

Chapter 5. Least Squares Estimation – Large-Sample Properties*

In Section 3 of Chapter 3, we derived the mean and variance of the least-squares estimator in the context of the linear regression model, but this is not a complete description of the sampling distribution and is thus not sufficient for inference. Furthermore, the theory does not apply in the context of the linear projection model, which is more relevant for empirical applications. In Section 8 of Chapter 3, we assume $u|\mathbf{x} \sim N(0, \sigma^2)$ and study the conditional distribution of $\hat{\beta}$, t -statistics and F -statistics given \mathbf{X} . In general the distribution of $u|\mathbf{x}$ is unknown and even if it were known, the unconditional distribution of $\hat{\beta}$ is hard to derive since $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ is a complicated function of $\{\mathbf{x}_i\}_{i=1}^n$. Perhaps we can view the results in the normal regression model as some sort of approximation to the sampling distributions without requiring the assumption of normality, but how can we be precise about this?

Take a simple example for illustration. Let y_i and x_i be drawn from the joint density

$$f(x, y) = \frac{1}{2\pi xy} \exp\left\{-\frac{1}{2}(\log y - \log x)^2\right\} \exp\left\{-\frac{1}{2}(\log x)^2\right\}$$

and let $\hat{\beta}_2$ be the slope coefficient estimate from a least-squares regression of y_i on x_i and a constant. Using simulation methods, the density function of $\hat{\beta}_2$ was computed and plotted in Figure 1 for sample sizes of $n = 25, 100$ and 1000 . The vertical line marks the true projection coefficient. From the Figure we can see that the density functions are dispersed and highly non-normal when n is small. As the sample size increases the density becomes more concentrated about the population coefficient and more like normality.

The asymptotic (or large sample) approach rigorizes the intuition above and tries to approximate the sampling distribution of $\hat{\beta}$ based on the limiting experiment that the sample size n tends to infinity. It aims to assess the distribution of $\hat{\beta}$ in actual practical samples where n is finite. A preliminary step in this approach is the demonstration that the estimator converges in probability to the true parameter as the sample size gets large. The second step is to study the distributional properties of $\hat{\beta}$ in the neighborhood of the true value, that is, the asymptotic normality of $\hat{\beta}$. The final step is to estimate the asymptotic variance which is necessary in statistical inferences such as hypothesis testing and confidence interval (CI) construction. In hypothesis testing, it is necessary to construct test statistics and derive their asymptotic distributions under the null. We will study the t -test and three asymptotically equivalent tests under both homoskedasticity and heteroskedasticity. It is also standard to develop the local power function for illustrating the power

*Email: pingyu@hku.hk

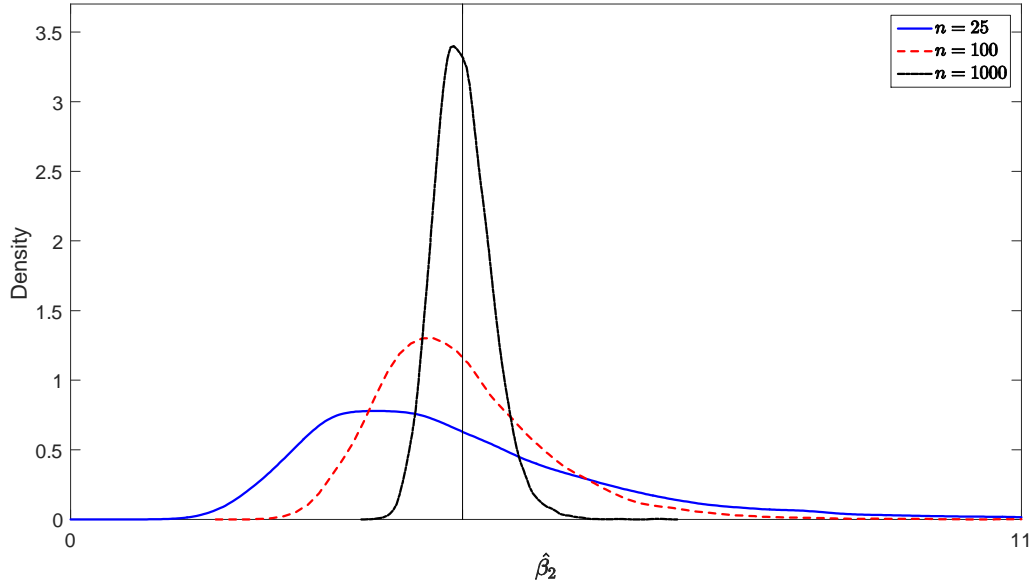


Figure 1: Finite Sample Density of $\hat{\beta}_2$ When $n = 25, 100$ and 1000

properties of the test.

This chapter concentrates on asymptotic properties related to the LSE. Related materials can be found in Chapter 2 of Hayashi (2000), Chapter 4 of Cameron and Trivedi (2005), Chapter 4 of Wooldridge (2010), and Chapters 7, 8 and 9 of Hansen (2022).

1 Asymptotics for the LSE

We first show that the LSE is CAN and then re-derive its asymptotic distribution by treating it as a MoM estimator.

1.1 Consistency

It is useful to express $\hat{\beta}$ as

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{u}) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}. \quad (1)$$

To show $\hat{\beta}$ is consistent, we impose the following additional assumptions.

Assumption OLS.1': $\text{rank}(E[\mathbf{x}\mathbf{x}']) = k$.

Assumption OLS.2': $y = \mathbf{x}'\beta + u$ with $E[\mathbf{x}u] = \mathbf{0}$.

Note that Assumption OLS.1' implicitly assumes that $E[\|\mathbf{x}\|^2] < \infty$. Assumption OLS.1' is the large-sample counterpart of Assumption OLS.1, and Assumption OLS.2' is weaker than Assumption OLS.2.

Theorem 1 Under Assumptions OLS.0, OLS.1', OLS.2' and OLS.3, $\hat{\beta} \xrightarrow{p} \beta$.

Proof. From (1), to show $\hat{\beta} \xrightarrow{p} \beta$, we need only to show that $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \xrightarrow{p} \mathbf{0}$. Note that

$$\begin{aligned} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i u_i \right) \\ &= g \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i', \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i u_i \right) \xrightarrow{p} E[\mathbf{x}_i \mathbf{x}_i']^{-1} E[\mathbf{x}_i u_i] = \mathbf{0}. \end{aligned}$$

Here, the convergence in probability is from (I) the WLLN which implies

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{p} E[\mathbf{x}_i \mathbf{x}_i'] \text{ and } \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i u_i \xrightarrow{p} E[\mathbf{x}_i u_i]; \quad (2)$$

(II) the fact that $g(\mathbf{A}, \mathbf{b}) = \mathbf{A}^{-1}\mathbf{b}$ is a continuous function at $(E[\mathbf{x}_i \mathbf{x}_i'], E[\mathbf{x}_i u_i])$. The last equality is from Assumption OLS.2'.

(I) To apply the WLLN, we require (i) $\mathbf{x}_i \mathbf{x}_i'$ and $\mathbf{x}_i u_i$ are i.i.d., which is implied by Assumption OLS.0 and that functions of i.i.d. data are also i.i.d.; (ii) $E[\|\mathbf{x}\|^2] < \infty$ (OLS.1') and $E[\|\mathbf{x}u\|] < \infty$. $E[\|\mathbf{x}u\|] < \infty$ is implied by the Cauchy-Schwarz inequality,¹

$$E[\|\mathbf{x}u\|] \leq E[\|\mathbf{x}\|^2]^{1/2} E[|u|^2]^{1/2},$$

which is finite by Assumption OLS.1' and OLS.3. (II) To guarantee $\mathbf{A}^{-1}\mathbf{b}$ to be a continuous function at $(E[\mathbf{x}_i \mathbf{x}_i'], E[\mathbf{x}_i u_i])$, we must assume that $E[\mathbf{x}_i \mathbf{x}_i']^{-1}$ exists which is implied by Assumption OLS.1'.² ■

Exercise 1 Take the model $y_i = \mathbf{x}_{1i}'\beta_1 + \mathbf{x}_{2i}'\beta_2 + u_i$ with $E[\mathbf{x}_i u_i] = \mathbf{0}$. Suppose that β_1 is estimated by regressing y_i on \mathbf{x}_{1i} only. Find the probability limit of this estimator. In general, is it consistent for β_1 ? If not, under what conditions is this estimator consistent for β_1 ?

We can similarly show that the estimators $\hat{\sigma}^2$ and s^2 are consistent for σ^2 .

Theorem 2 Under the assumptions of Theorem 1, $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$ and $s^2 \xrightarrow{p} \sigma^2$.

Proof. Note that

$$\hat{u}_i = y_i - \mathbf{x}_i' \hat{\beta} = u_i + \mathbf{x}_i' \beta - \mathbf{x}_i' \hat{\beta} = u_i - \mathbf{x}_i' (\hat{\beta} - \beta).$$

Thus

$$\hat{u}_i^2 = u_i^2 - 2u_i \mathbf{x}_i' (\hat{\beta} - \beta) + (\hat{\beta} - \beta)' \mathbf{x}_i \mathbf{x}_i' (\hat{\beta} - \beta) \quad (3)$$

¹Cauchy-Schwarz inequality: For any random $m \times n$ matrix \mathbf{X} and $m \times k$ matrix \mathbf{Y} , $E[\|\mathbf{X}'\mathbf{Y}\|] \leq E[\|\mathbf{X}\|^2]^{1/2} E[\|\mathbf{Y}\|^2]^{1/2}$, where $E[\|\mathbf{X}'\mathbf{Y}\|]$ can be interpreted as the absolute value of inner product, $|\langle \mathbf{X}, \mathbf{Y} \rangle|$.

²If $x_i \in \mathbb{R}$, $E[x_i x_i']^{-1} = E[x_i^2]^{-1}$ is the reciprocal of $E[x_i^2]$ which is a continuous function of $E[x_i^2]$ only if $E[x_i^2] \neq 0$.

and

$$\begin{aligned}
\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \\
&= \frac{1}{n} \sum_{i=1}^n u_i^2 - 2 \left(\frac{1}{n} \sum_{i=1}^n u_i \mathbf{x}_i' \right) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \\
&\xrightarrow{p} \sigma^2,
\end{aligned}$$

where the last line uses the WLLN, (2), Theorem 1 and the CMT.

Finally, since $n/(n-k) \rightarrow 1$, it follows that

$$s^2 = \frac{n}{n-k} \hat{\sigma}^2 \xrightarrow{p} \sigma^2$$

by the CMT. ■

One implication of this theorem is that multiple estimators can be consistent for the population parameter. While $\hat{\sigma}^2$ and s^2 are unequal in any given application, they are close in value when n is very large.

1.2 Asymptotic Normality

To study the asymptotic normality of $\hat{\boldsymbol{\beta}}$, we impose the following additional assumption.

Assumption OLS.5: $E[u^4] < \infty$ and $E[\|\mathbf{x}\|^4] < \infty$.

The assumption $E[u^4] < \infty$ is stronger than Assumption 3 in Chapter 3 which states $E[u^2] < \infty$.

Theorem 3 Under Assumptions OLS.0, OLS.1', OLS.2', and OLS.5,

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

where $\mathbf{V} = \mathbf{Q}^{-1} \boldsymbol{\Omega} \mathbf{Q}^{-1}$ with $\mathbf{Q} = E[\mathbf{x}_i \mathbf{x}_i']$ and $\boldsymbol{\Omega} = E[\mathbf{x}_i \mathbf{x}_i' u_i^2]$.

Proof. From (1),

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i u_i \right).$$

Note first that

$$E[\|\mathbf{x}_i \mathbf{x}_i' u_i^2\|] \leq E[\|\mathbf{x}_i \mathbf{x}_i'\|^2]^{1/2} E[u_i^4]^{1/2} \leq E[\|\mathbf{x}_i\|^4]^{1/2} E[u_i^4]^{1/2} < \infty, \quad (4)$$

where the first inequality is from the Cauchy-Schwarz inequality, the second inequality is from the

Schwarz matrix inequality,³ and the last inequality is from Assumption OLS.5. So by the CLT,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i u_i \xrightarrow{d} N(\mathbf{0}, \mathbf{\Omega}).$$

Given that $n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{p} \mathbf{Q}$,

$$\sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathbf{Q}^{-1} N(\mathbf{0}, \mathbf{\Omega}) = N(\mathbf{0}, \mathbf{V})$$

by Slutsky's theorem. ■

From Assumption OLS.2', $\mathbf{x}'\boldsymbol{\beta}$ is the best linear projector of y on $\text{span}(\mathbf{x})$ rather than the conditional mean of y given \mathbf{x} , so if we express $y = m(\mathbf{x}) + \varepsilon = \mathbf{x}'\boldsymbol{\beta} + u$ with $m(\mathbf{x}) = E[y|\mathbf{x}]$, then $u = m(\mathbf{x}) - \mathbf{x}'\boldsymbol{\beta} + \varepsilon$. Even if $E[\varepsilon^2|\mathbf{x}] = \sigma^2$, $E[u^2|\mathbf{x}] = E[(m(\mathbf{x}) - \mathbf{x}'\boldsymbol{\beta} + \varepsilon)^2|\mathbf{x}] = (m(\mathbf{x}) - \mathbf{x}'\boldsymbol{\beta})^2 + \sigma^2$ depends on \mathbf{x} if $m(\mathbf{x}) \neq \mathbf{x}'\boldsymbol{\beta}$; in other words, the model is heteroskedastic. Another case where heteroskedasticity would naturally arise is the random coefficient model discussed in Chapter 6.

If the model is indeed homoskedastic, \mathbf{V} reduces to $\mathbf{V}^0 = \sigma^2 \mathbf{Q}^{-1}$. We call \mathbf{V}^0 the **homoskedastic covariance matrix**. Sometimes, to state the asymptotic distribution of part of $\hat{\boldsymbol{\beta}}$ as in the residual regression, we partition \mathbf{Q} and $\mathbf{\Omega}$ as

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{pmatrix}, \mathbf{\Omega} = \begin{pmatrix} \mathbf{\Omega}_{11} & \mathbf{\Omega}_{12} \\ \mathbf{\Omega}_{21} & \mathbf{\Omega}_{22} \end{pmatrix}. \quad (5)$$

Recall from the proof of the FWL theorem,

$$\mathbf{Q}^{-1} = \begin{pmatrix} \mathbf{Q}_{11.2}^{-1} & -\mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22.1}^{-1} \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22.1}^{-1} \end{pmatrix},$$

where $\mathbf{Q}_{11.2} = \mathbf{Q}_{11} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21}$ and $\mathbf{Q}_{22.1} = \mathbf{Q}_{22} - \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} \mathbf{Q}_{12}$. Thus when the error is homoskedastic, $n \cdot \text{AVar}(\hat{\boldsymbol{\beta}}_1) = \sigma^2 \mathbf{Q}_{11.2}^{-1}$, and $n \cdot \text{ACov}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\beta}}_2) = -\sigma^2 \mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1}$. We can also derive the general formulas in the heteroskedastic case, but these formulas are not easily interpretable and so less useful.

To understand the form of $\mathbf{Q}_{11.2}$, think about the projection of \mathbf{x}_1 on $\text{span}(\mathbf{x}_2)$ in the L^2 space. The coefficient of the projection is

$$\boldsymbol{\Gamma} = E[\mathbf{x}_2 \mathbf{x}_2']^{-1} E[\mathbf{x}_2 \mathbf{x}_1'],$$

and the error term is

$$\mathbf{v} = \mathbf{x}_1 - \boldsymbol{\Gamma}' \mathbf{x}_2 \equiv \mathbf{x}_{1 \perp 2},$$

³Schwarz matrix inequality: For any random $m \times n$ matrices \mathbf{X} and \mathbf{Y} , $\|\mathbf{X}'\mathbf{Y}\| \leq \|\mathbf{X}\| \|\mathbf{Y}\|$. This is a special form of the Cauchy-Schwarz inequality, where $\|\mathbf{X}'\mathbf{Y}\|$ can be interpreted as the absolute value of inner product, $|\langle \mathbf{X}, \mathbf{Y} \rangle|$.

so

$$\begin{aligned} E[\mathbf{v}\mathbf{v}'] &= E[\mathbf{x}_1\mathbf{x}_1'] - \mathbf{\Gamma}' E[\mathbf{x}_2\mathbf{x}_1'] - E[\mathbf{x}_1\mathbf{x}_2'] \mathbf{\Gamma} + \mathbf{\Gamma}' E[\mathbf{x}_2\mathbf{x}_2'] \mathbf{\Gamma} \\ &= E[\mathbf{x}_1\mathbf{x}_1'] - E[\mathbf{x}_1\mathbf{x}_2'] E[\mathbf{x}_2\mathbf{x}_2']^{-1} E[\mathbf{x}_2\mathbf{x}_1'] = \mathbf{Q}_{11.2}, \end{aligned}$$

where we can get $E[\mathbf{v}\mathbf{v}'] = E[\mathbf{x}_1\mathbf{x}_1'] - E[\mathbf{x}_1\mathbf{x}_2'] E[\mathbf{x}_2\mathbf{x}_2']^{-1} E[\mathbf{x}_2\mathbf{x}_1']$ directly by the intuition of projection. In other words, $\mathbf{Q}_{11.2}$ is the error "length" in the projection of \mathbf{x}_1 on $\text{span}(\mathbf{x}_2)$. This is reminiscent of the FWL theorem where the influence of \mathbf{x}_2 is excluded from \mathbf{x}_1 to get $\hat{\beta}_1$, so $\mathbf{Q} = E[\mathbf{x}\mathbf{x}']$ in the asymptotic variance of $\hat{\beta}$ should be replaced by $\mathbf{Q}_{11.2} = E[\mathbf{x}_{1\perp 2}\mathbf{x}_{1\perp 2}']$ in the asymptotic variance of $\hat{\beta}_1$.

