

# Large-Sample Tests

Consider the linear regression model  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \beta_2\mathbf{x}_2 + \mathbf{u}$ . The  $u_i$  are assumed IID. Regressors may be predetermined and/or exogenous.

The  $t$  statistic for the hypothesis  $\beta_2 = \beta_2^0$  is

$$t_{\beta_2} = \frac{\hat{\beta}_2 - \beta_2^0}{\sqrt{s^2(\mathbf{X}^\top\mathbf{X})_{22}^{-1}}} = \frac{N^{1/2}(\hat{\beta}_2 - \beta_2^0)}{\sqrt{s^2N(\mathbf{X}^\top\mathbf{X})_{22}^{-1}}}, \quad (1)$$

where  $(\mathbf{X}^\top\mathbf{X})_{22}^{-1}$  is the last element on the main diagonal of  $(\mathbf{X}^\top\mathbf{X})^{-1}$ .

Theorem 4.3 tells us that:

- $N^{1/2}(\hat{\beta}_2 - \beta_2^0)$  is asymptotically normal with variance the last element on the main diagonal of  $\sigma_0^2(\mathbf{S}_{\mathbf{X}^\top\mathbf{X}})^{-1}$ ;
- $s^2$  times  $N(\mathbf{X}^\top\mathbf{X})_{22}^{-1}$  consistently estimates this variance.

Recall that, asymptotically,  $\mathbf{S}_{\mathbf{X}^\top\mathbf{X}} = N^{-1}\mathbf{X}^\top\mathbf{X}$ .

We conclude that

$$t_{\beta_2} \overset{a}{\sim} N(0, 1), \quad (2)$$

where “ $\overset{a}{\sim}$ ” means “**is asymptotically distributed as.**”

The result (2) justifies the use of  $t$  tests outside the confines of the classical normal linear model.

We can compute asymptotic  $P$  values or critical values using either the standard normal or  $t(N - k)$  distributions.

When the model is  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$ , the usual  $F$  statistic is asymptotically valid.

Rewriting  $F_{\beta_2}$  (under the null) in terms of quantities that are  $O_p(1)$ ,

$$F_{\beta_2} = \frac{N^{-1/2}\boldsymbol{\epsilon}^\top \mathbf{M}_1 \mathbf{X}_2 (N^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} N^{-1/2} \mathbf{X}_2^\top \mathbf{M}_1 \boldsymbol{\epsilon} / r}{\boldsymbol{\epsilon}^\top \mathbf{M}_X \boldsymbol{\epsilon} / (N - k)}, \quad (3)$$

where  $\boldsymbol{\epsilon} \equiv \mathbf{u} / \sigma_0$  and  $r = k_2$ , the dimension of  $\boldsymbol{\beta}_2$ .

We now show that  $rF_{\beta_2}$  is asymptotically distributed as  $\chi^2(r)$ , which follows from Theorem 4.3.

The denominator of (3) is  $\epsilon^\top \mathbf{M}_X \epsilon / (N - k)$ , which is just  $s^2 / \sigma_0^2$ .

Since  $s^2$  is consistent for  $\sigma_0^2$ , this denominator must tend to 1 asymptotically.

The numerator of  $rF_{\beta_2}$ , multiplied by  $r$ , is

$$N^{-1/2} \epsilon^\top \mathbf{M}_1 \mathbf{X}_2 (N^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} N^{-1/2} \mathbf{X}_2^\top \mathbf{M}_1 \epsilon. \quad (4)$$

Let  $\mathbf{v} = N^{-1/2} \mathbf{X}^\top \epsilon$ . Then a CLT shows that  $\mathbf{v} \stackrel{a}{\sim} N(\mathbf{0}, \mathbf{S}_{\mathbf{X}^\top \mathbf{X}})$ .

If we partition  $\mathbf{v}$  into two subvectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , we have

$$N^{-1/2} \mathbf{X}_2^\top \mathbf{M}_1 \epsilon = \mathbf{v}_2 - N^{-1} \mathbf{X}_2^\top \mathbf{X}_1 (N^{-1} \mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{v}_1. \quad (5)$$

This tends to  $\mathbf{v}_2 - \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{v}_1$  as  $N \rightarrow \infty$ . Here  $\mathbf{S}_{11}$  and  $\mathbf{S}_{21}$  are submatrices of  $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}$ , with  $N^{-1} \mathbf{X}_1^\top \mathbf{X}_1 \rightarrow \mathbf{S}_{11}$  and  $N^{-1} \mathbf{X}_2^\top \mathbf{X}_1 \rightarrow \mathbf{S}_{21}$ .

