

# Chapter 10

## The Bootstrap

### 10.1 Definition of the Bootstrap

Let  $F$  denote a distribution function for the population of observations  $(y_i, \mathbf{x}_i)$ . Let

$$T_n = T_n((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n), F)$$

be a statistic of interest, for example an estimator  $\hat{\theta}$  or a t-statistic  $(\hat{\theta} - \theta)/s(\hat{\theta})$ . Note that we write  $T_n$  as possibly a function of  $F$ . For example, the t-statistic is a function of the parameter  $\theta$  which itself is a function of  $F$ .

The exact CDF of  $T_n$  when the data are sampled from the distribution  $F$  is

$$G_n(u, F) = \Pr(T_n \leq u \mid F)$$

In general,  $G_n(u, F)$  depends on  $F$ , meaning that  $G$  changes as  $F$  changes.

Ideally, inference would be based on  $G_n(u, F)$ . This is generally impossible since  $F$  is unknown.

Asymptotic inference is based on approximating  $G_n(u, F)$  with  $G(u, F) = \lim_{n \rightarrow \infty} G_n(u, F)$ . When  $G(u, F) = G(u)$  does not depend on  $F$ , we say that  $T_n$  is asymptotically pivotal and use the distribution function  $G(u)$  for inferential purposes.

In a seminal contribution, Efron (1979) proposed the bootstrap, which makes a different approximation. The unknown  $F$  is replaced by a consistent estimate  $F_n$  (one choice is discussed in the next section). Plugged into  $G_n(u, F)$  we obtain

$$G_n^*(u) = G_n(u, F_n). \quad (10.1)$$

We call  $G_n^*$  the bootstrap distribution. Bootstrap inference is based on  $G_n^*(u)$ .

Let  $(y_i^*, \mathbf{x}_i^*)$  denote random variables with the distribution  $F_n$ . A random sample from this distribution is called the bootstrap data. The statistic  $T_n^* = T_n((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), F_n)$  constructed on this sample is a random variable with distribution  $G_n^*$ . That is,  $\Pr(T_n^* \leq u) = G_n^*(u)$ . We call  $T_n^*$  the bootstrap statistic. The distribution of  $T_n^*$  is identical to that of  $T_n$  when the true CDF is  $F_n$  rather than  $F$ .

The bootstrap distribution is itself random, as it depends on the sample through the estimator  $F_n$ .

In the next sections we describe computation of the bootstrap distribution.

### 10.2 The Empirical Distribution Function

Recall that  $F(y, \mathbf{x}) = \Pr(y_i \leq y, \mathbf{x}_i \leq \mathbf{x}) = \mathbb{E}(1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x}))$ , where  $1(\cdot)$  is the indicator function. This is a population moment. The method of moments estimator is the corresponding

sample moment:

$$F_n(y, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n 1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x}). \quad (10.2)$$

$F_n(y, \mathbf{x})$  is called the empirical distribution function (EDF).  $F_n$  is a nonparametric estimate of  $F$ . Note that while  $F$  may be either discrete or continuous,  $F_n$  is by construction a step function.

The EDF is a consistent estimator of the CDF. To see this, note that for any  $(y, \mathbf{x})$ ,  $1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x})$  is an iid random variable with expectation  $F(y, \mathbf{x})$ . Thus by the WLLN (Theorem 2.6.1),  $F_n(y, \mathbf{x}) \xrightarrow{p} F(y, \mathbf{x})$ . Furthermore, by the CLT (Theorem 2.8.1),

$$\sqrt{n}(F_n(y, \mathbf{x}) - F(y, \mathbf{x})) \xrightarrow{d} N(0, F(y, \mathbf{x})(1 - F(y, \mathbf{x}))).$$

To see the effect of sample size on the EDF, in the Figure below, I have plotted the EDF and true CDF for three random samples of size  $n = 25, 50, 100$ , and  $500$ . The random draws are from the  $N(0, 1)$  distribution. For  $n = 25$ , the EDF is only a crude approximation to the CDF, but the approximation appears to improve for the large  $n$ . In general, as the sample size gets larger, the EDF step function gets uniformly close to the true CDF.

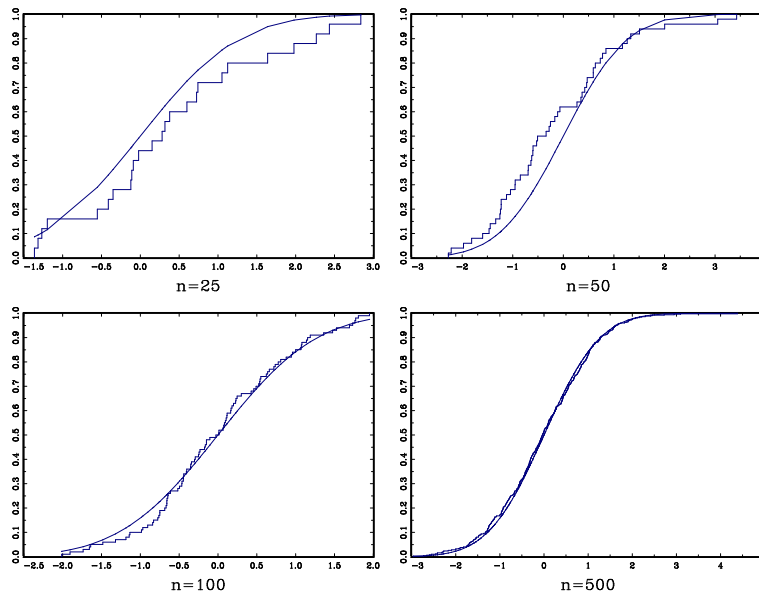


Figure 10.1: Empirical Distribution Functions

The EDF is a valid discrete probability distribution which puts probability mass  $1/n$  at each pair  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ . Notationally, it is helpful to think of a random pair  $(y_i^*, \mathbf{x}_i^*)$  with the distribution  $F_n$ . That is,

$$\Pr(y_i^* \leq y, \mathbf{x}_i^* \leq \mathbf{x}) = F_n(y, \mathbf{x}).$$

We can easily calculate the moments of functions of  $(y_i^*, \mathbf{x}_i^*)$ :

$$\begin{aligned} \mathbb{E}h(y_i^*, \mathbf{x}_i^*) &= \int h(y, \mathbf{x}) dF_n(y, \mathbf{x}) \\ &= \sum_{i=1}^n h(y_i, \mathbf{x}_i) \Pr(y_i^* = y_i, \mathbf{x}_i^* = \mathbf{x}_i) \\ &= \frac{1}{n} \sum_{i=1}^n h(y_i, \mathbf{x}_i), \end{aligned}$$

the empirical sample average.

