

Supplementary Notes for 8070 Review Quiz

These are supplementary notes that cover a handful of extra topics for the 8070 review quiz that were not provided elsewhere.

1. Bias and Variance of $\hat{\beta}$

We will consider the following assumptions throughout this part of the course:

1. Linear CEF: $Y = X'\beta + e$ and $\mathbb{E}[e|X] = 0$
2. Finite Moments: $\mathbb{E}[Y^2] < \infty$ and $\mathbb{E}[\|X\|^2] < \infty$
3. Positive definite design matrix: $\mathbb{E}[XX']$ is positive definite.

For some of the results below, we will also use the additional **homoskedasticity** condition: $\mathbb{E}[e^2|X] = \sigma^2$ (that is, the variance of the error term does not depend on X)

We'll continue to suppose that we have access to an i.i.d. sample. The main two properties that we'll consider are the **bias** of $\hat{\beta}$ and the **sampling variance** of $\hat{\beta}$. Before we consider those, let's start by defining what they are. Let $\hat{\theta}$ generically denote some estimator of a population parameter of interest θ . Then,

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta$$

$\hat{\theta}$ is said to be **unbiased** if $\text{Bias}(\hat{\theta}) = 0$, or, equivalently, if $\mathbb{E}[\hat{\theta}] = \theta$. It is worth pausing a moment to think conceptually about what is happening here. First, estimators are random — this point may not be immediately obvious though. In particular, given once you have access to a particular dataset, this typically pins down a value of $\hat{\theta}$. What it means that $\hat{\theta}$ is random is that we can carry out the thought experiment of repeatedly collecting n new observations from the same population and re-calculating $\hat{\theta}$ for the new data. In our thought experiment, given that we have new samples, the value of $\hat{\theta}$ would generally change with each new sample. If you were to carry this procedure out an extremely large number of times, this would give rise to a distribution of $\hat{\theta}$ in repeated samples; this distribution is called the **sampling distribution** of $\hat{\theta}$.

In practice, however, we only have one dataset and, therefore, only one value of $\hat{\theta}$. Given the above discussion, it is natural to consider the $\hat{\theta}$ that we have as a draw from the sampling distribution discussed above. Therefore, if an estimator is unbiased, what this means is that, on average (with respect to the sampling distribution), our estimator $\hat{\theta}$ is equal to the population parameter θ . Importantly, unbiasedness is generally a good property for an estimator to have, but, given that we only have one draw from the sampling distribution, even if our estimator is unbiased, it is still *possible* that our particular value of $\hat{\theta}$ could be far away from θ .

Practice: Show that $\bar{Y} := \frac{1}{n} \sum_{i=1}^n Y_i$ is unbiased for $\mathbb{E}[Y]$.

Next, the sampling variance of $\hat{\theta}$ is given by $\text{var}(\hat{\theta})$. You should think of this as the variance of $\hat{\theta}$ in the repeated sampling thought experiment mentioned above. All else equal, we would prefer estimators that have lower sampling variance.

Expectation of least squares estimator

H: 4.5, 4.7

Now, let's consider the bias of $\hat{\beta}$. To start with let's calculate $\mathbb{E}[\hat{\beta}|\mathbf{X}]$ (this sort of conditional expectation may feel a bit unusual as we are conditioning on the data matrix, but it is totally reasonable to do this)

$$\begin{aligned}\mathbb{E}[\hat{\beta}|\mathbf{X}] &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}[\mathbf{Y}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta \\ &= \beta\end{aligned}$$

To see the step that uses $\mathbb{E}[\mathbf{Y}|\mathbf{X}] = \mathbf{X}\beta$, let's point out a few things. First,

$$\mathbb{E}[Y_i|\mathbf{X}] = \mathbb{E}[Y_i|X_1, X_2, \dots, X_n] = \mathbb{E}[Y_i|X_i] = X'_i\beta$$

where the first equality holds immediately, the second equality holds by the independence in i.i.d. sampling, and the last equality holds by the linear CEF. Thus,

$$\mathbb{E}[\mathbf{Y}|\mathbf{X}] = \begin{pmatrix} \vdots \\ \mathbb{E}[Y_i|\mathbf{X}] \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ X'_i\beta \\ \vdots \end{pmatrix} = \mathbf{X}\beta$$

which is what we used above.

