

ECMA31000: Introduction to Empirical Analysis

Linear Regression; IV

Joe Hardwick

University of Chicago

Autumn 2021

Outline

- Last time:
 - OLS Estimation
 - Partitioned Regression
- Today:
 - Properties of the OLS estimator
 - Measures of Fit
 - Hypothesis testing in the linear regression model
 - Introduction to IV: Endogeneity.

Definitions

- Suppose (y, x, u) satisfy

$$y = x'\beta + u; \quad E(xu) = 0.$$

- Suppose we have an iid sample $\{y_i, x_i\}_{i=1}^n$.
- Assume a unique ordinary least squares estimator exists:

$$\hat{\beta}_n = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i.$$

Finite Sample properties of OLS

- Now make the additional assumption that $E(u_i|x_i) = 0$.
- Since (u_i, x_i) is independent of x_j for $j \neq i$, we obtain for each $i \leq n$:

$$0 = E(u_i|x_i) = E(u_i|x_1, \dots, x_n).$$

- This is equivalently written as

$$E(U|X) = 0.$$

- This is in turn equivalent to

$$E(Y|X) = X\beta.$$

Finite Sample properties of OLS: Bias

- We have

$$\begin{aligned} E(\hat{\beta}_n | X) &= E\left((X'X)^{-1} X' Y | X\right) \\ &= (X'X)^{-1} X' E(Y|X) \\ &= (X'X)^{-1} X' X \beta \\ &= \beta. \end{aligned}$$

- By the Law of Iterated Expectation:

$$E(\hat{\beta}_n) = E\left(E(\hat{\beta}_n | X)\right) = \beta,$$

so $\hat{\beta}_n$ is unbiased.

Finite Sample properties of OLS: Variance

- Suppose in addition that $\text{Var}(u_i|x_i) = \sigma^2$.
- Since the error variance does not depend on the value of x_i , it is called homoskedastic.
- If the error variance depends on x_i , u_i is heteroskedastic.
- Under homoskedasticity,

$$\text{Var}(U|X) = \sigma^2 I_n,$$

because the (i,j) entry is given by

$$\text{Var}(U|X)_{i,j} = \text{E}(u_i u_j | X) = \begin{cases} \text{E}(u_i^2 | x_i) & i = j \\ \text{E}(u_i u_j | x_i, x_j) & i \neq j \end{cases}$$

Finite Sample properties of OLS: Variance

- We have $E(u_i^2|x_i) = \sigma^2$, and

$$\begin{aligned} E(u_i u_j | x_i, x_j) &= E(u_i E(u_j | u_i, x_i, x_j) | x_i, x_j) \\ &= E(u_i E(u_j | x_j) | x_i, x_j) = 0. \end{aligned}$$

- Under heteroskedasticity, the off diagonal elements will still be 0 but $E(u_i^2|x_i) = \sigma(x_i)^2$, so

$$Var(U|X) = \begin{pmatrix} \sigma(x_1)^2 & 0 & 0 & 0 \\ 0 & \sigma(x_2)^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \sigma(x_n)^2 \end{pmatrix} := \Omega.$$

Finite Sample properties of OLS: Variance

- We now conclude that

$$\begin{aligned} \text{Var}(\hat{\beta}_n | X) &= \text{Var}\left(\beta + (X'X)^{-1}X'U | X\right) \\ &= (X'X)^{-1}X'\Omega X(X'X)^{-1}. \end{aligned}$$

- Under homoskedasticity, $\Omega = \sigma^2 I_n$, and so the variance becomes

$$\begin{aligned} \text{Var}(\hat{\beta}_n | X) &= \sigma^2 (X'X)^{-1} X' I_n X (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1}. \end{aligned}$$

Finite Sample properties of OLS: Variance

- Suppose $E(u_i|x_i) = 0$ and the error term u_i is homoskedastic.
- Under these assumptions, the OLS estimator is the “best linear unbiased estimator”, which means it has the “smallest” variance in the class of linear estimators that are also unbiased conditional on X . This result is known as the *Gauss-Markov Theorem*.
- Our goal is to show that

$$\sigma^2 (X'X)^{-1}$$

is “smaller” than the conditional variance of any other estimator $\tilde{\beta} = A(X) Y$ which also satisfies $E(\tilde{\beta}|X) = \beta$.

- Precisely: $Var(\tilde{\beta}|X) - Var(\hat{\beta}_n|X) = D$ is a positive-semidefinite matrix.

