

Section: Module 4

OLS

Kentaro Nakamura

GOV 2003

February 20th, 2026

Today's agenda

- OLS
 - QR Decomposition / OLS Computation (PS1)
- Hypothesis testing
 - Basic Principle
 - Derivation of finite-sample test (T-test / F-test)
 - Asymptotic inference
- Due to time constraints, I put some materials I cannot cover as Appendix
 - OLS as Best Linear Approximation
 - OLS estimates without i th observation
 - FWL theorem and within transformation for fixed effect regression

Orthonormal Matrix

Definition (Orthonormal Matrix)

An $m \times m$ matrix $\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_m) = (\mathbf{q}_1, \dots, \mathbf{q}_m)^\top$ is called an **orthonormal matrix** if

$$\mathbf{P}^\top \mathbf{P} = \mathbf{P} \mathbf{P}^\top = \mathbf{I}_m$$

- As is clear from the definition, $\mathbf{P}^{-1} = \mathbf{P}^\top$

Example

$$\mathbf{Q} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

is orthogonal matrix, satisfying $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$.

Instability of Matrix Inversion

- While we derived $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$, this is NOT how lm function computes the regression.
 - This is because the computation of $(\mathbf{X}^\top \mathbf{X})^{-1}$ is numerically unstable when columns of \mathbf{X} are numerically colinear
 - To understand this, we need to understand **QR decomposition**

Theorem (QR Decomposition)

Let \mathbf{X} be an $n \times p$ matrix with full column rank ($p \leq n$). Then there exist

- an $n \times p$ matrix \mathbf{Q} with orthonormal columns ($\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_p$),
 - a $p \times p$ upper-triangular matrix \mathbf{R} with positive diagonal entries,
- such that

$$\mathbf{X} = \mathbf{QR}.$$

- NOTE: $\mathbf{Q} \in \mathbb{R}^{n \times p}$ is not a square matrix (unless $n = p$), and this satisfies $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ but in general $\mathbf{Q}\mathbf{Q}^\top \neq \mathbf{I}$

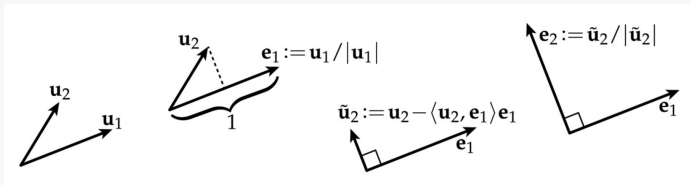
Gram–Schmidt Orthogonalization

Theorem (Gram–Schmidt)

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$ with linearly independent columns. Define recursively:

$$\mathbf{v}_1 = \mathbf{x}_1, \quad \mathbf{v}_k = \mathbf{x}_k - \sum_{j=1}^{k-1} \frac{\mathbf{v}_j^\top \mathbf{x}_k}{\mathbf{v}_j^\top \mathbf{v}_j} \mathbf{v}_j, \quad \mathbf{q}_k = \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|}.$$

Then $\{\mathbf{q}_1, \dots, \mathbf{q}_p\}$ is an orthonormal basis of $S(\mathbf{X})$.



QR Decomposition and Gram Schmidt

- Start with the columns of $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p]$.
- Gram-Schmidt step:** for each $k = 1, \dots, p$,

$$\mathbf{v}_k = \mathbf{x}_k - \sum_{j=1}^{k-1} \underbrace{(\mathbf{q}_j^\top \mathbf{x}_k)}_{\text{scalar}} \mathbf{q}_j, \quad \mathbf{q}_k = \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|}.$$

- Rearrange the equation for \mathbf{x}_k :

$$\mathbf{x}_k = (\mathbf{q}_1^\top \mathbf{x}_k) \mathbf{q}_1 + (\mathbf{q}_2^\top \mathbf{x}_k) \mathbf{q}_2 + \dots + (\mathbf{q}_k^\top \mathbf{x}_k) \mathbf{q}_k.$$

