sqrt{\frac{\boldsymbol{\epsilon} \cdot \boldsymbol{\epsilon}}{n}} \quad (22)$$

A simple formula exists relating SSE , SD_y and r_{xy} . To derive it, we introduce two quantities SST and SSM . The total sum square SST defined as

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = (\mathbf{Y} - \bar{y}\mathbf{X}_0) \cdot (\mathbf{Y} - \bar{y}\mathbf{X}_0) \quad (23)$$

It follows from the definition of standard deviation [see equation (5) or (15)] that

$$SST = nSD_y^2 \quad (24)$$

The sum square predicted by the linear model SSM is defined as

$$SSM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) \quad (25)$$

Using $\hat{\mathbf{Y}} = \beta_0\mathbf{X}_0 + \beta_1\mathbf{X}$ and $\beta_0 = \bar{y} - \beta_1\bar{x}$, we have

$$\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0 = (\bar{y} - \beta_1\bar{x})\mathbf{X}_0 + \beta_1\mathbf{X} - \bar{y}\mathbf{X}_0 = \beta_1(\mathbf{X} - \bar{x}\mathbf{X}_0) = \beta_1SD_x\mathbf{Z}_x, \quad (26)$$

where we have used equation (16) to obtain the last equality. Combining (25), (26), (17) (21) and (24), we obtain

$$SSM = n\beta_1^2SD_x^2 = n\left(r_{xy}\frac{SD_y}{SD_x}\right)^2 SD_x^2 = nr_{xy}^2SD_y^2 = r_{xy}^2SST \quad (27)$$

Next we want to prove an identity that $SST = SSM + SSE$. To see that, we start with the identity

$$\mathbf{Y} - \bar{y}\mathbf{X}_0 = (\mathbf{Y} - \hat{\mathbf{Y}}) + (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) = \boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0)$$

We then “square” both sides by taking the scalar product with itself:

$$(\mathbf{Y} - \bar{y}\mathbf{X}_0) \cdot (\mathbf{Y} - \bar{y}\mathbf{X}_0) = [\boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0)] \cdot [\boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0)]$$

The left hand side is SST , the right hand side is the sum of SSE , SSM and a cross term:

$$\begin{aligned} SST &= \boldsymbol{\epsilon} \cdot \boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) + 2\boldsymbol{\epsilon} \cdot (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) \\ &= SSE + SSM + 2\boldsymbol{\epsilon} \cdot (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) \end{aligned}$$

It follows from (26) and $\mathbf{Z}_x \cdot \boldsymbol{\epsilon} = 0$ that the cross term vanishes:

$$\boldsymbol{\epsilon} \cdot (\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0) = \beta_1SD_x\mathbf{Z}_x \cdot \boldsymbol{\epsilon} = 0$$

Another way of seeing this is to note that $\boldsymbol{\epsilon}$ is orthogonal to both \mathbf{X}_0 and \mathbf{X} , as required by the least square prescription (14). So $\boldsymbol{\epsilon}$ is orthogonal to any linear combination of \mathbf{X}_0 and \mathbf{X} . Since $\hat{\mathbf{Y}} - \bar{y}\mathbf{X}_0$ is a linear combination of \mathbf{X}_0 and \mathbf{X} , it is orthogonal to $\boldsymbol{\epsilon}$.

We have just proved that

$$\boxed{SST = SSM + SSE} \quad (28)$$

We have calculated above that $SSM = r_{xy}^2SST$ and $SST = nSD_y^2$. Using the identity $SST = SSM + SSE$, we obtain

$$SSE = SST - SSM = (1 - r_{xy}^2)SST = n(1 - r_{xy}^2)SD_y^2.$$

Combining this equation with the definition of $RMSE$ in (22), we find

$$\boxed{RMSE = \sqrt{1 - r_{xy}^2} SD_y} \quad (29)$$

2 Multiple Regression

Now we want to generalize the results of simple regression to multiple regression. We will first consider the case with two variables in Section 2.1, because simple analytic expressions exist in this case and the derivation is relatively straightforward. Then we will turn to the more general case with more than two variables in Section 2.2. Finally, we will derive a general result for the RMSE calculation in Section 2.4.