Gauss-Markov Theorem

- In particular, let r be a $(k+1) \times 1$ vector, and let $r'\beta$ denote a linear combination of the β_j .
- The Gauss-Markov theorem implies that $r'\hat{\beta}$ is the Best Linear Unbiased Estimator of $r'\beta$, since

$$\begin{aligned} \text{Var}(r'\tilde{\beta}|X) - \text{Var}(r'\hat{\beta}|X) &= r' \text{Var}(\tilde{\beta}|X) r - r' \text{Var}(\hat{\beta}|X) r \\ &= r' [\text{Var}(\tilde{\beta}|X) - \text{Var}(\hat{\beta}|X)] r \\ &= r'Dr \geq 0. \end{aligned}$$

- For example, if $r = \left(0, 0, \dots, 0, \underbrace{1}_{(j+1)\text{st position}}, 0, \dots, 0\right)$, then $r'\hat{\beta} = \hat{\beta}_j$, and so $\text{Var}(\hat{\beta}_j|X) \leq \text{Var}(\tilde{\beta}_j|X)$.

Gauss-Markov Theorem

- First note that a linear estimator $\tilde{\beta}$ of β satisfies

$$\tilde{\beta} = A(x_1, \dots, x_n) Y = AY,$$

for some $(k + 1) \times N$ matrix $A = A(x_1, \dots, x_n)$ depending only on the sample $\{x_i\}_{i=1}^n$.

- No matter what the true value of the unknown parameter β is, $\tilde{\beta}$ must also satisfy

$$E(\tilde{\beta}|X) = \beta.$$

- This implies

$$E(\tilde{\beta}|X) = E(AY|X) = AX\beta = \beta.$$

Gauss-Markov Theorem

- Since the final equality must hold for any β , it must be the case that $AX = I_{k+1}$.
- Next, we compute the variance of AY conditional on X :

$$\begin{aligned} \text{Var}(AY|X) &= A\text{Var}(Y|X)A' \\ &= A\text{Var}(U|X)A' \\ &= \sigma^2 AA'. \end{aligned}$$

- If $A = (X'X)^{-1}X'$, we obtain the OLS estimator, with variance $\sigma^2(X'X)^{-1}$.

Gauss-Markov Theorem

- It remains to show that $\text{Var}(\hat{\beta}_{OLS}|X) \leq \text{Var}(\tilde{\beta}|X)$, where \leq means that

$$\sigma^2 AA' - \sigma^2 (X'X)^{-1}$$

is a positive semi-definite matrix for any A such that $AX = I_{k+1}$.

- To this end, define

$$C = A - (X'X)^{-1} X'.$$

Gauss-Markov Theorem

- Note that

$$\begin{aligned}AA' - (X'X)^{-1} &= \left(C + (X'X)^{-1}X'\right)\left(C + (X'X)^{-1}X'\right)' \\&\quad - (X'X)^{-1} \\&= CC' + (X'X)^{-1}X'C' + CX(X'X)^{-1} \\&= CC',\end{aligned}$$

where the final equality holds since

$$CX = AX - (X'X)^{-1}X'X = I_{k+1} - I_{k+1} = 0.$$

- The conclusion follows because CC' is always positive semi-definite.

Questions?

Large Sample properties of OLS

- Now drop the assumptions that $E(u_i|x_i) = 0$ and $Var(u_i|x_i) = \sigma^2$. Rewrite the model as

$$y_i = x_i' \beta + u_i; \quad E(u_i x_i) = 0,$$

and suppose $E(xx')$ exists and is invertible.

- The OLS estimator is consistent, because

$$\frac{1}{n} \sum_{i=1}^n x_i x_i' \xrightarrow{\text{a.s.}} E(xx');$$

$$\frac{1}{n} \sum_{i=1}^n x_i y_i \xrightarrow{\text{a.s.}} E(xy),$$

by the SLLN. As we know, the convergence is joint, so

$$\left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i \xrightarrow{\text{a.s.}} E(xx')^{-1} E(xy) = \beta.$$

Asymptotic Normality of OLS

- Maintain the assumption that $E(xx')$ exists, and assume $Var(xu) = E(u^2xx')$ exists also. Then:

$$\sqrt{n} (\hat{\beta}_n - \beta) \xrightarrow{d} \mathcal{N}(0, \Sigma),$$

where $\Sigma = E(xx')^{-1} Var(xu) E(xx')^{-1}$.