- Notice that each column satisfies $\mathbf{x}_k = \sum_{j=1}^k r_{jk} \mathbf{q}_j$, where $r_{jk} := \mathbf{q}_j^\top \mathbf{x}_k$ for $j \leq k$ and $r_{jk} := 0$ for $j > k$.

$$\mathbf{X} = \begin{bmatrix} \mathbf{q}_1 & \dots & \mathbf{q}_p \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1p} \\ 0 & r_{22} & \dots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & r_{pp} \end{bmatrix} = \mathbf{QR}.$$

Gram–Schmidt Orthogonalization (Example 1/3)

- Let

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = [\mathbf{x}_1, \mathbf{x}_2], \quad \mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}.$$

- Step 1:**

$$\mathbf{v}_1 = \mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \|\mathbf{v}_1\| = \sqrt{1^2 + 0^2 + 0^2} = 1$$

and thus

$$\mathbf{q}_1 = \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

Gram–Schmidt Orthogonalization (Example 2/3)

- **Step 2:** compute the projection of \mathbf{x}_2 onto \mathbf{q}_1 . The projection is

$$\text{proj}_{\mathbf{q}_1}(\mathbf{x}_2) = (\mathbf{q}_1^\top \mathbf{x}_2) \mathbf{q}_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

- Subtract it:

$$\mathbf{v}_2 = \mathbf{x}_2 - \text{proj}_{\mathbf{q}_1}(\mathbf{x}_2) = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

- Normalize:

$$\|\mathbf{v}_2\| = \sqrt{0^2 + 1^2 + 0^2} = 1, \quad \mathbf{q}_2 = \frac{\mathbf{v}_2}{\|\mathbf{v}_2\|} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

Gram–Schmidt Orthogonalization (Example 3/3)

- Thus,

$$\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{Q}^\top \mathbf{Q} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{I}_2.$$

- Compute $\mathbf{R} = \mathbf{Q}^\top \mathbf{X}$. Since

$$\mathbf{R} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$

- Verify:

$$\mathbf{QR} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \mathbf{X}.$$

OLS via QR (Why $\text{lm}()$ does this)

- Recall the least squares problem:

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2.$$

- Suppose $\mathbf{X} = \mathbf{QR}$, where $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ and \mathbf{R} is upper triangular.

$$\begin{aligned}\|\mathbf{Y} - \mathbf{X}\beta\|_2^2 &= \|\mathbf{Y} - \mathbf{QR}\beta\|_2^2 \\ &= \|\mathbf{Y} - \mathbf{QQ}^\top \mathbf{Y} + \mathbf{QQ}^\top \mathbf{Y} - \mathbf{QR}\beta\|_2^2 \\ &= \|(\mathbf{I} - \mathbf{QQ}^\top)\mathbf{Y} + \mathbf{Q}(\mathbf{Q}^\top \mathbf{Y} - \mathbf{R}\beta)\|_2^2 \\ &= \|(\mathbf{I} - \mathbf{QQ}^\top)\mathbf{Y}\|_2^2 + \|\mathbf{Q}(\mathbf{Q}^\top \mathbf{Y} - \mathbf{R}\beta)\|_2^2 \\ &= \|(\mathbf{I} - \mathbf{QQ}^\top)\mathbf{Y}\|_2^2 + \|\mathbf{Q}^\top \mathbf{Y} - \mathbf{R}\beta\|_2^2\end{aligned}$$

where the second equality is by adding and subtracting the same quantity, the fourth equality is because of orthogonality (i.e., $(\mathbf{I} - \mathbf{QQ}^\top)\mathbf{Q} = \mathbf{0}$), and the fifth equality is because \mathbf{Q} is orthonormal (i.e., $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$).