2.1 Two Variables

Suppose we want to fit the data points $\{y_i\}$ with two variables $\{x_{1i}\}$ and $\{x_{2i}\}$ by a linear model

$$y_i = \hat{y}_i + \epsilon_i, \quad i = 1, 2, \dots, n$$

with

$$\hat{y}_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}.$$

The least square prescription is again to find β_0 , β_1 and β_2 to minimize the sum square of the residuals:

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

As in the case of simple regression, it is more convenient to rewrite the above equations in vector form as follows:

$$\mathbf{Y} = \hat{\mathbf{Y}} + \boldsymbol{\epsilon} = \beta_0 \mathbf{X}_0 + \beta_1 \mathbf{X}_1 + \beta_2 \mathbf{X}_2 + \boldsymbol{\epsilon} \quad (30)$$

$$\hat{\mathbf{Y}} = \beta_0 \mathbf{X}_0 + \beta_1 \mathbf{X}_1 + \beta_2 \mathbf{X}_2 \quad (31)$$

$$SSE = \boldsymbol{\epsilon} \cdot \boldsymbol{\epsilon} = (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}_1 - \beta_2 \mathbf{X}_2) \cdot (\mathbf{Y} - \beta_0 \mathbf{X}_0 - \beta_1 \mathbf{X}_1 - \beta_2 \mathbf{X}_2) \quad (32)$$

To employ the least square prescription, we set the derivatives of SSE with respect to β_0 , β_1 and β_2 to 0. The resulting equations are

$$\mathbf{X}_0 \cdot \boldsymbol{\epsilon} = \mathbf{X}_1 \cdot \boldsymbol{\epsilon} = \mathbf{X}_2 \cdot \boldsymbol{\epsilon} = 0. \quad (33)$$

That is to say that $\boldsymbol{\epsilon}$ is orthogonal to \mathbf{X}_0 , \mathbf{X}_1 and \mathbf{X}_2 . In other words, $\hat{\mathbf{Y}}$ is the vector \mathbf{Y} projected onto the vector space spanned by \mathbf{X}_0 , \mathbf{X}_1 and \mathbf{X}_2 .

For convenience, we denote \mathbf{Z}_1 and \mathbf{Z}_2 as the Z -score vectors associated with \mathbf{X}_1 and \mathbf{X}_2 , respectively. That is,

$$\mathbf{Z}_1 = \frac{\mathbf{X}_1 - \bar{x}_1 \mathbf{X}_0}{SD_1}, \quad \mathbf{Z}_2 = \frac{\mathbf{X}_2 - \bar{x}_2 \mathbf{X}_0}{SD_2}, \quad (34)$$

where SD_1 and SD_2 are the standard deviation of $\{x_{1i}\}$ and $\{x_{2i}\}$, respectively. We see that \mathbf{Z}_1 is a linear combination of \mathbf{X}_0 and \mathbf{X}_1 ; \mathbf{Z}_2 is a linear combination of \mathbf{X}_0 and \mathbf{X}_2 . Since $\boldsymbol{\epsilon}$ is orthogonal to \mathbf{X}_0 , \mathbf{X}_1 and \mathbf{X}_2 , it is orthogonal to \mathbf{Z}_1 and \mathbf{Z}_2 as well:

$$\mathbf{Z}_1 \cdot \boldsymbol{\epsilon} = \mathbf{Z}_2 \cdot \boldsymbol{\epsilon} = 0. \quad (35)$$

Recall that the mean of a vector \mathbf{U} is $\bar{u} = \mathbf{X}_0 \cdot \mathbf{U}/n$, and the correlation between \mathbf{U} and \mathbf{V} is $r_{uv} = \mathbf{Z}_u \cdot \mathbf{Z}_v/n$. Thus the orthogonality conditions of $\boldsymbol{\epsilon}$ mean that (1) $\boldsymbol{\epsilon}$ has zero mean, $\bar{\epsilon} = 0$, and (2) $\boldsymbol{\epsilon}$ is uncorrelated with both \mathbf{X}_1 and \mathbf{X}_2 : $r_{\epsilon 1} = r_{\epsilon 2} = 0$.