Since  $\mathbf{v}$  is asymptotically multivariate normal, so is  $\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1$ .

The vector  $\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1$  has mean  $\mathbf{0}$ , so its covariance matrix is the expectation of

$$\mathbf{v}_2\mathbf{v}_2^\top + \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1\mathbf{v}_1^\top\mathbf{S}_{11}^{-1}\mathbf{S}_{12} - 2\mathbf{v}_2\mathbf{v}_1^\top\mathbf{S}_{11}^{-1}\mathbf{S}_{12}. \quad (6)$$

Asymptotically, we can replace  $\mathbf{v}_1\mathbf{v}_1^\top$ ,  $\mathbf{v}_2\mathbf{v}_2^\top$ , and  $\mathbf{v}_2\mathbf{v}_1^\top$  by their plims, which are  $\mathbf{S}_{11}$ ,  $\mathbf{S}_{22}$ , and  $\mathbf{S}_{21}$ , respectively. To see this, observe that

$$\text{plim}(\mathbf{v}_1\mathbf{v}_1^\top) = \text{plim}(N^{-1}\mathbf{X}_1^\top\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top\mathbf{X}_1) = \text{plim}N^{-1}(\mathbf{X}_1^\top\mathbf{X}_1) = \mathbf{S}_{11}. \quad (7)$$

This allows us to conclude that

$$\text{Var}(\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1) = \mathbf{S}_{22} - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{S}_{12}. \quad (8)$$

Thus the numerator of the  $F$  statistic is asymptotically equal to

$$(\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1)^\top(\mathbf{S}_{22} - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{S}_{12})^{-1}(\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1). \quad (9)$$

(9) is a quadratic form in the  $r$ -vector  $\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1$ , which is asymptotically multivariate normal, and the inverse of  $\text{Var}(\mathbf{v}_2 - \mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{v}_1)$ , its covariance matrix.

By Theorem 4.1, expression (9) is asymptotically distributed as  $\chi^2(r)$ .

Because the denominator of the  $F$  statistic tends to 1, we conclude that

$$rF_{\beta_2} \stackrel{a}{\sim} \chi^2(r) \quad (10)$$

under  $H_0$  with predetermined regressors.

Since  $1/r$  times a  $\chi^2(r)$  random variable is distributed as  $F(r, \infty)$ , we may also conclude that  $F_{\beta_2} \stackrel{a}{\sim} F(r, N - k)$ .

(10) justifies **asymptotic  $F$  tests** when disturbances are not normally distributed and some regressors may be predetermined.

We can use either the  $\chi^2(r)$  or  $F(r, N - k)$  distributions. If we use the  $\chi^2(r)$  distribution, we have to multiply the  $F$  statistic by  $r$ .

# Wald Tests

A vector of  $r$  linear restrictions on a parameter vector  $\beta$  can always be written as

$$R\beta = r, \quad (11)$$

where  $R$  is an  $r \times k$  matrix and  $r$  is an  $r$ -vector.

For example, if  $k = 3$  and the restrictions were that  $\beta_1 = 0$  and  $\beta_2 = -1$ , equations (11) would be

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \beta = \begin{bmatrix} 0 \\ -1 \end{bmatrix}. \quad (12)$$

The elements of the matrix  $R$  and the vector  $r$  must be known. They are not functions of the data, and they are often integers.

Suppose that a  $k$ -vector of estimates  $\hat{\beta}$  has covariance matrix  $\text{Var}(\hat{\beta})$ .

It is easy to see that  $\text{Var}(R\hat{\beta})$  must be  $R \text{Var}(\hat{\beta}) R^\top$ .

We can test the restrictions (11) by calculating the **Wald statistic**

$$W(\hat{\beta}) = (R\hat{\beta} - r)^\top (R\widehat{\text{Var}}(\hat{\beta})R^\top)^{-1} (R\hat{\beta} - r), \quad (13)$$

where  $\widehat{\text{Var}}(\hat{\beta})$  estimates  $\text{Var}(\hat{\beta})$  consistently.

