

The previous section has described how to obtain parametric and non-parametric bootstrap samples in general, and illustrated the particular case of linear regression. In this section we assume that we have obtained  $B = 999$  bootstrap samples of  $\theta$ , the parameter of interest, and that we have sorted them into order. Let

$$\hat{\theta}_1^*, \dots, \hat{\theta}_{999}^* \tag{1}$$

denote this ordered set, so that  $\hat{\theta}_i^* < \hat{\theta}_j^*$ , for  $1 \leq i < j \leq 999$ . Of course, in the linear regression example the parameter of interest is  $\theta = \beta$ , the slope.

Asymptotic  
coverage error:

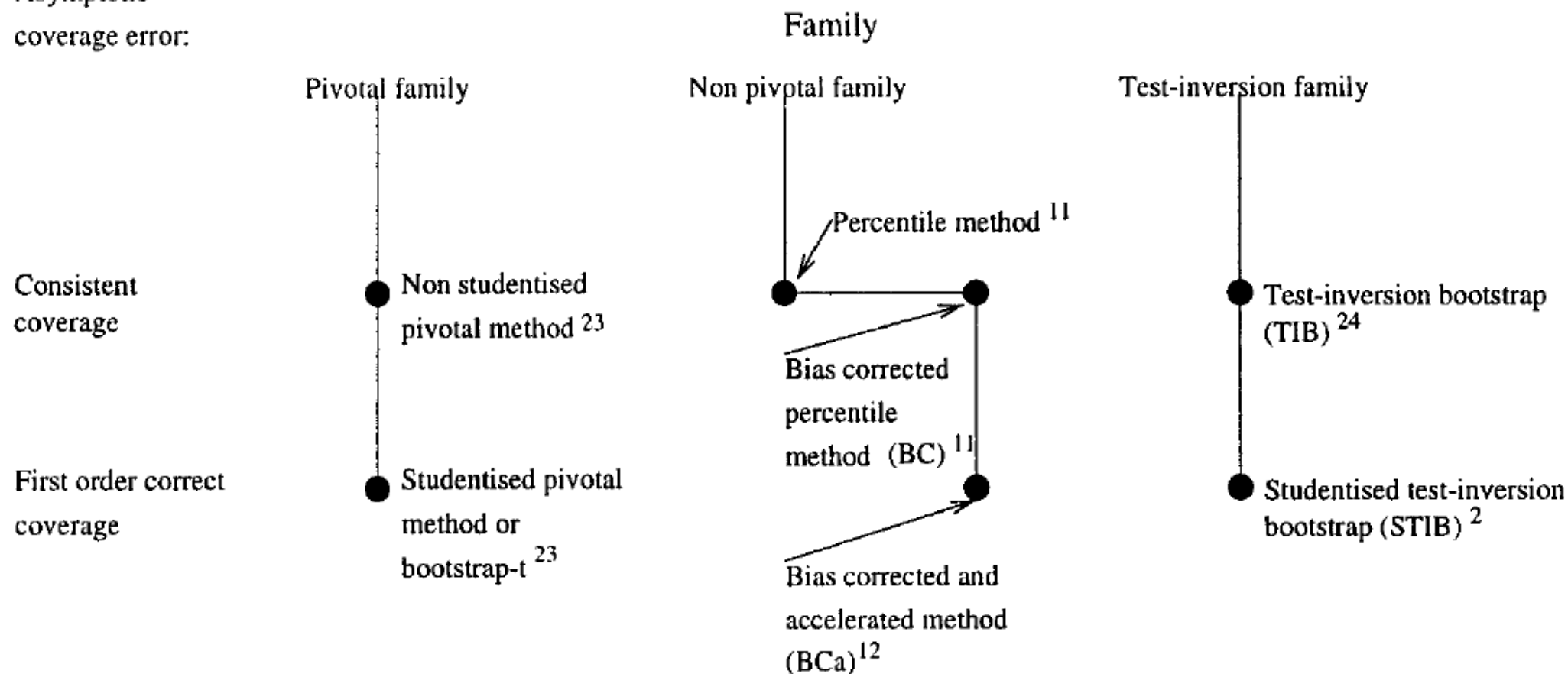


Figure 3. Schematic diagram describing relationship of bootstrap methods for constructing confidence intervals, with early references. *Consistent* coverage means the coverage error is  $O(n^{-1/2})$ ; *first-order correct* coverage means the coverage error is  $O(n^{-1})$ . See the introduction and Table III for more details.

