

Ordinary Least Squares and Variance

$$Y = X\beta + \varepsilon \quad E[\varepsilon | X] = 0$$

Estimate β by $\hat{\beta} = (X'X)^{-1}X'\hat{Y}$

$$E[\hat{\beta}] = \beta$$

$$\text{Var}(\hat{\beta}|X) = \sigma^2(X'X)^{-1} \quad \text{if } \text{Var}(\varepsilon|X) = \sigma^2 I_n$$

Heteroskedasticity

Suppose $\varepsilon_i \sim \text{Var}(\varepsilon|X)$ is something like

- Heteroskedasticity: ε_i 's are independent of each other but depend on X .

e.g. $\text{Var}(\varepsilon_i|X_i) = \sigma^2(X_i)$ for some nonnegative function

- Clustering: ε_j 's are representing groups with shared characteristics and so nonzero covariances

Shows up with "panel data" forms

Two indices $i=(j,t)$ $j=1 \dots J$, $t=1 \dots T$

each unit j comes in groups of size T

ε may have form $\left[\begin{array}{c} \varepsilon_{1,1} \\ \varepsilon_{1,2} \\ \vdots \\ \varepsilon_{1,T} \\ \hline \varepsilon_{2,1} \\ \varepsilon_{2,2} \\ \vdots \\ \varepsilon_{2,T} \\ \hline \vdots \\ \hline \varepsilon_{J,1} \\ \varepsilon_{J,2} \\ \vdots \\ \varepsilon_{J,T} \end{array} \right] := \text{cov}(\varepsilon_{ij}, \varepsilon_{is}) = \sigma_{ij,s} \neq 0$

- Sequential structure?

- Time series or spatial data have autocorrelation
- units arranged in ordered set $\{ \dots \}$
- $$\sigma_{ts+1} = \text{cov}(\varepsilon_t, \varepsilon_{t+1}) \quad \text{may depend on } |s-t|$$
- $$\varepsilon = \begin{bmatrix} \varepsilon_1 & \varepsilon_2 & \dots \\ \vdots & \vdots & \ddots \\ \varepsilon_s & \varepsilon_{s+1} & \dots \\ \vdots & \vdots & \ddots \\ \varepsilon_n & \varepsilon_{n-1} & \dots \end{bmatrix}$$

- What happens to OLS when $\varepsilon \neq \sigma^2 I$

Gauss-Markov no longer true. consequences OLS

ε is positive definite and symmetric

Let's further suppose ε^{-1} exists, will also be positive definite and symmetric

We can define square root of $\varepsilon^{-1} = C' C$

not unique, examples Cholesky decomposition

C is upper triangular, $C_{i,j} = 0$ if $i < j$
 C' is lower triangular

- Can use the spectral decomposition

If M is positive definite and symmetric on $\mathbb{R}^{N \times N}$

$$M = V \vartheta V' \quad \text{where } \vartheta \text{ is diagonal, } \vartheta_{ij} = 0 \text{ if } i \neq j$$

$$V \text{ is orthogonal} \quad V' V = I_N$$

$$V V' = I_N$$

ϑ is nonnegative and entries are the eigenvalues of M

columns of V are eigenvectors

$$\text{ie } MV_i - d_{ii} \mathbb{I}_N = 0 \quad V_i \neq 0,$$

square root ($\equiv D^{\frac{1}{2}}V$) takes square root of diagonal entries

$$d_{ii} = \max_{\substack{\|v\|=1 \\ v \in \mathbb{R}^N}} \sqrt{Mv} \quad \text{expression is } \\ \left\{ v \in \mathbb{R}^N : \|v\| = \sqrt{\sum v_i^2} = 1 \right\} \\ \text{is compact and } \sqrt{Mv} \text{ continuous}$$

$$f(v) = \sqrt{Mv} - \lambda(\sqrt{v} - 1)$$

$$\frac{\partial f}{\partial v} = 0 \text{ at optimum } \sqrt{M} - \lambda \sqrt{v} = 0 \quad \text{FOC} \\ Mv = \lambda v \quad \text{at } v = v_i$$