- The first term does not involve β , and since \mathbf{R} is invertible (as \mathbf{X} is full rank), the second component is minimized when $\mathbf{Q}^\top \mathbf{Y} = \mathbf{R}\beta$. 10

OLS via QR

- Suppose

$$\mathbf{R} = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}, \quad \mathbf{z} = \mathbf{Q}^\top \mathbf{Y} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}.$$

- Then $\mathbf{R}\boldsymbol{\beta} = \mathbf{z}$ means:

$$\begin{cases} r_{11}\beta_1 + r_{12}\beta_2 = z_1 \\ r_{22}\beta_2 = z_2 \end{cases}$$

- Start from the bottom equation:

$$r_{22}\beta_2 = z_2 \quad \Rightarrow \quad \beta_2 = \frac{z_2}{r_{22}}.$$

- Plug β_2 into the first equation:

$$r_{11}\beta_1 + r_{12}\beta_2 = z_1 \quad \Rightarrow \quad \beta_1 = \frac{z_1 - r_{12}\beta_2}{r_{11}}.$$

- This is called **backward substitution**
 - No matrix inversion. Just division and subtraction.

Numeric example of backward substitution

- Let $\mathbf{Y} = [3 \ 2 \ 0]^T$.
- Compute

$$\mathbf{z} = \mathbf{Q}^T \mathbf{Y} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \end{bmatrix}.$$

- Now solve $\mathbf{R}\hat{\boldsymbol{\beta}} = \mathbf{z}$:

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \end{bmatrix}.$$

- Bottom equation: $\hat{\beta}_2 = 2$.
- Top equation: $\hat{\beta}_1 + \hat{\beta}_2 = 3 \Rightarrow \hat{\beta}_1 = 1$.

Hypothesis testing: Basic principle (1)

- **Procedure**

- STEP 1: Assume the null hypothesis $H_0 : \theta \in \Theta_0$ ($\Leftrightarrow H_1 : \theta \in \Theta_1$, with $\Theta = \Theta_0 \cup \Theta_1$)
- STEP 2: Calculate the distribution of test statistics under the null (known as **reference distribution**)
- STEP 3: If the observed statistics is too extreme compared to the reference distribution, we reject the null hypothesis
 - Two decisions: (i) reject null or (ii) accept null

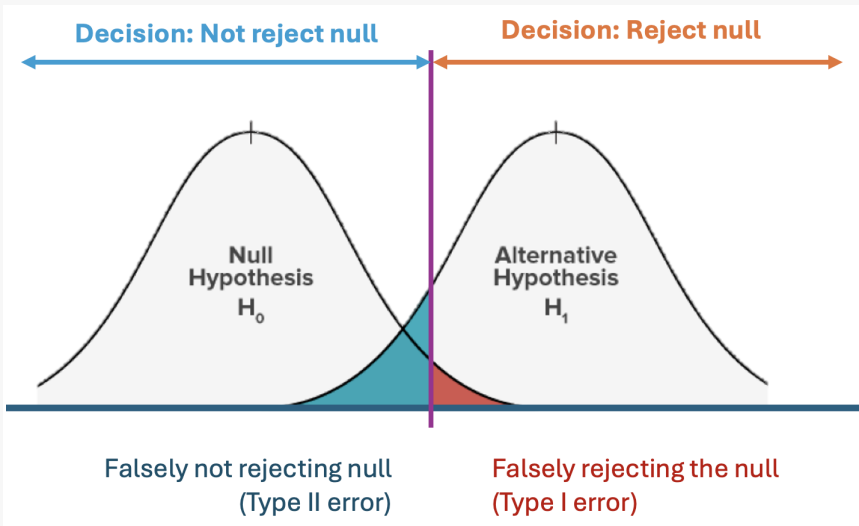
- Importantly, we need to know the distribution of test statistics
 - Thus, we typically either assume the assumption on our data or use asymptotic approach

Hypothesis testing: Basic principle (2)

	Accept H_0	Reject H_0
H_0 true	Correct Decision	Type I Error
H_1 true	Type II Error	Correct Decision

- The probability of Type I error is called the **size** of the test
- The rejection probability under the alternative hypothesis is called the **power** of the test
 - I.e., power equals 1 minus the probability of a Type II error
- There is a trade-off between size and power
 - Understand this trade-off from the next page's illustration