### 3.3. Percentile Method

We now move to the second arm of Figure 3. These methods are in some ways less intuitive than those described above, but have the advantage of not requiring  $\hat{\sigma}$ .

*3.3.1. Rationale.* Consider a monotonically increasing function  $g(\cdot)$ , and write  $\phi = g(\theta)$ ,  $\hat{\phi} = g(\hat{\theta})$  and  $\hat{\phi}^* = g(\hat{\theta}^*)$ . Choose (if possible)  $g(\cdot)$ , such that

$$\hat{\phi}^* - \hat{\phi} \sim \hat{\phi} - \phi \sim N(0, \sigma^2) \quad (5)$$

Then, since  $\hat{\phi} - \phi \sim N(0, \sigma^2)$ , the interval for  $\theta$  is

$$(-\infty, g^{-1}(\hat{\phi} - \sigma z_\alpha)) \quad (6)$$

where  $z_\alpha$  is the 100 $\alpha$  per cent point of the standard normal distribution. However, (5) implies that  $\hat{\phi} - \sigma z_\alpha = F_{\hat{\phi}^*}^{-1}(1 - \alpha)$ . Further, since  $g$  is monotonically increasing,  $F_{\hat{\phi}^*}^{-1}(1 - \alpha) = g(F_{\hat{\theta}^*}^{-1}(1 - \alpha))$ . Substituting in (6) gives the percentile interval

$$(-\infty, F_{\hat{\theta}^*}^{-1}(1 - \alpha)) \quad (7)$$

3.3.3. *Advantages.* Simplicity is the attraction of this method, and explains its continued popularity. Unlike the bootstrap- $t$ , no estimates of the  $\sigma$  are required. Further, no invalid parameter values can be included in the interval.

Another advantage of this group of methods over the pivotal methods is that they are transformation respecting.

3.3.4. *Disadvantages.* The coverage error is often substantial if the distribution of  $\hat{\theta}$  is not nearly symmetric (Efron and Tibshirani, Reference [6], p. 178 ff). The reason is that the justification of the method rests on the existence of a  $g(\cdot)$  such that (5) holds, and for many problems such a  $g$  does not exist.

3.3.5. *Coverage error* [9].

$$\mathbf{P}(\hat{\theta}_{1-\alpha}^* > \theta) = 1 - \alpha + O(n^{-1/2})$$

### 3.4. Bias corrected method

The quickly recognized shortcomings of the percentile method [11] led to the development of the bias corrected or BC method.

*3.4.1. Rationale.* Again, consider a monotonically increasing function  $g(\cdot)$ , and write  $\phi = g(\theta)$ ,  $\hat{\phi} = g(\hat{\theta})$  and  $\hat{\phi}^* = g(\hat{\theta}^*)$ . However, now (if possible) choose  $g(\cdot)$ , such that

$$\hat{\phi}^* - \hat{\phi} \sim \hat{\phi} - \phi \sim N(-b\sigma, \sigma^2) \quad (8)$$

for some constant  $b$ . An analogous (but slightly more complex) argument than that used in the case of the percentile interval then yields the BC interval

$$(-\infty, F_{\hat{\theta}^*}^{-1}(\Phi(2b - z_\alpha))) \quad (9)$$

where  $b$  is estimated by  $\Phi^{-1}(\mathbf{P}(\hat{\theta}^* \leq \hat{\theta}))$  and  $\Phi^{-1}$  is the inverse cumulative distribution function of the normal distribution.

