

# Efficiency of the OLS Estimator

For scalar parameters, one estimator is more **efficient** than another if the precision of the former is greater than that of the latter.

Let  $\check{\beta}$  and  $\hat{\beta}$  be unbiased estimators of a  $k$ -vector of parameters  $\beta$ , with covariance matrices  $\text{Var}(\check{\beta})$  and  $\text{Var}(\hat{\beta})$ , respectively.

- $\hat{\beta}$  is said to be more efficient than  $\check{\beta}$  if and only if  $\text{Var}(\hat{\beta})^{-1} - \text{Var}(\check{\beta})^{-1}$  is a non-zero positive semidefinite matrix.
- If  $A$  and  $B$  are positive definite matrices, then  $A - B$  is positive semidefinite if and only if  $B^{-1} - A^{-1}$  is positive semidefinite.
- Thus  $\hat{\beta}$  is more efficient than  $\check{\beta}$  if and only if  $\text{Var}(\check{\beta}) - \text{Var}(\hat{\beta})$  is a non-zero positive semidefinite matrix.
- If  $\hat{\beta}$  is more efficient than  $\check{\beta}$ , every element of  $\beta$ , and every linear combination of them, is estimated at least as efficiently by using  $\hat{\beta}$  as by using  $\check{\beta}$ .

Consider  $\gamma = \mathbf{w}^\top \boldsymbol{\beta}$ .

$$\begin{aligned} \text{Var}(\ddot{\gamma}) - \text{Var}(\hat{\gamma}) &= \mathbf{w}^\top \text{Var}(\ddot{\boldsymbol{\beta}}) \mathbf{w} - \mathbf{w}^\top \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{w} \\ &= \mathbf{w}^\top (\text{Var}(\ddot{\boldsymbol{\beta}}) - \text{Var}(\hat{\boldsymbol{\beta}})) \mathbf{w}. \end{aligned} \quad (1)$$

This must be either positive or zero if  $\text{Var}(\ddot{\boldsymbol{\beta}}) - \text{Var}(\hat{\boldsymbol{\beta}})$  is positive semidefinite. Thus  $\text{Var}(\hat{\gamma}) \leq \text{Var}(\ddot{\gamma})$  when  $\hat{\boldsymbol{\beta}}$  is more efficient than  $\ddot{\boldsymbol{\beta}}$ .

Let  $\hat{\boldsymbol{\beta}}$  denote the OLS estimator and  $\ddot{\boldsymbol{\beta}}$  some other linear estimator.

An estimator is **linear** if we can write it as a linear function of  $\mathbf{y}$ .  $\hat{\boldsymbol{\beta}}$  is linear, because it is equal to the matrix  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  times the vector  $\mathbf{y}$ .

For any linear estimator that is not the OLS estimator,

$$\ddot{\boldsymbol{\beta}} = \mathbf{A} \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} + \mathbf{C} \mathbf{y}, \quad (2)$$

where  $\mathbf{A}$  and  $\mathbf{C}$  are  $k \times N$  matrices that depend on  $\mathbf{X}$ .

To obtain the second equality in (2), define

$$\mathbf{C} \equiv \mathbf{A} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top. \quad (3)$$

Alternative linear estimators include **instrumental variables** and **generalized least squares**.

The **Gauss-Markov Theorem** says that  $\hat{\beta}$  is the **best linear unbiased estimator**, or **BLUE**. We are comparing it with every other unbiased estimator  $\check{\beta}$  of the form (2).

If  $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$  and  $E(\mathbf{u}\mathbf{u}^\top | \mathbf{X}) = \sigma^2 \mathbf{I}$ , then  $\text{Var}(\check{\beta}) - \text{Var}(\hat{\beta})$  is a positive semidefinite matrix.

Substituting for  $\mathbf{y}$  in (2), we find that

$$\check{\beta} = \mathbf{A}(\mathbf{X}\beta_0 + \mathbf{u}) = \mathbf{A}\mathbf{X}\beta_0 + \mathbf{A}\mathbf{u}. \quad (4)$$

Since we want  $\check{\beta}$  to be unbiased, the expectation of the rightmost expression in (4), conditional on  $\mathbf{X}$ , must be  $\beta_0$ .

Since  $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$ ,  $E(\mathbf{A}\mathbf{u} | \mathbf{x}) = \mathbf{0}$ , and, by the law of iterated expectations,  $E(\mathbf{A}\mathbf{u}) = \mathbf{0}$ .