S - λ is eigenvalue for eigenvector v_i

$$\text{So } v_i^T M v_i = v_i^T \lambda v_i = \lambda v_i^T v_i = \lambda$$

So optimum is $\lambda \equiv d_{ii}$

Defn. Q_i as $\begin{bmatrix} v_i & z \end{bmatrix}$ where z is

$z = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$ i^{th} orthogonal to v_i (defined by projecting basis vectors onto orthogonal complement of space V_i)

$$M Q_i = [d_{ii} v_i \quad Mz]$$

$$\text{So } Q_i^T M Q_i = \begin{bmatrix} d_{ii} - 0 \\ 0 \end{bmatrix} \quad 0's \text{ by orthogonality}$$

Can apply same proof to $Z^T M Z$ which is in $\mathbb{R}^{N \times N}$

and symmetric and positive definite and invertible (times

say) sequence $d_1, d_2, \dots, d_{N,N}$ of eigenvalues

and Q_1, Q_2, \dots mutually reach V

We can define estimators after this OLS

$$\text{Apply } E[C(Y - X\beta) | X] = C E[\varepsilon | X] = 0$$

$$\therefore C E[Y | X] = E[CX | X]\beta$$

can premultiply by $X^T C^T$ to get

$$E[X^T C^T C Y | X] = E[X^T C^T C X | X]\beta$$

so estimate by

$$\hat{\beta}^{OLS} = (X^T C^T C X)^{-1} X^T C^T C Y$$

$$= (X^T \varepsilon^{-1} X)^{-1} X^T \varepsilon^{-1} Y$$

Generalized Least squares

Properties: Unbiased

Conditional $\hat{\beta} = \beta + (X^T \varepsilon^{-1} X)^{-1} X^T \varepsilon^{-1} \varepsilon$

Variance $E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T | X] =$

$$E[(X^T \varepsilon^{-1} X)^{-1} (X^T \varepsilon^{-1} \varepsilon \varepsilon^T \varepsilon^{-1} X)(X^T \varepsilon^{-1} X)^{-1} | X]$$

$$= (X^T \varepsilon^{-1} X)^{-1}$$

Can apply same argument as in Gauss-Markov theorem fashion
 this is "best linear unbiased estimator"

Problem: Don't know Σ usually

Replace Σ by an estimate to get "feasible GLS"

Depending on properties of $\hat{\Sigma}$ estimator, make
 good or bad in finite samples.

Large sample properties & OLS

Using summation representation

$$X = \begin{bmatrix} x_{11} & \dots & x_{k1} & \dots & x_{N1} \\ x_{12} & \dots & x_{k2} & \dots & x_{N2} \\ \vdots & & \vdots & & \vdots \\ x_{1N} & \dots & x_{kN} & \dots & x_{NN} \end{bmatrix} = \begin{bmatrix} x_1' \\ \vdots \\ x_N' \end{bmatrix}$$

$$x_i' = (x_{1,i}, \dots, x_{k,i}, \dots, x_{N,i})$$

k variables for observation

$$\therefore X'X = [x_1 \dots x_N] \begin{bmatrix} x_1' \\ \vdots \\ x_N' \end{bmatrix} = \sum_{i=1}^N x_i x_i'$$

$$\text{Similarly } X'y = \sum_{i=1}^N x_i y_i \quad \text{since } y_i \text{ is scalar}$$

So OLS has

$$\hat{\beta} = (X'X)^{-1} X'y = \left(\sum_{i=1}^N x_i x_i' \right)^{-1} \left(\sum_{i=1}^N x_i y_i \right)$$

Usually rewrite as simple averages

$$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k, \bar{x}_N$$

$$= \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i y_i \right)$$

Assume now (y_i, x_i') are i.i.d and $E[\varepsilon_i | x_i] = 0$
Consistency

$$\hat{\beta} = \beta + \left[\frac{1}{N} \sum_{i=1}^N x_i x_i' \right]^{-1} \left[\frac{1}{N} \sum_{i=1}^N x_i \varepsilon_i \right]$$

