

# Advanced Econometrics

Finite Sample Properties of Least Squares (Hansen Chapters 4 and 5)

Instructor: Ma, Jun

Renmin University of China

October 25, 2021

# Random Sampling

- ▶ The simplest context is when the observations are mutually independent, in which case we say that they are independent and identically distributed, or i.i.d. It is also common to describe iid observations as a random sample.

## Assumption

*The observations  $\{(Y_1, \mathbf{X}_1), \dots, (Y_i, \mathbf{X}_i), \dots, (Y_n, \mathbf{X}_n)\}$  are independent and identically distributed.*

- ▶ If you take any two individuals  $i \neq j$  in a sample, the values  $(Y_i, \mathbf{X}_i)$  are independent of the values  $(Y_j, \mathbf{X}_j)$  yet have the same distribution.

# Linear Regression Model

## Assumption (*Linear Regression Model*)

*The observations  $(Y_i, X_i)$  satisfy the linear regression equation*

$$Y_i = X_i' \beta + e_i$$

$$\mathbb{E}(e_i | X_i) = 0.$$

- ▶ Heteroskedastic regression:  $\mathbb{E}(e_i^2 | X_i) = \sigma^2(X_i) = \sigma_i^2$ .
- ▶ Homoskedastic regression: the conditional variance is constant.

## Assumption

*The conditional variance of the error*

$$\mathbb{E}(e_i^2 | X_i) = \sigma^2(X_i) = \sigma^2$$

*is independent of  $X_i$ .*

# Mean of Least-Squares Estimator

- ▶ Since the observations are assumed to be i.i.d., then

$$\mathbb{E}(Y_i | \mathbf{X}) \stackrel{\text{i.i.d.}}{=} \mathbb{E}(Y_i | X_i) = \mathbf{X}'_i \boldsymbol{\beta}.$$

- ▶ By the conditioning theorem and the linearity of expectations,

$$\begin{aligned}\mathbb{E}(\hat{\boldsymbol{\beta}} | \mathbf{X}) &= \mathbb{E}\left(\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i\right)^{-1} \left(\sum_{i=1}^n \mathbf{X}_i Y_i\right) \mid \mathbf{X}\right) \\ &= \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i\right)^{-1} \mathbb{E}\left(\left(\sum_{i=1}^n \mathbf{X}_i Y_i\right) \mid \mathbf{X}\right) \\ &= \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i\right)^{-1} \sum_{i=1}^n \mathbb{E}(\mathbf{X}_i Y_i \mid \mathbf{X}) \\ &= \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i\right)^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbb{E}(Y_i \mid \mathbf{X}) \\ &= \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i\right)^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i \boldsymbol{\beta} = \boldsymbol{\beta}.\end{aligned}$$

- Using matrix notation,

$$\mathbb{E}(\mathbf{Y} | \mathbf{X}) = \begin{pmatrix} \vdots \\ \mathbb{E}(Y_i | \mathbf{X}) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbf{X}'_i \boldsymbol{\beta} \\ \vdots \end{pmatrix} = \mathbf{X}\boldsymbol{\beta}$$

$$\mathbb{E}(\mathbf{e} | \mathbf{X}) = \begin{pmatrix} \vdots \\ \mathbb{E}(e_i | \mathbf{X}) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbb{E}(e_i | \mathbf{X}_i) \\ \vdots \end{pmatrix} = \mathbf{0}.$$

- By the conditioning theorem and the linearity of expectations,

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

- Since  $Y = X\beta + e$ ,

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

- By the conditioning theorem and the linearity of expectations,

$$\begin{aligned}\mathbb{E}(\hat{\beta} - \beta \mid X) &= \mathbb{E}\left((X'X)^{-1} X'e \mid X\right) \\ &= (X'X)^{-1} X'\mathbb{E}(e \mid X) \\ &= \mathbf{0}.\end{aligned}$$

## Theorem

*In the linear regression model and i.i.d. sampling*

$$\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta$$

- ▶ The conditional distribution of  $\hat{\beta}$  is centered at  $\beta$ .
- ▶ Applying the law of iterated expectations,

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\mathbb{E}(\hat{\beta} | \mathbf{X})) = \beta.$$

# Variance of Least Squares Estimator

- ▶ For any  $r \times 1$  random vector  $\mathbf{Z}$ , define the  $r \times r$  covariance matrix

$$\begin{aligned}\text{Var}(\mathbf{Z}) &= \mathbb{E}((\mathbf{Z} - \mathbb{E}(\mathbf{Z}))(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))') \\ &= \mathbb{E}(\mathbf{Z}\mathbf{Z}') - (\mathbb{E}(\mathbf{Z}))(\mathbb{E}(\mathbf{Z}))' .\end{aligned}$$

- ▶ For any pair  $(\mathbf{Z}, \mathbf{X})$ , define the conditional covariance matrix

$$\text{Var}(\mathbf{Z} | \mathbf{X}) = \mathbb{E}((\mathbf{Z} - \mathbb{E}(\mathbf{Z} | \mathbf{X}))(\mathbf{Z} - \mathbb{E}(\mathbf{Z} | \mathbf{X}))' | \mathbf{X}) .$$

- ▶ Define

$$\mathbf{V}_{\hat{\beta}} = \text{Var}(\hat{\beta} | \mathbf{X}) ,$$

the conditional covariance matrix of the LS estimators.



- ▶ The conditional covariance matrix of the error  $\mathbf{e}$  is

$$\text{Var}(\mathbf{e} \mid \mathbf{X}) = \mathbb{E}(\mathbf{e}\mathbf{e}' \mid \mathbf{X}) = \mathbf{D}.$$

The  $i^{\text{th}}$  diagonal element of  $\mathbf{D}$  is

$$\mathbb{E}(e_i^2 \mid \mathbf{X}) = \mathbb{E}(e_i^2 \mid \mathbf{X}_i) = \sigma_i^2$$

while the  $ij^{\text{th}}$  off-diagonal element of  $\mathbf{D}$  is

$$\mathbb{E}(e_i e_j \mid \mathbf{X}) \stackrel{\text{i.i.d.}}{=} \mathbb{E}(e_i \mid \mathbf{X}_i) \mathbb{E}(e_j \mid \mathbf{X}_j) = 0.$$

The first equality holds because of independence of the observations.

