2. OLS Part II

The OLS residuals are orthogonal to the regressors. If the model includes an intercept, the orthogonality of the residuals and regressors gives rise to three results, which have limited practical usefulness but are good exercises for understanding the algebra of least squares.

Partitioned and partial regression is not treated as seriously as it might be in a graduate course, and is only given a cursory presentation. However, the residual maker matrix M_i is presented, and is used in to define R^2 , and in several other parts of the course.

2.1 Some basic properties of OLS

First, note that the LS residuals are "orthogonal" to the regressors -

$$X'X\mathbf{b} - X'\mathbf{y} = \mathbf{0} \qquad (\text{``normal equations''; } (\mathbf{k} \times \mathbf{1}))$$

So,

 $-X'(\boldsymbol{y}-X\boldsymbol{b})=-X'\boldsymbol{e}=0 ;$

 $X'\boldsymbol{e}=0$

or,

Orthogonality implies linear independence (but not *vice versa*). See "<u>Linearly Independent</u>, <u>Orthogonal</u>, and <u>Uncorrelated Variables</u>" (Rodgers et al., 1984), for the definition and relationship between the three terms.

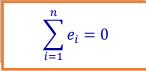
If the model includes an intercept term, then one regressor (say, the first column of X) is a unit vector.

In this case we get some further results:

1. The LS residuals sum to zero

$$X'\boldsymbol{e} = \begin{pmatrix} 1 & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ 1 & \cdots & x_{nk} \end{pmatrix}' \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ x_{1k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}$$
$$= \begin{pmatrix} \sum_i e_i \\ ? \\ ? \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

From the first element:



2. Fitted regression passes through sample mean

$$X' \mathbf{y} = X' X \mathbf{b} ,$$

or, $\begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ x_{1k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ x_{1k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} 1 & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ 1 & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix}$
So, $\begin{pmatrix} \sum_i y_i \\ ? \\ ? \end{pmatrix} = \begin{pmatrix} n & \sum_i x_{i2} & \cdots \\ ? & \cdots & ? \\ ? & \cdots & ? \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} .$

So,

From the first row of this vector equation –

$$\sum_{i} y_i = nb_1 + b_2 \sum_{i} x_{i2} + \dots + b_k \sum_{i} x_{ik}$$
$$\overline{y} = b_1 + b_2 \overline{x_2} + \dots + b_k \overline{x_k}$$

or,

3. Sample mean of the fitted y-values equals sample mean of actual y-values

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i = \mathbf{x}'_i \mathbf{b} + e_i = \widehat{y}_i + e_i$$

So,

$$\frac{1}{n}\sum_{i=1}^{n} y_i = \frac{1}{n}\sum_{i=1}^{n} \widehat{y}_i + \frac{1}{n}\sum_{i=1}^{n} e_i ,$$

 $\overline{y} = \overline{\hat{y}} + 0 = \overline{\hat{y}}$

or,

2.2 Partitioned and Partial Regression

We can solve for "blocks" of **b**. Suppose we partition our model into 2 blocks:

$$\mathbf{y} = X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

(n×1) (n×k_1)(k_1×1) (n×k_2)(k_2×1) (n×1)

We could solve for the OLS estimator of β_1 only:

$$\boldsymbol{b_1} = (X_1'X_1)^{-1}X_1'\boldsymbol{y} - (X_1'X_1)^{-1}X_1'X_2\boldsymbol{b_2}$$

When does $b_1 = (X_1'X_1)^{-1}X_1'y$?

Above we have a solution for b_1 in terms of b_2 . We can substitute in a similar solution for b_2 to obtain a solution for b_1 that is a function of only the X and y data.