To solve the regression equations (33), we find take the scalar product of equation (30) with \mathbf{X}_0 , resulting in the equation

$$\bar{y} = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 \quad (36)$$

This means that the regression line passes through the average point $(\bar{x}_1, \bar{x}_2, \bar{y})$. Multiplying equation (36) by \mathbf{X}_0 yields

$$\bar{y} \mathbf{X}_0 = \beta_0 \mathbf{X}_0 + \beta_1 \bar{x}_1 \mathbf{X}_0 + \beta_2 \bar{x}_2 \mathbf{X}_0$$

Subtracting the above equation from (30) gives

$$\mathbf{Y} - \bar{y} \mathbf{X}_0 = \beta_1 (\mathbf{X}_1 - \bar{x}_1 \mathbf{X}_0) + \beta_2 (\mathbf{X}_2 - \bar{x}_2 \mathbf{X}_0) + \boldsymbol{\epsilon}$$

Using the definition of the Z score vector, we can write the above equation as

$$SD_y \mathbf{Z}_y = \beta_1 SD_1 \mathbf{Z}_1 + \beta_2 SD_2 \mathbf{Z}_2 + \boldsymbol{\epsilon} \quad (37)$$

Taking the scalar product of the above equation with \mathbf{Z}_1 gives

$$SD_y r_{y1} = \beta_1 SD_1 + \beta_2 SD_2 r_{12}, \quad (38)$$

where $r_{y1} = \mathbf{Z}_y \cdot \mathbf{Z}_1/n$ is the correlation between \mathbf{Y} and \mathbf{X}_1 , $r_{12} = \mathbf{Z}_1 \cdot \mathbf{Z}_2/n$ is the correlation between \mathbf{X}_1 and \mathbf{X}_2 . Equation (38) can be written as

$$\beta_1 = r_{y1} \frac{SD_y}{SD_1} - \beta_2 \left(r_{12} \frac{SD_2}{SD_1} \right) = \beta_{y1} - \beta_2 \beta_{21}, \quad (39)$$

where $\beta_{y1} = r_{y1} SD_y / SD_1$ is the slope in the simple regression for predicting \mathbf{Y} from \mathbf{X}_1 , and $\beta_{21} = r_{12} SD_2 / SD_1$ is the slope in the simple regression for predicting \mathbf{X}_2 from \mathbf{X}_1 .

Before deriving the final expression for β_1 , let's point out two things. First, if $r_{12} = 0$ (or equivalently $\mathbf{X}_1 \cdot \mathbf{X}_2 = 0$) then $\beta_1 = \beta_{y1}$. That is to say that if \mathbf{X}_1 and \mathbf{X}_2 are uncorrelated (orthogonal), the slope β_1 in the multiple regression is exactly the same as the slope β_{y1} in the simple regression for prediction \mathbf{Y} from \mathbf{X}_1 . The adding of the \mathbf{X}_2 does not change the slope. Second, if $r_{12} > 0$ and $\beta_2 > 0$, then $\beta_1 < \beta_{y1}$. The slope β_1 decreases if \mathbf{X}_1 and \mathbf{X}_2 are positively correlated and if the slope of \mathbf{X}_2 in the multiple regression is positive.