### 10.3 Nonparametric Bootstrap

The **nonparametric bootstrap** is obtained when the bootstrap distribution (10.1) is defined using the EDF (10.2) as the estimate  $F_n$  of  $F$ .

Since the EDF  $F_n$  is a multinomial (with  $n$  support points), in principle the distribution  $G_n^*$  could be calculated by direct methods. However, as there are  $\binom{2n-1}{n}$  possible samples  $\{(y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*)\}$ , such a calculation is computationally infeasible. The popular alternative is to use simulation to approximate the distribution. The algorithm is identical to our discussion of Monte Carlo simulation, with the following points of clarification:

- The sample size  $n$  used for the simulation is the same as the sample size.
- The random vectors  $(y_i^*, \mathbf{x}_i^*)$  are drawn randomly from the empirical distribution. This is equivalent to sampling a pair  $(y_i, \mathbf{x}_i)$  randomly from the sample.

The bootstrap statistic  $T_n^* = T_n((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), F_n)$  is calculated for each bootstrap sample. This is repeated  $B$  times.  $B$  is known as the number of bootstrap replications. A theory for the determination of the number of bootstrap replications  $B$  has been developed by Andrews and Buchinsky (2000). It is desirable for  $B$  to be large, so long as the computational costs are reasonable.  $B = 1000$  typically suffices.

When the statistic  $T_n$  is a function of  $F$ , it is typically through dependence on a parameter. For example, the t-ratio  $(\hat{\theta} - \theta) / s(\hat{\theta})$  depends on  $\theta$ . As the bootstrap statistic replaces  $F$  with  $F_n$ , it similarly replaces  $\theta$  with  $\theta_n$ , the value of  $\theta$  implied by  $F_n$ . Typically  $\theta_n = \hat{\theta}$ , the parameter estimate. (When in doubt use  $\hat{\theta}$ .)

Sampling from the EDF is particularly easy. Since  $F_n$  is a discrete probability distribution putting probability mass  $1/n$  at each sample point, sampling from the EDF is equivalent to random sampling a pair  $(y_i, \mathbf{x}_i)$  from the observed data **with replacement**. In consequence, a bootstrap sample  $\{(y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*)\}$  will necessarily have some ties and multiple values, which is generally not a problem.

### 10.4 Bootstrap Estimation of Bias and Variance

The bias of  $\hat{\theta}$  is  $\tau_n = \mathbb{E}(\hat{\theta} - \theta_0)$ . Let  $T_n(\theta) = \hat{\theta} - \theta$ . Then  $\tau_n = \mathbb{E}(T_n(\theta_0))$ . The bootstrap counterparts are  $\hat{\theta}^* = \hat{\theta}((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*))$  and  $T_n^* = \hat{\theta}^* - \theta_n = \hat{\theta}^* - \hat{\theta}$ . The bootstrap estimate of  $\tau_n$  is

$$\tau_n^* = \mathbb{E}(T_n^*).$$

If this is calculated by the simulation described in the previous section, the estimate of  $\tau_n^*$  is

$$\begin{aligned} \hat{\tau}_n^* &= \frac{1}{B} \sum_{b=1}^B T_{nb}^* \\ &= \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^* - \hat{\theta} \\ &= \overline{\hat{\theta}^*} - \hat{\theta}. \end{aligned}$$

If  $\hat{\theta}$  is biased, it might be desirable to construct a biased-corrected estimator (one with reduced bias). Ideally, this would be

$$\tilde{\theta} = \hat{\theta} - \tau_n,$$

but  $\tau_n$  is unknown. The (estimated) bootstrap biased-corrected estimator is

$$\begin{aligned}\tilde{\theta}^* &= \hat{\theta} - \hat{\tau}_n^* \\ &= \hat{\theta} - (\overline{\hat{\theta}^*} - \hat{\theta}) \\ &= 2\hat{\theta} - \overline{\hat{\theta}^*}.\end{aligned}$$

Note, in particular, that the biased-corrected estimator is *not*  $\overline{\hat{\theta}^*}$ . Intuitively, the bootstrap makes the following experiment. Suppose that  $\hat{\theta}$  is the truth. Then what is the average value of  $\hat{\theta}$  calculated from such samples? The answer is  $\overline{\hat{\theta}^*}$ . If this is lower than  $\hat{\theta}$ , this suggests that the estimator is *downward-biased*, so a biased-corrected estimator of  $\theta$  should be *larger* than  $\hat{\theta}$ , and the best guess is the difference between  $\hat{\theta}$  and  $\overline{\hat{\theta}^*}$ . Similarly if  $\overline{\hat{\theta}^*}$  is higher than  $\hat{\theta}$ , then the estimator is upward-biased and the biased-corrected estimator should be lower than  $\hat{\theta}$ .

Let  $T_n = \hat{\theta}$ . The variance of  $\hat{\theta}$  is

$$V_n = \mathbb{E}(T_n - \mathbb{E}T_n)^2.$$

Let  $T_n^* = \hat{\theta}^*$ . It has variance

$$V_n^* = \mathbb{E}(T_n^* - \mathbb{E}T_n^*)^2.$$

The simulation estimate is

$$\hat{V}_n^* = \frac{1}{B} \sum_{b=1}^B \left( \hat{\theta}_b^* - \overline{\hat{\theta}^*} \right)^2.$$

A bootstrap standard error for  $\hat{\theta}$  is the square root of the bootstrap estimate of variance,  $s^*(\hat{\theta}) = \sqrt{\hat{V}_n^*}$ .