Define:

$$M_2 = (I - X_2(X_2'X_2)^{-1}X_2')$$

- *M*₂ is an "*idempotent*" matrix
- $M_2M_2 = M_2M_2' = M_2 = M_2'M_2.$

Then, we can write:

$$\boldsymbol{b_1} = (X_1'M_2X_1)^{-1}X_1'M_2\boldsymbol{y}$$

By the symmetry of the problem, we can interchange the "1" and "2" subscripts, and get:

$$\boldsymbol{b_2} = (X_2' M_1 X_2)^{-1} X_2' M_1 \boldsymbol{y}$$

So, finally, we can write:

$$\boldsymbol{b_1} = (X_1^* X_1^*)^{-1} X_1^* y_1^* \qquad \boldsymbol{b_2} = (X_2^* X_2^*)^{-1} X_2^* y_2^*$$

where:

$$X_1^* = M_2 X_1 \; ; \; X_2^* = M_1 X_2 \; ; \; \; \boldsymbol{y}_1^* = M_2 \boldsymbol{y} \; ; \; \; \boldsymbol{y}_2^* = M_1 \boldsymbol{y}$$

These results are important for several reasons:

- Historically, for computational reasons. See the "Frisch-Waugh-Lovell Theorem".
- In certain situations b_1 and b_2 may have different properties. This is difficult to show without have separate formulas.
- Having established the M_i "residual maker" matrix, it is used frequently to derive other results.

Why is M_i called a "residual maker" matrix?

2.3 Goodness-of-Fit

- One way of measuring the "quality" of fitted regression model is by the extent to which the model "explains" the *sample variation* for *y*.
- Sample variance of y is $\frac{1}{(n-1)}\sum_{i=1}^{n}(y_i-\bar{y})^2$.
- Or, we could just use $\sum_{i=1}^{n} (y_i \overline{y})^2$ to measure *variability*.
- Our "fitted" regression model, using LS, gives us

$$y = Xb + e = \hat{y} + e$$

 $\widehat{\mathbf{y}} = X\mathbf{b} = X(X'X)^{-1}X'\mathbf{y}$

where

• Recall that *if the model includes an intercept*, then the residuals sum to zero, and $\bar{y} = \bar{y}$.

To simplify things, introduce the following matrix:

$$M^{0} = [I_{n} - \frac{1}{n}ii']$$
$$i = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \qquad ; \qquad (n \times 1)$$

where:

Note that:

- M^0 is an idempotent matrix.
- $M^0 \boldsymbol{i} = \boldsymbol{0}$.
- M^0 transforms elements of a vector into deviations from sample mean.
- $y'M^0y = y'M^0M^0y = \sum_{i=1}^n (y_i \bar{y})^2$.

Let's check the third of these results:

$$M^{0}\boldsymbol{y} = \left\{ \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix} - \begin{bmatrix} 1/n & \cdots & 1/n \\ \vdots & \ddots & \vdots \\ 1/n & \cdots & 1/n \end{bmatrix} \right\} \begin{pmatrix} y_{1} \\ \vdots \\ y_{n} \end{pmatrix}$$

$$= \begin{bmatrix} y_1 - \frac{1}{n}y_1 - \frac{1}{n}y_2 \dots - \frac{1}{n}y_n \\ \vdots \\ y_n - \frac{1}{n}y_1 - \frac{1}{n}y_2 - \dots - \frac{1}{n}y_n \end{bmatrix} = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}.$$

Returning to our "fitted" model:

$$y = Xb + e = \hat{y} + e$$

So, we have:

$$M^0 \mathbf{y} = M^0 \widehat{\mathbf{y}} + M^0 \mathbf{e} = M^0 \widehat{\mathbf{y}} + \mathbf{e}$$
.

 $[M^0 e = e$; because the residuals sum to zero.]

Then –

$$\mathbf{y}' M^0 \mathbf{y} = \mathbf{y}' M^0 M^0 \mathbf{y} = (M^0 \hat{\mathbf{y}} + \mathbf{e})' (M^0 \hat{\mathbf{y}} + \mathbf{e})$$
$$= \hat{\mathbf{y}}' M^0 \hat{\mathbf{y}} + \mathbf{e}' \mathbf{e} + 2\mathbf{e}' M^0 \hat{\mathbf{y}}$$