Illustration: Type I & Type II error



Normal Distribution

Definition (Normal distribution)

A random variable Z is **normal** with mean μ and variance σ^2 , written $Z \sim \mathcal{N}(\mu, \sigma^2)$, if it has density

$$f(z) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(z - \mu)^2}{2\sigma^2}\right).$$

- If \mathbf{Z} is a vector, the distribution is called **multivariate normal**.
 - For example,

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \sim \text{MVN}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}\right)$$

where ρ is the correlation

- Important properties:
 - In the case of multivariate normal, no correlation is equivalent to independence (not true for other cases)
 - **reproductive property**: If $\mathbf{Z} \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then

$$\mathbf{AZ} + \mathbf{a} \sim \text{MVN}(\mathbf{A}\boldsymbol{\mu} + \mathbf{a}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

Chi-squared distribution

Definition (Chi-squared distribution)

If $Z_1, \dots, Z_\nu \stackrel{iid}{\sim} \mathcal{N}(0, 1)$, then

$$\sum_{j=1}^{\nu} Z_j^2 \sim \chi_\nu^2,$$

where ν is called the **degrees of freedom**.

t-distribution and F-distribution

Definition (Student t distribution)

Let $Z \sim \mathcal{N}(0, 1)$ and $U \sim \chi_{\nu}^2$ be independent. Then

$$T := \frac{Z}{\sqrt{U/\nu}} \sim t_{\nu}.$$

Definition (F distribution)

Let $U \sim \chi_{d_1}^2$ and $V \sim \chi_{d_2}^2$ be independent. Then

$$F := \frac{U/d_1}{V/d_2} \sim F_{d_1, d_2}.$$

- We will derive t-test and F-test statistic for OLS

Distribution of OLS Coefficient and Residual

- **Setup:** $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \mid \mathbf{X} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$

- Recall that

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}$$

- On the other hand, annihilator matrix $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is orthogonal to \mathbf{X} (i.e., $\mathbf{M}\mathbf{X} = 0$), and thus

$$\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \mathbf{P}\mathbf{Y} = \mathbf{M}\mathbf{Y} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = \mathbf{M}\boldsymbol{\epsilon}$$

- Thus,

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \hat{\boldsymbol{\epsilon}} \end{bmatrix} = \begin{bmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ \mathbf{M} \end{bmatrix} \boldsymbol{\epsilon}$$

- Therefore, (try to show this part by yourself!)

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \hat{\boldsymbol{\epsilon}} \end{bmatrix} \mid \mathbf{X} \sim \mathcal{N}\left(0, \begin{bmatrix} \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} & 0 \\ 0 & \sigma^2 \mathbf{M} \end{bmatrix}\right)$$

Extra: Derivation of OLS Variance Estimator

- So far, everything depends on the error variance $\sigma^2 = \mathbb{E}[\epsilon^2]$
 - We do not know the true error variance, so we need to propose the estimator $\hat{\sigma}^2$
- The natural estimator is to replace expectation with average; i.e.,

$$\hat{\sigma}^2 = \frac{1}{n} \hat{\epsilon}^\top \hat{\epsilon} = \frac{1}{n} (\mathbf{M}\epsilon)^\top \mathbf{M}\epsilon = \frac{1}{n} \epsilon^\top \mathbf{M}\epsilon$$

- This is actually not unbiased.
 - To see that, use $\text{tr}(\epsilon^\top \mathbf{M}\epsilon) = \text{tr}(\mathbf{M}\epsilon\epsilon^\top)$

$$\mathbb{E}[\hat{\sigma}^2 \mid \mathbf{X}] = \frac{1}{n} \text{tr}(\mathbb{E}[\mathbf{M}\epsilon\epsilon^\top \mid \mathbf{X}]) = \frac{1}{n} \text{tr}(\mathbf{M} \underbrace{\mathbb{E}[\epsilon\epsilon^\top \mid \mathbf{X}]}_{\sigma^2 \mathbf{I}_n}) = \frac{n-p}{n} \sigma^2$$

where the last line is by $\text{tr}(\mathbf{M}) = n - p$

- Thus, the unbiased estimator of variance is given by

$$s^2 = \frac{1}{n-p} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n-p} \epsilon^\top \epsilon$$