While this standard error may be calculated and reported, it is not clear if it is useful. The primary use of asymptotic standard errors is to construct asymptotic confidence intervals, which are based on the asymptotic normal approximation to the t-ratio. However, the use of the bootstrap presumes that such asymptotic approximations might be poor, in which case the normal approximation is suspected. It appears superior to calculate bootstrap confidence intervals, and we turn to this next.

## 10.5 Percentile Intervals

For a distribution function  $G_n(u, F)$ , let  $q_n(\alpha, F)$  denote its quantile function. This is the function which solves

$$G_n(q_n(\alpha, F), F) = \alpha.$$

[When  $G_n(u, F)$  is discrete,  $q_n(\alpha, F)$  may be non-unique, but we will ignore such complications.] Let  $q_n(\alpha)$  denote the quantile function of the true sampling distribution, and  $q_n^*(\alpha) = q_n(\alpha, F_n)$  denote the quantile function of the bootstrap distribution. Note that this function will change depending on the underlying statistic  $T_n$  whose distribution is  $G_n$ .

Let  $T_n = \hat{\theta}$ , an estimate of a parameter of interest. In  $(1 - \alpha)\%$  of samples,  $\hat{\theta}$  lies in the region  $[q_n(\alpha/2), q_n(1 - \alpha/2)]$ . This motivates a confidence interval proposed by Efron:

$$C_1 = [q_n^*(\alpha/2), q_n^*(1 - \alpha/2)].$$

This is often called the *percentile confidence interval*.

Computationally, the quantile  $q_n^*(\alpha)$  is estimated by  $\hat{q}_n^*(\alpha)$ , the  $\alpha$ 'th sample quantile of the simulated statistics  $\{T_{n1}^*, \dots, T_{nB}^*\}$ , as discussed in the section on Monte Carlo simulation. The  $(1 - \alpha)\%$  Efron percentile interval is then  $[\hat{q}_n^*(\alpha/2), \hat{q}_n^*(1 - \alpha/2)]$ .

The interval  $C_1$  is a popular bootstrap confidence interval often used in empirical practice. This is because it is easy to compute, simple to motivate, was popularized by Efron early in the history of the bootstrap, and also has the feature that it is translation invariant. That is, if we define  $\phi = f(\theta)$  as the parameter of interest for a monotonically increasing function  $f$ , then percentile method applied to this problem will produce the confidence interval  $[f(q_n^*(\alpha/2)), f(q_n^*(1 - \alpha/2))]$ , which is a naturally good property.

However, as we show now,  $C_1$  is in a deep sense very poorly motivated.

It will be useful if we introduce an alternative definition of  $C_1$ . Let  $T_n(\theta) = \hat{\theta} - \theta$  and let  $q_n(\alpha)$  be the quantile function of its distribution. (These are the original quantiles, with  $\theta$  subtracted.) Then  $C_1$  can alternatively be written as

$$C_1 = [\hat{\theta} + q_n^*(\alpha/2), \hat{\theta} + q_n^*(1 - \alpha/2)].$$

This is a bootstrap estimate of the “ideal” confidence interval

$$C_1^0 = [\hat{\theta} + q_n(\alpha/2), \hat{\theta} + q_n(1 - \alpha/2)].$$

The latter has coverage probability

$$\begin{aligned} \Pr(\theta_0 \in C_1^0) &= \Pr(\hat{\theta} + q_n(\alpha/2) \leq \theta_0 \leq \hat{\theta} + q_n(1 - \alpha/2)) \\ &= \Pr(-q_n(1 - \alpha/2) \leq \hat{\theta} - \theta_0 \leq -q_n(\alpha/2)) \\ &= G_n(-q_n(\alpha/2), F_0) - G_n(-q_n(1 - \alpha/2), F_0) \end{aligned}$$

which generally is not  $1 - \alpha$ ! There is one important exception. If  $\hat{\theta} - \theta_0$  has a symmetric distribution about 0, then  $G_n(-u, F_0) = 1 - G_n(u, F_0)$ , so

$$\begin{aligned} \Pr(\theta_0 \in C_1^0) &= G_n(-q_n(\alpha/2), F_0) - G_n(-q_n(1 - \alpha/2), F_0) \\ &= (1 - G_n(q_n(\alpha/2), F_0)) - (1 - G_n(q_n(1 - \alpha/2), F_0)) \\ &= \left(1 - \frac{\alpha}{2}\right) - \left(1 - \left(1 - \frac{\alpha}{2}\right)\right) \\ &= 1 - \alpha \end{aligned}$$

and this idealized confidence interval is accurate. Therefore,  $C_1^0$  and  $C_1$  are designed for the case that  $\hat{\theta}$  has a symmetric distribution about  $\theta_0$ .

When  $\hat{\theta}$  does not have a symmetric distribution,  $C_1$  may perform quite poorly.

However, by the translation invariance argument presented above, it also follows that if there exists some monotonically increasing transformation  $f(\cdot)$  such that  $f(\hat{\theta})$  is symmetrically distributed about  $f(\theta_0)$ , then the idealized percentile bootstrap method will be accurate.

Based on these arguments, many argue that the percentile interval should not be used unless the sampling distribution is close to unbiased and symmetric.

The problems with the percentile method can be circumvented, at least in principle, by an alternative method.