Let's go back to equation (39). If $r_{12} \neq 0$, the equation of β_1 involves β_2 . So β_1 and β_2 has to be solved together. Since \mathbf{X}_1 and \mathbf{X}_2 are symmetrical, the equation for β_2 can be obtained by simply exchanging the index between 1 and 2 of the β_1 equation:

$$\beta_2 = \beta_{y2} - \beta_1 \beta_{12}, \quad (40)$$

where

$$\beta_{y2} = r_{y2} \frac{SD_y}{SD_2} \quad \text{and} \quad \beta_{12} = r_{12} \frac{SD_1}{SD_2}. \quad (41)$$

Substituting β_2 from (40) into (39), we obtain

$$\beta_1 = \beta_{y1} - (\beta_{y2} - \beta_1 \beta_{12}) \beta_{21} = \beta_{y1} - \beta_{y2} \beta_{21} + \beta_1 \beta_{12} \beta_{21} \quad (42)$$

Note that the product

$$\beta_{12} \beta_{21} = \left(r_{12} \frac{SD_1}{SD_2} \right) \left(r_{12} \frac{SD_2}{SD_1} \right) = r_{12}^2.$$

Thus equation (42) becomes

$$\beta_1 = \beta_{y1} - \beta_{y2} \beta_{21} + r_{12}^2 \beta_1 \Rightarrow (1 - r_{12}^2) \beta_1 = \beta_{y1} - \beta_{y2} \beta_{21}$$

or

$$\beta_1 = \frac{\beta_{y1} - \beta_{y2} \beta_{21}}{1 - r_{12}^2}$$

Using the expressions for β_{y1} , β_{y2} and β_{12} , we find

$$\beta_1 = b_1 \frac{SD_y}{SD_1}, \quad (43)$$

where

$$b_1 = \frac{r_{y1} - r_{y2} r_{12}}{1 - r_{12}^2}$$

may be interpreted as the adjusted correlation between \mathbf{Y} and \mathbf{X}_1 taking into account the presence of \mathbf{X}_2 . The expression for β_2 is obtained by exchanging the index between 1 and 2:

$$\beta_2 = b_2 \frac{SD_y}{SD_2}, \quad b_2 = \frac{r_{y2} - r_{y1} r_{12}}{1 - r_{12}^2}.$$

Gathering all the results, we conclude that the regression coefficients are given by

$$\boxed{\beta_1 = b_1 \frac{SD_y}{SD_1} \quad , \quad \beta_2 = b_2 \frac{SD_y}{SD_2} \quad , \quad \beta_0 = \bar{y} - \beta_1 \bar{x}_1 - \beta_2 \bar{x}_2} \quad (44)$$

with

$$\boxed{b_1 = \frac{r_{y1} - r_{y2}r_{12}}{1 - r_{12}^2} \quad , \quad b_2 = \frac{r_{y2} - r_{y1}r_{12}}{1 - r_{12}^2}} \quad (45)$$

In Section 1.5, we see that the RMSE in simple regression is related to SD_y by $RMSE = \sqrt{1 - r_{xy}^2} SD_y$. In multiple regression, we will show in Section 2.4 that the formula is generalized to

$$RMSE = \sqrt{1 - R^2} SD_y,$$

where R is the correlation between \mathbf{Y} and $\hat{\mathbf{Y}}$:

$$R = \frac{\mathbf{Z}_y \cdot \mathbf{Z}_{\hat{y}}}{n}.$$

We will defer the calculation of R to Section 2.3. In the case of multiple regression with two variables considered here, R is given by

$$\boxed{R = \sqrt{\frac{r_{y1}^2 + r_{y2}^2 - 2r_{y1}r_{y2}r_{12}}{1 - r_{12}^2}}} \quad (46)$$