### 3.4.2. Calculation of 95 per cent interval.

1. Count the number of members of (1) that are less than  $\hat{\theta}$  (calculated from the original data). Call this number  $p$  and set  $b = \Phi^{-1}(p/B)$ .
2. Calculate  $Q = (B + 1)\Phi(2b - z_{0.05})$ , where  $z_{0.05} = -1.64$ .  $Q$  is the percentile of the bootstrap distribution required for the upper endpoint of the bias corrected confidence interval.
3. Estimate the endpoint of the interval by  $\hat{\theta}_{[Q]}^*$ , where  $[.]$  means ‘take the integer part’. If a more accurate estimate is required, interpolation can be used between the members of (1), as follows. Let the nearest integers to  $Q$  be  $a, b$ , so that  $a < Q < b$  and  $b = a + 1$ . Then the  $Q$ th percentile is estimated by

$$\hat{\theta}_Q^* \approx \hat{\theta}_a^* + \frac{\Phi^{-1}(\frac{Q}{B+1}) - \Phi^{-1}(\frac{a}{B+1})}{\Phi^{-1}(\frac{b}{B+1}) - \Phi^{-1}(\frac{a}{B+1})} (\hat{\theta}_b^* - \hat{\theta}_a^*) \quad (10)$$

The bias corrected interval, (7), is

$$(-\infty, \hat{\theta}_Q^*).$$



3.4.3. *Advantages.* The advantages are as for the percentile method, but see below.

3.4.4. *Disadvantages.* This method was devised as an improvement to the percentile method for non-symmetric problems. Hence, if the distribution of  $\hat{\theta}^*$  is symmetric about  $\hat{\theta}$ , then  $b = 0$  and the bias corrected and percentile intervals agree. However, the coverage error is still often substantial, because the validity of the method depends upon the existence of a  $g(\cdot)$  such that (8) holds, and for many problems such a  $g$  does not exist. In consequence, it has been omitted altogether from recent discussions [5, 6]. It is worth mentioning, though, as it is still the most accurate method implemented in the software package stata.

3.4.5. *Coverage error* [9].

$$\mathbf{P}(\hat{\theta}_Q^* > \theta) = 1 - \alpha + O(n^{-1/2})$$

### 3.5. Bias corrected and accelerated method

The shortcomings of the BC method in turn led [12] to the development of the bias corrected and accelerated or BCa method. The idea is to allow not only for the lack of symmetry of  $F_{\hat{\Theta}}(\cdot; \theta)$ , but also for the fact that its shape, or skewness, might change as  $\theta$  varies.

Note that the later abc method [13] is an analytic approximation to this method.

*3.5.1. Rationale.* Again, consider a monotonically increasing function  $g(\cdot)$ , and write  $\phi = g(\theta)$ ,  $\hat{\phi} = g(\hat{\theta})$  and  $\hat{\phi}^* = g(\hat{\theta}^*)$ . However, now (if possible) choose  $g(\cdot)$ , such that

$$\hat{\phi} \sim N(\phi - b\sigma(\phi), \sigma^2(\phi))$$

$$\hat{\phi}^* \sim N(\hat{\phi} - b\sigma(\hat{\phi}), \sigma^2(\hat{\phi}))$$

where  $\sigma(x) = 1 + ax$ . Again, an analogous argument to that used to justify the BC interval yields the BCa interval

$$\left( -\infty, F_{\hat{\Theta}^*}^{-1} \left( \Phi \left( b - \frac{z_{\alpha} - b}{1 + a(z_{\alpha} - b)} \right); \hat{\theta} \right) \right) \quad (11)$$

where  $b$  is defined as before and a formula for estimating  $a$  is given below.



### 3.5.2. Calculation of 95 per cent interval.

1. Calculate  $b$  as for the BC interval.
2. Next we need to calculate  $a$ . This calculation depends on whether the simulation is non-parametric or parametric, and in the latter case, whether nuisance parameters are present. For completeness we give a simple jack-knife estimate of  $a$ ; details about more sophisticated and accurate estimates can be found elsewhere [5, 6]. Let  $\mathbf{y}_{\text{obs}}^i$  represent the original data with the  $i$ th point omitted, and  $\hat{\theta}^i = \hat{\theta}(\mathbf{y}_{\text{obs}}^i)$  be the estimate of  $\theta$  constructed from this data. Let  $\tilde{\theta}$  be the mean of the  $\hat{\theta}^i$ 's. Then  $a$  is estimated by

$$\frac{\sum_{i=1}^n (\tilde{\theta} - \hat{\theta}^i)^3}{6 [\sum_{i=1}^n (\tilde{\theta} - \hat{\theta}^i)^2]^{3/2}}$$