Inserting appropriate powers of  $N$ , we can rewrite (13) as

$$W(\hat{\beta}) = (N^{1/2}(R\hat{\beta} - r))^\top (RN\widehat{\text{Var}}(\hat{\beta})R^\top)^{-1} (N^{1/2}(R\hat{\beta} - r)). \quad (14)$$

By Theorem 4.3,  $N^{1/2}(R\hat{\beta} - r)$  is asymptotically multivariate normal.

Therefore, we can apply Theorem 4.1 (asymptotically) to (14) and conclude that  $W(\hat{\beta}) \stackrel{a}{\sim} \chi^2(r)$  under the null hypothesis.

These results are not just for OLS. Equations (13) and (14) would still define a Wald statistic if  $\hat{\beta}$  were any root- $N$  consistent estimator and  $\widehat{\text{Var}}(\hat{\beta})$  were any valid estimator of its covariance matrix.

# Properties of Asymptotic Tests

Asymptotic tests are almost never exact in finite samples.

They may either **over-reject** (reject a true null more than  $100\alpha\%$  of the time) or **under-reject** (reject it less than  $100\alpha\%$  of the time).

Whether an asymptotic test over-rejects or under-rejects, and how severely, depends on many things, including:

- the sample size;
- the distribution of the disturbances;
- the number of explanatory variables and their properties;
- the number of restrictions;
- the covariance matrix estimator employed, which may depend on assumptions about the disturbances;

Wald tests of many restrictions are particularly problematic. They often over-reject severely.



# Performing Several Hypothesis Tests

When we perform two or more tests, we are engaged in **multiple testing**.  $P$  values cannot be interpreted in the usual way.

An unusually large test statistic (unusually small  $P$  value) is more likely to be obtained by accident when several tests are performed.

Suppose we perform  $m$  exact tests at level  $\alpha$ . Let the **familywise error rate**  $\alpha_m$  be the probability that at least one of the tests rejects.

For independent tests, the familywise error rate is one minus the probability that none of the tests rejects:

$$\alpha_m = 1 - (1 - \alpha)^m. \quad (15)$$

When  $m$  is large,  $\alpha_m$  can be much larger than  $\alpha$ . For example, if  $\alpha = 0.05$ , then  $\alpha_2 = 0.0975$ ,  $\alpha_4 = 0.18549$ ,  $\alpha_8 = 0.33658$ , and  $\alpha_{16} = 0.55987$ .

The simplest method to control familywise error rate is the **Bonferroni procedure**. It rejects whenever the smallest  $P$  value is less than  $\alpha/m$ .

This procedure is based on the **Bonferroni inequality**

$$\Pr \left( \cup_{j=1}^m (P_j \leq \alpha/m) \right) \leq \alpha, \quad (16)$$

where  $P_j$  is the  $P$  value for the  $j^{\text{th}}$  test.

For large values of  $m$ ,  $\alpha/m$  is very much smaller than  $\alpha$ .

It is even smaller than the value  $\alpha$  that solves (15) for a given  $\alpha_m$ , which would be appropriate if all the tests were independent.

When the  $P_j$  are positively correlated,  $\alpha/m$  can be much too small.

If there is perfect dependence, all the tests yield identical  $P$  values, and the familywise error rate is just  $\alpha$ .

One way to control the familywise error rate is to use the bootstrap.

The **Simes procedure** is less conservative than the Bonferroni one. Order the  $P$  values from the smallest,  $P_{(1)}$ , to the largest,  $P_{(m)}$ . Then reject the joint null hypothesis whenever

$$P_{(j)} \leq j\alpha/m \text{ for any } j = 1, \dots, m, \quad (17)$$

where  $\alpha$  is the desired familywise error rate.

If  $P_{(1)} < \alpha/m$ , Simes and Bonferroni both reject.

But Simes can also reject when the second-smallest  $P$  value is less than  $2\alpha/m$ , the third-smallest is less than  $3\alpha/m$ , and so on. It must reject if the largest  $P$  value is less than  $\alpha$ , i.e., if every one of the tests rejects.

When many coefficients in the same regression are of interest, or when we estimate many different regressions with just one coefficient of interest in each of them,  $m$  can be large.

In either case, a single “large”  $t$  statistic, with a small  $P$  value, may not mean very much.

# The Power of Hypothesis Tests

An  $F$  statistic for  $\beta_2 = \mathbf{0}$  is the ratio of  $\|\mathbf{P}_{M_1 X_2} \mathbf{y}\|^2$  to  $\|\mathbf{M}_X \mathbf{y}\|^2$ , each divided by its appropriate number of degrees of freedom.