So the key condition that  $\mathbf{A}$ , and thus  $\mathbf{C}$ , must satisfy for  $\check{\beta}$  to be unbiased is that  $E(\mathbf{A}\mathbf{X}\beta_0 | \mathbf{X}) = \beta_0$ .

This is the case for all  $\beta_0$  if and only if  $\mathbf{A}\mathbf{X} = \mathbf{I}$ .

Equivalently, it holds for all  $\beta_0$  whenever  $\mathbf{C}\mathbf{X} = \mathbf{O}$ , since

$$\mathbf{C}\mathbf{X} = \mathbf{A}\mathbf{X} - (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X} = \mathbf{A}\mathbf{X} - \mathbf{I}. \quad (5)$$

The condition that  $\mathbf{C}\mathbf{X} = \mathbf{O}$  implies that  $\mathbf{C}\mathbf{y} = \mathbf{C}\mathbf{u}$ . Since  $\mathbf{C}\mathbf{y} = \check{\beta} - \hat{\beta}$ , this makes it clear that  $\check{\beta} - \hat{\beta}$  has conditional expectation zero.

The unbiasedness condition on  $\check{\beta}$  also implies that the covariances of  $\check{\beta} - \hat{\beta}$  with  $\hat{\beta}$  are all zero.

In other words,  $\check{\beta} = \hat{\beta} + \mathbf{v}$ , where  $\mathbf{v}$  is uncorrelated with  $\hat{\beta}$ .

Since the covariances of  $\ddot{\beta} - \hat{\beta}$  with  $\hat{\beta}$  are zero,

$$E((\hat{\beta} - \beta_0)(\ddot{\beta} - \hat{\beta})^\top) = E((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \mathbf{u}^\top \mathbf{C}^\top) \quad (6)$$

$$= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma_0^2 \mathbf{I} \mathbf{C}^\top \quad (7)$$

$$= \sigma_0^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{C}^\top = \mathbf{O}. \quad (8)$$

Once again, we see that  $\mathbf{C}\mathbf{X} = \mathbf{O}$ , so that  $E(\mathbf{C}\mathbf{y}) = E(\mathbf{C}\mathbf{u}) = \mathbf{0}$ .

To complete the proof, note that

$$\begin{aligned} \text{Var}(\ddot{\beta}) &= \text{Var}(\hat{\beta} + (\ddot{\beta} - \hat{\beta})) \\ &= \text{Var}(\hat{\beta} + \mathbf{C}\mathbf{y}) \\ &= \text{Var}(\hat{\beta}) + \text{Var}(\mathbf{C}\mathbf{y}). \end{aligned} \quad (9)$$

The difference between  $\text{Var}(\ddot{\beta})$  and  $\text{Var}(\hat{\beta})$  is  $\text{Var}(\mathbf{C}\mathbf{y})$ , which must be a positive semidefinite matrix.

Any unbiased linear estimator  $\check{\beta}$  equals  $\hat{\beta}$  plus  $Cy$ , which has expectation zero and is uncorrelated with  $\hat{\beta}$ :

$$\check{\beta} = \hat{\beta} + Cy = \hat{\beta} + v. \quad (10)$$

- The Gauss-Markov theorem requires the disturbances to be independent and homoskedastic, but they do not have to be normally distributed.
- It applies only to a correctly specified model with exogenous regressors and disturbances with a scalar covariance matrix.
- It does *not* say that  $\hat{\beta}$  is more efficient than every imaginable estimator. Nonlinear and/or biased estimators may well perform better than OLS.

Something very similar to (10) holds asymptotically for estimators that are asymptotically efficient and asymptotically unbiased.

# Residuals and Disturbances

Least-squares residuals  $\hat{\mathbf{u}} \equiv \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  have both numerical and statistical properties.

Numerical:  $\hat{\mathbf{u}}$  is orthogonal to every vector that lies in  $\mathcal{S}(\mathbf{X})$ .