- This follows because

$$\sqrt{n} (\hat{\beta}_n - \beta) = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i.$$

Asymptotic Normality of OLS

- Since the sequence of vectors $\{(y_i, x_i')'\}_{i \geq 1}$ is iid, the sequence $\{x_i(y_i - x_i'\beta)\}_{i \geq 1}$ is iid. Therefore, the CLT implies:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i \xrightarrow{d} \mathcal{N}(0, \text{Var}(xu)).$$

- Applying Slutsky's Theorem gives

$$\left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i \xrightarrow{d} \mathcal{N}(0, \Sigma),$$

as desired.

Estimation of Σ

- Deriving asymptotic normality of $\hat{\beta}_n$ will enable us to test hypotheses about the unknown parameter β .
- Since we do not know Σ , we must construct a consistent estimator of it to yield an asymptotic distribution whose quantiles are known so we may conduct tests.
- First, suppose $E(u|x) = 0$ and $Var(u|x) = \sigma^2$. Then

$$\begin{aligned}Var(xu) &= E(u^2 xx') \\&= E(E(u^2|x) xx') \\&= \sigma^2 E(xx').\end{aligned}$$

- It follows that

$$\Sigma = \sigma^2 E(xx')^{-1}.$$

Estimation of Σ

- A natural estimator of $E(xx')^{-1}$ is $(\frac{1}{n} \sum_{i=1}^n x_i x'_i)^{-1}$.
- It remains to find a consistent estimator of σ^2 . We use

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 = \frac{1}{n} \|M_X Y\|^2,$$

since $M_X Y$ is the vector of residuals of a regression of y on x .

- Next, note that $M_X Y = M_X (X\beta + U) = M_X U$, so

$$\begin{aligned}\|M_X Y\|^2 &= \|M_X U\|^2 \\ &= U' M_X' M_X U \\ &= U' M_X U \\ &= U' U - U' P_X U.\end{aligned}$$

Estimation of Σ

- Finally, note that

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n u_i^2 - \left(\frac{1}{n} \sum_{i=1}^n u_i x'_i \right) \left(\frac{1}{n} \sum_{i=1}^n x_i x'_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i u_i \right) \\ &\xrightarrow{a.s.} \sigma^2 - 0 \cdot \mathbb{E}(xx')^{-1} \cdot 0 \\ &= \sigma^2,\end{aligned}$$

by the continuous mapping theorem.

- In summary, a consistent estimator of Σ is given by

$$\hat{\Sigma} = \hat{\sigma}^2 \left(\frac{1}{n} \sum_{i=1}^n x_i x'_i \right)^{-1}.$$

Estimation of Σ

- If we do not assume $E(U|X) = 0$ and $Var(U|X) = \sigma^2 I_n$, Σ does not simplify, and we use

$$\hat{\Sigma} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 x_i x_i' \right) \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1}.$$

- Proving consistency boils down to showing that

$$\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 x_i x_i' \xrightarrow{p} E(u^2 x x').$$

- First, decompose this quantity as

$$\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 x_i x_i' = \frac{1}{n} \sum_{i=1}^n u_i^2 x_i x_i' + \frac{1}{n} \sum_{i=1}^n (\hat{u}_i^2 - u_i^2) x_i x_i'.$$

Estimation of Σ

- We have

$$\frac{1}{n} \sum_{i=1}^n u_i^2 x_i x_i' \xrightarrow{\text{a.s.}} \mathbb{E}(u^2 x x') ,$$

so it remains to show that

$$\frac{1}{n} \sum_{i=1}^n (\hat{u}_i^2 - u_i^2) x_i x_i' = o_p(1) .$$

- We prove that every element of this matrix is $o_p(1)$.

Estimation of Σ

- Fix j, k element $\frac{1}{n} \sum_{i=1}^n (\hat{u}_i^2 - u_i^2) x_{i,j} x_{i,k}$, and observe that

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n (\hat{u}_i^2 - u_i^2) x_{i,j} x_{i,k} \right| &\leq \frac{1}{n} \sum_{i=1}^n |(\hat{u}_i^2 - u_i^2)| |x_{i,j} x_{i,k}| \\ &\leq \max_{1 \leq i \leq n} |\hat{u}_i^2 - u_i^2| \frac{1}{n} \sum_{i=1}^n |x_{i,j} x_{i,k}|. \end{aligned}$$