However,

$$\boldsymbol{e}'\boldsymbol{M}^{0}\boldsymbol{\hat{y}} = \boldsymbol{e}'\boldsymbol{M}^{0'}\boldsymbol{\hat{y}} = (\boldsymbol{M}^{0}\boldsymbol{e})'\boldsymbol{\hat{y}} = \boldsymbol{e}'\boldsymbol{\hat{y}} = \boldsymbol{e}'\boldsymbol{X}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y} = \boldsymbol{0} \; .$$

So, we have –

$$\mathbf{y}' M^0 \mathbf{y} = \mathbf{\hat{y}}' M^0 \mathbf{\hat{y}} + \mathbf{e}' \mathbf{e}$$
$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n e_i^2$$
$$\mathbf{SST} = \mathbf{SSR} + \mathbf{SSE}$$

Recall: $\overline{\hat{y}} = \overline{y}$.

This lets us define the "Coefficient of Determination" -

$$R^{2} = \left(\frac{SSR}{SST}\right) = 1 - \left(\frac{SSE}{SST}\right)$$

Note:

• The second equality in definition of R^2 holds only if model *includes an intercept*.

•
$$R^2 = \left(\frac{SSR}{SST}\right) \ge 0$$

•
$$R^2 = 1 - \left(\frac{SSE}{SST}\right) \le 1$$

- So, $0 \le R^2 \le 1$
- Interpretation of "0" and "1"?
- R^2 is *unitless*.

What happens if we add *any* regressor(s) to the model?

$$\mathbf{y} = X_1 \boldsymbol{\beta}_1 + \boldsymbol{\varepsilon} \qquad \qquad ; \qquad \qquad [1]$$

Then:

$$\boldsymbol{y} = X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2 + \boldsymbol{u} \quad ; \qquad [2]$$

(A) Applying LS to [2]:

min.
$$(\hat{\boldsymbol{u}}'\hat{\boldsymbol{u}})$$
 ; $\hat{\boldsymbol{u}} = \boldsymbol{y} - X_1\boldsymbol{b}_1 - X_2\boldsymbol{b}_2$

(B) Applying LS to [1]:

min.
$$(e'e)$$
 ; $e = y - X_1 \widehat{\beta}_1$

Problem (B) is just Problem (A), subject to restriction: $\beta_2 = 0$. Minimized value in (A) must be \leq minimized value in (B). So, $\hat{u}'\hat{u} \leq e'e$.

What does this imply?

- Adding *any* regressor(s) to the model *cannot increase* (and typically will *decrease*) the sum of squared residuals.
- So, adding *any* regressor(s) to the model *cannot decrease* (and typically will *increase*) the value of R^2 .
- Means that R^2 is not really a very interesting measure of the "quality" of the regression model, in terms of explaining sample variability of the dependent variable.
- For these reasons, we usually use the "adjusted" Coefficient of Determination.

We modify $R^2 = [1 - \frac{e'e}{y'M^0y}]$ to become:

$$\bar{R}^2 = \left[1 - \frac{e'e/(n-k)}{y'M^0y/(n-1)}\right].$$

• What are we doing here?

We're adjusting for "degrees of freedom" in numerator and denominator.

- "Degrees of freedom" = number of independent pieces of information.
- e = y Xb. We estimate k parameters from the n data-points. We have (n k) "degrees of freedom" associated with the fitted model.
- In denominator have constructed \overline{y} from sample. "Lost" one degree of freedom.
- Possible for $\overline{R}^2 < 0$ (even with intercept in the model).
- \overline{R}^2 can *increase or decrease* when we add regressors.
- When will it increase (decrease)?

In multiple regression, \overline{R}^2 will *increase* (decrease) if a variable is deleted, if and only if the associated t-statistic has *absolute value less than* (greater than) unity.

- If model *doesn't* include an intercept, then $SST \neq SSR + SSE$, and in this case no longer any guarantee that $0 \le R^2 \le 1$.
- Must be careful comparing R^2 and \overline{R}^2 values across models.

Example -

- (1) $\widehat{C}_{\iota} = 0.5 + 0.8Y_i$; $R^2 = 0.90$
- (2) $\log(\hat{C}_i) = 0.2 + 0.75Y_i$; $R^2 = 0.80$

Sample variation is in *different units*.