The book provides an alternative derivation for the same result which I think is also useful for quickly covering. Notice that we can alternatively write

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'(\mathbf{X}\beta + \mathbf{e})) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}\mathbf{e} \\ &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}\mathbf{e}\end{aligned}\tag{1}$$

The expression in Equation 1 is one that we will use a number of times throughout this semester, so I think it is worth highlighting.

Now, using this expression, notice that

$$\begin{aligned}\mathbb{E}[\hat{\beta}|\mathbf{X}] &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}\mathbb{E}[\mathbf{e}|\mathbf{X}] \\ &= \beta\end{aligned}$$

where the last equality holds because $\mathbb{E}[\mathbf{e}|\mathbf{X}] = \mathbf{0}$ which holds because $\mathbb{E}[e|X] = 0$ and by using similar arguments as for $\mathbb{E}[\mathbf{Y}|\mathbf{X}]$ above.

Given the result above, it then follows by the law of iterated expectations that

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[\mathbb{E}[\hat{\beta}|\mathbf{X}]] = \beta$$

and that, therefore, $\hat{\beta}$ is unbiased for β .

Variance of least squares estimator

H: 4.6, 4.7

Next, we'll calculate the sampling variance of $\hat{\beta}$. To this end, let's start by defining

$$\mathbf{D} := \text{var}(\mathbf{e}|\mathbf{X}) = \mathbb{E}[\mathbf{e}\mathbf{e}'|\mathbf{X}]$$

where the last equality holds because $\mathbb{E}[\mathbf{e}|\mathbf{X}] = \mathbf{0}$. It's worth momentarily thinking about some of the properties of \mathbf{D} . First, it is an $n \times n$ matrix. Second, its diagonal elements are given by $\mathbb{E}[e_i^2|\mathbf{X}] = \mathbb{E}[e_i^2|X_i] =: \sigma_i^2$. The off-diagonal elements are given by $\mathbb{E}[e_i e_j|\mathbf{X}] = \mathbb{E}[e_i|X_i]\mathbb{E}[e_j|X_j] = 0$ (here, the second equality holds by independence across observations). Thus, \mathbf{D} is a diagonal matrix. If we are willing to introduce the assumption of homoskedasticity, then $\mathbb{E}[e_i^2|X_i] = \sigma^2$ (and is therefore constant across i). In this case, $\mathbf{D} = \mathbf{I}_n \sigma^2$.

As a first step towards calculating $\text{var}(\hat{\beta})$, notice that

$$\begin{aligned}\text{var}(\mathbf{Y}|\mathbf{X}) &= \text{var}(\mathbf{X}\beta + \mathbf{e}|\mathbf{X}) \\ &= \text{var}(\mathbf{e}|\mathbf{X}) = \mathbf{D}\end{aligned}$$

where the first equality holds by plugging in for \mathbf{Y} , the second equality holds because we are conditioning on \mathbf{X} , and the last equality by the definition of \mathbf{D} .

Now, consider

$$\begin{aligned}\mathbf{V}_{\hat{\beta}} &:= \text{var}(\hat{\beta}|\mathbf{X}) \\ &= \text{var}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}|\mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{var}(\mathbf{Y}|\mathbf{X})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\end{aligned}$$

where the second equality holds by plugging in for $\hat{\beta}$, the third equality by the matrix version of $\text{var}(aZ) = a^2\text{var}(Z)$ when a is a constant and Z is a scalar random variable (and because $\mathbf{X}'\mathbf{X}$ is symmetric), and the last equality holds because $\text{var}(\mathbf{Y}|\mathbf{X}) = \mathbf{D}$ which we showed above. If we additionally invoke homoskedasticity, then this will simplify; in particular, in this case $\mathbf{X}'\mathbf{D}\mathbf{X} = \mathbf{X}'\mathbf{I}_n\sigma^2\mathbf{X} = \mathbf{X}'\mathbf{X}\sigma^2$. This implies that

$$\mathbf{V}_{\hat{\beta}}^0 = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

where I include the 0 superscript to indicate that this expression holds only under the additional condition of homoskedasticity.