3. Let  $\tilde{Q}$  be the integer part of  $(B + 1) \Phi(b - \frac{z_{0.05} - b}{1 + a(z_{0.05} - b)})$ , where  $z_{0.05} = -1.64$ .
4. Estimate the  $\tilde{Q}$ th percentile of the bootstrap distribution (1) as in step 3 of the bias corrected interval calculation. Then, the BCa interval, (11), is estimated by

$$(-\infty, \hat{\theta}_{\tilde{Q}}^*)$$

3.5.3. *Advantages.* The advantages are as for the percentile method, plus this method generally has a smaller coverage error than the percentile and BC intervals (Efron and Tibshirani, Reference [6], p. 184 ff), but see below.

3.5.4. *Disadvantages.* The calculation of  $a$  can be tortuous in complex parametric problems. The coverage error of this method increases as  $\alpha \rightarrow 0$ . To see why this is so, note that as this happens, the right hand endpoint of the interval should be estimated by ever larger elements of the set of ordered  $\hat{\theta}^*$ 's. However, this is not the case: as  $\alpha \rightarrow 0$

$$\Phi\left(b - \frac{z_\alpha - b}{1 + a(z_\alpha - b)}\right) \rightarrow \Phi(b - 1/a) \neq 1$$

This anomaly means that coverage can be erratic for small  $\alpha$ , typically  $\alpha < 0.025$ . (Davison and Hinkley, Reference [5], p. 205, p. 231 and Section 5 below).

3.5.5. *Coverage error* [9].

$$\mathbf{P}(\hat{\theta}_{\hat{Q}}^* > \theta) = 1 - \alpha + O(n^{-1})$$

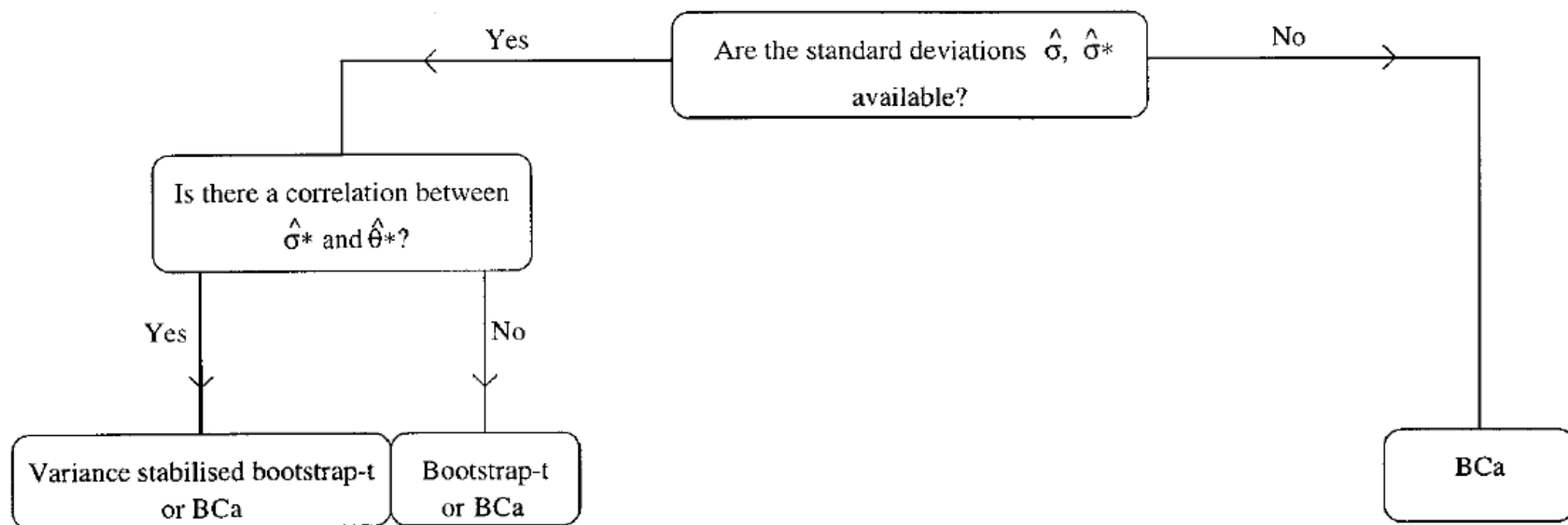


Figure 4. Guide to choosing a bootstrap confidence interval method when using non-parametric simulation.