Consistency of  $\hat{\boldsymbol{\beta}}$  implies that  $\hat{\mathbf{u}} \rightarrow \mathbf{u}$  as  $N \rightarrow \infty$ .

$$\mathbf{M}_X \mathbf{y} = \mathbf{M}_X \mathbf{X} \boldsymbol{\beta}_0 + \mathbf{M}_X \mathbf{u} = \mathbf{M}_X \mathbf{u} = \hat{\mathbf{u}}. \quad (11)$$

Each residual is a linear combination of every one of the disturbances:

$$\hat{u}_i = u_i - \mathbf{X}_i (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \quad (12)$$

$$= u_i - \mathbf{X}_i (\mathbf{X}^\top \mathbf{X})^{-1} \sum_{j=1}^N \mathbf{X}_j^\top u_j. \quad (13)$$

Even when the  $u_i$  are independent, the  $\hat{u}_i$  are not independent.

Now assume that  $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$ , so that  $E(u_i | \mathbf{X}) = 0$  for all  $i$ .

Since, by (13),  $\hat{u}_i$  is just a linear combination of all the  $u_i$ ,  $E(\hat{u}_i | \mathbf{X}) = 0$ . Therefore,  $\text{Var}(\hat{u}_i)$  is just  $E(\hat{u}_i^2)$ .

We know that  $\|\hat{\mathbf{u}}\| < \|\mathbf{u}(\boldsymbol{\beta})\| = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|$  for any  $\boldsymbol{\beta} \neq \hat{\boldsymbol{\beta}}$ . In particular,  $\|\hat{\mathbf{u}}\| \leq \|\mathbf{u}(\boldsymbol{\beta}_0)\|$ .

This implies that  $E(\|\hat{\mathbf{u}}\|^2) \leq E(\|\mathbf{u}\|^2)$ . If  $\text{Var}(u_i) = \sigma_0^2$ , then

$$\sum_{i=1}^N \text{Var}(\hat{u}_i) = \sum_{i=1}^N E(\hat{u}_i^2) = E\left(\sum_{i=1}^N \hat{u}_i^2\right) = E(\|\hat{\mathbf{u}}\|^2) \quad (14)$$

$$\leq E(\|\mathbf{u}\|^2) = E\left(\sum_{i=1}^N u_i^2\right) = \sum_{i=1}^N E(u_i^2) = N\sigma_0^2. \quad (15)$$

This suggests that, at least for most observations,  $\text{Var}(\hat{u}_i)$  must be less than  $\sigma_0^2$ . In fact,  $\text{Var}(\hat{u}_i)$  is less than  $\sigma_0^2$  for *every* observation.



The covariance matrix of the entire vector  $\hat{\mathbf{u}}$  is

$$\text{Var}(\hat{\mathbf{u}}) = \text{Var}(\mathbf{M}_X \mathbf{u}) = \text{E}(\mathbf{M}_X \mathbf{u} \mathbf{u}^\top \mathbf{M}_X) \quad (16)$$

$$= \mathbf{M}_X \text{E}(\mathbf{u} \mathbf{u}^\top) \mathbf{M}_X = \mathbf{M}_X \text{Var}(\mathbf{u}) \mathbf{M}_X \quad (17)$$

$$= \mathbf{M}_X (\sigma_0^2 \mathbf{I}) \mathbf{M}_X = \sigma_0^2 \mathbf{M}_X \mathbf{M}_X = \sigma_0^2 \mathbf{M}_X. \quad (18)$$

This uses the facts that  $\text{E}(\mathbf{M}_X \mathbf{u}) = \mathbf{0}$  and that  $\mathbf{M}_X$  is idempotent.

From (18),  $\text{E}(\hat{u}_i \hat{u}_j) \neq 0$  for  $i \neq j$ . Residuals are correlated even when disturbances are uncorrelated.

Residuals do not have constant variance, and the variance of every residual must always be smaller than  $\sigma_0^2$ .

A typical diagonal element of  $\mathbf{M}_X$  is  $1 - h_i$ . Therefore, it follows from (18) that

$$\text{Var}(\hat{u}_i) = \text{E}(\hat{u}_i^2) = (1 - h_i) \sigma_0^2. \quad (19)$$

Since  $0 \leq 1 - h_i < 1$ , we see that  $\text{E}(\hat{u}_i^2) < \sigma_0^2$ .

## Estimating the Variance of the Disturbances

(19) tells us that high-leverage observations, for which  $h_i$  is relatively large, must have residuals with unusually small variances.

The method of moments suggests that we can estimate  $\sigma^2$  by using the corresponding sample moment.