- Since

$$\frac{1}{n} \sum_{i=1}^n |x_{i,j} x_{i,k}| \xrightarrow{\text{a.s.}} \mathbb{E}(|x_{i,k} x_{i,k}|),$$

this term is $O_p(1)$. Therefore, it is sufficient to prove

$$\max_{1 \leq i \leq n} |\hat{u}_i^2 - u_i^2| = o_p(1).$$

Estimation of Σ

- To this end, note that

$$\hat{u}_i - u_i = x'_i (\beta - \hat{\beta}_n),$$

so since $\hat{u}_i = y_i - x'_i \hat{\beta}_n = x'_i (\beta - \hat{\beta}_n) + u_i$, we have:

$$\begin{aligned} |\hat{u}_i^2 - u_i^2| &= \left| x'_i (\beta - \hat{\beta}_n) (\hat{u}_i + u_i) \right| \\ &= \left| x'_i (\beta - \hat{\beta}_n) \left(x'_i (\beta - \hat{\beta}_n) + 2u_i \right) \right| \\ &= \left| \left(x'_i (\beta - \hat{\beta}_n) \right)^2 + 2u_i x'_i (\beta - \hat{\beta}_n) \right| \\ &\leq \left| \left(x'_i (\beta - \hat{\beta}_n) \right)^2 \right| + 2 \left| u_i x'_i (\beta - \hat{\beta}_n) \right|. \end{aligned}$$

Estimation of Σ

- Next, using the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned} \max_{1 \leq i \leq n} |\hat{u}_i^2 - u_i^2| &\leq \left\| \beta - \hat{\beta}_n \right\|^2 \max_{1 \leq i \leq n} \|x_i\|^2 \\ &\quad + 2 \left\| \beta - \hat{\beta}_n \right\| \max_{1 \leq i \leq n} \|x_i u_i\| \\ &= \left\| \sqrt{n} (\beta - \hat{\beta}_n) \right\|^2 \frac{\max_{1 \leq i \leq n} \|x_i\|^2}{n} \\ &\quad + 2 \left\| \sqrt{n} (\beta - \hat{\beta}_n) \right\| \frac{\max_{1 \leq i \leq n} \|x_i u_i\|}{\sqrt{n}}. \end{aligned}$$

Estimation of Σ

- Since $\sqrt{n} (\beta - \hat{\beta}_n) = O_p(1)$, we must show:

$$\frac{\max_{1 \leq i \leq n} \|x_i\|^2}{n} = o_p(1);$$
$$\frac{\max_{1 \leq i \leq n} \|x_i u_i\|}{\sqrt{n}} = o_p(1).$$

Lemma

Let $\{Z_i\}_{i \geq 1}$ be a sequence of identically distributed random vectors such that $E(\|Z_i\|^r) < \infty$. Then

$$\frac{\max_{1 \leq i \leq n} \|Z_i\|}{n^{1/r}} \xrightarrow{p} 0.$$

Estimation of Σ

Proof.

Fix $\epsilon > 0$ and note that

$$\begin{aligned} \text{P} \left(\max_{1 \leq i \leq n} \|Z_i\| > \epsilon n^{1/r} \right) &= \text{P} \left(\cup_{i=1}^n \{\|Z_i\|^r > \epsilon^r n\} \right) \\ &\leq \sum_{i=1}^n \text{P} (\|Z_i\|^r > \epsilon^r n) \\ &= \sum_{i=1}^n \text{P} (\|Z_i\|^r \mathbf{1}(\|Z_i\|^r > \epsilon^r n) > \epsilon^r n) \\ &\leq \frac{1}{n\epsilon^r} \sum_{i=1}^n \text{E} (\|Z_i\|^r \mathbf{1}(\|Z_i\|^r > \epsilon^r n)) \\ &= \frac{1}{\epsilon^r} \text{E} (\|Z_i\|^r \mathbf{1}(\|Z_i\|^r > \epsilon^r n)) \\ &\rightarrow 0. \end{aligned}$$

□

Estimation of Σ

- The second equality holds because

$$\{\|Z_i\|^r > \epsilon^r n\} = \{\|Z_i\|^r \mathbf{1}(\|Z_i\|^r > \epsilon^r n) > \epsilon^r n\}.$$

- The second inequality follows by Markov's inequality and the third because the Z_i have identical distribution.
- Convergence to 0 holds because $E(\|Z_i\|^r) < \infty$.
- Finally, since $E(\|x\|^2) < \infty$ and $E(\|ux\|^2) < \infty$, we conclude that

$$\max_{1 \leq i \leq n} |\hat{u}_i^2 - u_i^2| = o_p(1),$$

so

$$\frac{1}{n} \sum_{i=1}^n (\hat{u}_i^2 - u_i^2) x_i x'_i = o_p(1).$$

Questions?