Notice also that $\mathbf{Q}_{11} \geq \mathbf{Q}_{11.2}$, so in the following two homoskedastic linear regression models,

$$\begin{aligned} \mathcal{M}_1 &: y^1 = \mathbf{x}_1'\beta_1 + \mathbf{x}_2'\beta_2 + u, \\ \mathcal{M}_2 &: y^2 = \mathbf{x}_1'\beta_1 + u, \end{aligned}$$

where the same u appears in y^1 and y^2 ,⁴ and $E[u|\mathbf{x}] = 0$, $E[u^2|\mathbf{x}] = \sigma^2$, the asymptotic variance of LSE of β_1 in \mathcal{M}_1 is larger than that in \mathcal{M}_2 . Intuitively, β_1 can only be less precisely estimated in \mathcal{M}_1 due to the multicollinearity problem between \mathbf{x}_1 and \mathbf{x}_2 .

Exercise 2 *Of the variables $(y_i^*, y_i, \mathbf{x}_i)$ only the pair (y_i, \mathbf{x}_i) are observed. In this case, we say that y_i^* is a latent variable. Suppose*

$$\begin{aligned} y_i^* &= \mathbf{x}_i'\beta + u_i, \\ E[\mathbf{x}_i u_i] &= \mathbf{0}, \\ y_i &= y_i^* + v_i, \end{aligned}$$

where v_i is a measurement error satisfying $E[\mathbf{x}_i v_i] = \mathbf{0}$ and $E[y_i^* v_i] = 0$. Let $\hat{\beta}$ denote the OLS coefficient from the regression of y_i on \mathbf{x}_i .

- (i) Is β the coefficient from the linear projection of y_i on \mathbf{x}_i ?
- (ii) Is $\hat{\beta}$ consistent for β ?
- (iii) Find the asymptotic distribution of $\sqrt{n}(\hat{\beta} - \beta)$.

1.3 LSE as a MoM Estimator

The LSE is a MoM estimator. The corresponding moment conditions are the orthogonal conditions

$$E[\mathbf{x}u] = \mathbf{0},$$

⁴Note that y^1 and y^2 are different!

where $u = y - \mathbf{x}'\boldsymbol{\beta}$. So the sample analog is the normal equation

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \left(y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} \right) = E_n [\mathbf{x} \hat{u}] = \mathbf{0},$$

the solution of which is exactly the LSE. Now, $\mathbf{M} = -E [\mathbf{x}_i \mathbf{x}_i'] = -\mathbf{Q}$, and $\boldsymbol{\Omega} = E [\mathbf{x}_i \mathbf{x}_i' u_i^2]$, so

$$\sqrt{n} \left(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}),$$

the same as in Theorem 3. Note that the asymptotic variance \mathbf{V} takes the *sandwich* form. The larger the $E [\mathbf{x}_i \mathbf{x}_i']$, the smaller the \mathbf{V} . Although the LSE is a MoM estimator, it is a *special* MoM estimator because it can be treated as a "projection" estimator.

We provide more intuition on the asymptotic variance of $\hat{\boldsymbol{\beta}}$ below. Consider a simple linear regression model

$$y_i = \beta x_i + u_i,$$

where $E[x_i]$ is normalized to be 0. From introductory econometrics,

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{\widehat{Cov}(x, y)}{\widehat{Var}(x)},$$

and under homoskedasticity,

$$AVar(\hat{\beta}) = \frac{\sigma^2}{nVar(x)}.$$

So the larger the $Var(x)$, the smaller the $AVar(\hat{\beta})$. Actually, $Var(x) = \left| \frac{\partial E[xu]}{\partial \beta} \right|$, so the intuition in introductory courses matches the general results of the MoM estimator.

Similarly, we can derive the asymptotic distribution of the **weighted least squares** (WLS) **estimator**, a special GLS estimator with a diagonal weight matrix. Recall that

$$\hat{\boldsymbol{\beta}}_{\text{GLS}} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{y},$$

which reduces to

$$\hat{\boldsymbol{\beta}}_{\text{WLS}} = \left(\sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n w_i \mathbf{x}_i y_i \right)$$

when $\mathbf{W} = \text{diag}\{w_1, \dots, w_n\}$. Note that this estimator is a MoM estimator under the moment condition

$$E[w_i \mathbf{x}_i u_i] = \mathbf{0},$$

so

$$\sqrt{n} \left(\hat{\boldsymbol{\beta}}_{\text{WLS}} - \boldsymbol{\beta} \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\mathbf{W}}),$$

where $\mathbf{V}_W = E[w_i \mathbf{x}_i \mathbf{x}_i']^{-1} E[w_i^2 \mathbf{x}_i \mathbf{x}_i' u_i^2] E[w_i \mathbf{x}_i \mathbf{x}_i']^{-1}$.

Exercise 3 Suppose $w_i = \sigma_i^{-2}$, where $\sigma_i^2 = E[u_i^2 | \mathbf{x}_i]$. Derive the asymptotic distribution of $\sqrt{n}(\hat{\beta}_{WLS} - \beta)$.

2 Covariance Matrix Estimators

Since $\mathbf{Q} = E[\mathbf{x}_i \mathbf{x}_i']$ and $\mathbf{\Omega} = E[\mathbf{x}_i \mathbf{x}_i' u_i^2]$,

$$\hat{\mathbf{Q}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' = \frac{1}{n} \mathbf{X}' \mathbf{X},$$

and

$$\hat{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{u}_i^2 = \frac{1}{n} \mathbf{X}' \text{diag}\{\hat{u}_1^2, \dots, \hat{u}_n^2\} \mathbf{X} \equiv \frac{1}{n} \mathbf{X}' \hat{\mathbf{D}} \mathbf{X} \quad (6)$$

are the MoM estimators for \mathbf{Q} and $\mathbf{\Omega}$, where $\{\hat{u}_i\}_{i=1}^n$ are the OLS residuals. Given that $\mathbf{V} = \mathbf{Q}^{-1} \mathbf{\Omega} \mathbf{Q}^{-1}$, it can be estimated by

$$\hat{\mathbf{V}} = \hat{\mathbf{Q}}^{-1} \hat{\mathbf{\Omega}} \hat{\mathbf{Q}}^{-1},$$

and $AVar(\hat{\beta})$ is estimated by $\hat{\mathbf{V}}/n$. As in (5), we can partition $\hat{\mathbf{Q}}$ and $\hat{\mathbf{\Omega}}$ accordingly and the corresponding notations are just put a hat on. Recall from Exercise 8 of Chapter 3, $Var(\hat{\beta}_j | \mathbf{X}) = \sum_{i=1}^n w_{ij} \sigma_i^2 / SSR_j$. Since $\hat{\mathbf{V}}/n = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{D}} \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1}$ just replaces σ_i^2 in $Var(\hat{\beta} | \mathbf{X})$ by \hat{u}_i^2 , $\widehat{AVar}(\hat{\beta}_j) = \sum_{i=1}^n w_{ij} \hat{u}_i^2 / SSR_j$, $j = 1, \dots, k$, where $w_{ij} > 0$, $\sum_{i=1}^n w_{ij} = 1$, and SSR_j is the SSR in the regression of x_j on all other regressors.

Although this estimator is natural nowadays, it took some time to come into existence because $\mathbf{\Omega}$ is usually expressed as $E[\mathbf{x}_i \mathbf{x}_i' \sigma_i^2]$ whose estimation requires estimating n conditional variances. $\hat{\mathbf{V}}$ appeared first in the statistical literature Eicker (1967) and Huber (1967), and was introduced into econometrics by White (1980c). So this estimator is often called the "Eicker-Huber-White formula" or something of the kind. Other popular names for this estimator include the "heteroskedasticity-consistent (or robust) covariance matrix estimator" or the "sandwich-form covariance matrix estimator". The following theorem provides a direct proof of the consistency of $\hat{\mathbf{V}}$ (although it is a trivial corollary of the consistency of the MoM estimator).

Exercise 4 In the model

$$\begin{aligned} y_i &= \mathbf{x}_i' \beta + u_i, \\ E[\mathbf{x}_i u_i] &= 0, \mathbf{\Omega} = E[\mathbf{x}_i \mathbf{x}_i' u_i^2], \end{aligned}$$

find the MoM estimators of $(\beta, \mathbf{\Omega})$.

Theorem 4 Under the assumptions of Theorem 3, $\hat{\mathbf{V}} \xrightarrow{p} \mathbf{V}$.

Proof. From the WLLN, $\widehat{\mathbf{Q}}$ is consistent. As long as we can show $\widehat{\mathbf{\Omega}}$ is consistent, by the CMT $\widehat{\mathbf{V}}$ is consistent. Using (3)

$$\widehat{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \widehat{u}_i^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' u_i^2 - \frac{2}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i u_i + \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \left((\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \right)^2.$$

From (4), $E[\|\mathbf{x}_i \mathbf{x}_i' u_i^2\|] < \infty$, so by the WLLN, $n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' u_i^2 \xrightarrow{p} \mathbf{\Omega}$. We need only to prove the remaining two terms are $o_p(1)$.

Note that the second term satisfies

$$\begin{aligned} \left\| \frac{2}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i u_i \right\| &\leq \frac{2}{n} \sum_{i=1}^n \left\| \mathbf{x}_i \mathbf{x}_i' (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i u_i \right\| \\ &\leq \frac{2}{n} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{x}_i'\| \left| (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \right| |u_i| \leq \left(\frac{2}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^3 |u_i| \right) \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|, \end{aligned}$$

where the first inequality is from the triangle inequality,⁵ and the second and third inequalities are from the Schwarz matrix inequality. By Hölder's inequality,⁶

$$E[\|\mathbf{x}_i\|^3 |u_i|] \leq E[\|\mathbf{x}_i\|^4]^{3/4} E[|u_i|^4]^{1/4} < \infty,$$

so by the WLLN, $n^{-1} \sum_{i=1}^n \|\mathbf{x}_i\|^3 |u_i| \xrightarrow{p} E[\|\mathbf{x}_i\|^3 |u_i|] < \infty$. Given that $\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta} = o_p(1)$, the second term is $o_p(1)O_p(1) = o_p(1)$. The third term satisfies

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \left((\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \right)^2 \right\| \leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{x}_i'\| \left((\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \right)^2 \leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^4 \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|^2 = o_p(1),$$

where the steps follow from similar arguments as in the second term. ■

In the homoskedastic case, we can estimate \mathbf{V} by $\widehat{\mathbf{V}}^0 = \widehat{\sigma}^2 \widehat{\mathbf{Q}}^{-1}$, and correspondingly, $\widehat{AVar}(\widehat{\beta}_j) = n^{-1} \sum_{i=1}^n \widehat{u}_i^2 / SSR_j$. In other words, the weights in the general formula take a special form, $w_{ij} = n^{-1}$. It is hard to judge which formula, homoskedasticity-only or heteroskedasticity-robust, is larger (why?). Although either way is possible in theory, the heteroskedasticity-robust formula is usually larger than the homoskedasticity-only one in practice; Kauermann and Carroll (2001) show that the former is actually more variable than the latter in the normal regression model, which is the price paid for robustness.

2.1 Alternative Covariance Matrix Estimators (*)

MacKinnon and White (1985) find that $\widehat{\mathbf{V}}$ can have a substantial small sample downward bias and suggest three small-sample corrected versions of $\widehat{\mathbf{V}}$. The first is suggested in Hinkley (1977) and

⁵Triangle inequality: For any $m \times n$ matrices \mathbf{X} and \mathbf{Y} , $\|\mathbf{X} + \mathbf{Y}\| \leq \|\mathbf{X}\| + \|\mathbf{Y}\|$.

⁶Hölder's inequality: If $p > 1$ and $q > 1$ and $\frac{1}{p} + \frac{1}{q} = 1$, then for any random $m \times n$ matrices \mathbf{X} and \mathbf{Y} , $E[\|\mathbf{X}'\mathbf{Y}\|] \leq E[\|\mathbf{X}\|^p]^{1/p} E[\|\mathbf{Y}\|^q]^{1/q}$.

simply replaces $\hat{\Omega}$ by $\frac{n}{n-k}\hat{\Omega}$. This estimator simply adjusts the degrees of freedom (dof) of $\{\hat{u}_i\}_{i=1}^n$, termed as HC1. The second is based on the observation that $E[\hat{u}_i^2|\mathbf{X}] = (1 - h_i)\sigma^2$ and replaces \hat{u}_i^2 by $\hat{u}_i^2/(1 - h_i)$, where $h_i = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i$ is the leverage of \mathbf{x}_i . This estimator is termed as HC2, denoted as

$$\hat{\mathbf{V}}_{\text{HC2}} = \hat{\mathbf{Q}}^{-1} \left(\frac{1}{n} \sum_{i=1}^n \frac{\hat{u}_i^2}{1 - h_i} \mathbf{x}_i \mathbf{x}_i' \right) \hat{\mathbf{Q}}^{-1}.$$

The third is based on the jackknife principle which was introduced by Quenouille (1949, 1956) and Tukey (1958) (see also Miller (1964, 1974)). Recall in Section 6 of Chapter 3 the definition of $\hat{\beta}_{(-i)}$ as the least-squares estimator with the i 'th observation deleted. From equation (3.13) of Efron (1982), the jackknife estimator of the variance matrix for $\hat{\beta}$ is

$$\hat{\mathbf{V}}^* = (n-1) \sum_{i=1}^n \left(\hat{\beta}_{(-i)} - \bar{\beta} \right) \left(\hat{\beta}_{(-i)} - \bar{\beta} \right)', \quad (7)$$

where

$$\bar{\beta} = \frac{1}{n} \sum_{i=1}^n \hat{\beta}_{(-i)}.$$

Using formula (5) in Chapter 3, we can show that

$$\hat{\mathbf{V}}^* = \left(\frac{n-1}{n} \right) \hat{\mathbf{Q}}^{-1} \hat{\Omega}^* \hat{\mathbf{Q}}^{-1}, \quad (8)$$

where

$$\hat{\Omega}^* = \frac{1}{n} \sum_{i=1}^n (1 - h_i)^{-2} \mathbf{x}_i \mathbf{x}_i' \hat{u}_i^2 - \left(\frac{1}{n} \sum_{i=1}^n (1 - h_i)^{-1} \mathbf{x}_i \hat{u}_i \right) \left(\frac{1}{n} \sum_{i=1}^n (1 - h_i)^{-1} \mathbf{x}_i \hat{u}_i \right)'.$$

This estimator is termed as HC3. MacKinnon and White (1985) present numerical (simulation) evidence that HC3 works the best among the four available estimators. They also suggest that the scaling factor $(n-1)/n$ in (8) can be omitted.

Exercise 5 Show that the two expressions of $\hat{\mathbf{V}}^*$ in (7) and (8) are equal.

To understand why $n-1$ appears in (7), let's consider a simple situation where the regressor is a constant. In this case, $\hat{\beta} = \bar{y}$, and $V = \text{Var}(y)$. It turns out that the jackknife estimator of $\text{Var}(y)$ is exactly the usual unbiased estimator. To be specific, the leave-one-out estimator is

$$\bar{y}_{(-i)} = \frac{1}{n-1} \sum_{j \neq i} y_j = \frac{n}{n-1} \bar{y} - \frac{1}{n-1} y_i.$$

The sample mean of the leave-one-out estimators is

$$\frac{1}{n} \sum_{i=1}^n \bar{y}_{(-i)} = \frac{n}{n-1} \bar{y} - \frac{1}{n-1} \bar{y} = \bar{y}.$$

The difference is

$$\bar{y}_{(-i)} - \bar{y} = \frac{1}{n-1} (\bar{y} - y_i).$$

The jackknife estimate of $Var(y)$ is then

$$(n-1) \sum_{i=1}^n \left(\frac{1}{n-1} \right)^2 (\bar{y} - y_i)^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{y} - y_i)^2,$$

which is exactly the conventionally unbiased estimator of $Var(y)$. The downside of jackknife is that it requires n separate estimations, which in some cases can be computationally costly. However, it can be used when an explicit asymptotic variance formula is not available and can be used as a check on the reliability of an asymptotic formula.

Andrews (1991) suggests a similar estimator based on cross-validation, which is defined by replacing the OLS residual \hat{u}_i in (6) with the leave-one-out estimator $\hat{u}_{i,-i} = (1 - h_i)^{-1} \hat{u}_i$. With this substitution, Andrews' proposed estimator is

$$\hat{\mathbf{V}}^{**} = \hat{\mathbf{Q}}^{-1} \hat{\mathbf{\Omega}}^{**} \hat{\mathbf{Q}}^{-1},$$

where

$$\hat{\mathbf{\Omega}}^{**} = \frac{1}{n} \sum_{i=1}^n (1 - h_i)^{-2} \mathbf{x}_i \mathbf{x}_i' \hat{u}_i^2.$$

It is similar to the MacKinnon-White HC3 estimator, but omits the mean correction. Andrews (1991) argues that simulation evidence indicates that $\hat{\mathbf{V}}^{**}$ is an improvement on $\hat{\mathbf{V}}^*$.

The jackknife represents a linear approximation of the bootstrap proposed in Efron (1979). See Hall (1994) and Horowitz (1997, 2001, 2003, 2019) for an introduction of the bootstrap in econometrics and Efron and Gong (1983) and Efron and Tibshirani (1986) for an introduction in statistics. Other popular books on the bootstrap and resampling methods include Efron (1982), Hall (1992), Efron and Tibshirani (1993), Shao and Tu (1995), and Davison and Hinkley (1997).