If we observed the  $u_i$ , this sample moment would be

$$\frac{1}{N} \sum_{i=1}^N u_i^2. \quad (20)$$

In fact, we only observe the  $\hat{u}_i$ , so a natural MM estimator is

$$\hat{\sigma}^2 \equiv \frac{1}{N} \sum_{i=1}^N \hat{u}_i^2. \quad (21)$$

Because  $E(\hat{u}_i^2) < \sigma_0^2$ , by (19),  $\hat{\sigma}^2$  must be biased downward.

We know that  $\sum_{i=1}^N h_i = k$ . Therefore,

$$E(\hat{\sigma}^2) = \frac{1}{N} \sum_{i=1}^N E(\hat{u}_i^2) = \frac{1}{N} \sum_{i=1}^N (1 - h_i) \sigma_0^2 = \frac{N - k}{N} \sigma_0^2. \quad (22)$$

Since  $\hat{u} = M_X u$  and  $M_X$  is idempotent, the SSR is just  $u^\top M_X u$ , and

$$E(u^\top M_X u) = E(\text{SSR}(\hat{\beta})) = E\left(\sum_{i=1}^N \hat{u}_i^2\right) = (N - k) \sigma_0^2. \quad (23)$$

Adding one more regressor has exactly the same effect on the expectation of the SSR as taking away one observation.

The result (23) suggests another MM estimator which is unbiased whenever  $\hat{\beta}$  is:

$$s^2 \equiv \frac{1}{N - k} \sum_{i=1}^N \hat{u}_i^2. \quad (24)$$

The square root of  $s^2$  is  $s$ , the **standard error of the regression** (or **regression standard error**). Even though  $s^2$  provides an unbiased estimate of  $\sigma^2$ ,  $s$  does not provide an unbiased estimate of  $\sigma$ .

An unbiased estimator for  $\text{Var}(\hat{\boldsymbol{\beta}})$  is

$$\widehat{\text{Var}}(\hat{\boldsymbol{\beta}}) = s^2(\mathbf{X}^\top\mathbf{X})^{-1}. \quad (25)$$

This is the usual estimator of  $\text{Var}(\hat{\boldsymbol{\beta}})$  under the assumption of IID disturbances.

- Suppose that  $\mathbf{X}$  is fixed across samples. Then  $\sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$  is the same for every sample, but  $s^2(\mathbf{X}^\top\mathbf{X})^{-1}$  varies.
- However,  $s^2(\mathbf{X}^\top\mathbf{X})^{-1}$  is always proportional to  $\sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$ . Like the latter, it is normally  $O(N^{-1})$
- Evidently,  $s^2 = O_p(1)$ , and  $s^2/\sigma_0^2 \rightarrow 1$  as  $N \rightarrow \infty$ .

# Over-specification

A regression model is **over-specified** if some variables that belong to  $\Omega_i$  but do not appear in the DGP, are mistakenly included.

Since the DGP remains a special case of the model, there is no actual misspecification, merely a failure to incorporate zero restrictions that are true.

Consider the over-specified linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (26)$$

where  $\mathbf{X}_i$  and  $\mathbf{Z}_i$  belong to  $\Omega_i$ . The data are actually generated by

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma_0^2 \mathbf{I}). \quad (27)$$

The DGP (27) is a special case of (26), with  $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ ,  $\boldsymbol{\gamma} = \mathbf{0}$ , and  $\sigma^2 = \sigma_0^2$ .

The estimates  $\hat{\beta}$  from (26) are the same as those from the FWL regression

$$M_Z \mathbf{y} = M_Z \mathbf{X} \boldsymbol{\beta} + \text{residuals}, \quad (28)$$

where, as usual,  $M_Z = \mathbf{I} - \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ .

Thus we see that

$$\hat{\beta} = (\mathbf{X}^\top M_Z \mathbf{X})^{-1} \mathbf{X}^\top M_Z \mathbf{y}. \quad (29)$$

$\hat{\beta}$  must be unbiased if  $\mathbf{X}$  and  $\mathbf{Z}$  are exogenous.

If we replace  $\mathbf{y}$  by  $\mathbf{X} \boldsymbol{\beta}_0 + \mathbf{u}$ , we find that

$$\hat{\beta} = \boldsymbol{\beta}_0 + (\mathbf{X}^\top M_Z \mathbf{X})^{-1} \mathbf{X}^\top M_Z \mathbf{u}. \quad (30)$$

The second term on the r.h.s. has conditional mean  $\mathbf{0}$ , provided we take expectations conditional on  $\mathbf{Z}$  as well as on  $\mathbf{X}$ .