2.2 More Than Two Variables

Suppose we now want to fit $\{y_i\}$ with p variables $\{x_{1i}, x_{2i}, \dots, x_{pi}\}$. The regression equation has $p+1$ parameters $\beta_0, \beta_1, \beta_2, \dots, \beta_p$. It can be written in the vector form as

$$\mathbf{Y} = \hat{\mathbf{Y}} + \boldsymbol{\epsilon} = \sum_{j=0}^p \beta_j \mathbf{X}_j + \boldsymbol{\epsilon} \quad (47)$$

$$\hat{\mathbf{Y}} = \sum_{j=0}^p \beta_j \mathbf{X}_j \quad (48)$$

$$SSE = \boldsymbol{\epsilon} \cdot \boldsymbol{\epsilon} = \left(\mathbf{Y} - \sum_{j=0}^p \beta_j \mathbf{X}_j \right) \cdot \left(\mathbf{Y} - \sum_{j=0}^p \beta_j \mathbf{X}_j \right) \quad (49)$$

To minimize SSE , we set the derivatives of SSE with respect to β_j ($j = 0, 1, \dots, p$) to 0. The resulting equations can be written as

$$\mathbf{X}_j \cdot \boldsymbol{\epsilon} = 0 \quad , \quad j = 0, 1, \dots, p. \quad (50)$$

This means that $\boldsymbol{\epsilon}$ is orthogonal all $\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_p$, or equivalently, the mean of $\boldsymbol{\epsilon}$ is 0 and $\boldsymbol{\epsilon}$ is uncorrelated with any of the variables we are trying to fit.

To find the regression coefficients, we follow the same procedures as before. First, take the scalar product of equation (47) with \mathbf{X}_0 . The result is

$$\bar{y} = \beta_0 + \sum_{j=1}^p \beta_j \bar{x}_j \quad (51)$$

As before, this means that the regression line passes through the point of average. Next, we multiple the above equation by \mathbf{X}_0 :

$$\bar{y} \mathbf{X}_0 = \beta_0 \mathbf{X}_0 + \sum_{j=1}^p \beta_j \bar{x}_j \mathbf{X}_0$$

and then subtract it from equation (47):

$$\mathbf{Y} - \bar{y}\mathbf{X}_0 = \sum_{j=1}^p \beta_j (\mathbf{X}_j - \bar{x}_j \mathbf{X}_0) + \boldsymbol{\epsilon}.$$

Using the definition of the Z score vector (16), we can write the above equation as

$$SD_y \mathbf{Z}_y = \sum_{j=1}^p \beta_j SD_j \mathbf{Z}_j + \boldsymbol{\epsilon}.$$

Dividing both sides by SD_y results in

$$\mathbf{Z}_y = \sum_{j=1}^p b_j \mathbf{Z}_j + \frac{\boldsymbol{\epsilon}}{SD_y}, \quad (52)$$

where

$$b_j = \beta_j \frac{SD_j}{SD_y}. \quad (53)$$

Taking the scalar product of equation (52) with \mathbf{Z}_i ($i = 0, 1, 2, \dots, p$), we obtain

$$\sum_{j=1}^p r_{ij} b_j = r_{yi} \quad , \quad i = 1, 2, \dots, p. \quad (54)$$

This is a system of linear equations for b_1, b_2, \dots, b_p . Written them out, they look like

$$\begin{aligned} b_1 + r_{12}b_2 + r_{13}b_3 + \dots + r_{1p}b_p &= r_{y1} \\ b_1r_{21} + b_2 + r_{23}b_3 + \dots + r_{2p}b_p &= r_{y2} \\ b_1r_{31} + b_2r_{32} + b_3 + \dots + r_{3p}b_p &= r_{y3} \\ &\vdots \\ b_1r_{p1} + b_2r_{p2} + b_3r_{p3} + \dots + b_p &= r_{yp} \end{aligned}$$

If all of the variables \mathbf{X}_j are uncorrelated (i.e. $r_{ij} = 0$ if $i \neq j$), the solution is $b_j = r_{yj}$ and the slopes are $\beta_j = r_{yj}SD_y/SD_j$. This is exactly the same as the slopes in simple regression for predicting \mathbf{Y} from \mathbf{X}_j .