Imbens and Kolesár (2016) use simulations to show that the confidence intervals based on a degrees-of-freedom correction suggested by Bell and McCaffrey (2002, BM hereafter) perform the best in coverage and length. BM approximate the distribution of

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\mathbf{V}}_{\text{HC2}}^{jj}}}$$

by a t -distribution with the dof K_{BM} , where $\hat{\mathbf{V}}_{\text{HC2}}^{jj}$ is the j th diagonal element of $\hat{\mathbf{V}}_{\text{HC2}}$. In effect, the BM standard error $\sqrt{\hat{\mathbf{V}}_{\text{BM}}^{jj}} = \sqrt{\hat{\mathbf{V}}_{\text{HC2}}^{jj}} \frac{t_{K_{\text{BM}}, 0.025}}{1.96}$, where $t_{K_{\text{BM}}, 0.025}$ is the upper 0.025 quantile of the $t_{K_{\text{BM}}}$ distribution. How to choose K_{BM} ? It is chosen so that under homoskedasticity and normality, the distribution of $K_{\text{BM}} \cdot \hat{\mathbf{V}}_{\text{HC2}}^{jj} / \mathbf{V}_{jj}$ has the first two moments in common with $\chi_{K_{\text{BM}}}^2$. Because $E[\hat{\mathbf{V}}_{\text{HC2}}^{jj} | \mathbf{X}] = \mathbf{V}_{jj}$, the means automatically match. It remains to match the variances.

It turns out that

$$\widehat{\mathbf{V}}_{\text{HC2}}^{jj} = \sum_{i=1}^n \lambda_i Z_i,$$

where $Z_i \sim \chi_1^2$, and all Z_i 's are independent. The weights λ_i are eigenvalues of the $n \times n$ matrix $\sigma^2 \cdot \mathbf{G}'\mathbf{G}$, with the i th column of the $n \times n$ matrix \mathbf{G} , is equal to

$$\mathbf{G}_i = \frac{1}{\sqrt{1-h_i}} (e_{n,i} - \mathbf{P}_i) \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} e_{k,j},$$

where \mathbf{P}_i is the i th column of \mathbf{P} , and $e_{l,g}$ is the g th unit vector of \mathbb{R}^l . Solving

$$K_{\text{BM}}^2 \frac{\text{Var}(\widehat{\mathbf{V}}_{\text{HC2}}^{jj})}{\mathbf{V}_{jj}^2} = 2K_{\text{BM}},$$

we have

$$K_{\text{BM}} = \frac{2\mathbf{V}_{jj}^2}{\text{Var}(\widehat{\mathbf{V}}_{\text{HC2}}^{jj})} = \frac{(\sum_{i=1}^n \lambda_i)^2}{\sum_{i=1}^n \lambda_i^2},$$

which depends only on \mathbf{X} but not on σ^2 , and may be different for different j 's. Note that if without homoskedasticity, then λ_i should be the eigenvalues of $\mathbf{G}'\mathbf{\Omega}\mathbf{G}$. These weights are not feasible because $\mathbf{\Omega}$ is unknown in general. The feasible version of the Satterthwaite (1946) dof suggestion replaces $\mathbf{\Omega}$ by $\widehat{\mathbf{\Omega}} = \text{diag}\{\widehat{u}_i^2 / (1 - h_i)\}$. However, because $\widehat{\mathbf{\Omega}}$ is a noisy estimator of the conditional variance, the resulting confidence intervals are often substantially conservative.

Chesher and Jewitt (1987) show that the EWH estimator can be biased substantially even with moderately large samples if the distribution of the regressors is skewed, so it is the combination of the sample size and the distribution of the regressors that determines the accuracy of the EWH estimator in finite samples. The degrees-of-freedom correction of BM captures the skewness of regressors in \mathbf{G} ; when the distribution of regressors is skewed (e.g., the regressor of interest follows a log-normal distribution), the effective dof will be smaller.

3 Restricted Least Squares Revisited (*)

In Chapter 2, we derived the RLS estimator. We now study the asymptotic properties of the RLS estimator. Since the RLS estimator is a special MD estimator, we concentrate on the MD estimator under the constraints $\mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$ in this section. From Exercise 26 and the discussion in Section 5.1 of Chapter 2, we can show

$$\widehat{\boldsymbol{\beta}}_{MD} = \widehat{\boldsymbol{\beta}} - \mathbf{W}_n^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}_n^{-1}\mathbf{R})^{-1}(\mathbf{R}'\widehat{\boldsymbol{\beta}} - \mathbf{c}).$$

To derive its asymptotic distribution, we impose the following regularity conditions.

Assumption RLS.1: $\mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$ where \mathbf{R} is $k \times q$ with $\text{rank}(\mathbf{R}) = q$.

Assumption RLS.2: $\mathbf{W}_n \xrightarrow{p} \mathbf{W} > 0$.

Theorem 5 Under the Assumptions of Theorem 1, RLS.1 and RLS.2, $\hat{\boldsymbol{\beta}}_{MD} \xrightarrow{p} \boldsymbol{\beta}$. Under the Assumptions of Theorem 3, RLS.1 and RLS.2,

$$\sqrt{n} \left(\hat{\boldsymbol{\beta}}_{MD} - \boldsymbol{\beta} \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\mathbf{W}}),$$

where

$$\begin{aligned} \mathbf{V}_{\mathbf{W}} &= \mathbf{V} - \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V} - \mathbf{V} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W}^{-1} \\ &\quad + \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W}^{-1}, \end{aligned}$$

and $\mathbf{V} = \mathbf{Q}^{-1} \boldsymbol{\Omega} \mathbf{Q}^{-1}$.

From this theorem, the RLS estimator is CAN and its asymptotic variance is $\mathbf{V}_{\mathbf{Q}}$. Unless the model is homoskedastic, it is hard to compare $\mathbf{V}_{\mathbf{Q}}$ and \mathbf{V} .

Exercise 6 Prove the above theorem.

The asymptotic distribution of $\hat{\boldsymbol{\beta}}_{MD}$ depends on \mathbf{W} . A natural question is which \mathbf{W} is the best in the sense of minimizing $\mathbf{V}_{\mathbf{W}}$. This turns out to be \mathbf{V}^{-1} as shown in the following theorem. Since \mathbf{V}^{-1} is unknown, we can replace \mathbf{V}^{-1} with a consistent estimate $\hat{\mathbf{V}}^{-1}$ and the asymptotic distribution (and efficiency) are unchanged. We call the MD estimator setting $\mathbf{W}_n = \hat{\mathbf{V}}^{-1}$ the **efficient minimum distance (EMD) estimator**, which takes the form

$$\hat{\boldsymbol{\beta}}_{EMD} = \hat{\boldsymbol{\beta}} - \hat{\mathbf{V}} \mathbf{R} \left(\mathbf{R}' \hat{\mathbf{V}} \mathbf{R} \right)^{-1} \left(\mathbf{R}' \hat{\boldsymbol{\beta}} - \mathbf{c} \right). \quad (9)$$

Theorem 6 Under the Assumptions of Theorem 3 and RLS.1,

$$\sqrt{n} \left(\hat{\boldsymbol{\beta}}_{EMD} - \boldsymbol{\beta} \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}^*),$$

where

$$\mathbf{V}^* = \mathbf{V} - \mathbf{V} \mathbf{R} (\mathbf{R}' \mathbf{V} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}.$$

Since

$$\mathbf{V}^* \leq \mathbf{V},$$

$\hat{\boldsymbol{\beta}}_{EMD}$ has lower asymptotic variance than the unrestricted estimator. Furthermore, for any \mathbf{W} ,

$$\mathbf{V}^* \leq \mathbf{V}_{\mathbf{W}},$$

so $\hat{\boldsymbol{\beta}}_{EMD}$ is asymptotically efficient in the class of MD estimators.

Exercise 7 (i) Show that $\mathbf{V}_{\mathbf{V}^{-1}} = \mathbf{V}^*$. (ii*) Show that $\mathbf{V}^* \leq \mathbf{V}_{\mathbf{W}}$ for any \mathbf{W} .

Exercise 8 Consider the exclusion restriction $\beta_2 = \mathbf{0}$ in the linear regression model $y_i = \mathbf{x}'_1 \beta_1 + \mathbf{x}'_2 \beta_2 + u_i$.

- (i) Derive the asymptotic distribution of $\hat{\beta}_1$, $\hat{\beta}_{1R}$ and $\hat{\beta}_{1,EMD}$ and show that $AVar(\hat{\beta}_{1,EMD}) \leq AVar(\hat{\beta}_1)$ and $AVar(\hat{\beta}_{1,EMD}) \leq AVar(\hat{\beta}_{1R})$.
- (ii) When the model is homoskedastic, show that $AVar(\hat{\beta}_{1,EMD}) = AVar(\hat{\beta}_{1R}) \leq AVar(\hat{\beta}_1)$. When will the equality hold? (Hint: $\mathbf{Q}_{12} = \mathbf{0}$)
- (iii) When the model is heteroskedastic, provide an example where $AVar(\hat{\beta}_{1R}) > AVar(\hat{\beta}_1)$.

This theorem shows that the MD estimator with the smallest asymptotic variance is $\hat{\beta}_{EMD}$. One implication is that the RLS estimator is generally inefficient. The interesting exception is the case of conditional homoskedasticity, in which the optimal weight matrix is $\mathbf{W} = \mathbf{V}_\beta^{0-1}$ and thus the RLS estimator is an efficient MD estimator. When the error is conditionally heteroskedastic, there are asymptotic efficiency gains by using minimum distance rather than least squares.

The fact that the RLS estimator is generally inefficient appears counter-intuitive and requires some reflection to understand. Standard intuition suggests to apply the same estimation method (least squares) to the unconstrained and constrained models, and this is the most common empirical practice. But the above theorem shows that this is not the efficient estimation method. Instead, the EMD estimator has a smaller asymptotic variance. Why? Consider the RLS estimator with the exclusion restrictions. In this case, the least squares estimation does not make use of the regressor \mathbf{x}_{2i} . It ignores the information $E[\mathbf{x}_{2i}u_i] = \mathbf{0}$. This information is relevant when the error is heteroskedastic and the excluded regressors are correlated with the included regressors.

Finally, note that all asymptotic variances can be consistently estimated by their sample analogs, e.g.,

$$\hat{\mathbf{V}}^* = \hat{\mathbf{V}} - \hat{\mathbf{V}}\mathbf{R} \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \mathbf{R}'\hat{\mathbf{V}},$$

where $\hat{\mathbf{V}}$ is a consistent estimator of \mathbf{V} .

3.1 Orthogonality of Efficient Estimators

One important property of an efficient estimator is the following orthogonality property popularized by Hausman (1978).⁷

Theorem 7 (Orthogonality of Efficient Estimators) Let $\hat{\beta}$ and $\tilde{\beta}$ be two CAN estimators of β , and $\tilde{\beta}$ is efficient. Then the limiting distributions of $\sqrt{n}\hat{\Delta}$ and $\sqrt{n}(\tilde{\beta} - \beta)$ have zero covariance, where $\hat{\Delta} \equiv \hat{\beta} - \tilde{\beta}$.

⁷See also Lehmann and Casella (1998, Theorem 1.7, p. 85) and Rao (1973, Theorem 51.2.(i), p. 317) for a similar result. Lehmann and Casella cite Barankin (1949), Stein (1950), Bahadur (1957), and Luenberger (1969) as early references.

Proof. Suppose $\hat{\Delta}$ and $\tilde{\beta}$ are not orthogonal. Since $\text{plim}(\hat{\Delta}) = \mathbf{0}$, consider a new estimator $\bar{\beta} = \tilde{\beta} + a\mathbf{A}\hat{\Delta}$, where a is a scalar and \mathbf{A} is an arbitrary matrix to be chosen. The new estimator is CAN with asymptotic variance

$$\text{Var}(\bar{\beta}) = \text{Var}(\tilde{\beta}) + a\mathbf{A}\text{Cov}(\tilde{\beta}, \hat{\Delta}) + a\text{Cov}(\tilde{\beta}, \hat{\Delta})'\mathbf{A}' + a^2\mathbf{A}\text{Var}(\hat{\Delta})\mathbf{A}'.$$

Since $\tilde{\beta}$ is efficient, the minimizer of $\text{Var}(\bar{\beta})$ with respect to a is achieved at $a = 0$. Taking the first-order derivative with respect to a yields

$$\mathbf{AC} + \mathbf{C}'\mathbf{A}' + 2a\mathbf{A}\text{Var}(\hat{\Delta})\mathbf{A}',$$

where $\mathbf{C} \equiv \text{Cov}(\tilde{\beta}, \hat{\Delta})$. Choosing $\mathbf{A} = -\mathbf{C}'$, we have

$$-2\mathbf{C}'\mathbf{C} + 2a\mathbf{C}'\text{Var}(\hat{\Delta})\mathbf{C}.$$

Therefore, at $a = 0$, this derivative is equal to $-2\mathbf{C}'\mathbf{C} \leq 0$. Unless $\mathbf{C} = \mathbf{0}$, we may have a better estimator than $\tilde{\beta}$ by choosing a a little bit deviating from 0. ■

Exercise 9 (Finite-Sample Orthogonality of Efficient Estimators) *In the homoskedastic linear regression model, check $\text{Cov}(\hat{\beta} - \hat{\beta}_R, \hat{\beta}_R | \mathbf{X}) = \mathbf{0}$ directly. (Hint: recall that $\hat{\beta}_R = \mathbf{P}_{\mathbf{S} \perp \mathbf{X}'\mathbf{X}}\hat{\beta} + (\mathbf{I} - \mathbf{P}_{\mathbf{S} \perp \mathbf{X}'\mathbf{X}})\mathbf{s}$.)*

A direct corollary of the above theorem is as follows.

Corollary 1 *Let $\hat{\beta}$ and $\tilde{\beta}$ be two CAN estimators of β , and $\tilde{\beta}$ is efficient. Then $A\text{Var}(\hat{\beta} - \tilde{\beta}) = A\text{Var}(\hat{\beta}) - A\text{Var}(\tilde{\beta})$.*

3.2 Misspecification

We usually have two methods to study the effects of misspecification of the restrictions. In Method I, we assume the truth is $\mathbf{R}'\beta = \mathbf{c}^*$ with $\mathbf{c}^* \neq \mathbf{c}$; in Method II, we assume the truth is $\mathbf{R}'\beta_n = \mathbf{c} + n^{-1/2}\delta$. The specification in Method II need some explanation. In this specification, the constraint is "close" to correct, as the difference $\mathbf{R}'\beta_n - \mathbf{c} = n^{-1/2}\delta$ is "small" in the sense that it decreases with sample size n . We call such a misspecification as **local misspecification**. The reason why the deviation is proportional to $n^{-1/2}$ is because this is the only choice under which the localizing parameter δ appears in the asymptotic distribution but does not dominate it. We will give more discussions on this choice of rate in studying the local power of tests.

First discuss Method I. From the expression of $\hat{\beta}_{MD}$, it is not hard to see that

$$\hat{\beta}_{MD} \xrightarrow{p} \beta_{MD}^* = \beta - \mathbf{W}^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1}(\mathbf{c}^* - \mathbf{c}).$$

The second term, $\mathbf{W}^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1}(\mathbf{c}^* - \mathbf{c})$, shows that imposing an incorrect constraint leads to inconsistency - an asymptotic bias. We call the limiting value β_{MD}^* the minimum-distance

projection coefficient or the pseudo-true value implied by the restriction. There are more to say. Define

$$\beta_n^* = \beta - \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} (\mathbf{c}^* - \mathbf{c}).$$

(Note that β_n^* is different from β_{MD}^* .) Then

$$\begin{aligned} \sqrt{n} (\hat{\beta}_{MD} - \beta_n^*) &= \sqrt{n} (\hat{\beta} - \beta) - \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \sqrt{n} (\mathbf{R}' \hat{\beta} - \mathbf{c}^*) \\ &= (\mathbf{I} - \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \mathbf{R}') \sqrt{n} (\hat{\beta} - \beta) \\ &\xrightarrow{d} (\mathbf{I} - \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}') N(\mathbf{0}, \mathbf{V}) = N(\mathbf{0}, \mathbf{V}_{\mathbf{W}}). \end{aligned} \quad (10)$$

In particular,

$$\sqrt{n} (\hat{\beta}_{EMD} - \beta_n^*) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}^*).$$

This means that even when the constraint $\mathbf{R}'\beta = \mathbf{c}$ is misspecified, the conventional covariance matrix estimator is still an appropriate measure of the sampling variance, though the distribution is centered at the pseudo-true values (or projections) β_n^* rather than β . The fact that the estimators are biased is an unavoidable consequence of misspecification.

In Method II, since the true model is $y_i = \mathbf{x}_i' \beta_n + u_i$, it is not hard to show that

$$\sqrt{n} (\hat{\beta} - \beta_n) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}) \quad (11)$$

which is the same as when β is fixed. A difference arises in the constrained estimator. Since $\mathbf{c} = \mathbf{R}'\beta_n - n^{-1/2}\delta$,

$$\mathbf{R}'\hat{\beta} - \mathbf{c} = \mathbf{R}'(\hat{\beta} - \beta_n) + n^{-1/2}\delta,$$

and

$$\begin{aligned} \hat{\beta}_{MD} &= \hat{\beta} - \mathbf{W}_n^{-1} \mathbf{R} [\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R}]^{-1} (\mathbf{R}' \hat{\beta} - \mathbf{c}) \\ &= \hat{\beta} - \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \mathbf{R}' (\hat{\beta} - \beta_n) + n^{-1/2} \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \delta. \end{aligned}$$

It follows that

$$\sqrt{n} (\hat{\beta}_{MD} - \beta_n) \xrightarrow{d} (\mathbf{I} - \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \mathbf{R}') \sqrt{n} (\hat{\beta} - \beta_n) + \mathbf{W}_n^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}_n^{-1} \mathbf{R})^{-1} \delta.$$

The first term is asymptotically normal by (11). The second term converges in probability to a constant. This is because the $n^{-1/2}$ local scaling is exactly balanced by the \sqrt{n} scaling of the estimator. No alternative rate would have produced this result. Consequently, we find that the asymptotic distribution equals

$$\sqrt{n} (\hat{\beta}_{MD} - \beta_n) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\mathbf{W}}) + \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \delta = N(\delta^*, \mathbf{V}_{\mathbf{W}}), \quad (12)$$

where

$$\boldsymbol{\delta}^* = \mathbf{W}^{-1} \mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1} \boldsymbol{\delta}.$$

The asymptotic distribution (12) is an approximation of the sampling distribution of the restricted estimator under misspecification. The distribution (12) contains an asymptotic bias component $\boldsymbol{\delta}^*$. The approximation is not fundamentally different from (10) - they both have the same asymptotic variances, and both reflect the bias due to misspecification. The difference is that (10) puts the bias on the left-side of the convergence arrow, while (12) has the bias on the right-side. There is no substantive difference between the two, but (12) is more convenient for some purposes, such as the analysis of the power of tests, as we will explore in the last section of this chapter.