- ▶ Thus  $\mathbf{D}$  is a diagonal matrix with  $i^{\text{th}}$  diagonal element  $\sigma_i^2$ :

$$\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix}.$$

- ▶ In the special case of homoskedasticity,

$$\mathbb{E}(e_i^2 | \mathbf{X}_i) = \sigma_i^2 = \sigma^2$$

and we have  $\mathbf{D} = \mathbf{I}_n \sigma^2$ .

- ▶ For any  $n \times r$  matrix  $\mathbf{A} = \mathbf{A}(\mathbf{X})$ ,

$$\text{Var}(\mathbf{A}'\mathbf{Y} | \mathbf{X}) = \text{Var}(\mathbf{A}'\mathbf{e} | \mathbf{X}) = \mathbf{A}'\mathbf{D}\mathbf{A}.$$

- ▶ We write  $\hat{\boldsymbol{\beta}} = \mathbf{A}'\mathbf{Y}$  where  $\mathbf{A} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$  and thus

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = \text{Var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = \mathbf{A}'\mathbf{D}\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{D}\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}.$$

$$\mathbf{X}'\mathbf{D}\mathbf{X} = \sum_{i=1}^n X_i X_i' \sigma_i^2.$$

- ▶ In the special case of homoskedasticity,  $\mathbf{D} = \mathbf{I}_n \sigma^2$ , so

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

## Theorem

*In the linear regression model and i.i.d. sampling*

$$\begin{aligned} V_{\hat{\beta}} &= \text{Var}(\hat{\beta} | X) \\ &= (X'X)^{-1} X'DX (X'X)^{-1}. \end{aligned}$$

*In the homoskedastic linear regression model and i.i.d. sampling*

$$V_{\hat{\beta}} = \sigma^2 (X'X)^{-1}.$$

# Gauss-Markov Theorem

- ▶ Now consider the class of estimators that can be written as  $\tilde{\beta} = A'Y$ , where  $A$  is an  $n \times k$  matrix depending only on  $X$ .
- ▶ The LS estimator is the special case:  $A = X (X'X)^{-1}$ .
- ▶ The Gauss-Markov theorem says that the LS estimator is the best choice among linear unbiased estimators when the errors are homoskedastic, in the sense that the least-squares estimator has the smallest variance.

- ▶ For any linear estimator  $\tilde{\beta} = A'Y$  we have

$$\mathbb{E}(\tilde{\beta} | X) = A'\mathbb{E}(Y | X) = A'X\beta$$

so that  $\tilde{\beta}$  is unbiased if and only if  $A'X = I_k$ . Furthermore,

$$\text{Var}(\tilde{\beta} | X) = \text{Var}(A'Y | X) = A'DA = A'A\sigma^2.$$

- ▶ The best unbiased linear estimator is obtained by finding the matrix  $A_0$  satisfying  $A_0'X = I_k$  such that for any other matrix  $A$  satisfying  $A'X = I_k$  then  $A'A - A_0'A_0$  is positive semi-definite.

## Theorem

*In the homoskedastic linear regression model and i.i.d sampling, if  $\tilde{\beta}$  is a linear unbiased estimator of  $\beta$  then*

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

- ▶ The theorem is limited because the class of models is restricted to homoskedastic linear regression and the class of potential estimators is restricted to linear unbiased estimators.

# Residuals

- ▶ The residuals:

$$\hat{\mathbf{e}} = \mathbf{M}\mathbf{e}$$

where  $\mathbf{M} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ .

- ▶ We compute

$$\mathbb{E}(\hat{\mathbf{e}} | \mathbf{X}) = \mathbb{E}(\mathbf{M}\mathbf{e} | \mathbf{X}) = \mathbf{M}\mathbb{E}(\mathbf{e} | \mathbf{X}) = \mathbf{0}$$

$$\text{Var}(\hat{\mathbf{e}} | \mathbf{X}) = \text{Var}(\mathbf{M}\mathbf{e} | \mathbf{X}) = \mathbf{M}\text{Var}(\mathbf{e} | \mathbf{X})\mathbf{M} = \mathbf{M}\mathbf{D}\mathbf{M}.$$

- ▶ Under the assumption of conditional homoskedasticity,

$$\mathbb{E}(e_i^2 | \mathbf{X}_i) = \sigma^2 \text{ and } \text{Var}(\hat{\mathbf{e}} | \mathbf{X}) = \mathbf{M}\sigma^2.$$

## Estimation of Error Variance

- ▶ The method of moments estimator (MME) of  $\sigma^2 = \mathbb{E}(e_i^2)$  is the sample average of the squared residuals:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2.$$

- ▶ Observe

$$\hat{\sigma}^2 = \frac{1}{n} \mathbf{e}' \mathbf{M} \mathbf{e} = \frac{1}{n} \text{tr}(\mathbf{e}' \mathbf{M} \mathbf{e}) = \frac{1}{n} \text{tr}(\mathbf{M} \mathbf{e} \mathbf{e}')$$

and

$$\begin{aligned} \mathbb{E}(\hat{\sigma}^2 \mid \mathbf{X}) &= \frac{1}{n} \text{tr}(\mathbb{E}(\mathbf{M} \mathbf{e} \mathbf{e}' \mid \mathbf{X})) \\ &= \frac{1}{n} \text{tr}(\mathbf{M} \mathbb{E}(\mathbf{e} \mathbf{e}' \mid \mathbf{X})) \\ &= \frac{1}{n} \text{tr}(\mathbf{M} \mathbf{D}). \end{aligned}$$



- ▶ Under the homoskedasticity assumption  $\mathbb{E}(e_i^2 | \mathbf{X}_i) = \sigma^2$  so that  $\mathbf{D} = \mathbf{I}_n \sigma^2$ ,

$$\begin{aligned}\mathbb{E}(\hat{\sigma}^2 | \mathbf{X}) &= \frac{1}{n} \text{tr}(\mathbf{M}\sigma^2) \\ &= \sigma^2 \left( \frac{n-k}{n} \right).\end{aligned}$$

- ▶ To obtain an unbiased estimator is by rescaling the estimator:

$$s^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{e}_i^2.$$

Now  $\mathbb{E}(s^2 | \mathbf{X}) = \sigma^2$  and  $\mathbb{E}(s^2) = \sigma^2$ .