For the classical normal linear model, under both  $H_0$  and  $H_1$ , the numerator and denominator of this ratio are independent  $\chi^2$  variables, divided by their degrees of freedom.

Under either hypothesis,  $\mathbf{M}_X \mathbf{y} = \mathbf{M}_X \mathbf{u}$ . The difference in distribution under  $H_0$  and  $H_1$  must come from the numerator alone.

The numerator of  $F$ , times  $r/\sigma^2$ , is

$$\frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{M}_1 \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y}. \quad (18)$$

The vector  $\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{y}$  is normal under both hypotheses.

Its mean is  $\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 \beta_2$ , which vanishes when  $\beta_2 = \mathbf{0}$ , and its covariance matrix is  $\sigma^2 \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2$ .

Under the alternative, (18) follows the **noncentral chi-squared distribution**.

- If the  $m$ -vector  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$ , then  $\|\mathbf{z}\|^2 = \mathbf{z}^\top \mathbf{z}$  is distributed as (central) chi-squared with  $m$  degrees of freedom.
- Similarly, if  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{\Omega})$ , then  $\mathbf{x}^\top \mathbf{\Omega}^{-1} \mathbf{x} \sim \chi^2(m)$ .
- If instead  $\mathbf{z} \sim N(\boldsymbol{\mu}, \mathbf{I})$ , then  $\mathbf{z}^\top \mathbf{z}$  follows the noncentral chi-squared distribution with  $m$  degrees of freedom and **noncentrality parameter**, or **NCP**,  $\Lambda \equiv \boldsymbol{\mu}^\top \boldsymbol{\mu}$ .
- This distribution is written as  $\chi^2(m, \Lambda)$ . The expectation of  $\mathbf{z}^\top \mathbf{z}$  is  $m + \Lambda$ .
- When  $\boldsymbol{\mu} = \mathbf{0}$ , the  $\chi^2(m, 0)$  distribution is just the central  $\chi^2(m)$  distribution.

If  $\mathbf{x} \sim N(\boldsymbol{\nu}, \mathbf{\Omega})$ , then  $\mathbf{x}^\top \mathbf{\Omega}^{-1} \mathbf{x} \sim \chi^2(m, \boldsymbol{\nu}^\top \mathbf{\Omega}^{-1} \boldsymbol{\nu})$ . The distribution depends on  $\boldsymbol{\nu}$  and  $\mathbf{\Omega}$  only through the NCP,  $\Lambda = \boldsymbol{\nu}^\top \mathbf{\Omega}^{-1} \boldsymbol{\nu}$ .

In general, (18) is noncentral chi-squared, with  $r$  degrees of freedom and NCP

$$\Lambda \equiv \frac{1}{\sigma^2} \boldsymbol{\beta}_2^\top \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 \quad (19)$$

$$= \frac{1}{\sigma^2} \boldsymbol{\beta}_2^\top \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2. \quad (20)$$

Under the null hypothesis,  $\Lambda = 0$ .

Under either hypothesis, the distribution of the denominator of the  $F$  statistic, divided by  $\sigma^2$ , is central  $\chi^2(N - k)$ , and it is independent of the numerator.

The  $F$  statistic is thus the ratio of two random variables, divided by their degrees of freedom. The numerator is noncentral  $\chi^2(r)$ , and the denominator is central  $\chi^2(N - k)$ .

Under the normality assumption, these two random variables are independent of each other

We can write the distribution of the  $F$  statistic as

$$\frac{\chi^2(r, \Lambda)/r}{\chi^2(N-k)/(N-k)}, \quad (21)$$

with numerator and denominator mutually independent.

This is called the (singly) **noncentral  $F$  distribution**, with  $r$  and  $N - k$  degrees of freedom and NCP  $\Lambda$ .

- The difference between the distributions of the  $F$  statistic under  $H_0$  and  $H_1$  depends only on  $\Lambda$ .
- The numerator of the  $F$  statistic in (18), multiplied by  $r/\sigma^2$ , is (perhaps only asymptotically) distributed as  $\chi^2(r, \Lambda)$ . The NCP  $\Lambda$  depends on  $\sigma$ ,  $\mathbf{X}$ , and  $\beta_2$ .
- As  $\Lambda$  increases, the distribution of the asymptotic test statistic moves to the right and becomes more spread out, causing the power of the test to increase. Power also depends on  $r$ .