3.3 Nonlinear and Inequality Restrictions

In some cases it is desirable to impose nonlinear constraints on the parameter vector $\boldsymbol{\beta}$. They can be written as

$$\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}, \tag{13}$$

where $\mathbf{r}: \mathbb{R}^k \rightarrow \mathbb{R}^q$. This includes the linear constraints as a special case. An example of (13) which cannot be written in a linear form is $\beta_1\beta_2 = 1$, which is (13) with $r(\boldsymbol{\beta}) = \beta_1\beta_2 - 1$.

The RLS and MD estimators of $\boldsymbol{\beta}$ subject to (13) solve the minimization problems

$$\begin{aligned} \hat{\boldsymbol{\beta}}_R &= \arg \min_{\mathbf{r}(\boldsymbol{\beta})=0} SSR(\boldsymbol{\beta}), \\ \hat{\boldsymbol{\beta}}_{MD} &= \arg \min_{\mathbf{r}(\boldsymbol{\beta})=0} J_n(\boldsymbol{\beta}). \end{aligned}$$

The solutions can be achieved by the Lagrangian method. Computationally, there is in general no explicit expression for the solutions so they must be found numerically. Algorithms to numerically solve such Lagrangian problems are known as *constrained optimization* methods, and are available in programming languages such as Matlab.

The asymptotic distributions of $\hat{\boldsymbol{\beta}}_R$ and $\hat{\boldsymbol{\beta}}_{MD}$ are the same as in the linear constraints case except that \mathbf{R} is replaced by $\partial \mathbf{r}(\boldsymbol{\beta})'/\partial \boldsymbol{\beta}$, but the proof is more delicate.

We sometimes impose inequality constraints on $\boldsymbol{\beta}$,

$$\mathbf{r}(\boldsymbol{\beta}) \geq \mathbf{0}.$$

The most common example is a non-negative constraint $\beta_1 \geq 0$. The RLS and MD estimators of $\boldsymbol{\beta}$ can be written as

$$\begin{aligned} \hat{\boldsymbol{\beta}}_R &= \arg \min_{\mathbf{r}(\boldsymbol{\beta}) \geq \mathbf{0}} SSR(\boldsymbol{\beta}), \\ \hat{\boldsymbol{\beta}}_{MD} &= \arg \min_{\mathbf{r}(\boldsymbol{\beta}) \geq \mathbf{0}} J_n(\boldsymbol{\beta}). \end{aligned}$$

Except in special cases the constrained estimators do not have simple algebraic solutions. An

important exception is when there is a single non-negativity constraint, e.g., $\beta_1 \geq 0$ with $q = 1$. In this case the constrained estimator can be found by a two-step approach. First compute the unconstrained estimator $\hat{\beta}$. If $\hat{\beta}_1 \geq 0$ then $\hat{\beta}_R = \hat{\beta}_{MD} = \hat{\beta}$. Second, if $\hat{\beta}_1 < 0$ then impose $\beta_1 = 0$ (eliminate the regressor x_1) and re-estimate. This yields the constrained least-squares estimator. While this method works when there is a single non-negativity constraint, it does not immediately generalize to other contexts. The computational problems with inequality constraints are examples of *quadratic programming* problems. Quick and easy computer algorithms are available in programming languages such as Matlab.

Exercise 10 (Ridge Regression) Suppose the nonlinear constraints are $\beta' \Lambda \beta \leq B$, where $\Lambda > 0$ and $B > 0$. Show that

$$\hat{\beta}_R = \left(\mathbf{X}'\mathbf{X} + \hat{\lambda}\Lambda \right)^{-1} \mathbf{X}'\mathbf{y},$$

where $\hat{\lambda}$ is the Lagrange multiplier for the constraint $\beta' \Lambda \beta \leq B$.

Inference on inequality-constrained estimators is unfortunately quite challenging. The conventional asymptotic theory gives rise to the following dichotomy. If the true parameter satisfies the strict inequality $\mathbf{r}(\beta) > \mathbf{0}$, then asymptotically the estimator is not subject to the constraint and the inequality-constrained estimator has an asymptotic distribution equal to the unconstrained one. However if the true parameter is on the boundary, e.g., $\mathbf{r}(\beta) = \mathbf{0}$, then the estimator has a truncated structure. This is easiest to see in the one-dimensional case. If we have an estimator $\hat{\beta}$ which satisfies $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} Z = N(0, V)$ and $\beta = 0$, then the constrained estimator $\hat{\beta}_R (= \hat{\beta}_{MD}) = \max\{\hat{\beta}, 0\}$ will have the asymptotic distribution $\sqrt{n}\hat{\beta}_R \xrightarrow{d} \max\{Z, 0\}$, a “half-normal” distribution.

4 Functions of Parameters

Sometimes we are interested in some lower-dimensional function of the parameter vector $\beta = (\beta_1, \dots, \beta_k)'$. For example, we may be interested in a single coefficient β_j or a ratio β_j/β_l . In these cases we can write the parameter of interest as a function of β . Let $\mathbf{r}: \mathbb{R}^k \rightarrow \mathbb{R}^q$ denote this function and let

$$\theta = \mathbf{r}(\beta)$$

denote the parameter of interest. A natural estimate of θ is $\hat{\theta} = \mathbf{r}(\hat{\beta})$. To derive the asymptotic distribution of $\hat{\theta}$, we impose the following assumption.

Assumption RLS.1': $\mathbf{r}(\cdot)$ is continuously differentiable at the true value β and $\mathbf{R} = \frac{\partial}{\partial \beta} \mathbf{r}(\beta)'$ has rank q .

This assumption is an extension of Assumption RLS1.

Theorem 8 Under the assumptions of Theorem 3 and Assumption RLS.1',

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\theta).$$

where $\mathbf{V}_\theta = \mathbf{R}'\mathbf{V}\mathbf{R}$.

Proof. By the CMT, $\hat{\theta}$ is consistent for θ . By the Delta method, if $\mathbf{r}(\cdot)$ is differentiable at the true value β ,

$$\sqrt{n}(\hat{\theta} - \theta) = \sqrt{n}(\mathbf{r}(\hat{\beta}) - \mathbf{r}(\beta)) \xrightarrow{d} \mathbf{R}'N(\mathbf{0}, \mathbf{V}) = N(\mathbf{0}, \mathbf{V}_\theta).$$

■

A natural estimator of \mathbf{V}_θ is

$$\hat{\mathbf{V}}_\theta = \hat{\mathbf{R}}'\hat{\mathbf{V}}\hat{\mathbf{R}}, \quad (14)$$

where $\hat{\mathbf{R}} = \partial \mathbf{r}(\hat{\beta})' / \partial \beta$. If $\mathbf{r}(\cdot)$ is a $C^{(1)}$ function, then by the CMT, $\hat{\mathbf{V}}_\theta \xrightarrow{p} \mathbf{V}_\theta$ (why?).

In many cases, the function $\mathbf{r}(\beta)$ is linear:

$$\mathbf{r}(\beta) = \mathbf{R}'\beta$$

for some $k \times q$ matrix \mathbf{R} . In this case, $\frac{\partial}{\partial \beta} \mathbf{r}(\beta)' = \mathbf{R}$ and $\hat{\mathbf{R}} = \mathbf{R}$, so $\hat{\mathbf{V}}_\theta = \mathbf{R}'\hat{\mathbf{V}}\mathbf{R}$. For example, if \mathbf{R} is a "selector matrix"

$$\mathbf{R} = \begin{pmatrix} \mathbf{I}_{q \times q} \\ \mathbf{0}_{(k-q) \times q} \end{pmatrix},$$

so that if $\beta = (\beta'_1, \beta'_2)'$, then $\theta = \mathbf{R}'\beta = \beta_1$ and

$$\hat{\mathbf{V}}_\theta = (\mathbf{I}, \mathbf{0}) \hat{\mathbf{V}} \begin{pmatrix} \mathbf{I} \\ \mathbf{0} \end{pmatrix} = \hat{\mathbf{V}}_{11},$$

the upper-left block of $\hat{\mathbf{V}}$. When $q = 1$ (so $r(\beta)$ is real-valued), the standard error for $\hat{\theta}$ is the square root of $n^{-1}\hat{V}_\theta$, that is, $s(\hat{\theta}) = n^{-1/2}\sqrt{\hat{\mathbf{R}}'\hat{\mathbf{V}}\hat{\mathbf{R}}}$.

5 The t Test

Let $\theta = r(\beta): \mathbb{R}^k \rightarrow \mathbb{R}$ be any parameter of interest (for example, θ could be a single element of β), $\hat{\theta}$ its estimate and $s(\hat{\theta})$ its asymptotic standard error. Consider the studentized statistic

$$t_n(\theta) = \frac{\hat{\theta} - \theta}{s(\hat{\theta})}.$$

Since $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, V_\theta)$ and $\sqrt{n}s(\hat{\theta}) \xrightarrow{p} \sqrt{V_\theta}$, by Slutsky's theorem, we have

Theorem 9 *Under the assumptions of Theorem 8, if $V_\theta = \mathbf{R}'\mathbf{V}\mathbf{R} > 0$, then $t_n(\theta) \xrightarrow{d} N(0, 1)$.*

Given that \mathbf{R} has full rank 1, $\mathbf{Q} > 0$ and $\mathbf{V} = \mathbf{Q}^{-1}\mathbf{\Omega}\mathbf{Q}^{-1}$, a sufficient condition for $V_\theta > 0$ is $\mathbf{\Omega} > 0$. Under this condition, the asymptotic distribution of the t -ratio $t_n(\theta)$ is the standard

normal. Since the standard normal distribution does not depend on the parameters, we say that $t_n(\theta)$ is **asymptotically pivotal**. In special cases (such as the normal regression model), the statistic t_n has an exact t distribution, and is therefore exactly free of unknowns. In this case, we say that t_n is an **exactly pivotal** statistic. In general, however, pivotal statistics are unavailable and so we must rely on asymptotically pivotal statistics.

The most common one-dimensional hypotheses are the null

$$H_0 : \theta = \theta_0, \quad (15)$$

against the alternative

$$H_1 : \theta \neq \theta_0, \quad (16)$$

where θ_0 is some pre-specified value. The standard test for H_0 against H_1 is based on the absolute value of the t -statistic,

$$t_n = t_n(\theta_0) = \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})}.$$

Under H_0 , $t_n \xrightarrow{d} Z \sim N(0, 1)$, so $|t_n| \xrightarrow{d} |Z|$ by the CMT. $G(u) = P(|Z| \leq u) = \Phi(u) - (1 - \Phi(u)) = 2\Phi(u) - 1 \equiv \bar{\Phi}(u)$ is called the **asymptotic null distribution**.

The **asymptotic size** of the test is defined as the asymptotic probability of a Type I error:

$$\lim_{n \rightarrow \infty} P(|t_n| > c | H_0 \text{ true}) = P(|Z| > c) = 1 - \bar{\Phi}(c).$$

We see that the asymptotic size of the test is a simple function of the asymptotic null distribution G and the critical value c . As mentioned in Chapter 3, in the dominant approach to hypothesis testing, the researcher pre-selects a significance level $\alpha \in (0, 1)$ and then selects c so that the (asymptotic) size is no larger than α . We call c the **asymptotic critical value** because it has been selected from the asymptotic null distribution. Let $z_{\alpha/2}$ be the upper $\alpha/2$ quantile of the standard normal distribution. That is, if $Z \sim N(0, 1)$, then $P(Z > z_{\alpha/2}) = \alpha/2$ and $P(|Z| > z_{\alpha/2}) = \alpha$. For example, $z_{0.025} = 1.96$ and $z_{0.05} = 1.645$. A test of asymptotic significance α rejects H_0 if $|t_n| > z_{\alpha/2}$. Otherwise the test does not reject, or "accepts" H_0 .

The alternative hypothesis (16) is sometimes called a "two-sided" alternative. Sometimes we are interested in testing for one-sided alternatives such as

$$H_1 : \theta > \theta_0 \quad (17)$$

or

$$H_1 : \theta < \theta_0. \quad (18)$$

Tests of (15) against (17) or (18) are based on the signed t -statistic t_n . The hypothesis (15) is rejected in favor of (17) if $t_n > c$ where c satisfies $\alpha = 1 - \Phi(c)$. Negative values of t_n are not taken as evidence against H_0 , as point estimates $\hat{\theta}$ less than θ_0 do not point to (17). Since the critical

values are taken from the single tail of the normal distribution, they are smaller than for two-sided tests. Specifically, the asymptotic 5% critical value is $c = 1.645$. Thus, we reject (15) in favor of (17) if $t_n > 1.645$. Testing against (18) can be conducted in a similar way.

There seems to be an ambiguity. Should we use the two-sided critical value 1.96 or the one-sided critical value 1.645? The answer is that we should use one-sided tests and critical values only when the parameter space is known to satisfy a one-sided restriction such as $\theta \geq \theta_0$. This is when the test of (15) against (17) makes sense. If the restriction $\theta \geq \theta_0$ is not known a priori, then imposing this restriction to test (15) against (17) does not make sense. Since linear regression coefficients typically do not have a priori sign restrictions, we conclude that two-sided tests are generally appropriate.

Exercise 11 *Prove that if an additional regressor \mathbf{X}_{k+1} is added to \mathbf{X} , Theil's adjusted \bar{R}^2 increases if and only if $|t_{k+1}| > 1$, where $t_{k+1} = \hat{\beta}_{k+1}/s\left(\hat{\beta}_{k+1}\right)$ is the t -ratio for $\hat{\beta}_{k+1}$ and*

$$s\left(\hat{\beta}_{k+1}\right) = \left(s^2 \left[(\mathbf{X}'\mathbf{X})^{-1}\right]_{k+1,k+1}\right)^{1/2}$$

is the homoskedasticity-formula standard error. (Hint: Use the FWL theorem)

6 p -Value

An alternative approach, associated with R.A. Fisher, is to report an asymptotic p -value. The asymptotic p -value for $|t_n|$ is constructed as follows. Define the tail probability, or asymptotic p -value function

$$p(t) = P(|Z| > t) = 1 - G(t) = 2(1 - \Phi(t)),$$

where $G(\cdot)$ is the cdf of $|Z|$. Then the asymptotic p -value of the statistic $|t_n|$ is

$$p_n = p(|t_n|).$$

So the p -value is the probability of obtaining a test statistic result at least as extreme as the one that was actually observed or the smallest significance level at which the null would be rejected, assuming that the null is true. Since the distribution function G is monotonically increasing, the p -value is a monotonically decreasing function of t_n and is an equivalent test statistic. Figure 2 shows how to find p_n when $|t_n| = 1.85$ (the left panel) and p_n as a function of $|t_n|$ (the right panel). An important caveat is that the p -value p_n should not be interpreted as the probability that either hypothesis is true. For example, a common mis-interpretation is that p_n is the probability “that the null hypothesis is true.” This is incorrect. Rather, p_n is a measure of the strength of information against the null hypothesis.

A researcher will often “reject the null hypothesis” when the p -value turns out to be less than a predetermined significance level, often 0.05 or 0.01. Such a result indicates that the observed

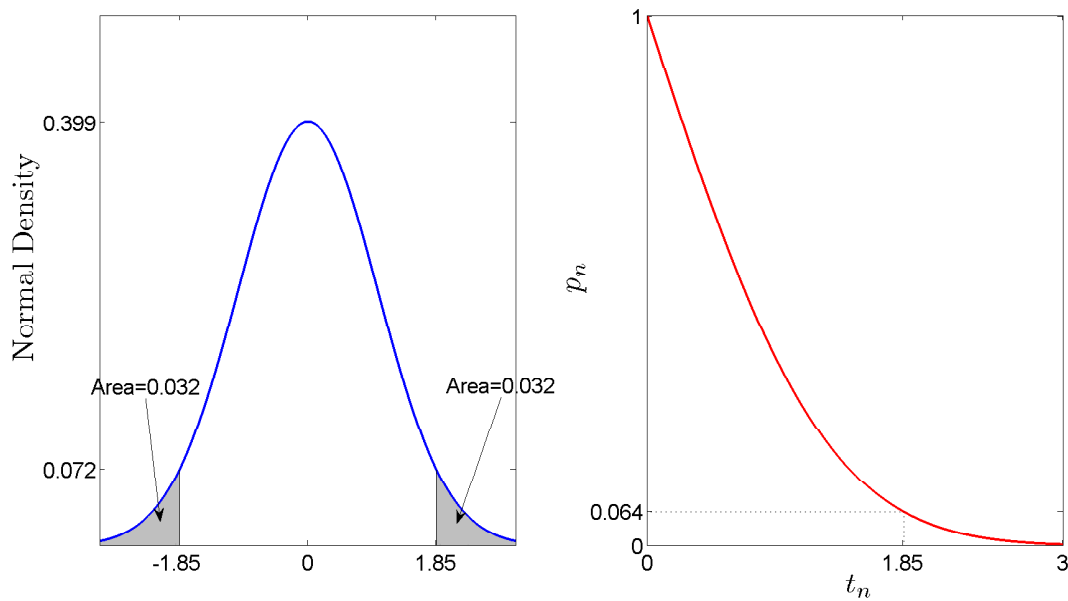


Figure 2: Obtaining the p -Value in a Two-Sided t -Test

result would be highly unlikely under the null hypothesis. In a sense, p -values and hypothesis tests are equivalent since $p_n < \alpha$ if and only if $|t_n| > z_{\alpha/2}$. Thus an equivalent statement of a Neyman-Pearson test is to reject at the α level if and only if $p_n < \alpha$. The p -value is more general, however, in that the reader is allowed to pick the level of significance α , in contrast to Neyman-Pearson rejection/acceptance reporting where the researcher picks the level.