If we want to calculate the unconditional variance of $\hat{\beta}$, then we can use the law of total variance. This is given in Theorem 2.8 in the textbook; in particular, as long as $\mathbb{E}[Y^2] < \infty$, then $\text{var}(Y) = \mathbb{E}[\text{var}(Y|X)] + \text{var}(\mathbb{E}[Y|X])$. Applying this to the present context, we have that

$$\begin{aligned}\text{var}(\hat{\beta}) &= \mathbb{E}[\text{var}(\hat{\beta}|\mathbf{X})] + \text{var}(\mathbb{E}[\hat{\beta}|\mathbf{X}]) \\ &= \mathbb{E}[\text{var}(\hat{\beta}|\mathbf{X})] + 0 \\ &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]\end{aligned}$$

as above, this can simplify under homoskedasticity.

Side-comment: It is worth briefly comparing the above results to similar results in the very simple case where we estimate $\mu := \mathbb{E}[Y]$ by \bar{Y} (the sample average of Y_i). In this case, recall that $\mathbb{E}[\bar{Y}] = \mu$, so that \bar{Y} is unbiased for μ , just like $\hat{\beta}$ is for β .

Further, recall that $\text{var}(\bar{Y}) = \frac{\text{var}(Y)}{n}$, which says that the sampling variance of \bar{Y} depends on the variance of Y , and it also tends to decrease for larger values of n . From the above discussion, it may not be immediately obvious whether or not the sampling variance of $\hat{\beta}$ decreases with n — it turns out that it does. To see this, recall that $(\mathbf{X}'\mathbf{X}) = \sum_{i=1}^n X_i X'_i$ which grows with n . Now, for simplicity, suppose that homoskedasticity holds (similar arguments will hold for the case without homoskedasticity), notice that we can rewrite

$$\mathbf{V}_{\hat{\beta}}^0 = \frac{\sigma^2}{n} \left(\frac{1}{n} \mathbf{X}' \mathbf{X} \right)^{-1}$$

which just multiplies and divides by n . Notice that, here, $\frac{1}{n} \mathbf{X}' \mathbf{X} = \frac{1}{n} \sum_{i=1}^n X_i X'_i$ is now an average that does not systematically grow with n . On the other hand, there is now an n in the denominator so that it is easier to see that the sampling variance of $\hat{\beta}$ *does* decrease with the sample size, just like for \bar{Y} .

Gauss-Markov Theorem

H: 4.8

The Gauss-Markov theorem says that, given the linear regression assumptions + homoskedasticity, $\hat{\beta}$ is **efficient** (has the smallest variance) among all possible *linear, unbiased* estimators (side-comment: Bruce Hansen has a recent paper showing that $\hat{\beta}$ is efficient among unbiased estimators; I am not sure that I fully understand his arguments, so I'm just going to teach the “classical” version of the Gauss-Markov theorem).

More specifically, the Gauss-Markov theorem says: Given the linear regression assumptions and homoskedasticity, for any possible linear, unbiased estimator of β , which we'll denote as $\tilde{\beta}$, $\text{var}(\tilde{\beta}|\mathbf{X}) \geq \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$

Efficiency is a very good property for an estimator to have, and, therefore, this kind of result provides a strong justification for using $\hat{\beta}$ as an estimate of β .

To prove this result, let's first see what linearity and unbiasedness “buys us”.

1. A linear estimator is one that we can write as $\tilde{\beta} = \mathbf{A}'\mathbf{Y}$ where \mathbf{A} is an $n \times k$ matrix that is a function of \mathbf{X}
2. Unbiasedness means that $\mathbb{E}[\tilde{\beta}|\mathbf{X}] = \beta$. If $\tilde{\beta}$ is also linear, notice that $\mathbb{E}[\mathbf{A}'\mathbf{Y}|\mathbf{X}] = \mathbf{A}'\mathbb{E}[\mathbf{Y}|\mathbf{X}] = \mathbf{A}'\mathbf{X}\beta$; then, unbiasedness therefore implies that $\mathbf{A}'\mathbf{X} = \mathbf{I}_k$.