Coefficient of Determination

- The coefficient of determination, or R^2 , is a measure of how well a linear model estimated by OLS fits a given dataset.
- R^2 is defined by

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - x'_i \hat{\beta})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\|M_{xy}\|^2}{\|M_c y\|^2},$$

where M_c is the residual maker of the matrix c consisting of a single column of ones. (Regressing y on a constant yields the OLS estimate \bar{y}).

- R^2 measures how well the model fits the data relative to a model with only a constant and no other regressors.

Coefficient of Determination

- We may equivalently write

$$R^2 := 1 - \frac{\|M_{Xy}\|^2}{\|M_{cy}\|^2} = 1 - \frac{\|M_X M_{cy}\|^2}{\|M_{cy}\|^2},$$

since $M_{Xy} = M_X M_{cy}$. To see this note that

$$\begin{aligned} M_X M_{cy} &= M_X (I - P_c) y \\ &= M_{Xy} - M_X P_c y \\ &= M_{Xy}, \end{aligned}$$

where the last equality follows because $P_c y$ is just the projection of y onto the span of the first column of X (which is a column of ones), it follows that $P_c y \in s(X)$. Therefore, $M_X P_c y = 0$, since the residual of a projection of $P_c y \in s(X)$ onto $s(X)$ is 0.

Coefficient of Determination

- It follows that we can also write

$$R^2 = 1 - \frac{\|M_X M_c y\|^2}{\|M_c y\|^2} = \frac{\|P_X M_c y\|^2}{\|M_c y\|^2}$$

- The second equality follows by the pythagorean theorem since $P_X M_c y$ is orthogonal to $M_X M_c y$:

$$\begin{aligned}\|P_X M_c y\|^2 + \|M_X M_c y\|^2 &= \|P_X M_c y + M_X M_c y\|^2 \\ &= \|M_c y\|^2.\end{aligned}$$

- These formulae show that $0 \leq R^2 \leq 1$.

Coefficient of Determination

- $\|M_c y\|^2 = \sum_{i=1}^n (y_i - \bar{y})^2 := TSS$ is called the total sum of squares.
- $\|M_{xy} y\|^2 = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 := SSR$ is the sum of squared residuals.
- $\|P_X M_c y\|^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 := ESS$ is the explained sum of squares.
- $R^2 = 1$ iff $SSR = 0$, which means $\hat{U} = 0$.
- $R^2 = 0$ iff $SSR = SST$, which says $\hat{\beta} = (\bar{y}, 0, \dots, 0)$, and the model fits the data no better than a constant would.

Coefficient of Determination

- High R^2 does not mean the linear model is “correct” in the sense that it provides a causal explanation of the relationship between x and y .
- On the other hand, low R^2 does not preclude a causal relationship between x and y in a linear model.
- R^2 is an estimator of the population R^2 , defined by

$$R_{pop}^2 = 1 - \frac{Var(u)}{Var(y)},$$

since $\frac{SSR}{n}$ is an estimate of $Var(u)$ and $\frac{TSS}{n}$ is an estimate of $Var(y)$.

Coefficient of Determination

- Note that R^2 and R_{pop}^2 always weakly increase when a regressor is added.
- Adjusted R^2 penalizes additional regressors, and is defined by

$$\bar{R}^2 = 1 - \frac{n-1}{n-k-1} \frac{SSR}{TSS} \leq R^2.$$

- Adding in a regressor which does not improve in-sample fit will cause \bar{R}^2 to decrease.

Questions?

Finite Sample Inference

- Suppose that $Y|X \sim \mathcal{N}(X\beta, \sigma^2 I_n)$, so that

$$Y = X\beta + U; \quad E(U|X) = 0, \quad \text{Var}(U|X) = \sigma^2 I_n.$$

- These assumptions (linear conditional mean, homoskedasticity) would be satisfied if the iid observations of (y_i, x_i) are multivariate normal.
- We showed in PSET 5 that

$$\hat{\beta}_n | X \sim \mathcal{N}\left(\beta, \sigma^2 (X'X)^{-1}\right).$$

Finite Sample Inference

- Let e_j be a $(k + 1) \times 1$ column vector consisting of zeroes, except in position j , where there is a 1.
- Then

$$e_j' \hat{\beta}_n | X = \hat{\beta}_{n,j} | X \sim \mathcal{N} \left(\beta_j, \left[\sigma^2 (X'X)^{-1} \right]_{j,j} \right),$$

where the subscript j, j indicates the (j, j) entry of the associated matrix.