Let  $T_n(\theta) = \hat{\theta} - \theta$ . Then

$$\begin{aligned} 1 - \alpha &= \Pr(q_n(\alpha/2) \leq T_n(\theta_0) \leq q_n(1 - \alpha/2)) \\ &= \Pr(\hat{\theta} - q_n(1 - \alpha/2) \leq \theta_0 \leq \hat{\theta} - q_n(\alpha/2)), \end{aligned}$$

so an exact  $(1 - \alpha)\%$  confidence interval for  $\theta_0$  would be

$$C_2^0 = [\hat{\theta} - q_n(1 - \alpha/2), \hat{\theta} - q_n(\alpha/2)].$$

This motivates a bootstrap analog

$$C_2 = [\hat{\theta} - q_n^*(1 - \alpha/2), \hat{\theta} - q_n^*(\alpha/2)].$$

Notice that generally this is very different from the Efron interval  $C_1$ ! They coincide in the special case that  $G_n^*(u)$  is symmetric about  $\hat{\theta}$ , but otherwise they differ.

Computationally, this interval can be estimated from a bootstrap simulation by sorting the bootstrap statistics  $T_n^* = (\hat{\theta}^* - \hat{\theta})$ , which are centered at the sample estimate  $\hat{\theta}$ . These are sorted to yield the quantile estimates  $\hat{q}_n^*(.025)$  and  $\hat{q}_n^*(.975)$ . The 95% confidence interval is then  $[\hat{\theta} - \hat{q}_n^*(.975), \hat{\theta} - \hat{q}_n^*(.025)]$ .

This confidence interval is discussed in most theoretical treatments of the bootstrap, but is not widely used in practice.

## 10.6 Percentile-t Equal-Tailed Interval

Suppose we want to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta < \theta_0$  at size  $\alpha$ . We would set  $T_n(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$  and reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $T_n(\theta_0) < c$ , where  $c$  would be selected so that

$$\Pr(T_n(\theta_0) < c) = \alpha.$$

Thus  $c = q_n(\alpha)$ . Since this is unknown, a bootstrap test replaces  $q_n(\alpha)$  with the bootstrap estimate  $\hat{q}_n^*(\alpha)$ , and the test rejects if  $T_n(\theta_0) < \hat{q}_n^*(\alpha)$ .

Similarly, if the alternative is  $\mathbb{H}_1 : \theta > \theta_0$ , the bootstrap test rejects if  $T_n(\theta_0) > \hat{q}_n^*(1 - \alpha)$ .

Computationally, these critical values can be estimated from a bootstrap simulation by sorting the bootstrap t-statistics  $T_n^* = (\hat{\theta}^* - \hat{\theta}) / s(\hat{\theta}^*)$ . Note, and this is important, that the bootstrap test statistic is centered at the estimate  $\hat{\theta}$ , and the standard error  $s(\hat{\theta}^*)$  is calculated on the bootstrap sample. These t-statistics are sorted to find the estimated quantiles  $\hat{q}_n^*(\alpha)$  and/or  $\hat{q}_n^*(1 - \alpha)$ .

Let  $T_n(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$ . Then taking the intersection of two one-sided intervals,

$$\begin{aligned} 1 - \alpha &= \Pr(q_n(\alpha/2) \leq T_n(\theta_0) \leq q_n(1 - \alpha/2)) \\ &= \Pr\left(q_n(\alpha/2) \leq (\hat{\theta} - \theta_0) / s(\hat{\theta}) \leq q_n(1 - \alpha/2)\right) \\ &= \Pr\left(\hat{\theta} - s(\hat{\theta})q_n(1 - \alpha/2) \leq \theta_0 \leq \hat{\theta} - s(\hat{\theta})q_n(\alpha/2)\right), \end{aligned}$$

so an exact  $(1 - \alpha)\%$  confidence interval for  $\theta_0$  would be

$$C_3^0 = [\hat{\theta} - s(\hat{\theta})q_n(1 - \alpha/2), \hat{\theta} - s(\hat{\theta})q_n(\alpha/2)].$$

This motivates a bootstrap analog

$$C_3 = [\hat{\theta} - s(\hat{\theta})\hat{q}_n^*(1 - \alpha/2), \hat{\theta} - s(\hat{\theta})\hat{q}_n^*(\alpha/2)].$$

This is often called a *percentile-t confidence interval*. It is *equal-tailed* or *central* since the probability that  $\theta_0$  is below the left endpoint approximately equals the probability that  $\theta_0$  is above the right endpoint, each  $\alpha/2$ .

Computationally, this is based on the critical values from the one-sided hypothesis tests, discussed above.

## 10.7 Symmetric Percentile-t Intervals

Suppose we want to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta \neq \theta_0$  at size  $\alpha$ . We would set  $T_n(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$  and reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $|T_n(\theta_0)| > c$ , where  $c$  would be selected so that

$$\Pr(|T_n(\theta_0)| > c) = \alpha.$$

Note that

$$\begin{aligned}\Pr(|T_n(\theta_0)| < c) &= \Pr(-c < T_n(\theta_0) < c) \\ &= G_n(c) - G_n(-c) \\ &\equiv \overline{G}_n(c),\end{aligned}$$

which is a symmetric distribution function. The ideal critical value  $c = q_n(\alpha)$  solves the equation

$$\overline{G}_n(q_n(\alpha)) = 1 - \alpha.$$

Equivalently,  $q_n(\alpha)$  is the  $1 - \alpha$  quantile of the distribution of  $|T_n(\theta_0)|$ .

The bootstrap estimate is  $q_n^*(\alpha)$ , the  $1 - \alpha$  quantile of the distribution of  $|T_n^*|$ , or the number which solves the equation

$$\overline{G}_n^*(q_n^*(\alpha)) = G_n^*(q_n^*(\alpha)) - G_n^*(-q_n^*(\alpha)) = 1 - \alpha.$$

Computationally,  $q_n^*(\alpha)$  is estimated from a bootstrap simulation by sorting the bootstrap t-statistics  $|T_n^*| = |\hat{\theta}^* - \hat{\theta}|/s(\hat{\theta}^*)$ , and taking the upper  $\alpha\%$  quantile. The bootstrap test rejects if  $|T_n(\theta_0)| > q_n^*(\alpha)$ .