# Covariance Matrix Estimation Under Homoskedasticity

- ▶ Under homoskedasticity, the covariance matrix takes the relatively simple form  $V_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2$  which is known up to the unknown scale  $\sigma^2$ .
- ▶ The classic covariance matrix estimator:

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

- ▶  $\hat{V}_{\hat{\beta}}^0$  is conditionally unbiased for  $V_{\hat{\beta}}$  under homoskedasticity:

$$\begin{aligned}\mathbb{E}(\hat{V}_{\hat{\beta}}^0 \mid \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \mathbb{E}(s^2 \mid \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \\ &= V_{\hat{\beta}}^0.\end{aligned}$$

- ▶ This was the dominant covariance matrix estimator in applied econometrics for many years, and is still the default method in most regression packages.

# Covariance Matrix Estimation Under Heteroskedasticity

- ▶ If the estimator  $\hat{V}_{\hat{\beta}}^0$  is used, but the regression error is heteroskedastic, it is possible for  $\hat{V}_{\hat{\beta}}^0$  to be quite biased for

$$V_{\hat{\beta}} = (X'X)^{-1} (X'DX) (X'X)^{-1}$$

where

$$\begin{aligned} D &= \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \\ &= \mathbb{E}(\mathbf{e}\mathbf{e}' \mid \mathbf{X}). \end{aligned}$$

- ▶ If  $e_i^2, i = 1, \dots, n$  are observed, we can construct an unbiased estimator for  $V_{\hat{\beta}}$ :

$$\hat{V}_{\hat{\beta}}^{ideal} = (X'X)^{-1} \left( \sum_{i=1}^n X_i X_i' e_i^2 \right) (X'X)^{-1}.$$

► Compute

$$\begin{aligned}\mathbb{E}\left(\hat{\mathbf{V}}_{\hat{\beta}}^{ideal} \mid \mathbf{X}\right) &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \mathbb{E}\left(e_i^2 \mid \mathbf{X}\right) \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \sigma_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{V}_{\hat{\beta}}.\end{aligned}$$

► A feasible version:

$$\hat{\mathbf{V}}_{\hat{\beta}}^W = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \hat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}.$$

This is known as the White covariance matrix estimator. It is biased:  $\mathbb{E}\left(\hat{\mathbf{V}}_{\hat{\beta}}^{ideal} \mid \mathbf{X}\right) \neq \mathbf{V}_{\hat{\beta}}$ .

# Measures of Fit

- ▶ A commonly reported measure of regression fit is the regression  $R^2$ :

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\hat{\sigma}^2}{\hat{\sigma}_Y^2}.$$

- ▶  $R^2$  can be viewed as an estimator of the population parameter

$$\rho^2 = \frac{\text{Var}(X_i' \beta)}{\text{Var}(Y_i)} = 1 - \frac{\sigma^2}{\sigma_Y^2}.$$

- ▶  $\hat{\sigma}^2$  and  $\hat{\sigma}_Y^2$  are biased estimators. The adjusted  $R^2$  uses unbiased versions:

$$\bar{R}^2 = 1 - \frac{s^2}{\tilde{\sigma}_Y^2} = 1 - \frac{(n-k)^{-1} \sum_{i=1}^n \hat{e}_i^2}{(n-1)^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2}.$$

- ▶  $R^2$  cannot be used for model selection, as it necessarily increases when regressors are added to a regression model.

# Normal Regression Model

- ▶ The normal regression model is the linear regression model with an independent normal error

$$Y = X'\beta + e$$
$$e \sim N(0, \sigma^2).$$

- ▶ The likelihood is the name for the joint probability density of the data, evaluated at the observed sample, and viewed as a function of the parameters.
- ▶ The maximum likelihood estimator is the value which maximizes this likelihood function.

- ▶ The conditional density of  $Y$  given  $\mathbf{X}$ :

$$f(Y | \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2} (Y - \mathbf{X}'\boldsymbol{\beta})^2\right).$$

- ▶ The conditional density of  $\mathbf{Y}$  given  $\mathbf{X}$ :

$$\begin{aligned} f_{\mathbf{Y}|\mathbf{X}}(\mathbf{Y} | \mathbf{X}) &= \prod_{i=1}^n f_{Y_i|\mathbf{X}_i}(Y_i | \mathbf{X}_i) \\ &= \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2} (Y_i - \mathbf{X}_i'\boldsymbol{\beta})^2\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{X}_i'\boldsymbol{\beta})^2\right) \\ &= L(\boldsymbol{\beta}, \sigma^2). \end{aligned}$$

$L(\boldsymbol{\beta}, \sigma^2)$  is called the likelihood function.

- ▶ Work with the natural logarithm:

$$\begin{aligned}\log f(\mathbf{Y} | \mathbf{X}) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{X}'_i \boldsymbol{\beta})^2 \\ &= \log L(\boldsymbol{\beta}, \sigma^2).\end{aligned}$$

- ▶ The MLE:

$$\left( \hat{\boldsymbol{\beta}}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2 \right) = \underset{\boldsymbol{\beta} \in \mathbb{R}^k, \sigma^2 > 0}{\operatorname{argmax}} \log L(\boldsymbol{\beta}, \sigma^2).$$



- ▶ In most applications of maximum likelihood, the MLE must be found by numerical methods. However, in the case of the normal regression model we can find an explicit expression.
- ▶ FOC:

$$0 = \frac{\partial \log L(\boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} \bigg|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{mle}}, \sigma^2=\hat{\sigma}_{\text{mle}}^2} = \frac{1}{\hat{\sigma}_{\text{mle}}^2} \sum_{i=1}^n \mathbf{X}_i (Y_i - \mathbf{X}_i' \hat{\boldsymbol{\beta}}_{\text{mle}})$$

$$0 = \frac{\partial \log L(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} \bigg|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{mle}}, \sigma^2=\hat{\sigma}_{\text{mle}}^2} = -\frac{n}{2\hat{\sigma}_{\text{mle}}^2} + \frac{1}{\hat{\sigma}_{\text{mle}}^4} \sum_{i=1}^n (Y_i - \mathbf{X}_i' \hat{\boldsymbol{\beta}}_{\text{mle}}).$$

- ▶ The MLE:

$$\hat{\beta}_{\text{mle}} = \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i Y_i \right) = \hat{\beta}_{\text{ols}}.$$