Another helpful observation is that the p -value function has simply made a unit-free transformation of the test statistic. That is, under H_0 , $p_n \xrightarrow{d} U[0, 1]$, so the "unusualness" of the test statistic can be compared to the easy-to-understand uniform distribution, regardless of the complication of the distribution of the original test statistic. To see this fact, note that the asymptotic distribution of $|t_n|$ is $G(x) = 1 - p(x)$. Thus

$$\begin{aligned} P(1 - p_n \leq u) &= P(1 - p(|t_n|) \leq u) = P(G(|t_n|) \leq u) \\ &= P(|t_n| \leq G^{-1}(u)) \rightarrow G(G^{-1}(u)) = u, \end{aligned}$$

establishing that $1 - p_n \xrightarrow{d} U[0, 1]$, from which it follows that $p_n \xrightarrow{d} U[0, 1]$.

7 Confidence Interval

A confidence interval (CI) C_n is an interval estimate of $\theta \in \mathbb{R}$ which is assumed to be *fixed*. It is a function of the data and hence is random. So it is not correct to say that " θ will *fall* in C_n with high probability", rather, C_n is designed to *cover* θ with high probability. Either $\theta \in C_n$ or $\theta \notin C_n$. The coverage probability is $P(\theta \in C_n)$.

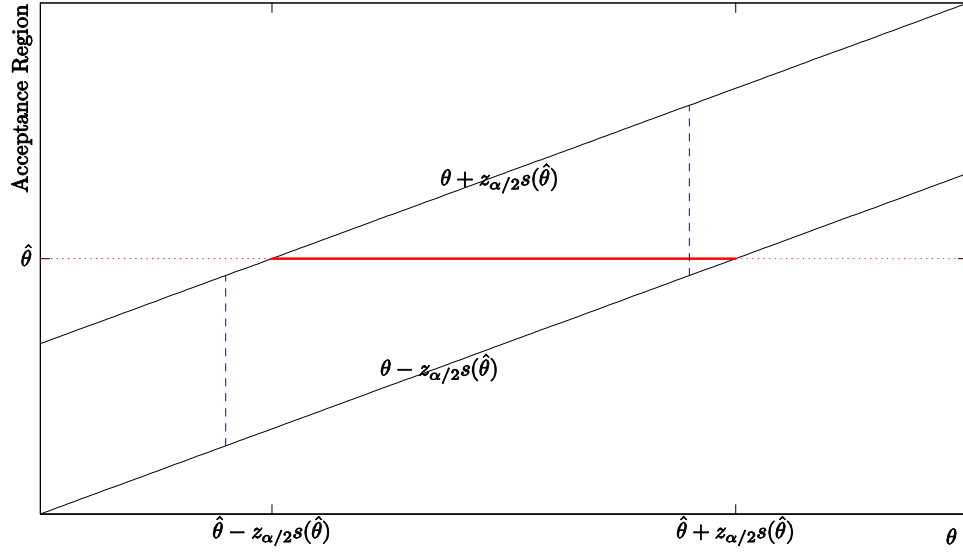


Figure 3: Test Statistic Inversion

We typically cannot calculate the exact coverage probability $P(\theta \in C_n)$. However we often can calculate the asymptotic coverage probability $\lim_{n \rightarrow \infty} P(\theta \in C_n)$. We say that C_n has asymptotic $(1 - \alpha)$ coverage for θ if $P(\theta \in C_n) \rightarrow 1 - \alpha$ as $n \rightarrow \infty$.

A good method for constructing a confidence interval is collecting parameter values which are not rejected by a statistical test, so-called "*test statistic inversion*" method. The t -test in Section 5 rejects $H_0: \theta = \theta_0$ if $|t_n(\theta_0)| > z_{\alpha/2}$. A confidence interval is then constructed using the values for which this test does not reject:

$$C_n = \{\theta \mid |t_n(\theta)| \leq z_{\alpha/2}\} = \left\{ \theta \mid -z_{\alpha/2} \leq \frac{\hat{\theta} - \theta}{s(\hat{\theta})} \leq z_{\alpha/2} \right\} = \left[\hat{\theta} - z_{\alpha/2}s(\hat{\theta}), \hat{\theta} + z_{\alpha/2}s(\hat{\theta}) \right].$$

Figure 3 illustrates the idea of inverting a test statistic. In Figure 3, the acceptance region for $\hat{\theta}$ at θ is $\left[\theta - z_{\alpha/2}s(\hat{\theta}), \theta + z_{\alpha/2}s(\hat{\theta}) \right]$, which is the region for $\hat{\theta}$ such that the hypothesis that the true value is θ cannot be rejected or is "accepted".

While there is no hard-and-fast guideline for choosing the coverage probability $1 - \alpha$, the most common professional choice is 95%, or $\alpha = .05$. This corresponds to selecting the confidence interval $\left[\hat{\theta} \pm 1.96s(\hat{\theta}) \right] \approx \left[\hat{\theta} \pm 2s(\hat{\theta}) \right]$. Thus values of θ within two standard errors of the estimated $\hat{\theta}$ are considered "reasonable" candidates for the true value θ , and values of θ outside two standard errors of the estimated $\hat{\theta}$ are considered unlikely or unreasonable candidates for the true value.

Finally, the interval has been constructed so that as $n \rightarrow \infty$,

$$P(\theta \in C_n) = P(|t_n(\theta)| \leq z_{\alpha/2}) \rightarrow P(|Z| \leq z_{\alpha/2}) = 1 - \alpha,$$

so C_n is indeed an asymptotic $(1 - \alpha)$ confidence interval. This justifies the duality between the confidence interval and inverting tests of each possible parameter value.

(*) Coverage accuracy is a basic requirement for a CI. Another property of a CI is its length (if it is fixed) or expected length (if it is random). Since at most one point of the interval is the true value, the expected length of the interval is a measure of the "average extent" of the false values included. In Appendix A, we discuss this issue following Pratt (1961).

8 Edgeworth Expansion (*)

Theorem 9 showed that the t -ratio $t_n(\theta)$ is asymptotically normal. In practice this means that we use the normal distribution to approximate the finite sample distribution of $t_n(\theta)$. How good is this approximation? Some insight into the accuracy of the normal approximation can be obtained by an Edgeworth expansion which is a higher-order approximation to the distribution of $t_n(\theta)$. The following result is an application of Theorem 13 of Chapter 4.

Theorem 10 *Under the assumptions of Theorem 9, if $E[u^{16}] < \infty$, $E[\|\mathbf{x}\|^{16}] < \infty$, $r(u)$ have five continuous derivatives in a neighborhood of β , and $\lim_{t \rightarrow \infty} \left| E \left[\exp \left(it \left(u^4 + \|\mathbf{x}\|^4 \right) \right) \right] \right| \leq B < 1$, as $n \rightarrow \infty$,*

$$P(t_n(\theta) \leq x) = \Phi(x) + n^{-1/2}p_1(x)\phi(x) + n^{-1}p_2(x)\phi(x) + o(n^{-1})$$

uniformly in x , where $p_1(x)$ is an even polynomial of order 2, and $p_2(x)$ is an odd polynomial of degree 5, with coefficients depending on the moments of $g(X)$ up to order 16.

Theorem 10 shows that the finite sample distribution of the t -ratio can be approximated up to $o(n^{-1})$ by the sum of three terms, the first being the standard normal distribution, the second a $O(n^{-1/2})$ adjustment, and the third a $O(n^{-1})$ adjustment.

First consider a one-sided confidence interval $C_n = [\hat{\theta} - z_\alpha s(\hat{\theta}), \infty)$. It has coverage

$$\begin{aligned} P(\theta \in C_n) &= P(t_n(\theta) \leq z_\alpha) \\ &= \Phi(z_\alpha) + n^{-1/2}p_1(z_\alpha)\phi(z_\alpha) + O(n^{-1}) \\ &= 1 - \alpha + O(n^{-1/2}). \end{aligned}$$

This means that the actual coverage is within $O(n^{-1/2})$ of the desired $1 - \alpha$ level. Next consider a two-sided interval $C_n = [\hat{\theta} - z_{\alpha/2} s(\hat{\theta}), \hat{\theta} + z_{\alpha/2} s(\hat{\theta})]$. It has coverage

$$\begin{aligned} P(\theta \in C_n) &= P(|t_n(\theta)| \leq z_{\alpha/2}) \\ &= 2\Phi(z_{\alpha/2}) + n^{-1}2p_2(z_{\alpha/2})\phi(z_{\alpha/2}) + o(n^{-1}) \\ &= 1 - \alpha + O(n^{-1}). \end{aligned}$$

This means that the actual coverage is within $O(n^{-1})$ of the desired $1 - \alpha$ level. The accuracy is better than the one-sided interval because the $O(n^{-1/2})$ term in the Edgeworth expansion has

offsetting effects in the two tails of the distribution.

9 The Wald Test

Sometimes $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta})$ is a $q \times 1$ vector, and it is desired to test the joint restrictions simultaneously. In this case the t -statistic approach does not work. We have the null and alternative

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0 \text{ vs } H_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0.^8$$

The natural estimate of $\boldsymbol{\theta}$ is $\hat{\boldsymbol{\theta}} = \mathbf{r}(\hat{\boldsymbol{\beta}})$. Suppose $\hat{\mathbf{V}}_{\boldsymbol{\theta}}$ is an estimate of the asymptotic covariance matrix of $\hat{\boldsymbol{\theta}}$, e.g., $\hat{\mathbf{V}}_{\boldsymbol{\theta}}$ in (14); then the Wald statistic⁹ for H_0 against H_1 is

$$W_n = n \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right)' \hat{\mathbf{V}}_{\boldsymbol{\theta}}^{-1} \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right).$$

We have known that $\sqrt{n} \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\boldsymbol{\theta}})$, and $\hat{\mathbf{V}}_{\boldsymbol{\theta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\theta}}$ under the null. So by Example 2 of Chapter 4, $W_n \xrightarrow{d} \chi_q^2$ under the null. We have established:

Theorem 11 *Under the assumptions of Theorem 8, $W_n \xrightarrow{d} \chi_q^2$ under H_0 .*

When \mathbf{r} is a linear function of $\boldsymbol{\beta}$, i.e., $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{R}'\boldsymbol{\beta}$, the Wald statistic takes the form

$$W_n = n \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)' \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right).$$

When $q = 1$, $W_n = t_n^2$. Correspondingly, the asymptotic distribution $\chi_1^2 = N(0, 1)^2$.

An asymptotic Wald test rejects H_0 in favor of H_1 if W_n exceeds $\chi_{q,\alpha}^2$, the upper- α quantile of the χ_q^2 distribution. For example, $\chi_{1,.05}^2 = 3.84 = z_{.025}^2$. The Wald test fails to reject if W_n is less than $\chi_{q,\alpha}^2$. The asymptotic p -value for W_n is $p_n = p(W_n)$, where $p(x) = P(\chi_q^2 \geq x)$ is the tail probability function of the χ_q^2 distribution. As before, the test rejects at the α level if $p_n < \alpha$, and p_n is asymptotically $U[0, 1]$ under H_0 .

(*) The Wald test can only determine whether $\boldsymbol{\theta}$ is equal to $\boldsymbol{\theta}_0$ completely or some elements of $\boldsymbol{\theta}$ are not equal to the counterparts of $\boldsymbol{\theta}_0$. In the latter case, one may want to know which element of $\boldsymbol{\theta}$ is not equal to that of $\boldsymbol{\theta}_0$. This is related to the **multiple testing** or **multiple comparisons** problem. See Chapter 9 of Lehmann and Romano (2005), Chapter 19 of Gourieroux and Monfort (1995) and the Handbook of Econometrics Chapter by Savin (1984) for an introduction.

⁸Different from t -tests, Wald tests are hard to apply for one-sided alternatives although not impossible.

⁹Abraham Wald (1902-1950) was a statistician at Columbia University. He is most famous for his work on decision theory and statistical sequential analysis. He and his wife died in an airplane crash while on an extensive lecture tour at the invitation of the Indian government.

10 Confidence Region

Similarly, we can construct confidence regions for multiple parameters, e.g., $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta}) \in \mathbb{R}^q$. By the test statistic inversion method, an asymptotic $(1 - \alpha)$ confidence region for $\boldsymbol{\theta}$ is

$$\mathbf{C}_n = \{\boldsymbol{\theta} | W_n(\boldsymbol{\theta}) \leq \chi_{q,\alpha}^2\},$$

where $W_n(\boldsymbol{\theta}) = n(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' \hat{\mathbf{V}}_{\boldsymbol{\theta}}^{-1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$. Since $\hat{\mathbf{V}}_{\boldsymbol{\theta}} > 0$, \mathbf{C}_n is an ellipsoid in the $\boldsymbol{\theta}$ plane.

To show \mathbf{C}_n intuitively, assume $q = 2$ and $\boldsymbol{\theta} = (\beta_1, \beta_2)'$. In this case, \mathbf{C}_n is an ellipse in the (β_1, β_2) plane as shown in Figure 4. In Figure 4, $\hat{\beta}_1$ and $\hat{\beta}_2$ are positively correlated. For comparison, we also show the $(1 - \alpha)$ CIs for β_1 and β_2 in Figure 4. It is tempting to use the rectangular region, say \mathbf{C}'_n , as a confidence region for (β_1, β_2) . However, $P((\beta_1, \beta_2) \in \mathbf{C}'_n)$ may not converge to $1 - \alpha$. For example, suppose $\hat{\beta}_1$ and $\hat{\beta}_2$ are asymptotically independent; then $P((\beta_1, \beta_2) \in \mathbf{C}'_n) \rightarrow (1 - \alpha)^2 < (1 - \alpha)$. Note also that the ellipse \mathbf{C}_n cannot be completely contained in \mathbf{C}'_n . This is because if we reverse the roles of $\boldsymbol{\beta}$ and $\hat{\boldsymbol{\beta}}$ in the figure, then $P(\hat{\boldsymbol{\beta}} \in \mathbf{C}_n) = 1 - \alpha$, while $P(\hat{\beta}_1 \in AB, \hat{\beta}_2 \in \mathbb{R}) = 1 - \alpha$.

Exercise 12 (*) (i) Show that $1 - 2\alpha \leq \lim_{n \rightarrow \infty} P((\beta_1, \beta_2) \in \mathbf{C}'_n) \leq 1 - \alpha$. (ii) When will $\lim_{n \rightarrow \infty} P((\beta_1, \beta_2) \in \mathbf{C}'_n) = 1 - \alpha$? (iii) Show that when $\hat{\beta}_1$ and $\hat{\beta}_2$ are asymptotically independent, $\lim_{n \rightarrow \infty} P((\beta_1, \beta_2) \in \mathbf{C}'_n) = (1 - \alpha)^2$. (iv) Is $\lim_{n \rightarrow \infty} P((\beta_1, \beta_2) \in \mathbf{C}'_n) = 1 - 2\alpha$ achievable? What is the exact lower bound of $\lim_{n \rightarrow \infty} P((\beta_1, \beta_2) \in \mathbf{C}'_n)$? (v) Show that $CD/AB = \sqrt{\chi_2^2(\alpha)/\chi_1^2(\alpha)}$ in Figure 4. (vi) Is it possible that $\mathbf{C}'_n \subseteq \mathbf{C}_n$?

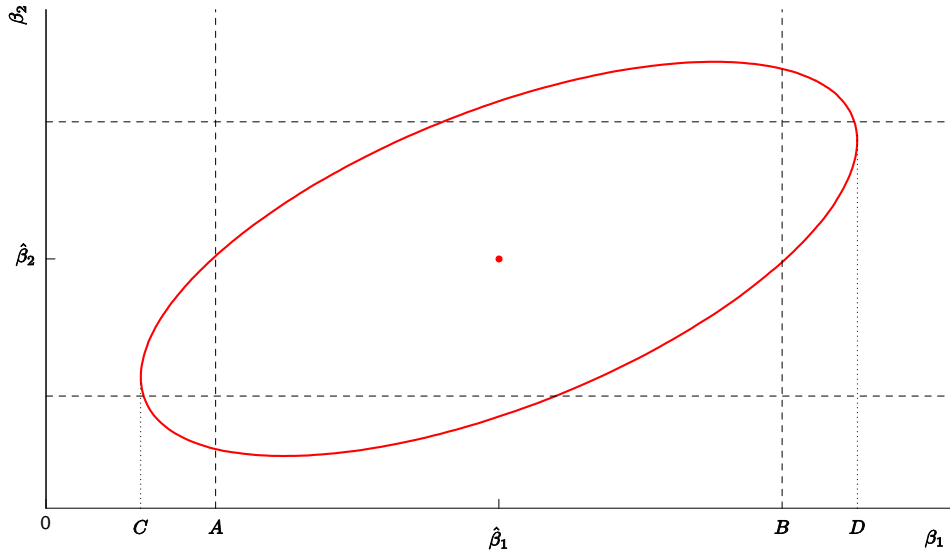


Figure 4: Confidence Region for (β_1, β_2)

11 Problems with Tests of Nonlinear Hypotheses

While the t test and Wald tests work well when the hypothesis is a linear restriction on β , they can work quite poorly when the restrictions are nonlinear. This can be seen in a simple example introduced by Lafontaine and White (1986). Take the model

$$y_i = \beta + u_i, u_i \sim N(0, \sigma^2)$$

and consider the hypothesis

$$H_0: \beta = 1.$$

Let $\hat{\beta}$ and $\hat{\sigma}^2$ be the sample mean and variance of y_i . The standard Wald test for H_0 is

$$W_n = n \frac{(\hat{\beta} - 1)^2}{\hat{\sigma}^2}.$$

Now notice that H_0 is equivalent to the hypothesis

$$H_0(s): \beta^s = 1,$$

for any positive integer s . Letting $r(\beta) = \beta^s$, and noting $R = s\beta^{s-1}$, we find that the standard Wald test for $H_0(s)$ is

$$W_n(s) = n \frac{(\hat{\beta}^s - 1)^2}{\hat{\sigma}^2 s^2 \hat{\beta}^{2s-2}}.$$

While the hypothesis $\beta^s = 1$ is unaffected by the choice of s , the statistic $W_n(s)$ varies with s . This is an unfortunate feature of the Wald statistic.