- We can also derive the distribution of

$$\hat{\sigma}^2 = \frac{SSR}{n - k - 1}.$$

This estimate is unbiased and consistent (see PSET 7).

Finite Sample Inference

- We will also see that $\hat{\beta}_n$ and $\hat{\sigma}^2$ are independent conditional on X , and:

$$\frac{(n - k - 1) \hat{\sigma}^2}{\sigma^2} \sim \chi_{n-k-1}^2.$$

- If σ^2 were known, we could write

$$\frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{e'_j (X'X)^{-1} e_j}} \sim \mathcal{N}(0, 1),$$

though typically we must replace it with $\hat{\sigma}^2$, which yields:

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{e'_j (X'X)^{-1} e_j}} \sim t_{n-k-1}.$$

- $\hat{\sigma} \sqrt{e'_j (X'X)^{-1} e_j}$ is called the standard error of $\hat{\beta}_j$, denoted $se(\hat{\beta}_j)$.

Finite Sample Inference

- To test $H_0 : \beta_j = \beta_j^0$, we can form the t-statistic:

$$T_n = \left| \frac{\hat{\beta}_j - \beta_j^0}{se(\hat{\beta}_j)} \right|,$$

which has a t_{n-k-1} distribution under H_0 .

- Let $t_{n-k-1, 1-\alpha/2}$ be the $1 - \alpha/2$ quantile of the t_{n-k-1} distribution. The test

$$\phi_n = \mathbf{1}(T_n > t_{n-k-1, 1-\alpha/2})$$

has rejection probability equal to α when H_0 is true, and rejection probability strictly greater than α when H_0 is false.

- Therefore ϕ_n is of size α .

Finite Sample Inference

- Let F denote the CDF of t_{n-k-1} . Conditional on the data, the p -value is given by

$$\begin{aligned}\hat{p}_n &= \inf_{\alpha \in (0,1)} \{ T_n > t_{n-k-1, 1-\alpha/2} \} \\ &= \inf_{\alpha \in (0,1)} \{ T_n > F^{-1}(1 - \alpha/2) \}.\end{aligned}$$

Since F^{-1} is strictly increasing and continuous, we obtain:

$$\begin{aligned}\alpha &= 2[1 - F(T_n)] \\ &= 2F(-T_n).\end{aligned}$$

Questions?

Large sample inference

- Joint normality allows us to obtain exact finite sample distributions for these statistics, but we may also appeal to asymptotic normality.
- A test is called asymptotically of size α if

$$\lim_{n \rightarrow \infty} \beta_n(\theta) \leq \alpha \quad \text{for all } \theta \in \Theta_0.$$

- Now suppose

$$\sqrt{n} (\hat{\beta}_n - \beta) \xrightarrow{d} \mathcal{N}(0, V),$$

where $V = E(xx')^{-1} E(u^2 xx') E(xx')^{-1}$. Suppose also that V is non-singular and $\hat{V}_n \xrightarrow{P} V$ is a consistent estimator of V .

Testing a single linear restriction

- Consider testing

$$H_0 : r'\beta = c \quad \text{vs.} \quad H_1 : r'\beta \neq c,$$

where r is some specified vector in \mathbb{R}^{k+1} , and c is a scalar.

- By the CMT:

$$\sqrt{n} \left(r' \hat{\beta}_n - r' \beta \right) \xrightarrow{d} \mathcal{N}(0, r' V r),$$

and so by Slutsky's theorem:

$$\frac{\sqrt{n} \left(r' \hat{\beta}_n - r' \beta \right)}{\sqrt{r' \hat{V}_n r}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Testing a single linear restriction

- It follows that under H_0 , the test statistic

$$T_n = \frac{\sqrt{n} (r' \hat{\beta}_n - c)}{\sqrt{r' \hat{V}_n r}} \xrightarrow{d} \mathcal{N}(0, 1).$$

- The test we use is $\phi_n = \mathbf{1}(|T_n| > z_{1-\alpha/2})$. It is of asymptotic size α , because under H_0 :

$$\begin{aligned} Pr(\text{Reject } H_0) &= P(|T_n| > z_{1-\alpha/2}) \\ &= P(T_n < -z_{1-\alpha/2}) + P(T_n > z_{1-\alpha/2}) \\ &\rightarrow \Phi(-z_{1-\alpha/2}) + 1 - \Phi(z_{1-\alpha/2}) \\ &= \frac{\alpha}{2} + 1 - \left(1 - \frac{\alpha}{2}\right) = \alpha. \end{aligned}$$

- Use $\phi_n = \mathbf{1}(T_n > z_{1-\alpha})$ for testing $H_0 : r'\beta \leq c$ vs. $H_1 : r'\beta > c$.