- ▶ The MLE for  $\sigma^2$ :

$$\hat{\sigma}_{\text{mle}}^2 = \frac{1}{n} \sum_{i=1}^n \left( Y_i - \mathbf{X}_i' \hat{\beta}_{\text{mle}} \right)^2 = \frac{1}{n} \sum_{i=1}^n \left( Y_i - \mathbf{X}_i' \hat{\beta}_{\text{ols}} \right)^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2.$$

- ▶ Maximized log-likelihood is a measure of goodness of fit:

$$\log L \left( \hat{\beta}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2 \right) = -\frac{n}{2} \log \left( 2\pi \hat{\sigma}_{\text{mle}}^2 \right) - \frac{n}{2}.$$

## Distribution of OLS Coefficient Vector

- ▶ The normality assumption  $e_i | X_i \sim N(0, \sigma^2)$  and iid assumption imply

$$\mathbf{e} | \mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2).$$

- ▶ The OLS estimator satisfies

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e},$$

which is a linear function of  $\mathbf{e}$ .

- ▶ Conditional on  $\mathbf{X}$ ,

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} | \mathbf{X} &\sim (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' N(0, \mathbf{I}_n \sigma^2) \\ &\sim N\left(0, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}\right) \\ &= N\left(0, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}\right) \end{aligned}$$

or

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

- ▶ This shows that under the assumption of normal errors, the OLS estimate has an exact normal distribution.

### Theorem

*In the linear regression model,*

$$\hat{\beta} \mid \mathbf{X} \sim \mathbf{N} \left( \beta, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \right)$$

- ▶ Any linear function of the OLS estimate is also normally distributed, including individual estimates:

$$\hat{\beta}_j \mid \mathbf{X} \sim \mathbf{N} \left( \beta_j, \sigma^2 \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]_{jj} \right).$$

## Distribution of OLS Residual Vector

- ▶ The OLS residual vector:  $\hat{e} = \mathbf{M}e$ .  $\hat{e}$  is linear in  $e$ .
- ▶ Conditional on  $X$ ,

$$\hat{e} = \mathbf{M}e \mid X \sim N\left(0, \sigma^2 \mathbf{M} \mathbf{M}'\right) = N\left(0, \sigma^2 \mathbf{M}\right).$$

- ▶ The joint distribution of  $\hat{\beta}$  and  $\hat{e}$ :

$$\begin{pmatrix} \hat{\beta} - \beta \\ \hat{e} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'e \\ \mathbf{M}e \end{pmatrix} = \begin{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\ \mathbf{M} \end{pmatrix} e.$$

- ▶ So

$$\begin{pmatrix} \hat{\beta} - \beta \\ \hat{e} \end{pmatrix} \mid X \sim N\left(\mathbf{0}, \begin{pmatrix} \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} & 0 \\ 0 & \sigma^2 \mathbf{M} \end{pmatrix}\right).$$

### Theorem

*In the linear regression model,  $\hat{e} \mid X \sim N(0, \sigma^2 \mathbf{M})$  and is independent of  $\hat{\beta}$ .*

## Distribution of Variance Estimate

- ▶  $s^2 = \hat{\mathbf{e}}'\hat{\mathbf{e}}/(n - k) = \mathbf{e}'\mathbf{M}\mathbf{e}/(n - k)$ .
- ▶ The spectral decomposition of  $\mathbf{M}$ :  $\mathbf{M} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  with  $\mathbf{H}'\mathbf{H} = \mathbf{I}_n$  and  $\mathbf{\Lambda}$  is diagonal with the eigenvalues of  $\mathbf{M}$  on the diagonal.
- ▶ Since  $\mathbf{M}$  is idempotent with rank  $n - k$ , it has  $n - k$  eigenvalues equalling 1 and  $k$  eigenvalues equalling 0:

$$\mathbf{\Lambda} = \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_k \end{bmatrix}.$$

- ▶  $\mathbf{U} = \mathbf{H}'\mathbf{e} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_n\sigma^2)$ .

$$\begin{aligned} (n - k) \frac{s^2}{\sigma^2} &= \frac{\mathbf{e}'\mathbf{M}\mathbf{e}}{\sigma^2} \\ &= \frac{\mathbf{e}'\mathbf{H}}{\sigma} \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_k \end{bmatrix} \frac{\mathbf{H}'\mathbf{e}}{\sigma} \\ &= \frac{\mathbf{U}'}{\sigma} \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_k \end{bmatrix} \frac{\mathbf{U}}{\sigma} \\ &\sim \chi_{n-k}^2. \end{aligned}$$

## Theorem

*In the linear regression model, conditional on  $\mathbf{X}$ ,*

$$\frac{(n - k) s^2}{\sigma^2} \sim \chi_{n-k}^2$$

*and is independent of  $\hat{\boldsymbol{\beta}}$ .*

## $t$ -statistic

- ▶ The “ $z$ -statistic”:

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]_{jj}}} \sim \text{N}(0, 1).$$

- ▶ Replace the unknown variance  $\sigma^2$  with its estimate  $s^2$ :

$$T = \frac{\hat{\beta}_j - \beta_j}{\sqrt{s^2 \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]_{jj}}} = \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)}.$$

- ▶ Write the  $t$ -statistic as the ratio of the standardized statistic and the square root of the scaled variance estimate:

$$\begin{aligned} T &= \frac{\hat{\beta}_j - \beta_j}{\sqrt{s^2 \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]_{jj}}} / \sqrt{\frac{(n-k)s^2}{\sigma^2} / (n-k)} \\ &\sim \frac{\text{N}(0, 1)}{\sqrt{\chi_{n-k}^2 / (n-k)}} \\ &\sim t_{n-k}. \end{aligned}$$



## Theorem

*In the normal regression model,  $T \sim t_{n-k}$ .*

- ▶ This derivation shows that the  $t$ -statistic has a sampling distribution which depends only on the quantity  $n - k$ . The distribution does not depend on any other features of the data.
- ▶ In this context, we say that the distribution of the  $t$ -statistic is pivotal, meaning that it does not depend on unknowns.
- ▶ The theorem only applies to the  $t$ -statistic constructed with the homoskedastic standard error estimate. It does not apply to a  $t$ -statistic constructed with the robust standard error estimates.