To demonstrate this effect, we plot in Figure 5 the Wald statistic $W_n(s)$ as a function of s , setting $n/\hat{\sigma}^2 = 10$. The increasing solid line is for the case $\hat{\beta} = 0.8$ and the decreasing dashed line is for the case $\hat{\beta} = 1.6$. It is easy to see that in each case there are values of s for which the test statistic is significant relative to asymptotic critical values, while there are other values of s for which the test statistic is insignificant.¹⁰ This is distressing since the choice of s is arbitrary and irrelevant to the actual hypothesis.

Our first-order asymptotic theory is not useful to help pick s , as $W_n(s) \xrightarrow{d} \chi_1^2$ under H_0 for any s . This is a context where **Monte Carlo simulation** can be quite useful as a tool to study and compare the exact distributions of statistical procedures in finite samples. The method uses random simulation to create artificial datasets, to which we apply the statistical tools of interest. This produces random draws from the statistic's sampling distribution. Through repetition, features of this distribution can be calculated. In the present context of the Wald statistic, one feature of importance is the Type I error of the test using the asymptotic 5% critical value 3.84 - the

¹⁰Breusch and Schmidt (1988) show that any positive value for the Wald test statistic is possible by rewriting H_0 in an algebraically equivalent form.

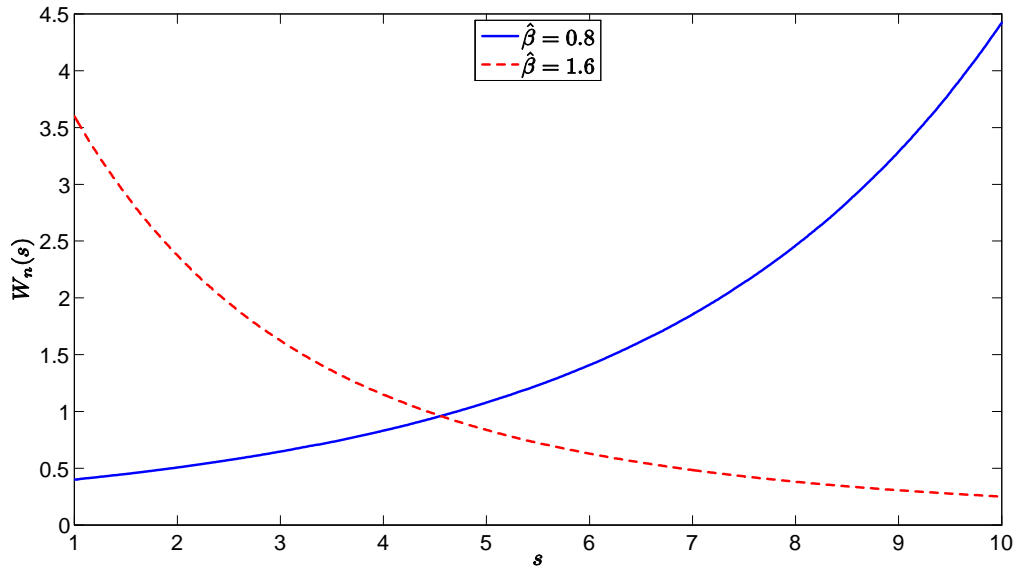


Figure 5: Wald Statistic as a function of s

probability of a false rejection, $P(W_n(s) > 3.84 | \beta = 1)$. Given the simplicity of the model, this probability depends only on s , n , and σ^2 . In Table 1 we report the results of a Monte Carlo simulation where we vary these three parameters: the value of s is varied from 1 to 10, n is varied among 20, 100 and 500, and σ is varied among 1 and 3. The table reports the simulation estimate of the Type I error probability from 50,000 random samples. Each row of the table corresponds to a different value of s - and thus corresponds to a particular choice of test statistic. The second through seventh columns contain the Type I error probabilities for different combinations of n and σ . These probabilities are calculated as the percentage of the 50,000 simulated Wald statistics $W_n(s)$ which are larger than 3.84. The null hypothesis $\beta^s = 1$ is true, so these probabilities are Type I error.

s	$\sigma = 1$			$\sigma = 3$		
	$n = 20$	$n = 100$	$n = 500$	$n = 20$	$n = 100$	$n = 500$
1	.06	.05	.05	.07	.05	.05
2	.08	.06	.05	.15	.08	.06
3	.10	.06	.05	.21	.12	.07
4	.13	.07	.06	.25	.15	.08
5	.15	.08	.06	.28	.18	.10
6	.17	.09	.06	.30	.20	.11
7	.19	.10	.06	.31	.22	.13
8	.20	.12	.07	.33	.24	.14
9	.22	.13	.07	.34	.25	.15
10	.23	.14	.08	.35	.26	.16

Table 1: Type I Error Probability of Asymptotic 5% $W_n(s)$ Test

Note: Rejection frequencies from 50,000 simulated random samples

To interpret the table, remember that the ideal Type I error probability is 5%(.05) with deviations indicating distortion. Type I error rates between 3% and 8% are considered reasonable. Error rates above 10% are considered excessive. Rates above 20% are unacceptable. When comparing statistical procedures, we compare the rates row by row, looking for tests for which rejection rates are close to 5% and rarely fall outside of the 3% – 8% range. For this particular example the only test which meets this criterion is the conventional $W_n = W_n(1)$ test. Any other choice of s leads to a test with unacceptable Type I error probabilities.

In Table 1 you can also see the impact of variation in sample size. In each case, the Type I error probability improves towards 5% as the sample size n increases. There is, however, no magic choice of n for which all tests perform uniformly well. Test performance deteriorates as s increases, which is not surprising given the dependence of $W_n(s)$ on s as shown in Figure 5.

In this example it is not surprising that the choice $s = 1$ yields the best test statistic. Other choices are arbitrary and would not be used in practice. While this is clear in this particular example, in other examples natural choices are not always obvious and the best choices may in fact appear counter-intuitive at first. This point can be illustrated through another example which is similar to one developed in Gregory and Veall (1985). Take the model

$$y_i = \beta_0 + x_{1i}\beta_1 + x_{2i}\beta_2 + u_i, E[\mathbf{x}_i u_i] = \mathbf{0} \quad (19)$$

and the hypothesis

$$H_0: \frac{\beta_1}{\beta_2} = \theta_0$$

where θ_0 is a known constant. Equivalently, define $\theta = \beta_1/\beta_2$, so the hypothesis can be stated as $H_0: \theta = \theta_0$. Let $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)'$ be the least-squares estimates of (19), $\hat{\mathbf{V}}_{\hat{\beta}}$ be an estimate of the

covariance matrix for $\hat{\beta}$ and $\hat{\theta} = \hat{\beta}_1/\hat{\beta}_2$.¹¹ Define

$$\hat{\mathbf{R}}_1 = \left(0, \frac{1}{\hat{\beta}_2}, -\frac{\hat{\beta}_1}{\hat{\beta}_2^2} \right)'$$

so that the standard error for $\hat{\theta}$ is $s(\hat{\theta}) = \left(\hat{\mathbf{R}}_1' \hat{\mathbf{V}} \hat{\mathbf{R}}_1 \right)^{1/2}$. In this case, a t -statistic for H_0 is

$$t_{1n} = \frac{\hat{\beta}_1/\hat{\beta}_2 - \theta_0}{s(\hat{\theta})}.$$

An alternative statistic can be constructed through reformulating the null hypothesis as

$$H_0: \beta_1 - \theta_0 \beta_2 = 0.$$

A t -statistic based on this formulation of the hypothesis is

$$t_{2n} = \frac{\hat{\beta}_1 - \theta_0 \hat{\beta}_2}{\left(\mathbf{R}_2' \hat{\mathbf{V}}_{\hat{\beta}} \mathbf{R}_2 \right)^{1/2}},$$

where $\mathbf{R}_2 = (0, 1, -\theta_0)'$.

	$n = 100$				$n = 500$			
	$P(t_n < -1.645)$		$P(t_n > 1.645)$		$P(t_n < -1.645)$		$P(t_n > 1.645)$	
β_2	t_{1n}	t_{2n}	t_{1n}	t_{2n}	t_{1n}	t_{2n}	t_{1n}	t_{2n}
.10	.47	.06	.00	.06	.28	.05	.00	.05
.25	.26	.06	.00	.06	.15	.05	.00	.05
.50	.15	.06	.00	.06	.10	.05	.00	.05
.75	.12	.06	.00	.06	.09	.05	.00	.05
1.00	.10	.06	.00	.06	.07	.05	.02	.05

Table 2: Type I Error Probability of Asymptotic 5% t -tests

To compare t_{1n} and t_{2n} we perform another simple Monte Carlo simulation. We let x_{1i} and x_{2i} be mutually independent $N(0, 1)$ variables, u_i be an independent $N(0, \sigma^2)$ draw with $\sigma = 3$, and normalize $\beta_0 = 1$ and $\beta_1 = 1$. This leaves β_2 as a free parameter, along with sample size n . We vary β_2 among .1, .25, .50, .75, and 1.0 and n among 100 and 500. The one-sided Type I error probabilities $P(t_n < -1.645)$ and $P(t_n > 1.645)$ are calculated from 50,000 simulated samples. The results are presented in Table 2. Ideally, the entries in the table should be 0.05. However, the rejection rates for the t_{1n} statistic diverge greatly from this value, especially for small values of β_2 . The left tail probabilities $P(t_{1n} < -1.645)$ greatly exceed 5%, while the right tail probabilities $P(t_{1n} > 1.645)$ are close to zero in most cases. In contrast, the rejection rates for the linear t_{2n}

¹¹If $\hat{\mathbf{V}}$ is used to estimate the asymptotic variance of $\sqrt{n}(\hat{\beta} - \beta)$, then $\hat{\mathbf{V}}_{\hat{\beta}} = \hat{\mathbf{V}}/n$.

statistic are invariant to the value of β_2 , and are close to the ideal 5% rate for both sample sizes. The implication of Table 2 is that the two t -ratios have dramatically different sampling behaviors.

The common message from both examples is that Wald statistics are sensitive to the algebraic formulation of the null hypothesis. In all cases, if the hypothesis can be expressed as a linear restriction on the model parameters, this formulation should be used. If no linear formulation is feasible, then the "most linear" formulation should be selected (as suggested by the theory of Phillips and Park (1988)), and alternatives to asymptotic critical values should be considered. It is also prudent to consider alternative tests to the Wald statistic, such as the minimum distance statistic developed in Section 13.

12 Monte Carlo Simulation

In Section 11 we introduced the method of Monte Carlo simulation to illustrate the small sample problems with tests of nonlinear hypotheses. In this section we describe the method in more detail.

Recall, our data consist of observations (y_i, \mathbf{x}_i) which are random draws from a population distribution F . Let $\boldsymbol{\theta}$ be a parameter and let $T_n = T_n((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n), \boldsymbol{\theta})$ be a statistic of interest, for example an estimator $\hat{\boldsymbol{\theta}}$ or a t -statistic $(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) / s(\hat{\boldsymbol{\theta}})$. The exact distribution of T_n is

$$G_n(u, F) = P(T_n \leq u | F).$$

While the asymptotic distribution of T_n might be known, the exact (finite sample, i.e., n is fixed) distribution G_n is generally unknown.

Monte Carlo simulation uses numerical simulation to compute $G_n(u, F)$ for selected choices of F . This is useful to investigate the performance of the statistic T_n in reasonable situations and sample sizes. The basic idea is that for any given F , the distribution function $G_n(u, F)$ can be calculated numerically through simulation. The name Monte Carlo derives from the famous Mediterranean gambling resort where games of chance are played.

The method of Monte Carlo is quite simple to describe. The researcher chooses F (the distribution of the data) and the sample size n . A "true" value of $\boldsymbol{\theta}$ is implied by this choice, or equivalently the value $\boldsymbol{\theta}$ is selected directly by the researcher which implies restrictions on F . Then the following experiment is conducted by computer simulation:

1. n independent random pairs (y_i^*, \mathbf{x}_i^*) , $i = 1, \dots, n$, are drawn from the distribution F using the computer's random number generator.
2. The statistic $T_n = T_n((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), \boldsymbol{\theta})$ is calculated on this pseudo data.

For step 1, most computer packages have built-in procedures for generating $U[0, 1]$ and $N(0, 1)$ random numbers, and from these most random variables can be constructed. (For example, a chi-square can be generated by sums of squares of normals.) For step 2, it is important that the statistic be evaluated at the "true" value of $\boldsymbol{\theta}$ corresponding to the choice of F .

The above experiment creates one random draw from the distribution $G_n(u, F)$. This is one observation from an unknown distribution. Clearly, from one observation very little can be said. So the researcher repeats the experiment B times, where B is a large number. Typically, we set $B = 1000$ or $B = 5000$. We will discuss this choice later. Notationally, let the b th experiment result in the draw T_{nb} , $b = 1, \dots, B$. These results are stored. After all B experiments have been calculated, these results constitute a random sample of size B from the distribution of $G_n(u, F) = P(T_{nb} \leq u) = P(T_n \leq u|F)$.

From a random sample, we can estimate any feature of interest using (typically) a method of moments estimator. We now describe some specific examples.

Suppose we are interested in the bias, mean-squared error (MSE), and/or variance of the distribution $\hat{\theta} - \theta$. We then set $T_n = \hat{\theta} - \theta$, run the above experiment, and calculate

$$\begin{aligned}\widehat{Bias}(\hat{\theta}) &= \frac{1}{B} \sum_{b=1}^B T_{nb} = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b - \theta, \\ \widehat{MSE}(\hat{\theta}) &= \frac{1}{B} \sum_{b=1}^B T_{nb}^2 = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}_b - \theta)^2, \\ \widehat{Var}(\hat{\theta}) &= \widehat{MSE}(\hat{\theta}) - \widehat{Bias}(\hat{\theta})^2.\end{aligned}$$

Suppose we are interested in the Type I error associated with an asymptotic 5% two-sided t -test. We would then set $T_n = |\hat{\theta} - \theta|/s(\hat{\theta})$ and calculate

$$\hat{p} = \frac{1}{B} \sum_{b=1}^B 1(T_{nb} \geq 1.96), \quad (20)$$

the percentage of the simulated t -ratios which exceed the asymptotic 5% critical value.

Suppose we are interested in the 5% and 95% quantile of $T_n = \hat{\theta}$ or $T_n = (\hat{\theta} - \theta)/s(\hat{\theta})$, we then compute the 5% and 95% sample quantiles of the sample $\{T_{nb}\}$. The α th sample quantile is a number q_α such that 100 α percent of the sample are less than q_α . A simple way to compute sample quantiles is to sort the sample $\{T_{nb}\}$ from low to high. Then q_α is the N 'th number in this ordered sequence, where $N = (B + 1)\alpha$. It is therefore convenient to pick B so that N is an integer. For example, if we set $B = 999$, then the 5% sample quantile is 50'th sorted value and the 95% sample quantile is the 950'th sorted value.

The typical purpose of a Monte Carlo simulation is to investigate the performance of a statistical procedure (estimator or test) in realistic settings. Generally, the performance will depend on n and F . In many cases, an estimator or test may perform wonderfully for some values, and poorly for others. It is therefore useful to conduct a variety of experiments, for a selection of choices of n and F .

As discussed above, the researcher must select the number of experiments, B . Often this is called the number of **replications**. Quite simply, a larger B results in more precise estimates of

the features of interest of G_n , but requires more computational time. In practice, therefore, the choice of B is often guided by the computational demands of the statistical procedure. Since the results of a Monte Carlo experiment are estimates computed from a random sample of size B , it is straightforward to calculate standard errors for any quantity of interest. If the standard error is too large to make a reliable inference, then B will have to be increased.

In particular, it is simple to make inferences about rejection probabilities from statistical tests, such as the percentage estimate reported in (20). The random variable $1(T_{nb} \geq 1.96)$ is iid Bernoulli, equalling 1 with probability $p = E[1(T_{nb} \geq 1.96)]$. The average (20) is therefore an unbiased estimator of p with standard deviation $s(\hat{p}) = \sqrt{p(1-p)/B}$. As p is unknown, this may be approximated by replacing p with \hat{p} or with an hypothesized value. For example, if we are assessing an asymptotic 5% test, then we can set $s(\hat{p}) = \sqrt{(.05)(.95)/B}$. Hence, standard errors for $B = 100, 1000$, and 5000 , are, respectively, $s(\hat{p}) = .022, .007$, and $.003$. See also Davidson and MacKinnon (1981a) for a more efficient estimator of p and its standard deviation.

Most papers in econometric methods, and some empirical papers, include the results of Monte Carlo simulations to illustrate the performance of their methods. When extending existing results, it is good practice to start by replicating existing (published) results. This is not exactly possible in the case of simulation results, as they are inherently random. For example suppose a paper investigates a statistical test, and reports a simulated rejection probability of 0.07 based on a simulation with $B = 100$ replications. Suppose you attempt to replicate this result, and find a rejection probability of 0.03 (again using $B = 100$ simulation replications). Should you conclude that you have failed in your attempt? Absolutely not! Under the hypothesis that both simulations are identical, you have two independent estimates, $\hat{p}_1 = 0.07$ and $\hat{p}_2 = 0.03$, of a common probability p . The asymptotic (as $B \rightarrow \infty$) distribution of their difference is $\sqrt{B}(\hat{p}_1 - \hat{p}_2) \rightarrow N(0, 2p(1-p))$, so a standard error for $\hat{p}_1 - \hat{p}_2 = 0.04$ is $\hat{s} = \sqrt{2p(1-p)/B} \approx 0.03$ using the estimate $p = (\hat{p}_1 + \hat{p}_2)/2$. Since the t -ratio $0.04/0.03 = 1.3$ is not statistically significant, it is incorrect to reject the null hypothesis that the two simulations are identical. The difference between the results $\hat{p}_1 = 0.07$ and $\hat{p}_1 = 0.03$ is consistent with random variation.