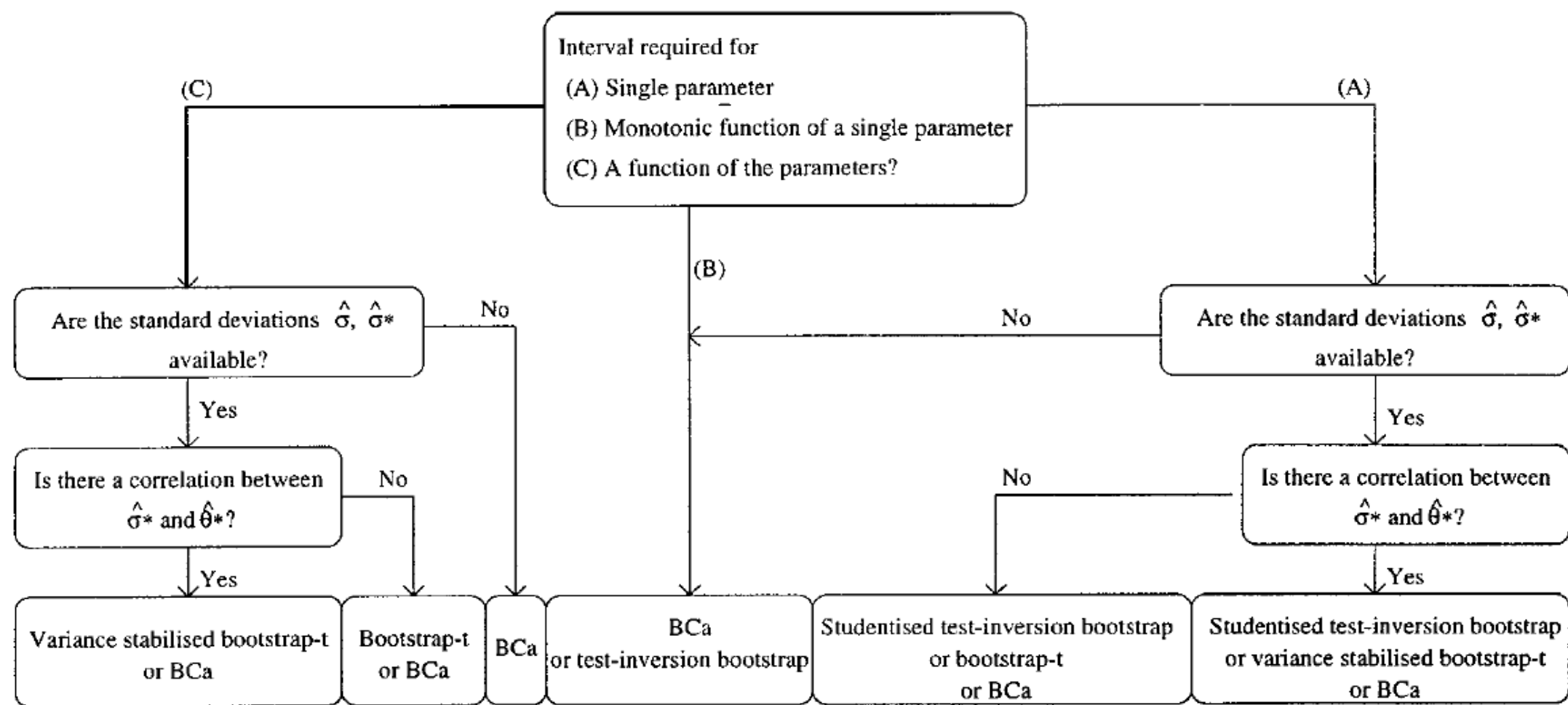


Figure 5. Guide to choosing a bootstrap confidence interval method when using parametric simulation, or simulating in a regression framework, as described in Section 2.3.