Distribution of Variance Estimator

- Since $\boldsymbol{\epsilon} = \hat{\boldsymbol{\epsilon}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$,

$$\begin{aligned}\boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} &= \hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}} + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + 2(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \hat{\boldsymbol{\epsilon}} \\ &= \hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}} + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\end{aligned}$$

- $\frac{\boldsymbol{\epsilon}^\top \boldsymbol{\epsilon}}{\sigma^2} \sim \chi_n^2$ because $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$
- $\frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\sigma^2} \sim \chi_p^2$ because $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$

- Therefore,

$$\frac{\hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}}}{\sigma^2} = \frac{(n-p)s^2}{\sigma^2} \sim \chi_{n-p}^2$$

t-test: Derivation

- **Setup:**

- $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ with $\epsilon | \mathbf{X} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$

- **t-statistic** is

$$T = \frac{\hat{\beta}_j - \beta_j}{\sqrt{s^2[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}}} = \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}}} / \sqrt{\frac{(n-p)s^2}{\sigma^2} / (n-p)}$$
$$\sim \frac{\mathcal{N}(0, 1)}{\sqrt{\chi_{n-p}^2 / (n-p)}} \sim t_{n-p}$$

- **Note:** Even though we start with conditional distribution $\epsilon | \mathbf{X} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$, the resulting t-distribution does not depend on \mathbf{X} (and also σ^2)

- This means that the reference distribution does not depend on any other features of the data, except the degrees of freedom $n - p$
 - Such test statistic is called **pivotal**, meaning that it does not depend on unknowns.

F-test: Derivation

- **Setup:**

- $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} | \mathbf{X} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$
- We are testing joint restriction $H_0 : \mathbf{A}\boldsymbol{\beta} = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{q \times p}$ with $\text{rank}(\mathbf{A}) = q$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

- As $\hat{\boldsymbol{\beta}}$ is normal, under H_0 ,

$$\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{b} \sim \mathcal{N}(0, \sigma^2 \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)$$

thus by normalization,

$$\frac{(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{b})^\top \{\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top\}^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{b})}{\sigma^2} \sim \chi_q^2$$

- You can easily derive denominator and see why the test statistic follows F-distribution using definition (try by yourself!)

Asymptotic Inference

- However, the assumption that error is distributed according to normal distribution is a strong assumption
 - Both t-test and F-test depends on this assumption
 - We want to relax it

- Instead, we want to resort to **asymptotic statistics**
 - **Idea:** Assume we have infinite amount of data (i.e., $N \rightarrow \infty$)
 - Use tools in asymptotic statistics to derive the distribution at the limiting (called **limiting distribution**)

Tools in Asymptotic Statistics

- **Law of Large Numbers (LLN):** If X_1, \dots, X_n are i.i.d.,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mathbb{E}[X_i]$$

- **Central Limit Theorem (CLT):** If X_1, \dots, X_n are i.i.d.,

$$\sqrt{n}(\bar{X} - \mathbb{E}[X_i]) \xrightarrow{d} \mathcal{N}(0, \mathbb{V}[X_i])$$

- **Slutsky's Lemma:** If $X_n \xrightarrow{d} X$ for some random variable X and $Y_n \xrightarrow{P} c$ for some constant c ,

$$X_n Y_n \xrightarrow{d} cX$$

Model-based asymptotic inference (1)

- Notice that by law of large numbers, the first term converges in probability to

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \xrightarrow{p} \mathbb{E}[\mathbf{X}^\top \mathbf{X}]$$

- On the other hand, because each $\mathbf{X}_i \epsilon_i$ are i.i.d., by central limit theorem, the second term converges in distribution to