Asymptotic Confidence Set for $r'\beta$

- It follows by definition of convergence in distribution that for any value of $r'\beta$:

$$\lim_{n \rightarrow \infty} P_{r'\beta} \left(\left| \frac{\sqrt{n} (r' \hat{\beta}_n - r'\beta)}{\sqrt{r' \hat{V}_n r}} \right| \leq z_{1-\alpha/2} \right) = 1 - \alpha.$$

- Since $z_{\alpha/2} = -z_{1-\alpha/2}$ by symmetry of the standard normal about 0, rearranging yields that

$$C_n = \left[r' \hat{\beta}_n - z_{1-\alpha/2} \sqrt{\frac{r' \hat{V}_n r}{n}}, r' \hat{\beta}_n + z_{1-\alpha/2} \sqrt{\frac{r' \hat{V}_n r}{n}} \right]$$

is an asymptotic $1 - \alpha$ confidence interval for $r'\beta$.

Testing Multiple Linear Restrictions

- Consider testing

$$H_0 : R\beta = c \quad \text{vs.} \quad R\beta \neq c,$$

where R is a $p \times (k + 1)$ -dimensional matrix of full row rank and c is a $p \times 1$ vector.

- The full rank condition means none of our restrictions are redundant.
- By the CMT:

$$\sqrt{n} (R\hat{\beta}_n - R\beta) \xrightarrow{d} \mathcal{N}(0, RVR'),$$

where RVR' is full rank (because R and V are), and hence positive definite (because V is).

Testing Multiple Linear Restrictions

- To see that RVR' is positive definite, note that if $a \neq 0$, $R'a \neq 0$, so

$$(R'a)' V (R'a) > 0,$$

because V is positive definite.

- A positive definite and symmetric matrix A has a square root $A^{1/2}$ with inverse $A^{-1/2} = (A^{-1})^{1/2}$.
- It follows by Slutsky's Theorem that

$$\left(R\hat{V}_n R' \right)^{-1/2} \sqrt{n} \left(R\hat{\beta}_n - R\beta \right) \xrightarrow{d} \mathcal{N}(0, I_p).$$

Testing Multiple Linear Restrictions

- It follows that

$$n \cdot \left(R\hat{\beta}_n - R\beta \right)' \left(R\hat{V}_n R' \right)^{-1} \left(R\hat{\beta}_n - R\beta \right) \xrightarrow{d} \chi_p^2.$$

- Under H_0 ,

$$T_n = n \cdot \left(R\hat{\beta}_n - c \right)' \left(R\hat{V}_n R' \right)^{-1} \left(R\hat{\beta}_n - c \right) \xrightarrow{d} \chi_p^2,$$

and so we reject iff $T_n > \chi_{p,1-\alpha}^2$.

Asymptotic Confidence Set for $R\beta$

- It follows that

$$C_n = \left\{ c \in \mathbb{R}^p : n \cdot (R\hat{\beta}_n - c)' (R\hat{V}_n R')^{-1} (R\hat{\beta}_n - c) \leq \chi_{p,1-\alpha}^2 \right\}$$

is an asymptotic $1 - \alpha$ confidence set for $R\beta$.

- This set is an ellipsoid centered at $R\hat{\beta}_n$, and satisfies

$$P_{R\beta} (R\beta \in C_n) \rightarrow 1 - \alpha.$$

- Taking $R = I_{k+1}$ yields an asymptotic $1 - \alpha$ confidence set for β .

Tests of Non-Linear Restrictions

- Finally, consider testing

$$H_0 : f(\beta) = 0 \quad \text{vs.} \quad H_1 : f(\beta) \neq 0,$$

where $f : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^p$ is continuously differentiable at β .

- Let $D_\beta f(\beta)$ denote the $p \times (k + 1)$ dimensional matrix of partial derivatives of f evaluated at β .
- The Delta Method implies

$$\sqrt{n} \left(f(\hat{\beta}_n) - f(\beta) \right) \xrightarrow{d} \mathcal{N}(0, D_\beta f(\beta) V D_\beta f(\beta)')$$

Tests of Non-Linear Restrictions

- The continuous mapping theorem implies that

$$D_\beta f(\hat{\beta}_n) \hat{V}_n D_\beta f(\hat{\beta}_n)' \xrightarrow{P} D_\beta f(\beta) V D_\beta f(\beta)'$$

- Now assume $D_\beta f(\beta)$ is full row rank. We can construct a statistic with asymptotic χ_p^2 distribution as before.
- Note that $f(\beta) = R\beta$ yields our previous analysis as a special case, since $D_\beta f(\beta) = R$.