Now, let's calculate the conditional variance of some generic linear, unbiased estimator of β

$$\begin{aligned}\text{var}(\tilde{\beta}|\mathbf{X}) &= \text{var}(\mathbf{A}'\mathbf{Y}|\mathbf{X}) \\ &= \text{var}(\mathbf{A}'(\mathbf{X}\beta + \mathbf{e})|\mathbf{X}) \\ &= \text{var}(\mathbf{A}'\mathbf{e}|\mathbf{X}) \\ &= \mathbf{A}'\text{var}(\mathbf{e}|\mathbf{X})\mathbf{A} \\ &= \mathbf{A}'\mathbf{A}\sigma^2\end{aligned}$$

where the first equality holds by linearity, the second equality substitutes for \mathbf{Y} , the third equality holds because the variance of the term involving $\mathbf{X}\beta$ is equal to 0 conditional on \mathbf{X} , the fourth equality holds by the property of variance that we used above (and because \mathbf{A} is a function of \mathbf{X}), and the last equality holds because $\text{var}(\mathbf{e}|\mathbf{X}) = \mathbf{I}_n\sigma^2$ under homoskedasticity.

Since, from earlier, we know that $\text{var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$, to complete the proof, we need to show that $\mathbf{A}'\mathbf{A} \geq (\mathbf{X}'\mathbf{X})^{-1}$. Towards this end, notice that

$$\begin{aligned}\mathbf{A}'\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1} &= \mathbf{A}'\mathbf{A} - \mathbf{A}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{A} \\ &= \mathbf{A}'(\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{A} \\ &= \mathbf{A}'\mathbf{M}\mathbf{A} \\ &= \mathbf{A}'\mathbf{M}'\mathbf{M}\mathbf{A} \\ &= (\mathbf{M}\mathbf{A})'\mathbf{M}\mathbf{A} \\ &\geq 0\end{aligned}$$

where the first equality uses $\mathbf{A}'\mathbf{X} = \mathbf{I}_k$, the second equality factors out \mathbf{A} , the third equality holds by the definition of \mathbf{M} , the fourth and fifth equalities hold because \mathbf{M} is idempotent and symmetric, the term in the last equality is positive semi-definite because it is a quadratic form.

Generalized least squares

H: 4.9

The Gauss-Markov theorem relied on the homoskedasticity condition. This begs the question of whether or not these efficiency results for $\hat{\beta}$ go through without this condition. Section 4.9 of the book considers this case. In fact, it considers a more general case than we have been considering so far where $\text{var}(\mathbf{e}|\mathbf{X}) = \Sigma\sigma^2$ where Σ is an $n \times n$ symmetric and positive semi-definite matrix (what's more general here is that this allows for relaxing the independence condition so that Σ can be non-diagonal).

Using similar arguments as above, we can show that, in this case

$$\text{var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\Sigma\mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}$$

However, Theorem 4.5 in the textbook shows that, under the linear regression assumptions (but not requiring homoskedasticity), for any possible linear, unbiased estimator of β (again, we'll denote it $\tilde{\beta}$),

$$\text{var}(\tilde{\beta}|\mathbf{X}) \geq \sigma^2(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}$$

Since $\text{var}(\hat{\beta}|\mathbf{X}) \neq \sigma^2(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}$, this suggests that we might ought to consider alternative estimators in this case. In particular, when Σ is known, consider pre-multiplying the regression by $\Sigma^{-1/2}$ to get

$$\tilde{\mathbf{Y}} = \tilde{\mathbf{X}}\beta + \tilde{\mathbf{e}}$$

where $\tilde{\mathbf{Y}} := \Sigma^{-1/2}\mathbf{Y}$, $\tilde{\mathbf{X}} := \Sigma^{-1/2}\mathbf{X}$, and $\tilde{\mathbf{e}} := \Sigma^{-1/2}\mathbf{e}$, and consider estimating this by OLS, so that

$$\begin{aligned}\tilde{\beta}_{gls} &= (\tilde{\mathbf{X}}'\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}'\tilde{\mathbf{Y}} \\ &= ((\Sigma^{-1/2}\mathbf{X})'\Sigma^{-1/2}\mathbf{X})^{-1}(\Sigma^{-1/2}\mathbf{X})'\Sigma^{-1/2}\mathbf{Y} \\ &= (\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{Y}\end{aligned}$$

Using the same sorts of arguments as we have been making above, you can show the following two results

$$\begin{aligned}\mathbb{E}[\tilde{\beta}_{gls}|\mathbf{X}] &= \beta \\ \text{var}(\tilde{\beta}_{gls}|\mathbf{X}) &= \sigma^2(\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1}\end{aligned}$$

This suggests that $\tilde{\beta}_{gls}$ is both unbiased and more efficient than $\hat{\beta}$ under heteroskedasticity.