What should be done? The first mistake was to copy the previous paper's choice of $B = 100$. Instead, suppose you set $B = 5000$. Suppose you now obtain $\hat{p}_2 = 0.04$. Then $\hat{p}_1 - \hat{p}_2 = 0.03$ and a standard error is $\hat{s} = \sqrt{p(1-p)(1/100 + 1/5000)} \approx 0.02$. Still we cannot reject the hypothesis that the two simulations are different. Even though the estimates (0.07 and 0.04) appear to be quite different, the difficulty is that the original simulation used a very small number of replications ($B = 100$) so the reported estimate is quite imprecise. In this case, it is appropriate to conclude that your results "replicate" the previous study, as there is no statistical evidence to reject the hypothesis that they are equivalent.

Most journals have policies requiring authors to make available their data sets and computer programs required for empirical results. They do not have similar policies regarding simulations. Nevertheless, it is good professional practice to make your simulations available. The best practice is to post your simulation code on your webpage. This invites others to build on and use your results, leading to possible collaboration, citation, and/or advancement.

13 Minimum Distance Test (*)

The likelihood ratio (LR) test is valid only under homoskedasticity. The counterpart of the LR test in the heteroskedastic environment is the minimum distance test. Based on the idea of the LR test, the minimum distance statistic is defined as

$$J_n = \min_{\mathbf{r}(\boldsymbol{\beta})=\mathbf{0}} J_n(\boldsymbol{\beta}) - \min_{\boldsymbol{\beta}} J_n(\boldsymbol{\beta}) = \min_{\mathbf{r}(\boldsymbol{\beta})=\mathbf{0}} J_n(\hat{\boldsymbol{\beta}}_{MD}),$$

where $\min_{\boldsymbol{\beta}} J_n(\boldsymbol{\beta}) = 0$ if no restrictions are imposed (why?). $J_n \geq 0$ measures the cost (on $J_n(\boldsymbol{\beta})$) of imposing the null restriction $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$. Usually, \mathbf{W}_n in $J_n(\boldsymbol{\beta})$ is chosen to be the efficient weight matrix $\hat{\mathbf{V}}^{-1}$, and the corresponding J_n is denoted as J_n^* with

$$J_n^* = n \left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{EMD} \right)' \hat{\mathbf{V}}^{-1} \left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{EMD} \right).$$

Consider the class of linear hypotheses $H_0: \mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$. In this case, we know from (9) that

$$\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{MD} = \hat{\mathbf{V}}\mathbf{R} \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} \right),$$

so

$$\begin{aligned} J_n^* &= n \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} \right)' \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \mathbf{R}'\hat{\mathbf{V}}\hat{\mathbf{V}}^{-1}\hat{\mathbf{V}}\mathbf{R} \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} \right) \\ &= n \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} \right)' \left(\mathbf{R}'\hat{\mathbf{V}}\mathbf{R} \right)^{-1} \left(\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} \right) = W_n. \end{aligned}$$

Thus for linear hypotheses, the efficient minimum distance statistic J_n^* is identical to the Wald statistic W_n which is heteroskedastic-robust.

Exercise 13 Show that $J_n^* = W$ under homoskedasticity when the null hypothesis is $H_0: \mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$, where W is the homoskedastic form of the Wald statistic defined in Chapter 4.

For nonlinear hypotheses, however, the Wald and minimum distance statistics are different. We know from Section 11 that the Wald statistic is not robust to the formulation of the null hypothesis. However, like the LR test statistic, the minimum distance statistic is invariant to the algebraic formulation of the null hypothesis, so is immune to this problem. Consequently, a simple solution to the problem associated with W_n in Section 11 is to use the minimum distance statistic J_n , which equals W_n with $s = 1$ in the first example, and $|t_{2n}|$ in the second example there. Whenever possible, the Wald statistic should not be used to test nonlinear hypotheses.

Exercise 14 Show that $J_n^* = W_n(1)$ in the first example and $J_n^* = t_{2n}^2$ in the second example of Section 11.

Newey and West (1987a) established the asymptotic null distribution of J_n^* for linear and non-linear hypotheses.

Theorem 12 Under the assumptions of Theorem 8, $J_n^* \xrightarrow{d} \chi_q^2$ under H_0 .

14 The Heteroskedasticity-Robust LM Test (*)

The validity of the LM test in the normal regression model depends on the assumption that the error is homoskedastic. In this section, we extend the homoskedasticity-only LM test to the heteroskedasticity-robust form. Suppose the null hypothesis is $H_0: \beta_2 = 0$, where β is decomposed as $(\beta_1', \beta_2')'$, β_1 and β_2 are $k_1 \times 1$ and $k_2 \times 1$ vectors, respectively, $k = k_1 + k_2$, and \mathbf{x} is decomposed as $(\mathbf{x}_1', \mathbf{x}_2')'$ correspondingly with \mathbf{x}_1 including the constant.

Exercise 15 Show that testing any linear constraints $\mathbf{R}'\beta = \mathbf{c}$ is equivalent to testing some original coefficients being zero in a new regression with redefined \mathbf{X} and \mathbf{y} , where $\mathbf{R} \in \mathbb{R}^{k_2 \times k}$ is full rank.

After some algebra we can write

$$LM = \left(n^{-1/2} \sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right)' \left(\tilde{\sigma}^2 n^{-1} \sum_{i=1}^n \hat{\mathbf{r}}_i \hat{\mathbf{r}}_i' \right)^{-1} \left(n^{-1/2} \sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right),$$

where $\tilde{\sigma}^2 = n^{-1} \sum_{i=1}^n \tilde{u}_i^2$ and each $\hat{\mathbf{r}}_i$ is a $k_2 \times 1$ vector of OLS residuals from the (multivariate) regression of \mathbf{x}_{i2} on \mathbf{x}_{i1} , $i = 1, \dots, n$. This statistic is not robust to heteroskedasticity because the matrix in the middle is not a consistent estimator of the asymptotic variance of $n^{-1/2} \sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i$ under heteroskedasticity. A heteroskedasticity-robust statistic is

$$\begin{aligned} LM_n &= \left(n^{-1/2} \sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right)' \left(n^{-1} \sum_{i=1}^n \tilde{u}_i^2 \hat{\mathbf{r}}_i \hat{\mathbf{r}}_i' \right)^{-1} \left(n^{-1/2} \sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right) \\ &= \left(\sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right)' \left(\sum_{i=1}^n \tilde{u}_i^2 \hat{\mathbf{r}}_i \hat{\mathbf{r}}_i' \right)^{-1} \left(\sum_{i=1}^n \hat{\mathbf{r}}_i \tilde{u}_i \right). \end{aligned}$$

Dropping the i subscript, this is easily obtained as $n - SSR_0$ from the OLS regression (without intercept)¹²

$$1 \text{ on } \tilde{u} \cdot \hat{\mathbf{r}}, \quad (21)$$

where $\tilde{u} \cdot \hat{\mathbf{r}} = (\tilde{u} \cdot \hat{r}_1, \dots, \tilde{u} \cdot \hat{r}_{k_2})'$ is the $k_2 \times 1$ vector obtained by multiplying \tilde{u} by each element of $\hat{\mathbf{r}}$ and SSR_0 is just the usual sum of squared residuals from the regression. Thus, we first regress each element of \mathbf{x}_2 onto all of \mathbf{x}_1 and collect the residuals in $\hat{\mathbf{r}}$. Then we form $\tilde{u} \cdot \hat{\mathbf{r}}$ (observation by observation) and run the regression in (21); $n - SSR_0$ from this regression is distributed asymptotically as $\chi_{k_2}^2$. For more details, see Davidson and MacKinnon (1985, 1993) or Wooldridge (1991, 1995).

¹²If there is an intercept, then the regression is trivial.

15 Test Consistency

We now define the test consistency against fixed alternatives. This concept was first introduced by Wald and Wolfowitz (1940).

Definition 1 A test of $H_0: \boldsymbol{\theta} \in \Theta_0$ is consistent against fixed alternatives if for all $\boldsymbol{\theta} \in \Theta_1$, $P(\text{Reject } H_0 | \boldsymbol{\theta}) \rightarrow 1$ as $n \rightarrow \infty$.

To understand this concept, consider the following simple example. Suppose that y_i is i.i.d. $N(\mu, 1)$. Consider the t -statistic $t_n(\mu) = \sqrt{n}(\bar{y} - \mu)$, and tests of $H_0: \mu = 0$ against $H_1: \mu > 0$. We reject H_0 if $t_n = t_n(0) > c$. Note that

$$t_n = t_n(\mu) + \sqrt{n}\mu$$

and $t_n(\mu) \equiv Z$ has an exact $N(0, 1)$ distribution. This is because $t_n(\mu)$ is centered at the true mean μ , while the test statistic $t_n(0)$ is centered at the (false) hypothesized mean of 0. The power of the test is

$$P(t_n > c | \mu) = P(Z + \sqrt{n}\mu > c) = 1 - \Phi(c - \sqrt{n}\mu).$$

This function is monotonically increasing in μ and n , and decreasing in c . Notice that for any c and $\mu \neq 0$, the power increases to 1 as $n \rightarrow \infty$. This means that for $\mu \in H_1$, the test will reject H_0 with probability approaching 1 as the sample size gets large. This is exactly test consistency.

For tests of the form “Reject H_0 if $T_n > c$ ”, a sufficient condition for test consistency is that T_n diverges to positive infinity with probability one for all $\boldsymbol{\theta} \in \Theta_1$. In general, the t -test and Wald test are consistent against fixed alternatives. For example, in testing $H_0: \theta = \theta_0$,

$$t_n = \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})} = \frac{\hat{\theta} - \theta}{s(\hat{\theta})} + \frac{\sqrt{n}(\theta - \theta_0)}{\sqrt{\hat{V}_\theta}} \quad (22)$$

since $s(\hat{\theta}) = \sqrt{\hat{V}_\theta/n}$. The first term on the right-hand-side converges in distribution to $N(0, 1)$. The second term on the right-hand-side equals zero if $\theta = \theta_0$, converges in probability to $+\infty$ if $\theta > \theta_0$, and converges in probability to $-\infty$ if $\theta < \theta_0$. Thus the two-sided t -test is consistent against $H_1: \theta \neq \theta_0$, and one-sided t -tests are consistent against the alternatives for which they are designed. For another example, the Wald statistic for $H_0: \boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$ against $H_1: \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ is

$$W_n = n(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \hat{\mathbf{V}}_\theta^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0).$$

Under H_1 , $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$. Thus $(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \hat{\mathbf{V}}_\theta^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{p} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{V}_\theta^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) > 0$. Hence under H_1 , $W_n \xrightarrow{p} \infty$. Again, this implies that Wald tests are consistent tests.

(*) Andrews (1986) introduces a testing analogue of estimator consistency, called **complete consistency**, which he shows is more appropriate than test consistency. It is shown that a sequence of estimators is consistent, if and only if certain tests based on the estimators (such as Wald or likelihood ratio tests) are completely consistent, for all simple null hypotheses.

16 Asymptotic Local Power

Consistency is a good property for a test, but does not give a useful approximation to the power of a test. To approximate the power function we need a distributional approximation. The standard asymptotic method for power analysis uses what are called **local alternatives**. This is similar to our analysis of restriction estimation under misspecification. The technique is to index the parameter by sample size so that the asymptotic distribution of the statistic is continuous in a localizing parameter. We first consider the t -test and then the Wald test.

In the t -test, we consider parameter vectors β_n which are indexed by sample size n and satisfy the real-valued relationship

$$\theta_n = r(\beta_n) = \theta_0 + n^{-1/2}h, \quad (23)$$

where the scalar h is called a **localizing parameter**. We index β_n and θ_n by sample size to indicate their dependence on n . The way to think of (23) is that the *true* value of the parameters are β_n and θ_n . The parameter θ_n is close to the hypothesized value θ_0 , with deviation $n^{-1/2}h$. Such a sequence of local alternatives θ_n is often called a **Pitman (1949) drift**¹³ or a **Pitman sequence**.¹⁴ We know for a fixed alternative, the power will converge to 1 as $n \rightarrow \infty$. To offset the effect of increasing n , we make the alternative harder to distinguish from H_0 as n gets larger. The rate $n^{-1/2}$ is the correct balance between these two forces. In the statistical literature, such alternatives are termed as "contiguous" local alternatives.

The specification (23) states that for any fixed h , θ_n approaches θ_0 as n gets large. Thus θ_n is "close" or "local" to θ_0 . The concept of a localizing sequence (23) might seem odd at first as in the actual world the sample size cannot mechanically affect the value of the parameter. Thus (23) should not be interpreted literally. Instead, it should be interpreted as a technical device which allows the asymptotic distribution of the test statistic to be continuous in the alternative hypothesis.

Similarly as in (22),

$$t_n = \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})} = \frac{\hat{\theta} - \theta_n}{s(\hat{\theta})} + \frac{\sqrt{n}(\theta_n - \theta_0)}{\sqrt{\hat{V}_\theta}} \xrightarrow{d} Z + \delta$$

under the local alternative (23), where $Z \sim N(0,1)$ and $\delta = h/\sqrt{V_\theta}$. In testing the one-sided alternative $H_1: \theta > \theta_0$, a t -test rejects H_0 for $t_n > z_\alpha$. The **asymptotic local power** of this test is the limit of the rejection probability under the local alternative (23),

$$\begin{aligned} \lim_{n \rightarrow \infty} P(\text{Reject } H_0 | \theta = \theta_n) &= \lim_{n \rightarrow \infty} P(t_n > z_\alpha | \theta = \theta_n) \\ &= P(Z + \delta > z_\alpha) = 1 - \Phi(z_\alpha - \delta) = \Phi(\delta - z_\alpha) \equiv \pi_\alpha(\delta). \end{aligned}$$

¹³Edwin J.G. Pitman (1897-1993) was an Australian mathematician who made significant contributions to statistics and probability theory. In particular, he is remembered primarily as the originator of the Pitman permutation test, Pitman nearness and Pitman efficiency, where Pitman efficiency is briefly discussed in Appendix B.

¹⁴See McManus (1991) for who invented local power analysis and Noether (1955) for exposition.

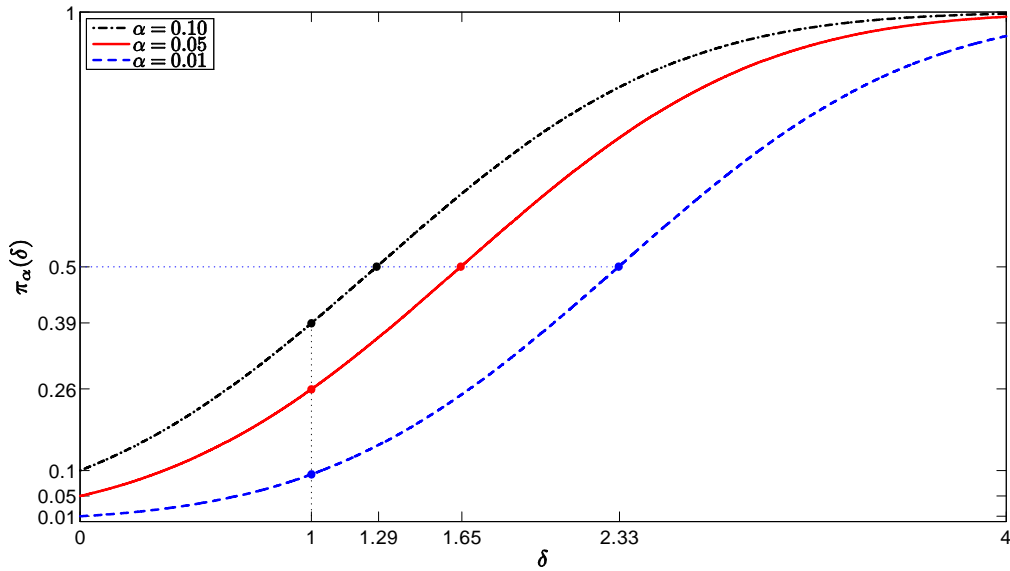


Figure 6: Asymptotic Local Power Function of One-Sided t -Test

We call $\pi_\alpha(\delta)$ the **local power function**.

Exercise 16 *Derive the local power function for the two-sided t -test.*

In Figure 6 we plot the local power function $\pi_\alpha(\delta)$ as a function of $\delta \in [0, 4]$ for tests of asymptotic size $\alpha = 0.10, \alpha = 0.05$, and $\alpha = 0.01$. We do not consider $\delta < 0$ since θ_n should be greater than θ_0 . $\delta = 0$ corresponds to the null hypothesis so $\pi_\alpha(0) = \alpha$. The power functions are monotonically increasing in both δ and α . The monotonicity with respect to α is due to the inherent trade-off between size and power. Decreasing size induces a decrease in power, and vice versa. The coefficient δ can be interpreted as the parameter deviation measured as a multiple of the standard error $s(\hat{\theta})$. To see this, recall that $s(\hat{\theta}) = n^{-1/2}\sqrt{\hat{V}_\theta} \approx n^{-1/2}\sqrt{V_\theta}$ and then note that

$$\delta = \frac{h}{\sqrt{V_\theta}} \approx \frac{n^{-1/2}h}{s(\hat{\theta})} = \frac{\theta_n - \theta_0}{s(\hat{\theta})},$$

meaning that δ equals the deviation $\theta_n - \theta_0$ expressed as multiples of the standard error $s(\hat{\theta})$. Thus as we examine Figure 6, we can interpret the power function at $\delta = 1$ (e.g., 26% for a 5% size test) as the power when the parameter θ_n is one standard error above the hypothesized value.