Under the alternative, (18) is equal to

$$\frac{1}{\sigma^2} \mathbf{u}^\top \mathbf{M}_1 \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{u} \quad (22)$$

$$+ \frac{2}{\sigma^2} \mathbf{u}^\top \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + \frac{1}{\sigma^2} \boldsymbol{\beta}_2^\top \mathbf{X}_2^\top \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2. \quad (23)$$

The first term follows the central  $\chi^2(r)$  distribution.

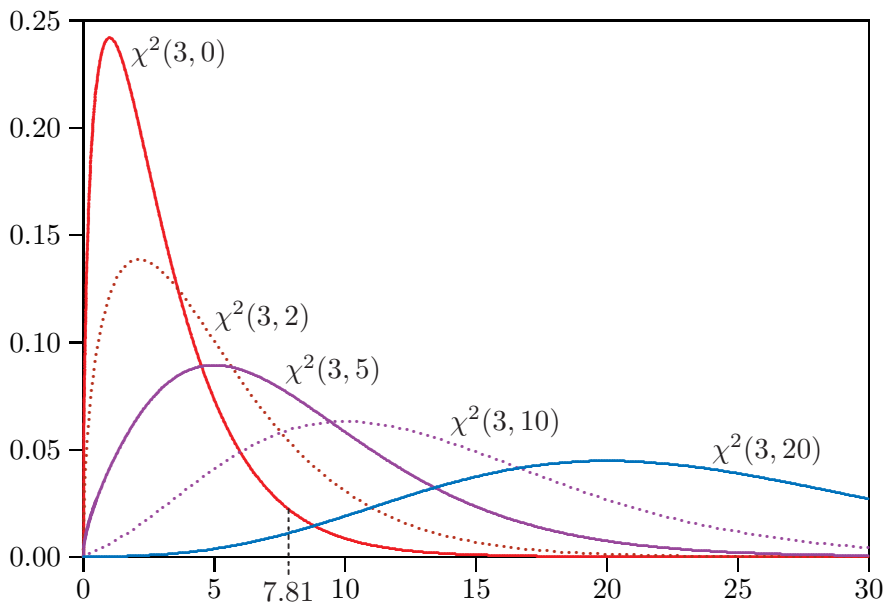
The second term is a random scalar which is normally distributed with mean zero and variance  $4\Lambda$ . It causes the  $\chi^2(r, \Lambda)$  distribution to become more spread out as  $\Lambda$  increases.

The third term is the NCP,  $\Lambda$ . It is what causes the  $\chi^2(r, \Lambda)$  distribution to move to the right as  $\Lambda$  increases.

Thus, informally, (18) equals  $\chi^2(r) + N(0, 4\Lambda) + \Lambda$ .

The next figure shows the density of the  $\chi^2(3, \Lambda)$  distribution for noncentrality parameters of 0, 2, 5, 10, and 20.





The signed square root of an  $F$  statistic for a single restriction is the  $t$  statistic:

$$t_{\beta_2} = \frac{\mathbf{x}_2^\top \mathbf{M}_1 \mathbf{y}}{s(\mathbf{x}_2^\top \mathbf{M}_1 \mathbf{x}_2)^{1/2}}. \quad (24)$$

The numerator,  $\mathbf{x}_2^\top \mathbf{M}_1 \mathbf{y}$ , is normally distributed under both  $H_0$  and  $H_1$ , with variance  $\sigma^2 \mathbf{x}_2^\top \mathbf{M}_1 \mathbf{x}_2$  and mean  $\mathbf{x}_2^\top \mathbf{M}_1 \mathbf{x}_2 \beta_2$ .

Thus  $s/\sigma$  times (24) is normal with variance 1 and mean

$$\lambda \equiv \frac{1}{\sigma} (\mathbf{x}_2^\top \mathbf{M}_1 \mathbf{x}_2)^{1/2} \beta_2. \quad (25)$$