# Confidence Intervals for Regression Coefficients

- ▶ An OLS estimate  $\hat{\beta}$  is a point estimate for the coefficients  $\beta$ .
- ▶ An interval estimate takes the form  $\hat{C} = [\hat{L}, \hat{U}]$ . The goal of an interval estimate  $\hat{C}$  is to contain the true value with high probability.
- ▶ The interval estimate  $\hat{C}$  is a function of the data and hence is random.
- ▶ An interval estimate  $\hat{C}$  is called a  $1 - \alpha$  confidence interval when  $\Pr(\beta \in \hat{C}) = 1 - \alpha$ .
- ▶ A good choice for a confidence interval is by adding and subtracting from the estimate  $\hat{\beta}$  a fixed multiple of the standard error:

$$\hat{C} = [\hat{\beta} - c \cdot s(\hat{\beta}), \hat{\beta} + c \cdot s(\hat{\beta})].$$

- ▶  $\widehat{C}$  is the set of parameter values for  $\beta$  such that the t-statistic  $T(\beta)$  is smaller than some constant  $c$ :

$$\widehat{C} = [\beta : |T(\beta)| \leq c] = \left\{ \beta : -c \leq \frac{\hat{\beta} - \beta}{s(\hat{\beta})} \leq c \right\}.$$

- ▶ The coverage probability is

$$\begin{aligned} \Pr(\beta \in \widehat{C}) &= \Pr(|T(\beta)| \leq c) \\ &= \Pr(-c \leq T(\beta) \leq c) \\ &= 2 \cdot F(c) - 1 \end{aligned}$$

where  $F$  is the  $t$  distribution with  $n - k$  degrees of freedom ( $F(-c) = 1 - F(c)$ ).

### Theorem

*In the normal regression model,  $\widehat{C}$  with  $c = F^{-1}(1 - \alpha/2)$  has coverage probability  $\Pr(\beta \in \widehat{C}) = 1 - \alpha$ .*

# Hypothesis Testing

- ▶ Let  $\theta \in \Theta \subset \mathbb{R}^d$  be a parameter of interest. Some examples of  $\theta$  include:
  - ▶ The coefficient of one of the regressors:  $\theta = \beta_1$ ,  $d = 1$ ,  $\Theta = \mathbb{R}$ .
  - ▶ A vector of coefficients:  $\theta = (\beta_1, \dots, \beta_l)'$ ,  $d = l$ ,  $\Theta = \mathbb{R}^l$ .
  - ▶ The variance of errors:  $\theta = \sigma^2$ ,  $d = 1$ ,  $\Theta = (0, \infty)$ .
- ▶ A statistical hypothesis is an assertion about  $\theta$ . Usually, we have two competing hypotheses, and we want to draw a conclusion, based on the data, as to which of the hypotheses is true. Let  $\Theta_0 \subset \Theta$  and  $\Theta_1 \subset \Theta$  such that  $\Theta_0 \cap \Theta_1 = \emptyset$  and  $\Theta_0 \cup \Theta_1 = \Theta$ . The two competing hypotheses are:
  - ▶ Null hypothesis  $\mathbb{H}_0 : \theta \in \Theta_0$ . This is a hypothesis that is held as true, unless data provides sufficient evidence against it.
  - ▶ Alternative hypothesis  $\mathbb{H}_1 : \theta \in \Theta_1$ . This is a hypothesis against which the null is tested. It is held to be true if the null is found false.

- ▶ The subsets  $\Theta_0$  and  $\Theta_1$  are chosen by the econometrician and therefore are known. Their union defines the maintained hypothesis, i.e. the space of values that  $\theta$  can take. For example, when  $\Theta = \mathbb{R}$ , one may consider  $\Theta_0 = \{0\}$ , and  $\Theta_1 = \mathbb{R} \setminus \{0\}$ . Another example is  $\Theta_0 = (-\infty, 0]$  and  $\Theta_1 = (0, \infty)$ .
- ▶ When  $\Theta_0$  has exactly one element ( $\Theta_0$  is a singleton), we say that  $\mathbb{H}_0 : \theta \in \Theta_0$  is a simple hypothesis. Otherwise, we say that  $\mathbb{H}_0$  is a composite hypothesis. Similarly,  $\mathbb{H}_1 : \theta \in \Theta_1$  can be simple or composite depending on whether  $\Theta_1$  is a singleton or not.
- ▶ Let  $S \in \mathcal{S}$  denote a statistic and the range of its values. A decision rule is defined by a partition of  $\mathcal{S}$  into acceptance region  $\mathcal{A}$  and rejection (critical) region  $\mathcal{R}$  ( $\mathcal{A} \cap \mathcal{R} = \emptyset$  and  $\mathcal{A} \cup \mathcal{R} = \mathcal{S}$ ).
- ▶  $\mathbb{H}_0$  is rejected when the test statistic falls in to the rejection region  $\mathcal{R}$ .

# Type I and Type II Errors

- ▶ There are two types of errors that the econometrician can make:

		Truth	
		$H_0$	$H_1$
Decision	$H_0$	✓	Type II error
	$H_1$	Type I error	✓

- ▶ Type I error is the error of rejecting  $H_0$  when  $H_0$  is true.
- ▶ Type II error is the error of accepting  $H_0$  when  $H_1$  is true.

# Power Function

- ▶ The probabilities of Type I and II errors can be described using the power function.
- ▶ Consider a test based on  $S$  that rejects  $\mathbb{H}_0$  when  $S \in \mathcal{R}$ . The power function of this test is defined as:

$$\pi(\boldsymbol{\theta}) = \Pr_{\boldsymbol{\theta}}(S \in \mathcal{R}),$$

where  $\Pr_{\boldsymbol{\theta}}(\cdot)$  denotes that the probability must be calculated under the assumption that the true value of the parameter is  $\boldsymbol{\theta}$ .

- ▶ The largest probability of Type I error (rejecting  $\mathbb{H}_0$  when it is true) is

$$\sup_{\boldsymbol{\theta} \in \Theta_0} \pi(\boldsymbol{\theta}) = \sup_{\boldsymbol{\theta} \in \Theta_0} \Pr_{\boldsymbol{\theta}}(S \in \mathcal{R}).$$

The expression above is also called the *size* of a test.