One issue, however, is that this estimator is generally infeasible because Σ is not typically known. Instead, in practice, you can replace Σ with a suitable estimate $\hat{\Sigma}$. This is called **feasible GLS**. My sense is that GLS/FGLS is not very common in applied work, especially relative to OLS combined with “heteroskedasticity robust” standard errors. I think there are several reasons for this. First, estimating Σ may be hard to do in practice. For example, if we return to the simpler case where $\text{var}(\mathbf{e}|\mathbf{X}) = \mathbf{D}$ and recalling that \mathbf{D} is diagonal with diagonal elements equal to $\mathbb{E}[e_i^2|X_i]$. To estimate \mathbf{D} then would require estimating $\mathbb{E}[e^2|X]$. In practice, you could write down a parametric model for $\mathbb{E}[e^2|X]$, but this might be difficult in practice. If the model is not correctly specified, then the efficiency arguments above may not hold anymore. Second, the arguments that rationalize FGLS typically require $n \rightarrow \infty$ and amount to showing that FGLS and GLS are equivalent in this case (I think the finite sample arguments that we have been considering above for OLS/GLS are not straightforward when $\text{var}(\mathbf{e}|\mathbf{X})$ has to be estimated). This somewhat weakens the positive results for GLS mentioned above. Finally, the arguments in this section have been for the case where the CEF is actually linear, so it is less clear if there is a gain to using FGLS when we view $\hat{\beta}$ as the

linear projection coefficient instead of the coefficient from a linear CEF model.

2. Omitted variable bias

H: 2.24

Partition X as follows $X = (X'_1, X'_2)'$ where X_1 is a k_1 dimensional vector, X_2 is a k_2 dimensional vector, and $k = k_1 + k_2$. Likewise partition β into $\beta = (\beta'_1, \beta'_2)'$. Suppose that we are interested in β_1 from the linear projection of Y onto X_1 and X_2 :

$$Y = X'_1\beta_1 + X'_2\beta_2 + e \quad (2)$$

Since this is a linear projection, it implies that $\mathbb{E}[Xe] = 0$.

However, let's suppose that X_2 is not observed, so that it is infeasible to run a regression of Y on X_1 and X_2 . In this section, we consider properties of the following **short regression**

$$Y = X'_1\gamma_1 + u$$

which is the linear projection of Y on X_1 only. Since this is a linear projection, we also have that $\mathbb{E}[X_1u] = 0$ and that

$$\begin{aligned} \gamma_1 &= \mathbb{E}[X_1X'_1]^{-1}\mathbb{E}[X_1Y] \\ &= \mathbb{E}[X_1X'_1]^{-1}\mathbb{E}[X_1(X'_1\beta_1 + X'_2\beta_2 + e)] \\ &= \beta_1 + \mathbb{E}[X_1X'_1]^{-1}\mathbb{E}[X_1X'_2]\beta_2 \\ &= \beta_1 + \Gamma_{21}\beta_2 \end{aligned}$$

where the first equality holds by the definition of linear projection of Y on X_1 , the second equality holds by substituting for Y , the third equality combines and cancels terms and also holds because $\mathbb{E}[X_1e] = 0$ (since $\mathbb{E}[Xe] = 0$), and the last equality holds because we define $\Gamma_{21} = \mathbb{E}[X_1X'_1]^{-1}\mathbb{E}[X_1X'_2]$, which is a $k_1 \times k_2$ matrix of coefficients from the projection of all k_2 elements of X_2 on X_1 .