Imposing the restriction that  $\boldsymbol{\gamma} = \mathbf{0}$  yields  $\tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ , which is also unbiased because the restriction is true.

We now have two unbiased linear estimators,  $\hat{\beta}$  and  $\tilde{\beta}$ .

When the restriction holds, the matrix  $\text{Var}(\hat{\beta}) - \text{Var}(\tilde{\beta})$  must be positive semidefinite. We know that

$$\text{Var}(\tilde{\beta}) = \sigma_0^2(\mathbf{X}^\top\mathbf{X})^{-1}. \quad (31)$$

It is easily shown that

$$\text{Var}(\hat{\beta}) = \text{E}((\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^\top) \quad (32)$$

$$= (\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{M}_Z\text{E}(\mathbf{u}\mathbf{u}^\top)\mathbf{M}_Z\mathbf{X}(\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1} \quad (33)$$

$$= \sigma_0^2(\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{M}_Z\mathbf{I}\mathbf{M}_Z\mathbf{X}(\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1} \quad (34)$$

$$= \sigma_0^2(\mathbf{X}^\top\mathbf{M}_Z\mathbf{X})^{-1}. \quad (35)$$

When there is only one regressor,  $x$ , and one parameter,  $\beta$ , it is easy to see that  $\tilde{\beta}$  is more efficient than  $\hat{\beta}$ .

Since  $M_Z$  is a projection matrix,  $\|M_Z x\|$  must be smaller (or at least, no larger) than  $\|x\|$ . Thus  $x^\top M_Z x \leq x^\top x$ , which implies that

$$\sigma_0^2 (x^\top M_Z x)^{-1} \geq \sigma_0^2 (x^\top x)^{-1}. \quad (36)$$

Notice that  $\text{Var}(\hat{\beta})$  is proportional to the inverse of the SSR from a regression of  $x$  on the other regressors.

Showing that  $\text{Var}(\hat{\beta}) - \text{Var}(\tilde{\beta})$  is positive semidefinite is equivalent to showing that  $\text{Var}(\tilde{\beta})^{-1} - \text{Var}(\hat{\beta})^{-1}$  is positive semidefinite:

$$\begin{aligned} X^\top X - X^\top M_Z X &= X^\top (I - M_Z) X \\ &= X^\top P_Z X = (P_Z X)^\top P_Z X. \end{aligned} \quad (37)$$

This is the transpose of a matrix times itself. Thus  $\sigma_0^2 (X^\top M_Z X)^{-1} - \sigma_0^2 (X^\top X)^{-1}$  must be positive semidefinite.

Adding additional variables that do not really belong in a model leads to less accurate estimates unless  $P_Z X = \mathbf{O}$ , so that  $M_Z X = X$ .



We used the FWL theorem to show that  $\hat{\beta} = (X^T M_Z X)^{-1} X^T M_Z y$ .

If we were unaware of this theorem, we would need to write

$$\begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix} = \begin{bmatrix} X^T X & X^T Z \\ Z^T X & Z^T Z \end{bmatrix}^{-1} \begin{bmatrix} X^T y \\ Z^T y \end{bmatrix}. \quad (38)$$

Then we would have to invert the matrix in (38) and multiply the first row of the inverse by the vector that follows it.

The inverse is (not easily!) seen to be

$$\begin{bmatrix} (X^T M_Z X)^{-1} & -(X^T M_Z X)^{-1} X^T Z (Z^T Z)^{-1} \\ -(Z^T Z)^{-1} Z^T X (X^T M_Z X)^{-1} & (Z^T M_X Z)^{-1} \end{bmatrix}. \quad (39)$$

This leads directly to  $\hat{\beta} = (X^T M_Z X)^{-1} X^T M_Z y$ .

To obtain  $\hat{\gamma} = (Z^T M_X Z)^{-1} Z^T M_X y$ , it is easier to write the lower left-hand submatrix as  $-(Z^T M_X Z)^{-1} Z^T X (X^T X)^{-1}$ .

# Under-specification

Now suppose the DGP is really

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma_0^2 \mathbf{I}). \quad (40)$$

The estimator  $\hat{\boldsymbol{\beta}}$  is now the “correct” one to use.