- ▶ When  $\mathbb{H}_0$  is simple, i.e.  $\Theta = \{\boldsymbol{\theta}_0\}$ , the size can be computed simply as  $\pi(\boldsymbol{\theta}_0) = \Pr_{\boldsymbol{\theta}_0}(S \in \mathcal{R})$ .

- ▶ The probability of Type II error is:

$$1 - \pi(\theta) = 1 - \Pr_{\theta}(S \in \mathcal{R}) \quad \text{for } \theta \in \Theta_1.$$

- ▶ Typically,  $\Theta_1$  has many elements, and therefore the probability of Type II error depends on the true value  $\theta$ .
- ▶ One would like to have the probabilities of Type I and II errors to be as small as possible, but unfortunately, they are inversely related.
- ▶ To reduce the probability of Type I error, one should make  $\mathcal{R}$  smaller. This, however, will increase the probability of Type II error.



## Definition

A test with power function  $\pi(\theta)$  is said to be a level  $\alpha$  test if  $\sup_{\theta \in \Theta_0} \pi(\theta) \leq \alpha$ . We say it is a size  $\alpha$  test if  $\sup_{\theta \in \Theta_0} \pi(\theta) = \alpha$ .

- ▶ Significance level of a test is the largest Type I error probability one tolerates. Typically, the significance level is chosen to be a small number close to zero: for example, 0.01, 0.05, 0.10.
- ▶ By convention, a valid test must control the probability of Type I error (level  $\alpha$  test, where  $\alpha$  is equal to the significance level).
- ▶ We want the probability of a Type II error probability to be as small as possible for given Type I error probability.

# Steps of Hypothesis Testing

1. Specify  $\mathbb{H}_0$  and  $\mathbb{H}_1$ .
2. Choose the significance level  $\alpha$ .
3. Define a decision rule (a test statistic  $S$  and a rejection region  $\mathcal{R}_\alpha$ ) so that the resulting test is a level  $\alpha$  test. Note that  $\mathcal{R}_\alpha$  typically depends on  $\alpha$ .
4. Perform the test.

# $p$ -Value

- ▶ The lowest significance level consistent with rejecting  $\mathbb{H}_0$  is called the  $p$ -value:

$$p\text{-value} = \min \{0 < \alpha < 1 : S \in \mathcal{R}_\alpha\}.$$

- ▶ Note that  $p$ -value is a statistic and a measure of the evidence against  $\mathbb{H}_0$ .
- ▶ If the  $p$ -value is smaller than our tolerance (significance level), then we reject  $\mathbb{H}_0$ .

# Power of a Test

- ▶ The power of a test with the power function  $\pi(\boldsymbol{\theta})$  is defined as

$$\pi(\boldsymbol{\theta}) \quad \text{for } \boldsymbol{\theta} \in \Theta_1.$$

- ▶ Given two level  $\alpha$  tests, we should prefer a more powerful test.
- ▶ We say that a level  $\alpha$  test with power function  $\pi_1(\boldsymbol{\theta})$  is uniformly more powerful than a level  $\alpha$  test with power function  $\pi_2(\boldsymbol{\theta})$  if  $\pi_1(\boldsymbol{\theta}) \geq \pi_2(\boldsymbol{\theta})$  for all  $\boldsymbol{\theta} \in \Theta_1$ .

## A Simple Example

- ▶ Consider a sample  $(X_1, X_2, \dots, X_n)$  from a normal population with mean  $\mu$  and variance  $\sigma^2$ . Suppose we know  $\sigma^2$  for now. Consider  $\mathbb{H}_0 : \mu = 0$  against  $\mathbb{H}_1 : \mu > 0$ .
- ▶ Consider  $T = \frac{\bar{X}}{\sigma/\sqrt{n}}$ , where  $\bar{X} = n^{-1} \sum_{i=1}^n X_i$ , which is a  $N(0, 1)$  random variable under  $\mathbb{H}_0$ .
- ▶ Consider  $\mathcal{R}_\alpha = [z_{1-\alpha}, \infty)$ , where  $\Pr(N(0, 1) > z_{1-\alpha}) = \alpha$ . Under  $\mathbb{H}_0$ ,  $\Pr[T \in \mathcal{R}_\alpha] = \alpha$ .
- ▶ For any given value of  $\mu$ , define

$$T_\mu = T - \frac{\mu}{\sigma/\sqrt{n}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}},$$

which is always  $N(0, 1)$  if the true mean is  $\mu$ .

- ▶ We reject the test when  $T \geq z_{1-\alpha}$ , which holds if and only if

$$T_\mu = T - \frac{\mu}{\sigma/\sqrt{n}} \geq z_{1-\alpha} - \frac{\mu}{\sigma/\sqrt{n}}.$$

It follows for this simple example, that the power function is

$$\pi(\mu) = 1 - \Phi\left(z_{1-\alpha} - \frac{\sqrt{n}\mu}{\sigma}\right).$$

- ▶ We notice that  $\pi(\mu)$  is smaller for all  $\mu$  if the significance level  $\alpha$  is smaller (and hence  $z_{1-\alpha}$  is larger). This reflects the trade-off between Type I error and Type II error probabilities: we cannot reduce both simultaneously.
- ▶  $\pi(\mu)$  is increasing in  $\mu$ . For  $\mu$ 's that are farther away from 0, the test can detect such deviation at a higher probability.
- ▶ As  $\mu \rightarrow \infty$ , the power converges to 1. The test is very likely to reject  $\mathbb{H}_0$  if the true mean is very large.
- ▶  $\pi(\mu)$  increases with the sample size  $n$ . The test can detect falseness of  $\mathbb{H}_0$  at a higher probability if our sample contains more information.

## $t$ Test

- ▶ The null hypothesis:

$$\mathbb{H}_0 : \beta_j = \beta_{j,0}.$$

- ▶ The alternative hypothesis:

$$\mathbb{H}_1 : \beta_j \neq \beta_{j,0}.$$

- ▶ The standard testing statistic is

$$|T| = \left| \frac{\hat{\beta}_j - \beta_{j,0}}{s(\hat{\beta}_j)} \right|.$$

- ▶ If  $\mathbb{H}_0$  is true, we expect  $|T|$  to be small, but if  $\mathbb{H}_1$  is true, then we would expect  $|T|$  to be large. Hence the standard rule is to reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  for large values of the t-statistic  $|T|$ :

Reject  $\mathbb{H}_0$  if  $|T| > c$ .