Importantly, the previous expression implies that γ_1 is not generally equal to β_1 ; that is, in general, we are not able to recover the parameter of interest β_1 from the feasible regression of Y on X_1 . This is probably not surprising — otherwise, our lives would be much easier! The difference between γ_1 and β_1 is called **omitted variable bias** and is a very important concern in many applications.

The only case where $\gamma_1 = \beta_1$ is when $\Gamma_{21}\beta_2 = 0$. The main cases where this can happen are when either $\Gamma_{21} = 0$ or $\beta_2 = 0$. $\Gamma_{21} = 0$ if $\mathbb{E}[X_1X'_2] = 0$ which would be the case if X_1 and X_2 are uncorrelated. $\beta_2 = 0$ occurs when the coefficient on X_2 in Equation 2 is equal to 0. In words, the cases where you can recover β_1 while only using the short regression are (i) if the omitted variables are uncorrelated with the included variables or (ii) if the omitted variables have no effect on the

outcome.

Side-Comment: There are a number of cases where you might be able to figure out the sign of the omitted variable bias. The textbook gives the following simple example. Consider the case where Y is a person's earnings, X_1 is a person's years of education, and X_2 is a person's "ability", and where you are interested in β_1 (the coefficient on years of education). However, suppose that ability is not observed. In this case, it might be reasonable to suppose that $\beta_2 > 0$ (i.e., that, conditional on years of education, individuals with higher ability tend to have higher earnings) and that $\Gamma_{21} > 0$ (i.e., that higher ability is positively correlated with more education). Under these conditions, it would be the case that $\gamma_1 > \beta_1$. This discussion suggests that a regression that only includes years of education would overestimate the effect of years of education relative to a model that included both education and ability. This sort of argument is quite common in applied work — something like: "even though we are not able to control for some important variable, its correlation with the observed variable of interest and likely sign in the long regression indicate that the estimate of our coefficient of interest is likely a lower (or upper) bound."

3. Frisch-Waugh-Lovell Theorem

H: 3.16

I will start with the classical discussion of FWL in terms of estimated regression coefficients, and then turn to the population version that is emphasized in the review. There are a large number of cases where we may be more interested in some of the regression parameters than others (e.g., the treatment effects discussion that we had earlier this semester), so it's useful to have some specific expressions for subsets of the parameters. For this, let's partition $\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2]$ and, likewise, $\beta = (\beta'_1, \beta'_2)'$. Using this notation, we can immediately write

$$\begin{aligned}\mathbf{Y} &= \mathbf{X}\hat{\beta} + \hat{\mathbf{e}} \\ &= [\mathbf{X}_1 \quad \mathbf{X}_2] \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \\ &= \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \hat{\mathbf{e}}\end{aligned}$$

Recall that, $\hat{\beta}_1$ and $\hat{\beta}_2$ minimize the sum of squared residuals

$$\begin{aligned}(\hat{\beta}'_1, \hat{\beta}'_2)' &= \underset{b_1, b_2}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - X'_{i1}b_1 - X'_{i2}b_2)^2 \\ &= \underset{b_1, b_2}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{X}_1b_1 - \mathbf{X}_2b_2)'(\mathbf{Y} - \mathbf{X}_1b_1 - \mathbf{X}_2b_2)\end{aligned}$$

If we are just focused on $\hat{\beta}_1$, we can alternatively express this as

$$\hat{\beta}_1 = \underset{b_1}{\operatorname{argmin}} \left\{ \min_{b_2} (\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}b_2)'(\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}b_2) \right\} \quad (3)$$

This sort of nested minimization is often referred to as “concentrating out” b_2 and is a fairly common estimation strategy (it doesn’t really apply here, but there are some cases where this sort of step may lead to estimators that are notably less computationally complex). The idea here is roughly that we can minimize the overall function by first minimizing it with respect to b_2 (treating b_1 as fixed). This results in our recovering $\hat{\beta}_2(b_1)$ (that is the value of b_2 that minimizes the objective function for a given value of b_1). Then, we can fully minimize the function by taking $\hat{\beta}_1$ to be the value of b_1 that minimizes the objective function taking into account $\hat{\beta}_2(b_1)$. [Also, notice that the inside minimization uses “min” rather than “argmin” because we are still interested in minimizing the objective function itself which is (obviously) quite different from minimizing $\hat{\beta}_2(b_1)$.]