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \epsilon_i = \sqrt{n} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i}_{\text{Form of Avg!}} \xrightarrow{d} \mathcal{N}(0, \mathbb{E}[(\mathbf{X}^\top \epsilon (\mathbf{X}^\top \epsilon)^\top)])$$

- Note that mean of normal distribution here is 0 because $\mathbb{E}[\mathbf{X}_i \epsilon_i] = 0$

Model-based asymptotic inference (2)

- Therefore,

$$\begin{aligned}\sqrt{n}(\hat{\beta} - \beta) &= \sqrt{n}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ &\rightarrow \mathcal{N}\left(0, \underbrace{\mathbb{E}[\mathbf{X}^\top \mathbf{X}]^{-1} \mathbb{E}[\mathbf{X}^\top \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top \mathbf{X}]}_{\text{Asymptotic Variance}} \mathbb{E}[\mathbf{X}^\top \mathbf{X}]^{-1}\right)\end{aligned}$$

- **Homoskedasticity:** $\mathbb{V}[\epsilon_i | \mathbf{X}_i] = \sigma^2$

- This gives $\mathbb{V}[\epsilon_i | \mathbf{X}_i] = \mathbb{E}[\epsilon_i^2 | \mathbf{X}_i] - \underbrace{\mathbb{E}[\epsilon_i | \mathbf{X}_i]^2}_{=0} = \mathbb{E}[\epsilon_i^2 | \mathbf{X}_i]$

- NOTE: Homoskedasticity is usually not plausible

$$\begin{aligned}\mathbb{E}[\mathbf{X}^\top \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top \mathbf{X}] &= \mathbb{E}[\mathbb{E}[\mathbf{X}^\top \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top \mathbf{X} | \mathbf{X}]] \quad (\because \text{Law of Iterated Expectation}) \\ &= \mathbb{E}[\mathbf{X}^\top \underbrace{\mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top | \mathbf{X}]}_{=\sigma^2 I_n} \mathbf{X}] = \sigma^2 \mathbb{E}[\mathbf{X}^\top \mathbf{X}]\end{aligned}$$

and thus the asymptotic variance is simplified to

$$\mathbb{E}[\mathbf{X}^\top \mathbf{X}]^{-1} \mathbb{E}[\mathbf{X}^\top \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top \mathbf{X}] \mathbb{E}[\mathbf{X}^\top \mathbf{X}]^{-1} = \sigma^2 \mathbb{E}[\mathbf{X}^\top \mathbf{X}]^{-1}$$

Wald Test Statistic

- **Null hypothesis:** $H_0 : \mathbf{A}\beta = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{q \times p}$ is full row rank
 - Any linear restriction can be written as $a_1\beta_1 + \dots + a_p\beta_p = b$
 - Can test any number of linearly independent restrictions (joint test)
- From asymptotic normality, assume

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

Then under H_0 ,

$$\sqrt{n}(\mathbf{A}\hat{\beta} - \mathbf{b}) \xrightarrow{d} \mathcal{N}(0, \mathbf{A}\Sigma\mathbf{A}^\top).$$

- **Wald statistic** (quadratic form):

$$\begin{aligned} W &:= (\mathbf{A}\hat{\beta} - \mathbf{b})^\top \left\{ \mathbf{A}\widehat{\mathbf{V}}(\hat{\beta})\mathbf{A}^\top \right\}^{-1} (\mathbf{A}\hat{\beta} - \mathbf{b}) \\ &= \left[\sqrt{n}(\mathbf{A}\hat{\beta} - \mathbf{b}) \right]^\top \left\{ \mathbf{A}\widehat{\Sigma}\mathbf{A}^\top \right\}^{-1} \left[\sqrt{n}(\mathbf{A}\hat{\beta} - \mathbf{b}) \right], \quad \widehat{\Sigma} := n\widehat{\mathbf{V}}(\hat{\beta}). \end{aligned}$$

- Under H_0 , $W \xrightarrow{d} \chi_{\text{rank}(\mathbf{A})}^2$.
- If W is larger than the critical value, reject H_0 .