Take $N \rightarrow \infty$

$$\frac{1}{N} \sum_{i=1}^N x_i x_i' \xrightarrow{P} E[x_i x_i'] \equiv \Sigma_{xx}$$

By Law of Large Numbers if there exists

Now, if Σ_{xx}^{-1} exists, by continuity

$$\begin{aligned} \lim \left[\frac{1}{N} \sum_{i=1}^N x_i x_i' \right]^{-1} &= \left[\lim \frac{1}{N} \sum_{i=1}^N x_i x_i' \right]^{-1} \\ &= \Sigma_{xx}^{-1} \end{aligned}$$

Next, by CLN

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N x_i \varepsilon_i &\xrightarrow{P} E[x_i \varepsilon_i] = E[x_i E[\varepsilon_i | x_i]] \\ &= E[x_i x_0] = 0 \end{aligned}$$

So $\hat{\beta} \xrightarrow{P} \beta$: consistency

In short, we didn't need $E[\varepsilon_i | x_i] = 0$ further

Argument just used $E[x_i \varepsilon_i] = 0$ - this is
entirely general, does not restrict distribution of data

$$\text{N.S. } E[x_i y_i] = E[x_i x_i' \beta] + \underbrace{E[x_i \varepsilon_i]}_{0 \text{ by assumption}}$$

$$\text{then } \beta = E[x_i x_i']^{-1} E[x_i y_i]$$

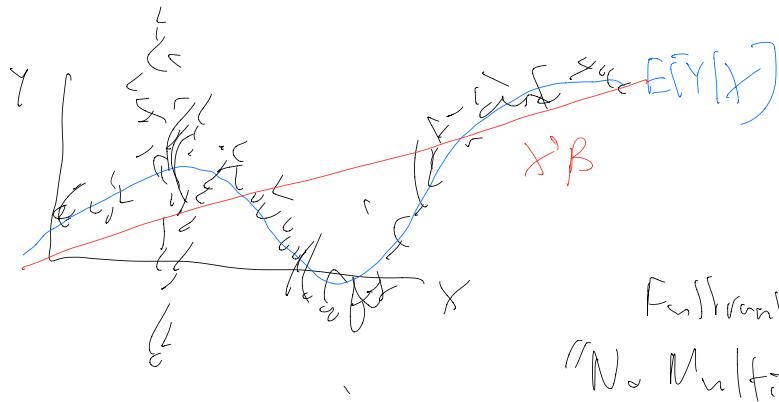
and we can define β as the best fit for any distribution with $E[x_i x_i']$ full rank

with β defined this way,

$$y_i = x_i \beta + \varepsilon_i \quad \text{because } \varepsilon_i = y_i - x_i \beta$$

$$\text{satisfy } E[x_i \varepsilon_i] = 0$$

Interpretation: OLS estimates linear best fit of y given x , even if conditional expectation not linear.



Fullrank $E[x_i x_i']$ assumption
"No Multicollinearity"
Gives a unique linear best fit

If fit fails, can choose non-linear linear best fit, e.g.,

Use a pseudo inverse instead of inverse,

Distribution Theory

- Need measure result about convergence in distribution
- Relationship of convergence in probability
- Suppose $\hat{\theta}_n \xrightarrow{P} \theta$ a constant

Then $\hat{\theta}_n \xrightarrow{d} \theta$

Proof:

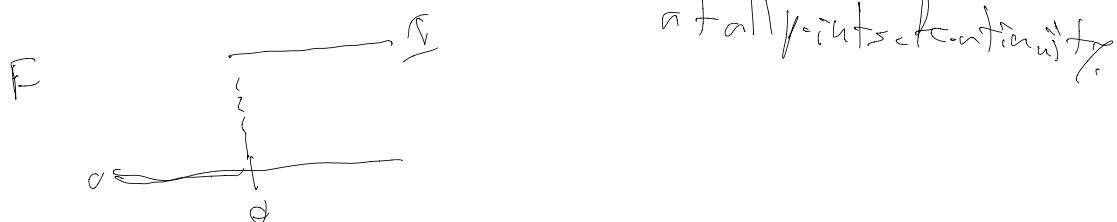
$$\begin{aligned} \Pr(\hat{\theta} < z) &= \Pr(\hat{\theta} - \theta < z - \theta) \\ &\geq \Pr(|\hat{\theta} - \theta| < |z - \theta|) \end{aligned}$$