There are no simple analytic expressions for b_j in general, but there are several well-known procedures to obtain the solution by successive algebraic operations. We will mention one in the Appendix for interested readers. Suppose all the b 's have been solved using one of those procedures, the slopes are given by equation (53) as

$$\beta_j = b_j \frac{SD_j}{SD_y} \quad , \quad j = 1, 2, \dots, p \quad (55)$$

and the intercept β_0 is given by equation (51) as

$$\beta_0 = \bar{y} - \sum_{j=1}^p \beta_j \bar{x}_j. \quad (56)$$

Finally, we should mention that the term “linear” in linear regression refers to be model being linear in the fitting parameters $\beta_0, \beta_1, \dots, \beta_p$. For example, we can fit y_i by the model

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 \sqrt{x_i} + \beta_3 x_i^2 + \beta_4 \ln x_i + \beta_5 \frac{x_i^3}{1 + 2x_i} \quad (57)$$

using the technique of multiple linear regression since the model is linear in $\beta_0, \beta_1, \dots, \beta_5$. We simply label

$$x_{1i} = x_i \quad , \quad x_{2i} = \sqrt{x_i} \quad , \quad x_{3i} = x_i^2 \quad , \quad x_{4i} = \ln x_i \quad , \quad x_{5i} = \frac{x_i^3}{1 + 2x_i}$$

and equation (57) can be written in vector form as

$$\mathbf{Y} = \beta_0 \mathbf{X}_0 + \beta_1 \mathbf{X}_1 + \beta_2 \mathbf{X}_2 + \beta_3 \mathbf{X}_3 + \beta_4 \mathbf{X}_4 + \beta_5 \mathbf{X}_5 ,$$

which is equation (47) with $p = 5$. The key is to note that the least square prescription is to minimize SSE by varying the parameters β 's, not the independent variables x 's.

2.3 Correlation Between \mathbf{Y} and $\hat{\mathbf{Y}}$

To generalize the $RMSE$ expression in Section 1.5 for multiple regression, we first consider the quantity R defined as the correlation between \mathbf{Y} and $\hat{\mathbf{Y}}$:

$$R = \frac{\mathbf{Z}_y \cdot \mathbf{Z}_{\hat{y}}}{n} . \quad (58)$$

We first calculate the average of $\hat{\mathbf{Y}}$:

$$\bar{\hat{y}} = \frac{1}{n} \hat{\mathbf{Y}} \cdot \mathbf{X}_0 = \frac{1}{n} \left(\beta_0 \mathbf{X}_0 + \sum_{j=1}^p \beta_j \mathbf{X}_j \right) \cdot \mathbf{X}_0 = \beta_0 + \sum_{j=1}^p \beta_j \bar{x}_j = \bar{y},$$

where we have used equation (48) for $\hat{\mathbf{Y}}$ and (51) for \bar{y} . From the definition of the Z score vector (16) and $\bar{\hat{y}} = \bar{y}$, we can write (58) as

$$\begin{aligned} R &= \frac{(\mathbf{Y} - \bar{y} \mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)}{n SD_y SD_{\hat{y}}} \\ &= \frac{(\hat{\mathbf{Y}} + \boldsymbol{\epsilon} - \bar{y} \mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)}{n SD_y SD_{\hat{y}}} \\ &= \frac{(\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)}{n SD_y SD_{\hat{y}}} \\ &= \frac{SD_{\hat{y}}}{SD_y} \end{aligned} \quad (59)$$

where we have used $\boldsymbol{\epsilon} \cdot \mathbf{X}_0 = 0$ and $\boldsymbol{\epsilon} \cdot \hat{\mathbf{Y}} = 0$ (since $\hat{\mathbf{Y}}$ is a linear combination of $\mathbf{X}_0, \dots, \mathbf{X}_p$ and all are orthogonal to $\boldsymbol{\epsilon}$). We have also used the definition of the standard deviation (15) to obtain the last line. Hence we have

$$R^2 = \frac{SD_{\hat{y}}^2}{SD_y^2} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} , \quad (60)$$

which is interpreted as the fraction of the variance of \mathbf{Y} explained by the linear model.