Appendix: OLS as Best Linear Approximation

- Linear regression is a **best linear approximation** of conditional expectation function $m(\mathbf{X}_i) = \mathbb{E}[Y_i | \mathbf{X}_i]$
- To see this, consider minimizing the mean squared loss between $m(\mathbf{X}_i)$ and $\mathbf{X}_i^\top \beta$

$$d(\beta) = \mathbb{E}[(m(\mathbf{X}_i) - \mathbf{X}_i^\top \beta)^2]$$

- Notice that this minimization problem is just replacing Y_i with $m(\mathbf{X}_i)$ for OLS. Thus, the best linear approximation $\tilde{\beta}$ is given by

$$\begin{aligned}\tilde{\beta} &= (\mathbb{E}[\mathbf{X}_i \mathbf{X}_i^\top])^{-1} \mathbb{E}[\mathbf{X}_i m(\mathbf{X}_i)] \\ &= (\mathbb{E}[\mathbf{X}_i \mathbf{X}_i^\top])^{-1} \mathbb{E}[\mathbf{X}_i \mathbb{E}[Y_i | \mathbf{X}_i]] \\ &= (\mathbb{E}[\mathbf{X}_i \mathbf{X}_i^\top])^{-1} \mathbb{E}[\mathbb{E}[\mathbf{X}_i Y_i | \mathbf{X}_i]] \\ &= (\mathbb{E}[\mathbf{X}_i \mathbf{X}_i^\top])^{-1} \mathbb{E}[\mathbf{X}_i Y_i] = \beta\end{aligned}$$

- OLS at the population (covered last week) is equal to the best linear approximation of $m(\mathbf{X}_i) = \mathbb{E}[Y_i | \mathbf{X}_i]$

Appendix: Sherman–Morrison Formula

- We derive the OLS estimate without the i th observation
- To do this, it is useful to know the following theorem.

Theorem (Sherman–Morrison formula Formula)

Suppose $\mathbf{A} \in \mathbb{R}^{n \times n}$ be an invertible square matrix and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ be vectors. If $\mathbf{A} + \mathbf{u}\mathbf{v}^\top$ is invertible, then

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^\top\mathbf{A}^{-1}}{1 + \mathbf{v}^\top\mathbf{A}^{-1}\mathbf{u}}$$

- See next page for the proof.

Appendix: Sherman–Morrison Formula: Proof

$$\begin{aligned} & \left(\mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} \mathbf{A}^{-1}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} \right) (\mathbf{A} + \mathbf{u} \mathbf{v}^{\top}) \\ &= \mathbf{A}^{-1} \mathbf{A} + \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{A}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} \\ &= \mathbf{I} + \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} - \frac{\mathbf{A}^{-1} \mathbf{u} (\mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}) \mathbf{v}^{\top}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} \\ &= \mathbf{I} + \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} - \frac{\mathbf{A}^{-1} \mathbf{u} (1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}) \mathbf{v}^{\top}}{1 + \mathbf{v}^{\top} \mathbf{A}^{-1} \mathbf{u}} \\ &= \mathbf{I} + \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} - \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^{\top} \\ &= \mathbf{I}. \end{aligned}$$

Appendix: OLS estimates without i th observation (1)

- Now, notice that $\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)}$ ($\mathbf{X}^\top \mathbf{X}$ removing i th observations) is written as

$$\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)} = \mathbf{X}^\top \mathbf{X} - \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top$$

- As a result, the OLS estimates without i th observation is given by

$$\hat{\beta}_{(i)} = (\mathbf{X}^\top \mathbf{X} - \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top)^{-1} (\mathbf{X}^\top \mathbf{Y} - \mathbf{x}_{(i)} \mathbf{y}_{(i)}^\top)$$

- Using Sherman-Morrison formula, we obtain

$$(\mathbf{X}^\top \mathbf{X} - \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top)^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top (\mathbf{X}^\top \mathbf{X})^{-1}}{1 - \mathbf{x}_{(i)} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)}^\top}$$