It follows that  $t_{\beta_2}$  has a distribution which can be written as

$$\frac{N(\lambda, 1)}{(\chi^2(N - k) / (N - k))^{1/2}}, \quad (26)$$

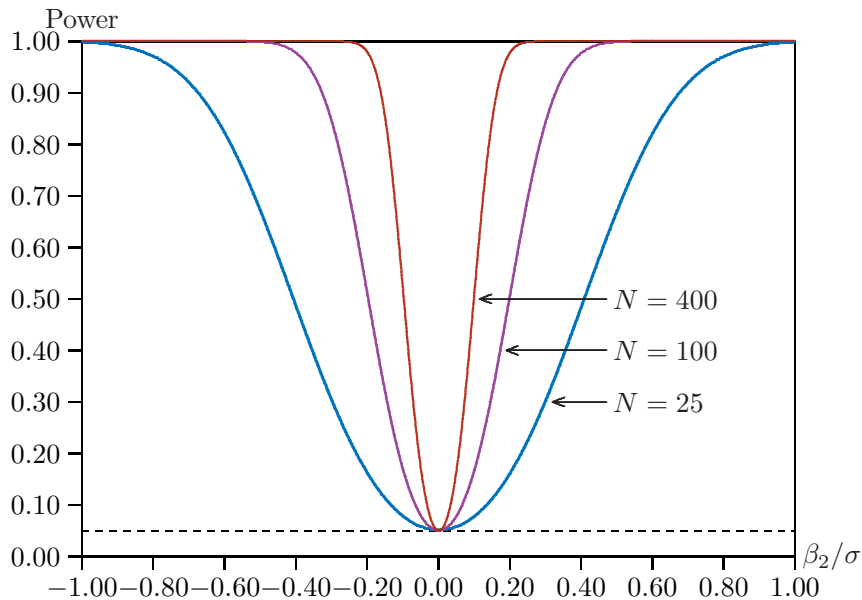
with independent numerator and denominator.

The distribution (26) is the **noncentral  $t$  distribution**, with  $N - k$  degrees of freedom and noncentrality parameter  $\lambda$ ; it is written as  $t(N - k, \lambda)$ .

- $\lambda^2 = \Lambda$ , where  $\Lambda$  is the NCP of the corresponding  $F$  statistic.
- The  $t(N - k, \lambda)$  distribution is similar to the  $N(\lambda, 1)$  distribution.
- It is also like an ordinary (**central**)  $t$  distribution with its mean shifted from the origin to (25), but it has a bit more variance.
- The distribution of  $t_{\beta_2}$  under  $H_1$  depends only on the NCP  $\lambda$ .
- For a given regressor matrix  $\mathbf{X}$  and sample size  $N$ ,  $\lambda$  in turn depends on the parameters only through the ratio  $\beta_2/\sigma$ ; see (25).
- Therefore, the power of the  $t$  test depends on  $\beta_2/\sigma$ .

The **power function** of a test gives the power of the test at any given level as a function of the parameters of the alternative hypothesis.

The next figure shows power functions for tests at the .05 level. The only regressor,  $x_2$ , is a constant.



- Since the test is exact, all the power functions equal .05 when  $\beta_2 = 0$ . Power then increases as  $\beta_2$  moves away from 0.
- Power when  $N = 400$  exceeds power when  $N = 100$ , which in turn exceeds power when  $N = 25$ , for every value of  $\beta_2 \neq 0$ .
- As  $N \rightarrow \infty$ , the power function converges to a T, with the foot of the vertical segment at .05 and the horizontal segment at 1.0.
- Asymptotically, the test rejects with probability 1 whenever the null is false. In finite samples, however,  $N^{1/2}\beta_2/\sigma$  must be sufficiently large for rejection to be likely.

In finite samples, asymptotic tests may have power functions that do not look like the figure.

- Power may be greater or less than .05 when  $H_0$  holds, and it may be minimized at a value that does not correspond to the null.
- Power may not even tend to unity as the parameter under test becomes infinitely far from the null hypothesis.