$\rightarrow 1$ as $n \rightarrow \infty$ if $z > \theta$
by convergence in probability

Can also show by symmetric argument

$$\Pr(\hat{\theta} < z) \rightarrow 0 \text{ if } z < \theta$$

So CDF converges to CDF of a constant θ



Saint Petersburg Theorem

Suppose $A_n \xrightarrow{d} Z$, $B_n \xrightarrow{P} b$ a constant. Then

$$(A_n, B_n) \xrightarrow{d} (Z, L)$$

→ Statistics can be replaced by their plim in asymptotic distributions.

Note: statement that $A_n \xrightarrow{d} Z, B_n \xrightarrow{d} W$
does not imply $(A_n, B_n) \xrightarrow{d} (Z, W)$
 Joint distribution is not defined for Z, W

Asymptotic Distribution of OLS

$$\sqrt{N} (\hat{\beta} - \beta) = \left(\underbrace{\frac{1}{N} \sum_{i=1}^N x_i x_i'}_{\xrightarrow{\text{CLT}} \Sigma_{xx}} \right)^{-1} \left(\underbrace{\frac{1}{N} \sum_{i=1}^N x_i \varepsilon_i}_{\xrightarrow{\text{CLT}} \Sigma_{x\varepsilon}} \right)$$

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N x_i \varepsilon_i \xrightarrow{d} N(0, E[x_i \varepsilon_i \varepsilon_i'])$$

So by Slutsky and continuous mapping

$$\sqrt{N} (\hat{\beta} - \beta) \xrightarrow{d} \Sigma_{xx}^{-1} Z$$

$$\sim N(0, \Sigma_{xx}^{-1} (E[x_i \varepsilon_i \varepsilon_i' x_i']) \Sigma_{xx}^{-1})$$

$$\text{Var}(Av) = A \text{Var}(v) A'$$

If $E[\varepsilon_i | x_i] = \sigma^2$ then variance simplifies

$$\sigma^2 \Sigma_{xx}^{-1} \quad \text{which can be estimated}$$

$$\text{by } \widehat{\Sigma}_{xx}^{-1} \text{ or before}$$

Otherwise, if $E(\varepsilon_i^2 | \mathbf{x}_i)$ is not always constant,

limits still holds, we can estimate the variance

$$\sum_{\mathbf{x}\mathbf{x}} E(\mathbf{x}_i \mathbf{\varepsilon}_i \mathbf{\varepsilon}_i' \mathbf{x}_i') \sum_{\mathbf{x}\mathbf{x}}$$

Called "Bartlett" or "Heteroscedasticity Consistent"

Variance estimation, or Eicker-White

Estimated by sample analogues. Replace expectedly

average, replace $\mathbf{\varepsilon}_i$ by $\hat{\mathbf{\varepsilon}}_i = \mathbf{y}_i - \hat{\mathbf{x}}_i'\hat{\beta}$

$(\frac{1}{N} \mathbf{x}_i \mathbf{x}_i')$ $\rightarrow \sum_{\mathbf{x}\mathbf{x}}$ so that's handled

check $\frac{1}{N} \sum \mathbf{x}_i \mathbf{x}_i' \hat{\mathbf{\varepsilon}}_i^2$

$$= \underbrace{\frac{1}{N} \sum \mathbf{x}_i \mathbf{x}_i' \hat{\mathbf{\varepsilon}}_i^2}_{\mathbb{E}[\mathbf{x}_i \mathbf{x}_i' \hat{\mathbf{\varepsilon}}_i^2]} + \underbrace{\frac{1}{N} \sum \mathbf{x}_i \mathbf{x}_i' (\hat{\mathbf{\varepsilon}}_i^2 - \mathbb{E}[\hat{\mathbf{\varepsilon}}_i^2])}_{\mathbb{E}[\mathbf{x}_i \mathbf{x}_i' (\hat{\mathbf{\varepsilon}}_i^2 - \mathbb{E}[\hat{\mathbf{\varepsilon}}_i^2])]}$$