Appendix: OLS estimates without i th observation (2)

- By multiplying $\mathbf{X}^\top \mathbf{Y} - \mathbf{x}_{(i)} \mathbf{y}_{(i)}^\top$, we obtain

$$\begin{aligned}\hat{\beta}_{(i)} &= \left(\mathbf{X}^\top \mathbf{X} - \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top \right)^{-1} \left(\mathbf{X}^\top \mathbf{Y} - \mathbf{x}_{(i)} \mathbf{y}_{(i)} \right) \\ &= \left[\left(\mathbf{X}^\top \mathbf{X} \right)^{-1} + \frac{\left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1}}{1 - \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)}} \right] \\ &\quad \times \left(\mathbf{X}^\top \mathbf{Y} - \mathbf{x}_{(i)} \mathbf{y}_{(i)} \right) \\ &= \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{Y} - \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)} \mathbf{y}_{(i)} \\ &\quad + \frac{\left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{Y}}{1 - \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)}} \\ &\quad - \frac{\left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)} \mathbf{y}_{(i)}}{1 - \mathbf{x}_{(i)}^\top \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{x}_{(i)}}\end{aligned}$$

Appendix: OLS estimates without i th observation (3)

- Thus,

$$\begin{aligned}\hat{\beta}_{(i)} &= \hat{\beta} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} y_{(i)} \\ &\quad + \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \mathbf{x}_{(i)}^\top \hat{\beta}}{1 - \ell_i} - \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \ell_i y_{(i)}}{1 - \ell_i} \\ &= \hat{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \left[\frac{\mathbf{x}_{(i)}^\top \hat{\beta} - \ell_i y_{(i)}}{1 - \ell_i} - y_{(i)} \right] \\ &= \hat{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \left[\frac{\mathbf{x}_{(i)}^\top \hat{\beta} - y_{(i)}}{1 - \ell_i} \right] \\ &= \hat{\beta} - \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)} \hat{\varepsilon}_{(i)}}{1 - \ell_i},\end{aligned}$$

where $\ell_i \equiv \mathbf{x}_{(i)}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{(i)}$ and $\hat{\varepsilon}_{(i)} \equiv y_{(i)} - \mathbf{x}_{(i)}^\top \hat{\beta}$.

Appendix: Fixed effect (1)

- Consider the fixed effect model:

$$Y = D\alpha + X\beta + \epsilon$$

where D is a $n \times G$ matrix for dummy variable for G groups

- One approach to estimate β in this context is to regress the regression model above
 - This approach is called the **least squares dummy variable (LSDV)**
- However, this is not scalable when G is large
 - Especially problematic when applied to panel data, where you want to use both time and unit fixed effects
- In such case, we can apply **within-transformation** approach to estimate it
 - This is actually well-explained from FWL theorem

Appendix: Fixed effect (2)

- Consider the partitioned regressor (\mathbf{D}, \mathbf{X}) . Then, the regression coefficient β is obtained by

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^\top \{I - P_D\} \mathbf{X})^{-1} \mathbf{X}^\top \{I - P_D\} \mathbf{Y} \\ &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{Y}}\end{aligned}$$

where $P_D = D(D^\top D)^{-1}D^\top$, $\tilde{\mathbf{X}} = \{I - P_D\}\mathbf{X}$, and $\tilde{\mathbf{Y}} = \{I - P_D\}\mathbf{Y}$

- Now, notice that

$$\tilde{\mathbf{Y}} = \{I - P_D\}\mathbf{Y} = \mathbf{Y} - P_D\mathbf{Y}$$

and you see that $P_D\mathbf{Y}$ is the regression of \mathbf{Y} on \mathbf{D} , which gives us the average of \mathbf{Y} on each group

- Now, $\tilde{\mathbf{Y}}$ is the demeaned outcome (demeaned by group averages)
 - You can also think $\tilde{\mathbf{X}}$ is also the demeaned regressor
 - This is basically the within-transformation