To compute R , we use equation (48) for $\hat{\mathbf{Y}}$ and (56) for β_0 , and write

$$\begin{aligned} SD_{\hat{y}}^2 &= \frac{1}{n} (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) \\ &= \frac{1}{n} \left[(\beta_0 - \bar{y}) \mathbf{X}_0 + \sum_{i=1}^p \beta_i \mathbf{X}_i \right] \cdot \left[(\beta_0 - \bar{y}) \mathbf{X}_0 + \sum_{j=1}^p \beta_j \mathbf{X}_j \right] \\ &= \frac{1}{n} \left[\sum_{i=1}^p \beta_i (\mathbf{X}_i - \bar{x}_i \mathbf{X}_0) \right] \cdot \left[\sum_{j=1}^p \beta_j (\mathbf{X}_j - \bar{x}_j \mathbf{X}_0) \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \left(\sum_{i=1}^p SD_i \beta_i \mathbf{Z}_i \right) \cdot \left(\sum_{j=1}^p SD_j \beta_j \mathbf{Z}_j \right) \\
&= \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^p SD_i SD_j \beta_i \beta_j \mathbf{Z}_i \cdot \mathbf{Z}_j \\
&= \sum_{i=1}^p \sum_{j=1}^p SD_i SD_j \beta_i \beta_j r_{ij}
\end{aligned}$$

$$R^2 = \frac{SD_y^2}{SD_y^2} = \sum_{i=1}^p \sum_{j=1}^p b_i b_j r_{ij},$$

where we have used the definition of b_j in equation (53). The sum can be simplified by noting that b_j satisfy equation (54), and thus we obtain

$$\boxed{R = \sqrt{\sum_{i=1}^p b_i r_{yi}}} \quad (61)$$

When the solution of b_i is obtained, R can be calculated using the above equation. In the case of multiple regression with two variables ($p = 2$), b_1 and b_2 are given by equation (45). Plugging them into equation (61), we obtain equation (46).

2.4 Root Mean Square Error

As in Section 1.5, we define SST and SSM as

$$\begin{aligned}
SST &= \sum_{i=1}^n (y_i - \bar{y})^2 = (\mathbf{Y} - \bar{y} \mathbf{X}_0) \cdot (\mathbf{Y} - \bar{y} \mathbf{X}_0) \\
SSM &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)
\end{aligned}$$

It follows from the definition of standard deviation that

$$SST = n SD_y^2 \quad (62)$$

and it follows from (60) that

$$SSM = R^2 SST. \quad (63)$$

The identity $SST = SSM + SSE$ still holds in multiple linear regression. The proof is almost exactly the same as in Section 1.5.

We start with the identity

$$\mathbf{Y} - \bar{y} \mathbf{X}_0 = (\mathbf{Y} - \hat{\mathbf{Y}}) + (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) = \boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0).$$

Take the scalar product with itself:

$$\begin{aligned}
(\mathbf{Y} - \bar{y} \mathbf{X}_0) \cdot (\mathbf{Y} - \bar{y} \mathbf{X}_0) &= [\boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)] \cdot [\boldsymbol{\epsilon} + (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)] \\
SST &= SSE + SSM + 2\boldsymbol{\epsilon} \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0)
\end{aligned}$$

Since $\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0$ is a linear combination of $\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_p$ and $\boldsymbol{\epsilon}$ is orthogonal to all these vectors, $\boldsymbol{\epsilon} \cdot (\hat{\mathbf{Y}} - \bar{y} \mathbf{X}_0) = 0$ and so we have

$$\boxed{SST = SSM + SSE} \quad (64)$$

Combining equation (62), (63) and (64), we have $SSE = n(1 - R^2)SD_y^2$. It follows from the definition of $RMSE = \sqrt{SSE/n}$ that

$$\boxed{RMSE = \sqrt{1 - R^2} SD_y} \tag{65}$$

which is the generalization of equation (29).

Appendix: Solving the Multiple Linear Regression Equations