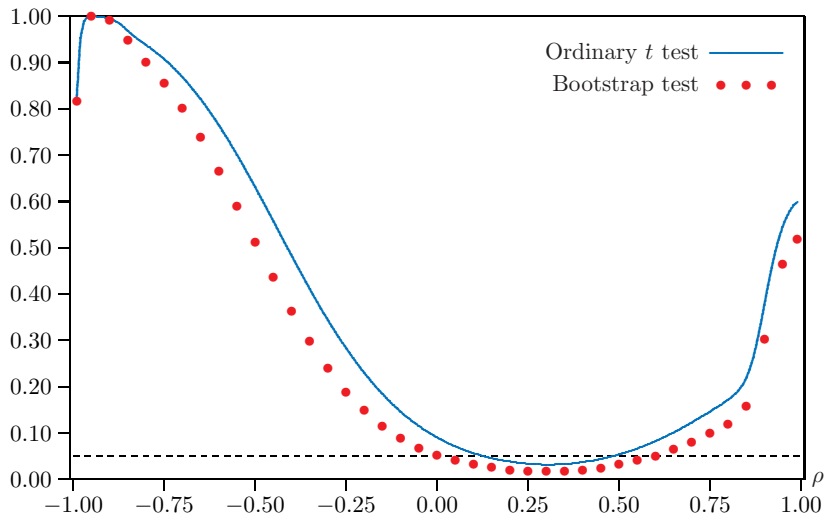
As an example of a power function that does not look the way asymptotic theory suggests, consider the model

$$y_t = \beta_1 + \sum_{j=2}^4 \beta_j X_{tj} + \delta y_{t-1} + u_t, \quad u_t = \rho u_{t-1} + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2), \quad (27)$$

where  $T = 20$ , the  $X_{tj}$  follow an AR(1) process with coefficient 0.75, all the  $\beta_j$  are equal to 1,  $\sigma = 0.1$ , and  $\delta = 0.9$ .

The next figure shows the power of asymptotic  $t$  and bootstrap tests for  $\rho = 0$  as a function of  $\rho$ , which varies between  $-0.99$  and  $0.99$ .

- For neither test does the power function achieve its minimum at  $\rho = 0$  or increase monotonically as  $|\rho|$  increases.
- Power actually declines sharply as  $\rho$  approaches  $-1$ .
- The asymptotic test has less power for values of  $\rho$  between 0 and about 0.62 than it does for  $\rho = 0$ .
- The bootstrap test rejects less than 5% of the time for values between 0 and about 0.61.



# Pre-Testing

Consider the linear regression model

$$y = X\beta + Z\gamma + u, \quad u \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}). \quad (28)$$

Here  $\beta$  is a  $k$ -vector,  $\gamma$  is an  $r$ -vector, and the regressors in  $X$  and  $Z$  are assumed, for simplicity, to be exogenous.

The parameters of interest are the  $k$  elements of  $\beta$ . We do not care about  $\gamma$ . They are **nuisance parameters**.

The unrestricted OLS estimator of  $\beta$  and its covariance matrix are

$$\hat{\beta} = (X^\top M_Z X)^{-1} X^\top M_Z y \quad \text{and} \quad \text{Var}(\hat{\beta}) = \sigma^2 (X^\top M_Z X)^{-1}. \quad (29)$$

Similarly, the restricted OLS estimator and its covariance matrix are

$$\tilde{\beta} = (X^\top X)^{-1} X^\top y \quad \text{and} \quad \text{Var}(\tilde{\beta}) = \sigma^2 (X^\top X)^{-1}. \quad (30)$$



Except when  $\mathbf{X}$  and  $\mathbf{Z}$  are orthogonal,  $\tilde{\boldsymbol{\beta}}$  is more efficient than  $\hat{\boldsymbol{\beta}}$ .

However, because  $\tilde{\boldsymbol{\beta}}$  is biased if  $\boldsymbol{\gamma} \neq \mathbf{0}$ , its MSE matrix is larger than  $\text{Var}(\tilde{\boldsymbol{\beta}})$  in that case.

$$\text{MSE}(\tilde{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \boldsymbol{\gamma} \boldsymbol{\gamma}^\top \mathbf{Z}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (31)$$

which is  $\text{Var}(\tilde{\boldsymbol{\beta}})$  plus the bias vector times itself transposed.

Since  $\tilde{\boldsymbol{\beta}}$  is more efficient than  $\hat{\boldsymbol{\beta}}$  when  $\boldsymbol{\gamma}$  is zero, it seems natural to test whether  $\boldsymbol{\gamma} = \mathbf{0}$  and use  $\hat{\boldsymbol{\beta}}$  when the test rejects and  $\tilde{\boldsymbol{\beta}}$  when it does not.