Need to show is that

$$\frac{1}{N} \sum \mathbf{x}_i \mathbf{x}_i' (\hat{\mathbf{\varepsilon}}_i^2 - \mathbb{E}[\hat{\mathbf{\varepsilon}}_i^2]) \xrightarrow{P} 0$$

$$\hat{\mathbf{\varepsilon}}_i = \mathbf{y}_i - \hat{\mathbf{x}}_i'\hat{\beta} = \mathbf{x}_i'\hat{\beta} + \mathbf{\varepsilon}_i - \mathbf{x}_i'\hat{\beta} = -\mathbf{x}_i'(\hat{\beta} - \beta) + \mathbf{\varepsilon}_i$$

$$\text{so } (\hat{\mathbf{\varepsilon}}_i^2) = (\hat{\beta} - \beta)' \mathbf{x}_i \mathbf{x}_i' (\hat{\beta} - \beta) + \mathbf{\varepsilon}_i^2 - 2\mathbf{\varepsilon}_i \mathbf{x}_i' (\hat{\beta} - \beta)$$

$$\text{so } \frac{1}{N} \sum \mathbf{x}_i \mathbf{x}_i' (\hat{\mathbf{\varepsilon}}_i^2 - \mathbb{E}[\hat{\mathbf{\varepsilon}}_i^2]) =$$

$$\begin{aligned}
& \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \sum_{k=1}^K \mathbb{E}[x_i x_j' / (\beta - \hat{\beta})] x_i x_k' / (\beta - \hat{\beta}) \\
& - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \sum_{k=1}^K x_i x_j' \mathbb{E}[x_i x_k' / (\hat{\beta} - \beta)] \\
& = \sum_{j=1}^K \sum_{k=1}^K \underbrace{(\beta_j - \hat{\beta}_j) (\beta_k - \hat{\beta}_k)}_{\rightarrow 0} \underbrace{\frac{1}{N} \sum_{i=1}^N x_i x_i' x_j x_k' / \mathbb{E}[x_i x_i' / x_j x_k]}_{\rightarrow \mathbb{E}[x_i x_i' / x_j x_k]} \\
& - \sum_{k=1}^K \underbrace{(\hat{\beta}_k - \beta_k)}_{\rightarrow 0} \underbrace{\frac{1}{N} \sum_{i=1}^N x_i x_i' \mathbb{E}[x_i x_k']}_{\text{if finite}} \xrightarrow{\rightarrow} \mathbb{E}[x_i x_i' / x_k x_k]
\end{aligned}$$

$$\xrightarrow{P} 0$$

$$\begin{aligned}
\text{So } \hat{\Sigma}_{\text{robust}} &= \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \mathbb{E}[x_i x_i'] \right) \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \\
&\xrightarrow{P} \hat{\Sigma}_{\beta} \text{ asymptotic variance of } \hat{\beta}
\end{aligned}$$

If data is iid, you can almost always work
 formulas if homoskedastic, it is constant, but
 also if not, you will see.

Reg Y | X, robust

Application of these formulas

Tests (and later, confidence intervals)

Goals to distinguish between "hypotheses"

H_0 (null hypothesis)

H_1 (alternative hypothesis)

Two disjoint sets of probability measures over data
usually corresponding to disjoint parameter sets.
We will use shorthand:

$$H_0 : \beta = \beta_0, \quad H_1 : \beta \neq \beta_0$$

is shorthand for $H_0 = \{ \text{probabilistic measures} \}$

under some set of assumptions, where population

regression coefficient = $\beta_0 \}$

and $H_1 = \{ \text{st. f. measures satisfying same assumptions} \}$

β in population $\neq \beta_0 \}$

Test: a function of data $\hat{\gamma} : \{\mathbf{z}_i\}_{i=1}^N \rightarrow \{0,1\}$

that picks H_0 or H_1

(contains a "randomized test" that is also a function
of additional random variables)