- ▶  $c$  is called the critical value. Its value is selected to control the probability of false rejections.
- ▶ When the null hypothesis is true,  $T$  has an exact student distribution. The probability of false rejection is

$$\begin{aligned}
 \Pr(\text{Reject } \mathbb{H}_0 \mid \mathbb{H}_0) &= \Pr(|T| > c \mid \mathbb{H}_0) \\
 &= \Pr(T > c \mid \mathbb{H}_0) + \Pr(T < -c \mid \mathbb{H}_0) \\
 &= 1 - F(c) + F(-c) \\
 &= 2(1 - F(c)).
 \end{aligned}$$

- ▶ We select the value  $c$  so that this probability equals the significance level:  $F(c) = 1 - \alpha/2$ .
- ▶ The  $p$ -value of a  $t$ -statistic is  $p = 2(1 - F(|T|))$ .

### Theorem

*In the normal regression model, if the null hypothesis is true, then  $|T| \sim t_{n-k}$ . If  $c$  is set so that  $\Pr(|t_{n-k}| \geq c) = \alpha$ , then the test “Reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $|T| > c$ ” has level  $\alpha$ .*



# Power

- ▶ Assume that the true value is given by  $\beta_j$ . Assume for simplicity that  $\sigma^2$  is known, so that  $s(\hat{\beta}_j) = \sqrt{\sigma^2 \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]_{jj}}$ .

- ▶ Write

$$T = \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} + \frac{\beta_j - \beta_{j,0}}{s(\hat{\beta}_j)}. \quad (1)$$

- ▶ We have that

$$Z = \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} \mid \mathbf{X} \sim \mathbf{N}(0, 1)$$

and

$$\frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} + \frac{\beta_j - \beta_{j,0}}{s(\hat{\beta}_j)} \mid \mathbf{X} \sim \mathbf{N}\left(\frac{\beta_j - \beta_{j,0}}{s(\hat{\beta}_j)}, 1\right).$$

- ▶ The critical value is  $c = \Phi^{-1}(1 - \alpha/2)$  ( $\Pr(Z > c) = \alpha/2$ ).

- ▶ If the null hypothesis is false, the distribution of the test statistic is not centered around zero, and we will see rejection rates higher than  $\alpha$ . The probability to reject is a function of the true value  $\beta_j$  and depends on the magnitude of  $|\beta_j - \beta_{j,0}| / s(\hat{\beta}_j)$ .
- ▶ Now

$$\begin{aligned} \pi(\beta_1) &= \Pr\left(\left|\frac{\hat{\beta}_j - \beta_{j,0}}{s(\hat{\beta}_j)}\right| > c\right) = \Pr\left(\left|\frac{\hat{\beta}_j - \beta_j + \beta_j - \beta_{j,0}}{s(\hat{\beta}_j)}\right| > c\right) \\ &= \Pr\left(\left|Z + \frac{\beta_j - \beta_{j,0}}{s(\hat{\beta}_j)}\right| > c\right). \end{aligned}$$

# One-sided Test

- ▶ In the case of one-sided tests, the null and alternative hypotheses may be specified as

$$\mathbb{H}_0 : \beta_j \leq \beta_{j,0},$$

$$\mathbb{H}_1 : \beta_j > \beta_{j,0}.$$

- ▶ In this case, a valid test should satisfy the following condition:

$$\sup_{\beta_j \leq \beta_{j,0}} \Pr(\text{reject } \mathbb{H}_0 \mid \beta_j) \leq \alpha, \quad (2)$$

i.e. the maximum probability to reject  $H_0$  when it is true should not exceed  $\alpha$ .

- ▶ Consider the following test (decision rule):

Reject  $\mathbb{H}_0$  when  $T > c$ .

where  $c$  is set so that  $\Pr(t_{n-k} \geq c) = \alpha$ .

- ▶ Under  $\mathbb{H}_0$ , we have:

$$\begin{aligned}\Pr(\text{reject } \mathbb{H}_0 \mid \beta_j \leq \beta_{j,0}) &= \Pr(T > c \mid \beta_j \leq \beta_{j,0}) \\ &= \Pr\left(\frac{\hat{\beta}_j - \beta_{j,0}}{s(\hat{\beta}_j)} > c \mid \beta_j \leq \beta_{j,0}\right) \\ &\leq \Pr\left(\frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} > c \mid \beta_j \leq \beta_{j,0}\right) \\ &\quad \text{since } \beta_j \leq \beta_{j,0} \\ &= \alpha \text{ (since } \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} \sim t_{n-k}).\end{aligned}$$

- ▶ The size control condition is satisfied.

# Testing a Single Linear Restriction

- ▶ Suppose we want to test

$$\mathbb{H}_0 : \mathbf{c}'\boldsymbol{\beta} = r,$$

$$\mathbb{H}_1 : \mathbf{c}'\boldsymbol{\beta} \neq r.$$

- ▶ In this case,  $\mathbf{c}$  is a  $k$ -vector,  $r$  is a scalar, and under the null hypothesis

$$c_1\beta_1 + \dots + c_k\beta_k - r = 0.$$

- ▶ We have that the LS estimator of  $\boldsymbol{\beta}$

$$\hat{\boldsymbol{\beta}} \mid \mathbf{X} \sim N\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}\right).$$

Then, under  $H_0$ ,

$$\frac{\mathbf{c}'\hat{\boldsymbol{\beta}} - r}{\sqrt{\sigma^2 \mathbf{c}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{c}}} \mid \mathbf{X} \sim N(0, 1).$$

- ▶ The  $t$ -statistic

$$\begin{aligned}
 T &= \frac{\mathbf{c}'\hat{\boldsymbol{\beta}} - r}{\sqrt{s^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}} \\
 &= \left( \frac{\mathbf{c}'\hat{\boldsymbol{\beta}} - r}{\sqrt{\sigma^2\mathbf{c}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{c}}} \right) / \sqrt{\frac{\mathbf{e}'\mathbf{M}\mathbf{e}}{\sigma^2} / (n - k)}.
 \end{aligned}$$

- ▶  $\mathbf{e}'\mathbf{M}\mathbf{e}/\sigma^2 \mid \mathbf{X} \sim \chi_{n-k}^2$  and is independent of  $\hat{\boldsymbol{\beta}}$ . Therefore, under  $\mathbb{H}_0$ ,  $T \mid \mathbf{X} \sim t_{n-k}$ .
- ▶ The significance level  $\alpha$  two-sided test of  $\mathbb{H}_0 : \mathbf{c}'\boldsymbol{\beta} = r$  is given by “reject  $\mathbb{H}_0$  if  $|T| > c$ ”, where  $\Pr(|t_{n-k}| > c) = \alpha$ .