Let’s focus on the inside minimization first. For the inside minimization, we treat b_1 as being fixed and the value of b_2 that minimizes this expression will be a function of b_1 ; I’ll call the value of b_2 that minimizes $\hat{\beta}_2(b_1)$. The inside minimization just amounts to just a regression of $\mathbf{Y} - \mathbf{X}_1\beta_1$ on \mathbf{X}_2 which implies that

$$\hat{\beta}_2(b_1) = (\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 (\mathbf{Y} - \mathbf{X}_1\beta_1)$$

Notice that, from the inside minimization problem, we are not directly interested in $\hat{\beta}_2(b_1)$, but rather the value of the function at $\hat{\beta}_2(b_1)$ (this is because of the “min” rather than “argmin”). This means that the term inside the large curly braces in Equation 3 comes from plugging in this value of $\hat{\beta}_2(b_1)$, i.e.,

$$\min_{b_2} (\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}b_2)'(\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}b_2) = (\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}\hat{\beta}_2(b_1))'(\mathbf{Y} - \mathbf{X}b_1 - \mathbf{X}\hat{\beta}_2(b_1))$$

Moreover, notice that

$$\begin{aligned} \mathbf{Y} - \mathbf{X}_1\beta_1 - \mathbf{X}_2\hat{\beta}_2(b_1) &= \mathbf{Y} - \mathbf{X}_1\beta_1 - \underbrace{\mathbf{X}_2(\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2}_{\mathbf{P}_2} (\mathbf{Y} - \mathbf{X}_1\beta_1) \\ &= (\mathbf{I}_n - \mathbf{P}_2)(\mathbf{Y} - \mathbf{X}_1\beta_1) \\ &= \mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1\beta_1) \end{aligned}$$

where $\mathbf{P}_2 := \mathbf{X}_2(\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2$ and $\mathbf{M}_2 := (\mathbf{I}_n - \mathbf{P}_2)$. Therefore, the inside term in Equation 3 can be written as

$$\begin{aligned} \min_{b_2} (\mathbf{Y} - \mathbf{X}\beta_1 - \mathbf{X}\beta_2)'(\mathbf{Y} - \mathbf{X}\beta_1 - \mathbf{X}\beta_2) &= (\mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1\beta_1))'(\mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1\beta_1)) \\ &= (\mathbf{Y} - \mathbf{X}_1\beta_1)' \mathbf{M}_2(\mathbf{Y} - \mathbf{X}_1\beta_1) \end{aligned}$$

where the first equality uses the expression from above and the last equality uses that \mathbf{M}_2 is symmetric and idempotent. Now, let's plug this into the outside minimization problem.

$$\begin{aligned}\hat{\beta}_1 &= \underset{b_1}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{X}_1 b_1)' \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 b_1) \\ &= \underset{b_1}{\operatorname{argmin}} \mathbf{Y}' \mathbf{M}_2 \mathbf{Y} - 2b_1' \mathbf{X}_1' \mathbf{M}_2 \mathbf{Y} + b_1' \mathbf{X}_1' \mathbf{M}_2 \mathbf{X}_1 b_1\end{aligned}$$

Taking the derivative of the right hand side and setting equal to 0, we have that

$$0 = -2\mathbf{X}_1' \mathbf{M}_2 \mathbf{Y} + 2\mathbf{X}_1' \mathbf{M}_2 \mathbf{X}_1 \hat{\beta}_1$$

which implies that

$$\hat{\beta}_1 = (\mathbf{X}_1' \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_2 \mathbf{Y}$$

The arguments above are symmetric, so you could make the same sorts of calculations and derive a similar result for $\hat{\beta}_2$.