Let

$$C_4 = [\hat{\theta} - s(\hat{\theta})q_n^*(\alpha), \quad \hat{\theta} + s(\hat{\theta})q_n^*(\alpha)],$$

where  $q_n^*(\alpha)$  is the bootstrap critical value for a two-sided hypothesis test.  $C_4$  is called the symmetric percentile-t interval. It is designed to work well since

$$\begin{aligned}\Pr(\theta_0 \in C_4) &= \Pr(\hat{\theta} - s(\hat{\theta})q_n^*(\alpha) \leq \theta_0 \leq \hat{\theta} + s(\hat{\theta})q_n^*(\alpha)) \\ &= \Pr(|T_n(\theta_0)| < q_n^*(\alpha)) \\ &\simeq \Pr(|T_n(\theta_0)| < q_n(\alpha)) \\ &= 1 - \alpha.\end{aligned}$$

If  $\theta$  is a vector, then to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta \neq \theta_0$  at size  $\alpha$ , we would use a Wald statistic

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

or some other asymptotically chi-square statistic. Thus here  $T_n(\theta) = W_n(\theta)$ . The ideal test rejects if  $W_n \geq q_n(\alpha)$ , where  $q_n(\alpha)$  is the  $(1 - \alpha)\%$  quantile of the distribution of  $W_n$ . The bootstrap test rejects if  $W_n \geq q_n^*(\alpha)$ , where  $q_n^*(\alpha)$  is the  $(1 - \alpha)\%$  quantile of the distribution of

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

Computationally, the critical value  $q_n^*(\alpha)$  is found as the quantile from simulated values of  $W_n^*$ . Note in the simulation that the Wald statistic is a quadratic form in  $(\hat{\theta}^* - \hat{\theta})$ , not  $(\hat{\theta}^* - \theta_0)$ . [This is a typical mistake made by practitioners.]

## 10.8 Asymptotic Expansions

Let  $T_n \in \mathbb{R}$  be a statistic such that

$$T_n \xrightarrow{d} N(0, \sigma^2). \quad (10.3)$$

In some cases, such as when  $T_n$  is a t-ratio, then  $\sigma^2 = 1$ . In other cases  $\sigma^2$  is unknown. Equivalently, writing  $T_n \sim G_n(u, F)$  then for each  $u$  and  $F$

$$\lim_{n \rightarrow \infty} G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right),$$

or

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + o(1). \quad (10.4)$$

While (10.4) says that  $G_n$  converges to  $\Phi\left(\frac{u}{\sigma}\right)$  as  $n \rightarrow \infty$ , it says nothing, however, about the rate of convergence, or the size of the divergence for any particular sample size  $n$ . A better asymptotic approximation may be obtained through an *asymptotic expansion*.

The following notation will be helpful. Let  $a_n$  be a sequence.

**Definition 10.8.1**  $a_n = o(1)$  if  $a_n \rightarrow 0$  as  $n \rightarrow \infty$

**Definition 10.8.2**  $a_n = O(1)$  if  $|a_n|$  is uniformly bounded.

**Definition 10.8.3**  $a_n = o(n^{-r})$  if  $n^r |a_n| \rightarrow 0$  as  $n \rightarrow \infty$ .

Basically,  $a_n = O(n^{-r})$  if it declines to zero like  $n^{-r}$ .

We say that a function  $g(u)$  is *even* if  $g(-u) = g(u)$ , and a function  $h(u)$  is *odd* if  $h(-u) = -h(u)$ . The derivative of an even function is odd, and vice-versa.

**Theorem 10.8.1** Under regularity conditions and (10.3),

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + \frac{1}{n^{1/2}}g_1(u, F) + \frac{1}{n}g_2(u, F) + O(n^{-3/2})$$

uniformly over  $u$ , where  $g_1$  is an even function of  $u$ , and  $g_2$  is an odd function of  $u$ . Moreover,  $g_1$  and  $g_2$  are differentiable functions of  $u$  and continuous in  $F$  relative to the supremum norm on the space of distribution functions.

The expansion in Theorem 10.8.1 is often called an **Edgeworth expansion**.

We can interpret Theorem 10.8.1 as follows. First,  $G_n(u, F)$  converges to the normal limit at rate  $n^{1/2}$ . To a second order of approximation,

$$G_n(u, F) \approx \Phi\left(\frac{u}{\sigma}\right) + n^{-1/2}g_1(u, F).$$

Since the derivative of  $g_1$  is odd, the density function is skewed. To a third order of approximation,

$$G_n(u, F) \approx \Phi\left(\frac{u}{\sigma}\right) + n^{-1/2}g_1(u, F) + n^{-1}g_2(u, F)$$

which adds a symmetric non-normal component to the approximate density (for example, adding leptokurtosis).



[Side Note: When  $T_n = \sqrt{n}(\bar{X}_n - \mu)/\sigma$ , a standardized sample mean, then

$$\begin{aligned} g_1(u) &= -\frac{1}{6}\kappa_3(u^2 - 1)\phi(u) \\ g_2(u) &= -\left(\frac{1}{24}\kappa_4(u^3 - 3u) + \frac{1}{72}\kappa_3^2(u^5 - 10u^3 + 15u)\right)\phi(u) \end{aligned}$$

where  $\phi(u)$  is the standard normal pdf, and

$$\begin{aligned} \kappa_3 &= \mathbb{E}(X - \mu)^3 / \sigma^3 \\ \kappa_4 &= \mathbb{E}(X - \mu)^4 / \sigma^4 - 3 \end{aligned}$$

the standardized skewness and excess kurtosis of the distribution of  $X$ . Note that when  $\kappa_3 = 0$  and  $\kappa_4 = 0$ , then  $g_1 = 0$  and  $g_2 = 0$ , so the second-order Edgeworth expansion corresponds to the normal distribution.]

### Francis Edgeworth

Francis Ysidro Edgeworth (1845-1926) of Ireland, founding editor of the *Economic Journal*, was a profound economic and statistical theorist, developing the theories of indifference curves and asymptotic expansions. He also could be viewed as the first econometrician due to his early use of mathematical statistics in the study of economic data.