If instead we use  $\tilde{\boldsymbol{\beta}}$ , there really is misspecification, and the restricted estimator  $\tilde{\boldsymbol{\beta}}$  is biased.

$$\mathbb{E}(\tilde{\boldsymbol{\beta}}) = \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u})) \quad (41)$$

$$= \boldsymbol{\beta}_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}) \quad (42)$$

$$= \boldsymbol{\beta}_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z}\boldsymbol{\gamma}_0. \quad (43)$$

The second term in (43) equals zero only when  $\mathbf{X}^\top \mathbf{Z} = \mathbf{O}$  or  $\boldsymbol{\gamma}_0 = \mathbf{0}$ .

In all other cases,  $\tilde{\beta}$  is biased. The magnitude of the bias depends on  $\gamma_0$  and on the  $X$  and  $Z$  matrices.

This bias does not vanish as  $N \rightarrow \infty$ , so  $\tilde{\beta}$  is also generally inconsistent.

Since  $\tilde{\beta}$  is biased, we cannot use  $\text{Var}(\tilde{\beta})$  to evaluate its accuracy.

Instead, we can use the **mean squared error matrix**, or **MSE matrix**:

$$\text{MSE}(\tilde{\beta}) \equiv E((\tilde{\beta} - \beta_0)(\tilde{\beta} - \beta_0)^\top). \quad (44)$$

The MSE matrix is equal to  $\text{Var}(\tilde{\beta})$  if  $\tilde{\beta}$  is unbiased, but not otherwise.

For a scalar parameter  $\beta$ ,

$$\text{MSE}(\tilde{\beta}) = \text{Var}(\tilde{\beta}) + (E(\tilde{\beta}) - \beta_0)^2. \quad (45)$$

In this case, it is common to report the **root mean squared error**, or **RMSE**, instead of the MSE.

It is easy to see that

$$\tilde{\beta} - \beta_0 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}. \quad (46)$$

Therefore,  $\tilde{\beta} - \beta_0$  times itself transposed is equal to

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0 \gamma_0^\top \mathbf{Z}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \mathbf{u}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (47)$$

$$+ (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0 \mathbf{u}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \gamma_0^\top \mathbf{Z}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (48)$$

The second term here has expectation  $\sigma_0^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ , and the third and fourth terms have expectation zero. Thus

$$\text{MSE}(\tilde{\beta}) = \sigma_0^2 (\mathbf{X}^\top \mathbf{X})^{-1} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0 \gamma_0^\top \mathbf{Z}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (49)$$

The first term here is  $O(N^{-1})$ , and the second is  $O(1)$ .

We would like to compare  $\text{MSE}(\tilde{\beta})$  with  $\text{MSE}(\hat{\beta}) = \text{Var}(\hat{\beta})$ .

If the bias is small, the second term in (49) must be small, and in that case  $\tilde{\beta}$  may well have smaller MSE than  $\hat{\beta}$ .

However, if the bias is large, the second term in (49) is necessarily large, and  $\tilde{\beta}$  must have larger MSE than  $\hat{\beta}$ .

$\tilde{\beta}$  may yield smaller MSE for some parameters and  $\hat{\beta}$  for others.

$\widehat{\text{Var}}(\tilde{\beta})$  calculated by a least-squares regression program attempts to estimate the first term in (49), but it ignores the second.

Because  $s^2$  is typically larger than  $\sigma_0^2$  if some regressors have been incorrectly omitted, this estimate is biased.

- Under-specification causes bias and inconsistency, but over-specification “merely” causes inefficiency. Which problem is more severe?
- It depends on how much information the sample contains. In sufficiently large samples, avoid under-specification at all costs.
- However, in samples of modest size, the gain in efficiency from omitting some variables may be very large relative to the bias that is caused by their omission. There is a **bias-variance tradeoff**.

# Measures of Goodness of Fit

The most commonly used (and misused) measure of goodness of fit is the **coefficient of determination**, or  $R^2$ .

The **uncentered**  $R^2$ , denoted  $R_u^2$ , is the ratio of the explained sum of squares (ESS) of the regression to the total sum of squares (TSS).

As a consequence of Pythagoras' Theorem,

$$\text{TSS} = \|\mathbf{y}\|^2 = \|\mathbf{P}_X\mathbf{y}\|^2 + \|\mathbf{M}_X\mathbf{y}\|^2 = \text{ESS} + \text{SSR}. \quad (50)$$

Therefore,

$$R_u^2 = \frac{\text{ESS}}{\text{TSS}} = \frac{\|\mathbf{P}_X\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = 1 - \frac{\|\mathbf{M}_X\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = 1 - \frac{\text{SSR}}{\text{TSS}} = \cos^2\theta, \quad (51)$$

where  $\theta$  is the angle between  $\mathbf{y}$  and  $\mathbf{P}_X\mathbf{y}$ .