Residual Regression

H: 3.18

The previous result is very closely related to a famous result in econometrics called the Frisch, Waugh, Lovell Theorem. In particular, from the previous expression for $\hat{\beta}_1$, we have that

$$\begin{aligned}\hat{\beta}_1 &= (\mathbf{X}_1' \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_2 \mathbf{Y} \\ &= (\mathbf{X}_1' \mathbf{M}_2' \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_2' \mathbf{M}_2 \mathbf{Y} \\ &= ((\mathbf{M}_2 \mathbf{X}_1)' \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{M}_2 \mathbf{X}_1)' \mathbf{M}_2 \mathbf{Y} \\ &= (\tilde{\mathbf{X}}_1' \tilde{\mathbf{X}}_1)^{-1} \tilde{\mathbf{X}}_1 \tilde{\mathbf{e}}_2\end{aligned}$$

which uses that \mathbf{M}_2 is symmetric and idempotent and where $\tilde{\mathbf{X}}_1 := \mathbf{M}_2 \mathbf{X}_1$ (i.e., the residuals from a regression of \mathbf{X}_1 on \mathbf{X}_2) and $\tilde{\mathbf{e}}_2 := \mathbf{M}_2 \mathbf{Y}$ (i.e., the residuals from the regression of \mathbf{Y} on \mathbf{X}_2).

This implies an algebraic equivalence between $\hat{\beta}_1$ from the regression of \mathbf{Y} on \mathbf{X}_1 and \mathbf{X}_2 and the following estimation procedure:

1. Regress \mathbf{Y} on \mathbf{X}_2 and recover the residuals $\tilde{\mathbf{e}}_2$.
2. Regress \mathbf{X}_1 on \mathbf{X}_2 and recover the residuals $\tilde{\mathbf{X}}_1$.
3. Regress $\tilde{\mathbf{e}}_2$ on $\tilde{\mathbf{X}}_1$.

This procedure delivers exactly the same estimate of $\hat{\beta}_1$. That this procedure recovers exactly the same estimate of $\hat{\beta}_1$ is called the Frisch-Waugh-Lovell Theorem.

This result gives a nice interpretation to the estimates of $\hat{\beta}_1$. It is equivalent to a regression of \mathbf{Y} on \mathbf{X}_1 after “partialling out” (i.e., removing the effect of \mathbf{X}_2 on both \mathbf{Y} and \mathbf{X}_1). Besides that, the FWL Theorem is computationally useful in some important cases too such as some of the panel data approaches that we’ll consider later in the semester.

Population Version of FWL

H: 2.23

A population version of FWL is given in H: 2.23. This is the version of FWL that we will use several times this semester, and I think the arguments are a little bit easier to follow. For simplicity (and because it is the leading case), let’s consider the case where X_1 is scalar and write

$$Y = X_1\beta_1 + X'_2\beta_2 + e$$

where $\mathbb{E}[Xe] = 0$. Now, consider the projection of X_1 on X_2 , that is,

$$X_1 = X'_2\lambda + v$$

where $\mathbb{E}[X_2v] = 0$. Now, notice that

$$\begin{aligned}\mathbb{E}[vY] &= \mathbb{E}[vX_1]\beta_1 + \underbrace{\mathbb{E}[vX'_2]\beta_2}_{=0} + \mathbb{E}[ve] \\ &= \mathbb{E}[v(X'_2\lambda + v)]\beta_1 + \mathbb{E}[(X_1 - X'_2\lambda)e] \\ &= \mathbb{E}[v^2]\beta_1\end{aligned}$$

where the second equality holds by substituting for X_1 in the first term and for v in the last term, and the last equality holds because $\mathbb{E}[X_2v] = 0$ and because $\mathbb{E}[X_1e] = 0$ and $\mathbb{E}[X_2e] = 0$. This implies that

$$\beta_1 = \frac{\mathbb{E}[vY]}{\mathbb{E}[v^2]}$$

which is one of the results that was emphasized in the 8070 review.

Next, notice that we can write the linear projection of Y on X_2

$$Y = X'_2\gamma + u$$

with $\mathbb{E}[X_2u] = 0$. Substituting this in to the previous expression for β_1 , we have that

$$\beta_1 = \frac{\mathbb{E}[v(X'_2\gamma + u)]}{\mathbb{E}[v^2]} = \frac{\mathbb{E}[vu]}{\mathbb{E}[v^2]}$$

where the second equality holds because $\mathbb{E}[X_2 v] = 0$. This is exactly a population version of FWL as it is the population version of the linear projection of u (the projection error from projecting Y on X_2) on v (the projection error from projecting X_1 on X_2).