## 10.9 One-Sided Tests

Using the expansion of Theorem 10.8.1, we can assess the accuracy of one-sided hypothesis tests and confidence regions based on an asymptotically normal t-ratio  $T_n$ . An asymptotic test is based on  $\Phi(u)$ .

To the second order, the exact distribution is

$$\Pr(T_n < u) = G_n(u, F_0) = \Phi(u) + \frac{1}{n^{1/2}}g_1(u, F_0) + O(n^{-1})$$

since  $\sigma = 1$ . The difference is

$$\begin{aligned} \Phi(u) - G_n(u, F_0) &= \frac{1}{n^{1/2}}g_1(u, F_0) + O(n^{-1}) \\ &= O(n^{-1/2}), \end{aligned}$$

so the order of the error is  $O(n^{-1/2})$ .

A bootstrap test is based on  $G_n^*(u)$ , which from Theorem 10.8.1 has the expansion

$$G_n^*(u) = G_n(u, F_n) = \Phi(u) + \frac{1}{n^{1/2}}g_1(u, F_n) + O(n^{-1}).$$

Because  $\Phi(u)$  appears in both expansions, the difference between the bootstrap distribution and the true distribution is

$$G_n^*(u) - G_n(u, F_0) = \frac{1}{n^{1/2}}(g_1(u, F_n) - g_1(u, F_0)) + O(n^{-1}).$$

Since  $F_n$  converges to  $F$  at rate  $\sqrt{n}$ , and  $g_1$  is continuous with respect to  $F$ , the difference  $(g_1(u, F_n) - g_1(u, F_0))$  converges to 0 at rate  $\sqrt{n}$ . Heuristically,

$$\begin{aligned} g_1(u, F_n) - g_1(u, F_0) &\approx \frac{\partial}{\partial F} g_1(u, F_0) (F_n - F_0) \\ &= O(n^{-1/2}), \end{aligned}$$

The “derivative”  $\frac{\partial}{\partial F} g_1(u, F)$  is only heuristic, as  $F$  is a function. We conclude that

$$G_n^*(u) - G_n(u, F_0) = O(n^{-1}),$$

or

$$\Pr(T_n^* \leq u) = \Pr(T_n \leq u) + O(n^{-1}),$$

which is an improved rate of convergence over the asymptotic test (which converged at rate  $O(n^{-1/2})$ ). This rate can be used to show that one-tailed bootstrap inference based on the t-ratio achieves a so-called *asymptotic refinement* – the Type I error of the test converges at a faster rate than an analogous asymptotic test.

## 10.10 Symmetric Two-Sided Tests

If a random variable  $y$  has distribution function  $H(u) = \Pr(y \leq u)$ , then the random variable  $|y|$  has distribution function

$$\overline{H}(u) = H(u) - H(-u)$$

since

$$\begin{aligned} \Pr(|y| \leq u) &= \Pr(-u \leq y \leq u) \\ &= \Pr(y \leq u) - \Pr(y \leq -u) \\ &= H(u) - H(-u). \end{aligned}$$

For example, if  $Z \sim N(0, 1)$ , then  $|Z|$  has distribution function

$$\overline{\Phi}(u) = \Phi(u) - \Phi(-u) = 2\Phi(u) - 1.$$

Similarly, if  $T_n$  has exact distribution  $G_n(u, F)$ , then  $|T_n|$  has the distribution function

$$\overline{G}_n(u, F) = G_n(u, F) - G_n(-u, F).$$

A two-sided hypothesis test rejects  $\mathbb{H}_0$  for large values of  $|T_n|$ . Since  $T_n \xrightarrow{d} Z$ , then  $|T_n| \xrightarrow{d} |Z| \sim \overline{\Phi}$ . Thus asymptotic critical values are taken from the  $\overline{\Phi}$  distribution, and exact critical values are taken from the  $\overline{G}_n(u, F_0)$  distribution. From Theorem 10.8.1, we can calculate that

$$\begin{aligned} \overline{G}_n(u, F) &= G_n(u, F) - G_n(-u, F) \\ &= \left( \Phi(u) + \frac{1}{n^{1/2}} g_1(u, F) + \frac{1}{n} g_2(u, F) \right) \\ &\quad - \left( \Phi(-u) + \frac{1}{n^{1/2}} g_1(-u, F) + \frac{1}{n} g_2(-u, F) \right) + O(n^{-3/2}) \\ &= \overline{\Phi}(u) + \frac{2}{n} g_2(u, F) + O(n^{-3/2}), \end{aligned} \tag{10.5}$$

where the simplifications are because  $g_1$  is even and  $g_2$  is odd. Hence the difference between the asymptotic distribution and the exact distribution is

$$\overline{\Phi}(u) - \overline{G}_n(u, F_0) = \frac{2}{n} g_2(u, F_0) + O(n^{-3/2}) = O(n^{-1}).$$

The order of the error is  $O(n^{-1})$ .

Interestingly, the asymptotic two-sided test has a better coverage rate than the asymptotic one-sided test. This is because the first term in the asymptotic expansion,  $g_1$ , is an even function, meaning that the errors in the two directions exactly cancel out.

Applying (10.5) to the bootstrap distribution, we find

$$\bar{G}_n^*(u) = \bar{G}_n(u, F_n) = \bar{\Phi}(u) + \frac{2}{n}g_2(u, F_n) + O(n^{-3/2}).$$

Thus the difference between the bootstrap and exact distributions is

$$\begin{aligned} \bar{G}_n^*(u) - \bar{G}_n(u, F_0) &= \frac{2}{n}(g_2(u, F_n) - g_2(u, F_0)) + O(n^{-3/2}) \\ &= O(n^{-3/2}), \end{aligned}$$

the last equality because  $F_n$  converges to  $F_0$  at rate  $\sqrt{n}$ , and  $g_2$  is continuous in  $F$ . Another way of writing this is

$$\Pr(|T_n^*| < u) = \Pr(|T_n| < u) + O(n^{-3/2})$$

so the error from using the bootstrap distribution (relative to the true unknown distribution) is  $O(n^{-3/2})$ . This is in contrast to the use of the asymptotic distribution, whose error is  $O(n^{-1})$ . Thus a two-sided bootstrap test also achieves an asymptotic refinement, similar to a one-sided test.