# Testing Multiple Linear Restrictions

- ▶ Suppose we want to test

$$\mathbb{H}_0 : \mathbf{R}\boldsymbol{\beta} = \mathbf{r},$$

$$\mathbb{H}_1 : \mathbf{R}\boldsymbol{\beta} \neq \mathbf{r},$$

where  $\mathbf{R}$  is a  $q \times k$  matrix and  $\mathbf{r}$  is a  $q$ -vector.

- ▶  $\mathbf{R} = \mathbf{I}_k, \mathbf{r} = \mathbf{0}$ . In this case, we test that  $\beta_1 = \dots = \beta_k = 0$ .
- ▶  $\mathbf{R} = \begin{pmatrix} 1 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \end{pmatrix}, \mathbf{r} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ . In this case,  
 $H_0 : \beta_1 + \beta_2 = 1, \beta_3 = 0$ .
- ▶ Consider the  $F$ -statistic

$$F = \frac{(RSS_r - RSS_{ur}) / q}{RSS_{ur} / (n - k)}.$$

- ▶  $RSS_r$ : the restricted Residual Sum of Squares.
- ▶  $RSS_{ur}$ : the unrestricted Residual Sum of Squares.

- ▶ Consider the restricted problem

$$\min_{\mathbf{b}} (\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b}) \quad \text{subject to } \mathbf{R}\mathbf{b} = \mathbf{r}.$$

- ▶ A Lagrangian function for this problem is

$$L(\mathbf{b}, \boldsymbol{\lambda}) = (\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b}) + 2\boldsymbol{\lambda}' (\mathbf{R}\mathbf{b} - \mathbf{r}),$$

where  $\boldsymbol{\lambda}$  is a  $q$ -vector.

- ▶ Let  $\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\lambda}}$  be the solution, where  $\tilde{\boldsymbol{\beta}}$  is the restricted LS estimator. It has to satisfy the first-order conditions

$$\frac{\partial L(\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\lambda}})}{\partial \mathbf{b}} = 2\mathbf{X}'\mathbf{X}\tilde{\boldsymbol{\beta}} - 2\mathbf{X}'\mathbf{Y} + 2\mathbf{R}'\tilde{\boldsymbol{\lambda}} = \mathbf{0}, \quad (3)$$

$$\frac{\partial L(\tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\lambda}})}{\partial \boldsymbol{\lambda}} = \mathbf{R}\tilde{\boldsymbol{\beta}} - \mathbf{r} = \mathbf{0}. \quad (4)$$

- ▶ The restricted LS estimator is given by

$$\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' \left( \mathbf{R} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' \right)^{-1} \left( \mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} \right),$$

where  $\hat{\boldsymbol{\beta}}$  is the LS estimator without the restriction  $\mathbf{R}\mathbf{b} = \mathbf{r}$ .



- Define the restricted residuals

$$\begin{aligned}\tilde{\mathbf{e}} &= \mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}} \\ &= (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\left(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \\ &= \hat{\mathbf{e}} + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\left(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}),\end{aligned}$$

- Then,

$$\begin{aligned}RSS_r &= \tilde{\mathbf{e}}'\tilde{\mathbf{e}} \\ &= \hat{\mathbf{e}}'\hat{\mathbf{e}} + (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})'\left(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \\ &\quad + 2\hat{\mathbf{e}}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\left(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \\ &= RSS_{ur} + (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})'\left(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}).\end{aligned}$$

- Since  $s^2 = \hat{\mathbf{e}}'\hat{\mathbf{e}}/(n - k) = RSS_{ur}/(n - k)$ ,

$$F = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})'\left(s^2\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'\right)^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})/q.$$

- ▶ We show next that under  $\mathbb{H}_0$ ,  $F | \mathbf{X} \sim F_{q,n-k}$ . First,

$$\mathbf{R}\hat{\boldsymbol{\beta}} | \mathbf{X} \sim \text{N}\left(\mathbf{R}\boldsymbol{\beta}, \sigma^2 \mathbf{R} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'\right).$$

- ▶ Then, under  $\mathbb{H}_0$ ,

$$\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} | \mathbf{X} \sim \text{N}\left(0, \sigma^2 \mathbf{R} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'\right).$$

It follows that

$$\left(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}\right)' \left(\sigma^2 \mathbf{R} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'\right)^{-1} \left(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}\right) | \mathbf{X} \sim \chi_q^2.$$

The result follows from

$\mathbf{e}'\mathbf{M}\mathbf{e}/\sigma^2 | \mathbf{X} \sim \chi_{n-k}^2$  and independent of  $\hat{\boldsymbol{\beta}}$  and the definition of  $F$ -distribution.

- ▶ Therefore, the test is given by “reject  $\mathbb{H}_0$  if  $F > c$ ”, where  $\Pr(F_{q,n-k} > c) = \alpha$ .

# Test of Model Significance

- ▶ Consider a model with the intercept

$$Y_i = \beta_1 + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + U_i,$$

- ▶ Consider the null hypothesis  $H_0 : \beta_2 = \dots \beta_k = 0$ . The restricted model is given by

$$Y_i = \beta_1 + U_i.$$

- ▶ In this case, the restricted LS estimator is  $\tilde{\beta}_1 = n^{-1} \sum_{i=1}^n Y_i = \bar{Y}$ , and  $RSS_r = TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2$ . In this case,

$$\begin{aligned} F &= \frac{(TSS - RSS_{ur}) / (k - 1)}{RSS_{ur} / (n - k)} \\ &= \frac{ESS / (k - 1)}{RSS_{ur} / (n - k)} \\ &= \frac{R^2 / (k - 1)}{(1 - R^2) / (n - k)} \\ &\sim F_{k-1, n-k}. \end{aligned}$$