Questions?

Introduction to IV

- Let (y, x, u) be a random vector such that y and u are scalar random variables and $x \in \mathbb{R}^{k+1}$.
- Assume the first component of x equals 1:

$$x = (x_0, x_1 \dots, x_k),$$

where $x_0 = 1$.

- Let $\beta = (\beta_0, \dots, \beta_k) \in \mathbb{R}^{k+1}$ be a constant vector of unknown parameters such that

$$y = x' \beta + u.$$

- We *no longer* assume $E(ux) = 0$, so β may not represent the best linear predictor, and therefore not the best predictor either.

Introduction to IV

- We are therefore interpreting this regression as a causal model.
- If $E(ux_j) = 0$ for some j , x_j is exogenous.
- If $E(ux_j) \neq 0$ for some j , x_j is endogenous.
- x_0 can always be made exogenous by shifting β_0 such that $E(x_0 u) = E(u) = 0$.
- Multiply the model by x and take expectations:

$$E(xy) = E(xx')\beta + E(xu).$$

Introduction to IV

- It follows that

$$E(xx')^{-1} E(xy) = \beta + E(xx')^{-1} E(xu).$$

- Therefore,

$$\hat{\beta}_n^{OLS} = \left(\frac{X'X}{n} \right)^{-1} \frac{X'Y}{n} \xrightarrow{a.s.} \beta + E(xx')^{-1} E(xu) \neq \beta.$$

- The OLS estimator is now an inconsistent estimator of β under endogeneity.

Instrumental Variables

- Our goal is to use a random vector $z \in \mathbb{R}^{I+1}$ such that $E(zu) = 0$ to identify β .
- The condition $E(zu) = 0$ is called instrument validity.
(Multivariate version of $Cov(z, u) = 0$)
- First, note that any exogenous component of x is included in z . These components of x are called included instruments.
- The constant 1 is included, since we can always set $E(u) = 0$. So, letting $z_0 = 1$:

$$z = (z_0, z_1, \dots, z_I) \in \mathbb{R}^{I+1}.$$

Instrumental Variables

- How to get β as a function of quantities we can estimate?
Model

$$y = x'\beta + u.$$

- Pre-multiply by z :

$$zy = zx'\beta + zu.$$

- Take expectations:

$$\begin{aligned} E(zy) &= E(zx')\beta + E(zu) \\ &= E(zx')\beta. \end{aligned}$$

- If $I = k$ (exactly as many instruments as regressors), $E(zx')$ is square, so

$$\beta = [E(zx')]^{-1} E(zy).$$

Instrumental Variables

- The components of z are called instrumental variables.
- We further assume that $E(zx')$ has rank $k + 1$. (Instrument relevance/rank condition) (Multivariate version of $\text{Cov}(z, x) \neq 0$).
- Finally, we assume $E(zz') < \infty$ and that there is no perfect collinearity in z .
- A necessary condition for the rank condition is $I \geq k$. This is called the order condition. In other words, we must have as many valid instruments as we have endogenous regressors.

Instrumental Variables: Order Condition

- If $I = k$, the system is exactly identified.
- If $I > k$, the system is overidentified, since we have more instruments than we need to identify β .
- Notice: If x_j is endogenous, it is not an included instrument.
- Given the order condition holds, the rank condition is necessary and sufficient to uniquely determine β .
- Later: What to do with extra instruments? Could throw them out and get an IV estimate, but this is inefficient.

IV Estimator

- We showed under validity and relevance assumptions:

$$\beta = E(zx')^{-1} E(zy).$$

- The sample analog principle yields

$$\frac{1}{n} \sum_{i=1}^n z_i (y_i - x'_i \hat{\beta}_{IV}) = 0,$$

or

$$\begin{aligned}\hat{\beta}_{IV} &= \left(\frac{1}{n} \sum_{i=1}^n z_i x'_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n z_i y_i \right) \\ &\xrightarrow{P} E(zx')^{-1} E(zy) = \beta.\end{aligned}$$

using the LLN and continuous mapping theorem, so the IV estimator is consistent.

Questions?