A reader might get confused between the two simultaneous effects. Two-sided tests have better rates of convergence than the one-sided tests, and bootstrap tests have better rates of convergence than asymptotic tests.

The analysis shows that there may be a trade-off between one-sided and two-sided tests. Two-sided tests will have more accurate size (Reported Type I error), but one-sided tests might have more power against alternatives of interest. Confidence intervals based on the bootstrap can be asymmetric if based on one-sided tests (equal-tailed intervals) and can therefore be more informative and have smaller length than symmetric intervals. Therefore, the choice between symmetric and equal-tailed confidence intervals is unclear, and needs to be determined on a case-by-case basis.

## 10.11 Percentile Confidence Intervals

To evaluate the coverage rate of the percentile interval, set  $T_n = \sqrt{n}(\hat{\theta} - \theta_0)$ . We know that  $T_n \xrightarrow{d} N(0, V)$ , which is not pivotal, as it depends on the unknown  $V$ . Theorem 10.8.1 shows that a first-order approximation

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + O(n^{-1/2}),$$

where  $\sigma = \sqrt{V}$ , and for the bootstrap

$$G_n^*(u) = G_n(u, F_n) = \Phi\left(\frac{u}{\hat{\sigma}}\right) + O(n^{-1/2}),$$

where  $\hat{\sigma} = \sqrt{V(F_n)}$  is the bootstrap estimate of  $\sigma$ . The difference is

$$\begin{aligned} G_n^*(u) - G_n(u, F_0) &= \Phi\left(\frac{u}{\hat{\sigma}}\right) - \Phi\left(\frac{u}{\sigma}\right) + O(n^{-1/2}) \\ &= -\phi\left(\frac{u}{\sigma}\right) \frac{u}{\sigma} (\hat{\sigma} - \sigma) + O(n^{-1/2}) \\ &= O(n^{-1/2}) \end{aligned}$$

Hence the order of the error is  $O(n^{-1/2})$ .

The good news is that the percentile-type methods (if appropriately used) can yield  $\sqrt{n}$ -convergent asymptotic inference. Yet these methods do not require the calculation of standard

errors! This means that in contexts where standard errors are not available or are difficult to calculate, the percentile bootstrap methods provide an attractive inference method.

The bad news is that the rate of convergence is disappointing. It is no better than the rate obtained from an asymptotic one-sided confidence region. Therefore if standard errors are available, it is unclear if there are any benefits from using the percentile bootstrap over simple asymptotic methods.

Based on these arguments, the theoretical literature (e.g. Hall, 1992, Horowitz, 2001) tends to advocate the use of the percentile-t bootstrap methods rather than percentile methods.

## 10.12 Bootstrap Methods for Regression Models

The bootstrap methods we have discussed have set  $G_n^*(u) = G_n(u, F_n)$ , where  $F_n$  is the EDF. Any other consistent estimate of  $F$  may be used to define a feasible bootstrap estimator. The advantage of the EDF is that it is fully nonparametric, it imposes no conditions, and works in nearly any context. But since it is fully nonparametric, it may be inefficient in contexts where more is known about  $F$ . We discuss bootstrap methods appropriate for the linear regression model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i | \mathbf{x}_i) &= 0. \end{aligned}$$

The non-parametric bootstrap resamples the observations  $(y_i^*, \mathbf{x}_i^*)$  from the EDF, which implies

$$\begin{aligned} y_i^* &= \mathbf{x}_i^{*'} \hat{\boldsymbol{\beta}} + e_i^* \\ \mathbb{E}(\mathbf{x}_i^* e_i^*) &= \mathbf{0} \end{aligned}$$

but generally

$$\mathbb{E}(e_i^* | \mathbf{x}_i^*) \neq 0.$$

The bootstrap distribution does not impose the regression assumption, and is thus an inefficient estimator of the true distribution (when in fact the regression assumption is true.)

One approach to this problem is to impose the very strong assumption that the error  $\varepsilon_i$  is independent of the regressor  $\mathbf{x}_i$ . The advantage is that in this case it is straightforward to construct bootstrap distributions. The disadvantage is that the bootstrap distribution may be a poor approximation when the error is not independent of the regressors.

To impose independence, it is sufficient to sample the  $\mathbf{x}_i^*$  and  $e_i^*$  independently, and then create  $y_i^* = \mathbf{x}_i^{*'} \hat{\boldsymbol{\beta}} + e_i^*$ . There are different ways to impose independence. A non-parametric method is to sample the bootstrap errors  $e_i^*$  randomly from the OLS residuals  $\{\hat{e}_1, \dots, \hat{e}_n\}$ . A parametric method is to generate the bootstrap errors  $e_i^*$  from a parametric distribution, such as the normal  $e_i^* \sim N(0, \hat{\sigma}^2)$ .

For the regressors  $\mathbf{x}_i^*$ , a nonparametric method is to sample the  $\mathbf{x}_i^*$  randomly from the EDF or sample values  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . A parametric method is to sample  $\mathbf{x}_i^*$  from an estimated parametric distribution. A third approach sets  $\mathbf{x}_i^* = \mathbf{x}_i$ . This is equivalent to treating the regressors as *fixed in repeated samples*. If this is done, then all inferential statements are made conditionally on the observed values of the regressors, which is a valid statistical approach. It does not really matter, however, whether or not the  $\mathbf{x}_i$  are really “fixed” or random.

The methods discussed above are unattractive for most applications in econometrics because they impose the stringent assumption that  $\mathbf{x}_i$  and  $e_i$  are independent. Typically what is desirable is to impose only the regression condition  $\mathbb{E}(e_i | \mathbf{x}_i) = 0$ . Unfortunately this is a harder problem.