This test is called a **preliminary test**, or **pre-test** for short. It implicitly defines a new estimator, which is called a **pre-test estimator**:

$$\hat{\boldsymbol{\beta}} = \mathbb{I}(F_{\boldsymbol{\gamma}=\mathbf{0}} > c_\alpha) \hat{\boldsymbol{\beta}} + \mathbb{I}(F_{\boldsymbol{\gamma}=\mathbf{0}} \leq c_\alpha) \tilde{\boldsymbol{\beta}}, \quad (32)$$

where  $F_{\boldsymbol{\gamma}=\mathbf{0}}$  is the  $F$  statistic for  $\boldsymbol{\gamma} = \mathbf{0}$ , and  $c_\alpha$  is the critical value for an  $F$  test with  $r$  and  $N - k - r$  degrees of freedom at level  $\alpha$ .

# Properties of Pre-Test Estimators

Consider the model

$$y = \beta x + \gamma z + u, \quad u \sim \text{NID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (33)$$

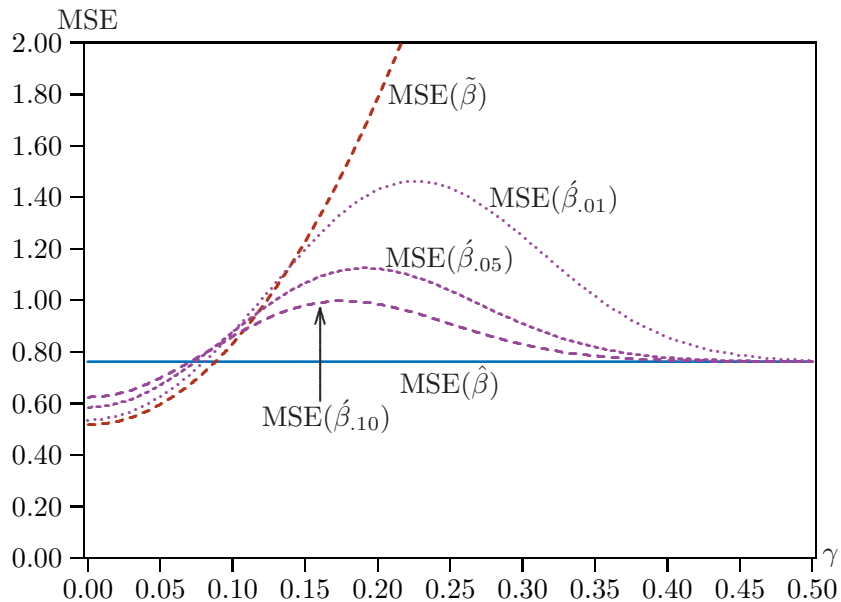
for which MSE is a scalar rather than a matrix

The two regressors are bivariate normal with correlation  $\rho = 0.5$ .

The potential reduction in variance from using  $\tilde{\beta}$  or  $\hat{\beta}$  rather than  $\hat{\beta}$  is evidently increasing in  $|\rho|$ , but so is the potential bias.

The figure shows MSE for five estimators of  $\beta$  as functions of  $\gamma$ . They are  $\hat{\beta}$ ,  $\tilde{\beta}$ , and three pre-test estimators,  $\hat{\beta}_{.01}$ ,  $\hat{\beta}_{.05}$ , and  $\hat{\beta}_{.10}$ .

The figure would look different if the NCP of the test were on the horizontal axis, but since the NCP is proportional to  $\gamma^2$ , both figures would contain the same information.



The MSE of  $\hat{\beta}$  is a horizontal line. It does not depend on  $\gamma$ .

$\text{MSE}(\tilde{\beta})$  is lower than  $\text{MSE}(\hat{\beta})$  when  $\gamma$  is sufficiently small, but it increases in proportion to  $\gamma^2$  (i.e., in proportion to the NCP).

- For small values of  $\gamma$ , the pre-test estimators have smaller MSE than  $\hat{\beta}$  but greater MSE than  $\tilde{\beta}$ . For very large values of  $\gamma$ , they perform essentially the same as  $\hat{\beta}$ .
- There is a large region in the middle of the figure where the pre-test estimators perform better than  $\tilde{\beta}$  but less well than  $\hat{\beta}$ .
- Near the point where  $\text{MSE}(\tilde{\beta})$  crosses  $\text{MSE}(\hat{\beta})$ , the pre-test estimator performs worse than either  $\tilde{\beta}$  or  $\hat{\beta}$ .

The level of the pre-test is important. When the level is high, the potential gain in efficiency for small values of  $\gamma$  is small, but so is the potential increase in MSE due to bias for intermediate values.

There is no reason to use a “conventional” significance level like .05 when pretesting. It is probably safer to use a higher level.