- For any angle  $\theta$ ,  $-1 \leq \cos \theta \leq 1$ . Consequently,  $0 \leq R_u^2 \leq 1$ .
- If  $\theta = 0$ ,  $\mathbf{y}$  and  $\mathbf{X}\hat{\boldsymbol{\beta}}$  would coincide, the residual vector  $\hat{\mathbf{u}}$  would vanish, and we would have a **perfect fit**, with  $R_u^2 = 1$ .
- If  $R_u^2 = 0$ , the fitted value vector would vanish, and  $\mathbf{y}$  would coincide with  $\hat{\mathbf{u}}$ .
- $R_u^2$  depends on the data only through residuals and fitted values. It is invariant to nonsingular linear transformations of  $\mathbf{X}$ .
- Because it is a ratio,  $R_u^2$  is invariant to changes in the scale of  $\mathbf{y}$ .

The **centered**  $R^2$ , denoted  $R_c^2$ , is much more commonly encountered than  $R_u^2$ . All variables are centered, that is, expressed as deviations from their means, before ESS and TSS are calculated.

By adding a large enough constant to all the  $y_i$ , we could make  $R_u^2$  become arbitrarily close to 1, since the SSR would stay the same and the TSS would increase without limit.

But  $R_c^2$  is invariant to changes in the mean of the regressand.

- Both versions of  $R^2$  are valid only if a regression model is estimated by least squares. Only then is  $TSS = ESS + SSR$ .
- The centered version is not valid if the regressors do not include a constant term (or equivalent), i.e. if  $\iota$  does not belong to  $\mathcal{S}(X)$ .
- Both  $R_u^2$  and  $R_c^2$  increase whenever more regressors are added.

Consider restricted and unrestricted models with the same dependent variable. They have the same TSS, so the regression with the larger ESS (smaller SSR) must also have the larger  $R^2$ .

The ESS from the unrestricted model is  $\|P_{X,Z}y\|^2$ , and the ESS from the restricted model is  $\|P_Xy\|^2$ . The difference between them is

$$y^\top P_{X,Z}y - y^\top P_Xy = y^\top (P_{X,Z} - P_X)y. \quad (52)$$

Since the matrix  $P_{X,Z} - P_X$  is an orthogonal projection matrix, (52) must be non-negative.



Why is  $P_{X,Z} - P_X$  an orthogonal projection matrix? Because  $\mathcal{S}(X) \subset \mathcal{S}(X, Z)$ . This implies that  $P_X$  projects on to a subspace of the image of  $P_{X,Z}$ .

We conclude that the ESS, and hence the  $R^2$ , from the unrestricted model can be no less than those from the restricted model.

The  $R^2$  can be modified so that adding additional regressors does not necessarily increase its value.

If  $\iota \in \mathcal{S}(X)$ , the centered  $R^2$  can be written as

$$R_c^2 = 1 - \frac{\sum_{i=1}^N \hat{u}_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2}. \quad (53)$$

The SSR has expectation  $(N - k)\sigma_0^2$  under standard assumptions. The denominator is  $N - 1$  times an unbiased estimator of the variance of  $y_i$  about its true mean. As such, it has expectation  $(N - 1) \text{Var}(y)$ .

Thus the second term of (53) can be thought of as the ratio of two biased estimators.

If we replace these by unbiased ones, we obtain the **adjusted**  $R^2$ ,

$$\bar{R}^2 \equiv 1 - \frac{\frac{1}{N-k} \sum_{i=1}^N \hat{u}_i^2}{\frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{(N-1) \mathbf{y}^\top \mathbf{M}_X \mathbf{y}}{(N-k) \mathbf{y}^\top \mathbf{M}_I \mathbf{y}}. \quad (54)$$

$\bar{R}^2$  and  $R_c^2$  are generally very similar, except when  $(N-k)/(N-1) \ll 1$ .

One nice feature of  $R_u^2$  and  $R_c^2$  is that they are constrained to lie between 0 and 1.

In contrast,  $\bar{R}^2$  can actually be negative when  $(N-1)/(N-k)$  is greater than TSS/SSR.

- Never compare any form of  $R^2$  for models that are estimated using different datasets!
- Models with high  $R^2$  can be complete nonsense.