One proposal which imposes the regression condition without independence is the *Wild Bootstrap*. The idea is to construct a conditional distribution for  $e_i^*$  so that

$$\begin{aligned} \mathbb{E}(e_i^* | \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^{*2} | \mathbf{x}_i) &= \hat{e}_i^2 \\ \mathbb{E}(e_i^{*3} | \mathbf{x}_i) &= \hat{e}_i^3. \end{aligned}$$

A conditional distribution with these features will preserve the main important features of the data. This can be achieved using a two-point distribution of the form

$$\begin{aligned}\Pr\left(e_i^* = \left(\frac{1 + \sqrt{5}}{2}\right) \hat{e}_i\right) &= \frac{\sqrt{5} - 1}{2\sqrt{5}} \\ \Pr\left(e_i^* = \left(\frac{1 - \sqrt{5}}{2}\right) \hat{e}_i\right) &= \frac{\sqrt{5} + 1}{2\sqrt{5}}\end{aligned}$$

For each  $\mathbf{x}_i$ , you sample  $e_i^*$  using this two-point distribution.

## Exercises

**Exercise 10.1** Let  $F_n(\mathbf{x})$  denote the EDF of a random sample. Show that

$$\sqrt{n}(F_n(\mathbf{x}) - F_0(\mathbf{x})) \xrightarrow{d} N(0, F_0(\mathbf{x})(1 - F_0(\mathbf{x}))).$$

**Exercise 10.2** Take a random sample  $\{y_1, \dots, y_n\}$  with  $\mu = \mathbb{E}y_i$  and  $\sigma^2 = \text{var}(y_i)$ . Let the statistic of interest be the sample mean  $T_n = \bar{y}_n$ . Find the population moments  $\mathbb{E}T_n$  and  $\text{var}(T_n)$ . Let  $\{y_1^*, \dots, y_n^*\}$  be a random sample from the empirical distribution function and let  $T_n^* = \bar{y}_n^*$  be its sample mean. Find the bootstrap moments  $\mathbb{E}T_n^*$  and  $\text{var}(T_n^*)$ .

**Exercise 10.3** Consider the following bootstrap procedure for a regression of  $y_i$  on  $\mathbf{x}_i$ . Let  $\hat{\beta}$  denote the OLS estimator from the regression of  $\mathbf{y}$  on  $\mathbf{X}$ , and  $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\beta}$  the OLS residuals.

- Draw a random vector  $(\mathbf{x}^*, e^*)$  from the pair  $\{(\mathbf{x}_i, \hat{e}_i) : i = 1, \dots, n\}$ . That is, draw a random integer  $i'$  from  $[1, 2, \dots, n]$ , and set  $\mathbf{x}^* = \mathbf{x}_{i'}$  and  $e^* = \hat{e}_{i'}$ . Set  $y^* = \mathbf{x}^{*\prime}\hat{\beta} + e^*$ . Draw (with replacement)  $n$  such vectors, creating a random bootstrap data set  $(\mathbf{y}^*, \mathbf{X}^*)$ .
- Regress  $\mathbf{y}^*$  on  $\mathbf{X}^*$ , yielding OLS estimates  $\hat{\beta}^*$  and any other statistic of interest.

Show that this bootstrap procedure is (numerically) *identical* to the non-parametric bootstrap.

**Exercise 10.4** Consider the following bootstrap procedure. Using the non-parametric bootstrap, generate bootstrap samples, calculate the estimate  $\hat{\theta}^*$  on these samples and then calculate

$$T_n^* = (\hat{\theta}^* - \hat{\theta})/s(\hat{\theta}),$$

where  $s(\hat{\theta})$  is the standard error in the original data. Let  $q_n^*(.05)$  and  $q_n^*(.95)$  denote the 5% and 95% quantiles of  $T_n^*$ , and define the bootstrap confidence interval

$$C = [\hat{\theta} - s(\hat{\theta})q_n^*(.95), \quad \hat{\theta} - s(\hat{\theta})q_n^*(.05)].$$

Show that  $C$  exactly equals the Alternative percentile interval (not the percentile-t interval).

**Exercise 10.5** You want to test  $\mathbb{H}_0 : \theta = 0$  against  $\mathbb{H}_1 : \theta > 0$ . The test for  $\mathbb{H}_0$  is to reject if  $T_n = \hat{\theta}/s(\hat{\theta}) > c$  where  $c$  is picked so that Type I error is  $\alpha$ . You do this as follows. Using the non-parametric bootstrap, you generate bootstrap samples, calculate the estimates  $\hat{\theta}^*$  on these samples and then calculate

$$T_n^* = \hat{\theta}^*/s(\hat{\theta}^*).$$

Let  $q_n^*(.95)$  denote the 95% quantile of  $T_n^*$ . You replace  $c$  with  $q_n^*(.95)$ , and thus reject  $\mathbb{H}_0$  if  $T_n = \hat{\theta}/s(\hat{\theta}) > q_n^*(.95)$ . What is wrong with this procedure?

**Exercise 10.6** Suppose that in an application,  $\hat{\theta} = 1.2$  and  $s(\hat{\theta}) = .2$ . Using the non-parametric bootstrap, 1000 samples are generated from the bootstrap distribution, and  $\hat{\theta}^*$  is calculated on each sample. The  $\hat{\theta}^*$  are sorted, and the 2.5% and 97.5% quantiles of the  $\hat{\theta}^*$  are .75 and 1.3, respectively.

- Report the 95% Efron Percentile interval for  $\theta$ .
- Report the 95% Alternative Percentile interval for  $\theta$ .
- With the given information, can you report the 95% Percentile-t interval for  $\theta$ ?

**Exercise 10.7** The datafile `hprice1.dat` contains data on house prices (sales), with variables listed in the file `hprice1.pdf`. Estimate a linear regression of price on the number of bedrooms, lot size, size of house, and the colonial dummy. Calculate 95% confidence intervals for the regression coefficients using both the asymptotic normal approximation and the percentile-t bootstrap.