Exercise 17 *Suppose we have n_0 data points, and we want to know the power when the true value of θ is ϑ . Which δ we should confer in Figure 6?*

The difference between power functions can be measured either vertically or horizontally. For example, in Figure 6 there is a vertical dotted line at $\delta = 1$, showing that the asymptotic local

power $\pi_\alpha(1)$ equals 39% for $\alpha = 0.10$, equals 26% for $\alpha = 0.05$ and equals 9% for $\alpha = 0.01$. This is the difference in power across tests of differing sizes, holding fixed the parameter in the alternative. A horizontal comparison can also be illuminating. To illustrate, in Figure 6 there is a horizontal dotted line at 50% power. 50% power is a useful benchmark, as it is the point where the test has equal odds of rejection and acceptance. The dotted line crosses the three power curves at $\delta = 1.29$ ($\alpha = 0.10$), $\delta = 1.65$ ($\alpha = 0.05$), and $\delta = 2.33$ ($\alpha = 0.01$). This means that the parameter θ must be at least 1.65 standard errors above the hypothesized value for the one-sided test to have 50% (approximate) power. The ratio of these values (e.g., $1.65/1.29 = 1.28$ for the asymptotic 5% versus 10% tests) measures the relative parameter magnitude needed to achieve the same power. (Thus, for a 5% size test to achieve 50% power, the deviation $\theta_n - \theta_0$ must be 28% larger than for a 10% size test.) Even more interesting, the square of this ratio (e.g., $(1.65/1.29)^2 = 1.64$) can be interpreted as the increase in sample size needed to achieve the same power under fixed parameters. That is, to achieve 50% power, a 5% size test needs 64% more observations than a 10% size test. This interpretation follows by the following informal argument. By definition and (23) $\delta = h/\sqrt{V_\theta} = \sqrt{n}(\theta_n - \theta_0)/\sqrt{V_\theta}$. Thus holding θ and V_θ fixed, we can see that δ^2 is proportional to n .

We next generalize the local power analysis to the case of vector-valued alternatives. Now the local parametrization takes the form

$$\theta_n = \mathbf{r}(\beta_n) = \theta_0 + n^{-1/2}\mathbf{h}, \quad (24)$$

where \mathbf{h} is a $q \times 1$ vector. Under (24),

$$\sqrt{n}(\hat{\theta} - \theta_0) = \sqrt{n}(\hat{\theta} - \theta_n) + \mathbf{h} \xrightarrow{d} \mathbf{Z}_h \sim N(\mathbf{h}, \mathbf{V}_\theta),$$

a normal random vector with mean \mathbf{h} and variance matrix \mathbf{V}_θ . Applied to the Wald statistic we find

$$W_n = n(\hat{\theta} - \theta_0)' \hat{\mathbf{V}}_\theta^{-1} (\hat{\theta} - \theta_0) \xrightarrow{d} \mathbf{Z}_h' \mathbf{V}_\theta^{-1} \mathbf{Z}_h \sim \chi_q^2(\lambda).$$

where $\chi_q^2(\lambda)$ is a **non-central chi-square distribution** with q degrees of freedom and non-central parameter (or noncentrality) $\lambda = \mathbf{h}' \mathbf{V}_\theta^{-1} \mathbf{h}$.¹⁵ Under the null, $\mathbf{h} = \mathbf{0}$, and the $\chi_q^2(\lambda)$ distribution then degenerates to the usual χ_q^2 distribution. In the case of $q = 1$, $|Z + \delta|^2 \sim \chi_1^2(\lambda)$ with $\lambda = \delta^2$. The asymptotic local power of the Wald test at the level α is

$$P(\chi_q^2(\lambda) > \chi_{q,\alpha}^2) \equiv \pi_{q,\alpha}(\lambda).$$

Figure 7 plots $\pi_{q,0.05}(\lambda)$ (the power of asymptotic 5% tests) as a function of λ for $q = 1, 2$ and 3 . The power functions are monotonically increasing in λ and asymptote to one. Figure 7 also shows the power loss for fixed non-centrality parameter λ as the dimensionality of the test increases. The power curves shift to the right as q increases, resulting in a decrease in power. This is illustrated by

¹⁵See formula (26.4.25) in Abramowitz and Stegun (1965) for the infinite series representation of the noncentral chi-square.

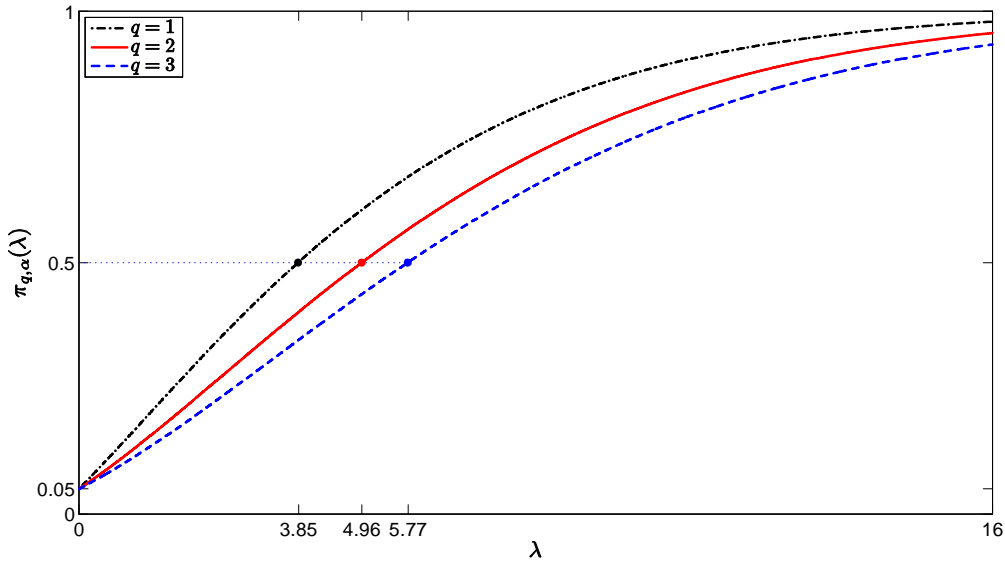


Figure 7: Asymptotic Local Power Function of the Wald Test

the dotted line at 50% power. The dotted line crosses the three power curves at $\lambda = 3.85$ ($q = 1$), $\lambda = 4.96$ ($q = 2$), and $\lambda = 5.77$ ($q = 3$). The ratio of these λ values correspond to the relative sample sizes needed to obtain the same power. Thus increasing the dimension of the test from $q = 1$ to $q = 2$ requires a 28% increase in sample size, or an increase from $q = 1$ to $q = 3$ requires a 50% increase in sample size, to obtain a test with 50% power. Intuitively, when testing more restrictions, we need more deviation from the null (or equivalently, more data points) to achieve the same power.

Exercise 18 (i) Show that $\widehat{\mathbf{V}}_{\boldsymbol{\theta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\theta}}$ under the local alternative $\boldsymbol{\beta}_n = \boldsymbol{\beta}_0 + n^{-1/2}\mathbf{b}$. (ii) If the local alternative $\boldsymbol{\beta}_n$ is specified as as in (i), what is the local power?

Exercise 19 (*) Derive the local power function of the minimum distance test with local alternatives (24). Is it the same as the Wald test?

(*) From the discussion above, we define the **asymptotic relative efficiency** (ARE) of any two tests.

Definition 2 (ARE) For any two tests of the same hypothesis, the same size, and same power with respect to the same alternative: if the first test requires n_1 observations and the second requires n_2 observations, then the relative efficiency of the second test with respect to the first is n_1/n_2 . The limit of this ratio as both numerator and denominator tend to infinity is the asymptotic relative efficiency of the two tests.

It turns out that for two t -tests with different $V_{\boldsymbol{\theta}}$'s, say $V_{\boldsymbol{\theta}}^1$ and $V_{\boldsymbol{\theta}}^2$, Pitman's ARE of the second test with respect to the first is $V_{\boldsymbol{\theta}}^1/V_{\boldsymbol{\theta}}^2$, independent of h . While for the Wald test, the corresponding

ARE is λ_1/λ_2 which generally depends on \mathbf{h} . See Pitman (1979) and Rao (1973) for additional details. In Appendix B, we review other approaches, besides Pitman's, toward assessment of the ARE of any two tests.

Exercise 20 (Empirical) *The data set invest.dat contains data on 565 U.S. firms extracted from Compustat for the year 1987. The variables, in order, are*

- I_i Investment to Capital Ratio (multiplied by 100).
- Q_i Total Market Value to Asset Ratio (Tobin's Q).
- C_i Cash Flow to Asset Ratio.
- D_i Long Term Debt to Asset Ratio.

The flow variables are annual sums for 1987. The stock variables are beginning of year.

- (a) *Estimate a linear regression of I_i on the other variables. Calculate appropriate standard errors.*
- (b) *Calculate asymptotic confidence intervals for the coefficients.*
- (c) *This regression is related to Tobin's q theory of investment, which suggests that investment should be predicted solely by Q_i . Thus the coefficient on Q_i should be positive and the others should be zero. Test the joint hypothesis that the coefficients on C_i and D_i are zero. Test the hypothesis that the coefficient on Q_i is zero. Are the results consistent with the predictions of the theory?*
- (d) *Now try a non-linear (quadratic) specification. Regress I_i on $Q_i, C_i, D_i, Q_i^2, C_i^2, D_i^2, Q_i C_i, Q_i D_i, C_i D_i$. Test the joint hypothesis that the six interaction and quadratic coefficients are zero.*

Exercise 21 (Empirical) *In a paper in 1963, Marc Nerlove analyzed a cost function for 145 American electric companies. The data file nerlov.dat contains his data. The variables are described as follows,*

- Column 1: total costs (call it TC) in millions of dollars
- Column 2: output (Q) in billions of kilowatt hours
- Column 3: price of labor (PL)
- Column 4: price of fuels (PF)
- Column 5: price of capital (PK)

Nerlove was interested in estimating a cost function: $TC = f(Q, PL, PF, PK)$.

(a) First estimate an unrestricted Cobb-Douglas specification

$$\log TC_i = \beta_1 + \beta_2 \log Q_i + \beta_3 \log PL_i + \beta_4 \log PK_i + \beta_5 \log PF_i + u_i. \quad (25)$$

Report parameter estimates and standard errors.

(b) Using a Wald statistic, test the hypothesis $H_0: \beta_3 + \beta_4 + \beta_5 = 1$.

(c) Estimate (25) by least-squares imposing this restriction by substitution. Report your parameter estimates and standard errors.

(d) Estimate (25) subject to $\beta_3 + \beta_4 + \beta_5 = 1$ using the RLS estimator. Do you obtain the same estimates as in part (c)?

Appendix A: Expected Length of a Confidence Interval

Suppose $L \leq \theta \leq U$ is a CI for θ , where L and U are random. The expected length of the interval is $E[U - L]$, which turns out equal to

$$\int_{\theta \neq \theta_0} P(L \leq \theta \leq U) d\theta,$$

the integrated probability of the interval covering the false value, where $P(\cdot)$ is the probability measure under the truth θ_0 . To see why, note that

$$\begin{aligned} E[U - L] &= \int \int 1(L \leq \theta \leq U) d\theta dP \\ &= \int \int 1(L \leq \theta \leq U) dP d\theta \\ &= \int P(L \leq \theta \leq U) d\theta = \int_{\theta \neq \theta_0} P(L \leq \theta \leq U) d\theta. \end{aligned}$$

So minimizing the expected length of a CI is equivalent to minimizing the coverage probability for each $\theta \neq \theta_0$. From the discussion in the main text,

$$P(L \leq \theta \leq U) = P(X \in A(\theta)),$$

where X is the data and $A(\theta)$ is the acceptance region for the null that θ is the true value. As a result, for $\theta \neq \theta_0$, $P(L \leq \theta \leq U)$ is the type II error in testing $H_0: \theta$ vs $H_1: \theta_0$, and minimizing $P(L \leq \theta \leq U)$ is equivalent to maximizing the power. By the Neyman-Pearson Lemma, the most powerful test is

$$1\left(\frac{f(X)}{f_{\theta_0}(X)} > c(\theta)\right), \quad (26)$$

where $f(X)$ is the true density of X (or density under θ_0), and $c(\theta)$ is the critical value for $H_0: \theta$. Collecting all θ 's such that $\frac{f(X)}{f_{\theta_0}(X)} \leq c(\theta)$ is the length minimizing CI.

The arguments above assume θ_0 were known, but if it were known, why do we need a CI for it? A more natural measure for the interval length is the average expected length $\int E_{\vartheta} [U - L] d\nu(\vartheta)$, where ϑ is the parameter in the model and $\nu(\cdot)$ is a measure for ϑ representing some prior information. Similarly, we can show

$$\int E_{\vartheta} [U - L] d\nu(\vartheta) = \int \int P_{\vartheta}(X \in A(\theta)) d\nu(\vartheta) d\theta.$$

Minimizing the average expected length is equivalent to maximizing the power in testing $H_0 : \theta$ vs $H_1 : P(\cdot)$, where $P(X \in A) = \int P_{\vartheta}(X \in A) d\nu(\vartheta)$ has the density $\int f_{\vartheta}(X) d\nu(\vartheta)$ with $f_{\vartheta}(\cdot)$ being the density associated with $P_{\vartheta}(\cdot)$. For any test φ , the average power

$$\int E_{\vartheta} [\varphi] d\nu(\vartheta) = \int \int \varphi f_{\vartheta}(X) dX d\nu(\vartheta) = \int \varphi \int f_{\vartheta}(X) d\nu(\vartheta) dX = \int \varphi dP(X)$$

is just the power in the above test, so the test maximizing the power against $P(\cdot)$ is equivalent to maximizing the average power $\int E_{\vartheta} [\varphi] d\nu(\vartheta)$. The corresponding test is the same as (26) but replaces $f(X)$ by $\int f_{\vartheta}(X) d\nu(\vartheta)$.

Appendix B: Approaches Toward ARE

Our discussions in this appendix follow Chapter 10 of Serfling (1980). Besides the earliest approach of ARE by Pitman (1949), there are five other approaches due to Chernoff (1952), Bahadur (1960), Hodges and Lehmann (1956), Hoeffding (1965), and Rubin and Sethuraman (1965). To understand the difference of these approaches, we first introduce some notations.

In Section 7 of Chapter 3, we represent H_0 and H_1 as Θ_0 and Θ_1 . For more general setups, we usually represent H_0 and H_1 as \mathcal{F}_0 and \mathcal{F}_1 for specified families of distributions for the data. For any test procedure T , define the function

$$\gamma_n(T, F) = P_F(T_n \text{ reject } H_0),$$

where $P_F(\cdot)$ means $P(\cdot | F \text{ is the truth})$, and the subscript n in T_n indicates the version of T based on a sample of size n . For $F \in \mathcal{F}_0$, $\gamma_n(T, F)$ represents the probability of a Type I error, and for $F \in \mathcal{F}_1$, it is the power function. The size of the test is

$$\alpha_n(T, \mathcal{F}_0) = \sup_{F \in \mathcal{F}_0} \gamma_n(T, F),$$

and the probability of a Type II error is

$$\beta_n(T, F) = 1 - \gamma_n(T, F).$$

We shall consider several performance criteria. Each entails specifications regarding

(a) $\alpha = \lim_{n \rightarrow \infty} \alpha_n(T, \mathcal{F}_0)$,

- (b) an alternative distribution $F^{(n)}$ allowed to depend on n ,
- (c) $\beta = \lim_{n \rightarrow \infty} \beta_n (T, F^{(n)})$.

With respect to (a), the cases $\alpha = 0$ and $\alpha > 0$ are distinguished. With respect to (c), the cases $\beta = 0$ and $\beta > 0$ are distinguished. With respect to (b), the cases $F^{(n)} \equiv F$ (fixed), and $F^{(n)} \rightarrow \mathcal{F}_0$ in some sense, are distinguished.

The following Table 3 summarizes relevant details and notation regarding the six approaches of ARE. Each of the approaches has its own special motivation and appeal. In an actual econometric problem, which approach to apply usually involves a trade-off between relevant intuitive considerations and availability of mathematical tools suitable for derivation of relevant $e(\cdot, \cdot)$. In the Pitman approach, the key mathematical tool is *central limit theory*. In the Rubin-Sethuraman approach, the theory of *moderate deviations* is used. In the other approaches, the theory of *large deviations* is employed.

Names of Contributors	Behavior of α_n	Behavior of β_n	Behavior of Alternatives	Notation for ARE
Pitman	$\alpha_n \rightarrow \alpha > 0$	$\beta_n \rightarrow \beta > 0$	$F^{(n)} \rightarrow \mathcal{F}_0$	$e_P(\cdot, \cdot)$
Chernoff	$\alpha_n \rightarrow 0$	$\beta_n \rightarrow 0$	$F^{(n)} = F$ (fixed)	$e_C(\cdot, \cdot)$
Bahadur	$\alpha_n \rightarrow 0$	$\beta_n \rightarrow \beta > 0$	$F^{(n)} = F$ (fixed)	$e_B(\cdot, \cdot)$
Hodges and Lehmann	$\alpha_n \rightarrow \alpha > 0$	$\beta_n \rightarrow 0$	$F^{(n)} = F$ (fixed)	$e_{HL}(\cdot, \cdot)$
Hoeffding	$\alpha_n \rightarrow 0$	$\beta_n \rightarrow 0$	$F^{(n)} = F$ (fixed)	$e_H(\cdot, \cdot)$
Rubin and Sethuraman	$\alpha_n \rightarrow 0$	$\beta_n \rightarrow 0$	$F^{(n)} \rightarrow \mathcal{F}_0$	$e_{RS}(\cdot, \cdot)$

Table 3: Summary of Six Approaches of ARE