

# Fast Double Bootstrap Tests of Nonnested Linear Regression Models

by

**Russell Davidson**

GREQAM  
Centre de la Vieille Charité  
2 rue de la Charité  
13002 Marseille, France

Department of Economics  
Queen's University  
Kingston, Ontario, Canada  
K7L 3N6

**russell@ehess.cnrs-mrs.fr**

and

**James G. MacKinnon**

Department of Economics  
Queen's University  
Kingston, Ontario, Canada  
K7L 3N6

**jgm@qed.econ.queensu.ca**

## Abstract

It has been shown in previous work that bootstrapping the  $J$  test for nonnested linear regression models dramatically improves its finite-sample performance. We provide evidence that a more sophisticated bootstrap procedure, which we call the fast double bootstrap, produces a very substantial further improvement in cases where the ordinary bootstrap does not work as well as it might. This FDB procedure is only about twice as expensive as the usual single bootstrap.

This research was supported, in part, by grants from the Social Sciences and Humanities Research Council of Canada.

August, 2001

## 1. Introduction

The  $J$  test proposed by Davidson and MacKinnon (1981) is the most widely-used procedure for testing nonnested regression models; see McAleer (1995). Its popularity stems from the fact that it is conceptually simple and easy to compute. However, its finite-sample distribution can be very far from the standard normal distribution that it follows asymptotically. As a consequence, it often overrejects severely. A natural way to improve the finite-sample properties of the test is to bootstrap it, as Fan and Li (1995) and Godfrey (1998) were among the first to suggest.

In Davidson and MacKinnon (2002), we developed a theoretical approach which enabled us to show precisely what determines the finite-sample distribution of the  $J$  test. By using our theoretical results to design simulation experiments, we showed that, in most cases, the  $J$  test will perform very reliably in finite samples if it is bootstrapped. However, there can still be some cases in which the bootstrapped  $J$  test rejects noticeably more often than it should.

In Davidson and MacKinnon (2001), we developed a simple and inexpensive technique which we called the “fast double bootstrap,” or FDB. Like the double bootstrap proposed by Beran (1988), it leads to a theoretical improvement in the performance of bootstrap tests. Unlike the double bootstrap, it requires only about twice as much computation as the ordinary, or single, bootstrap. Although the FDB is not as widely applicable as the double bootstrap, it can be applied to a broad range of econometric tests, including the  $J$  test. In this paper, we develop the FDB  $J$  test and demonstrate, by means of simulations, that it works extraordinarily well.

In the next section, we briefly describe the  $J$  test and discuss some standard ways in which it can be bootstrapped. In Section 3, we describe how the fast double bootstrap can be used to make the  $J$  test more reliable than the ordinary, or single, bootstrap  $J$  test. Finally, in Section 4, we present some simulation results on the performance of the single bootstrap and FDB  $J$  tests.

## 2. The $J$ Test

Consider two nonnested, linear regression models with IID errors:

$$\begin{aligned} H_1: \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, & \mathbf{u} &\sim \text{IID}(\mathbf{0}, \sigma_1^2 \mathbf{I}), \text{ and} \\ H_2: \mathbf{y} &= \mathbf{Z}\boldsymbol{\gamma} + \mathbf{v}, & \mathbf{v} &\sim \text{IID}(\mathbf{0}, \sigma_2^2 \mathbf{I}), \end{aligned}$$

where  $\mathbf{y}$ ,  $\mathbf{u}$ , and  $\mathbf{v}$  are  $n \times 1$  vectors,  $\mathbf{X}$  and  $\mathbf{Z}$  are, respectively,  $n \times k$  and  $n \times l$  matrices of regressors, and  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  are, respectively, a  $k$ -vector and an  $l$ -vector of unknown parameters. The  $J$  statistic for testing  $H_1$  is the ordinary  $t$  statistic for  $\alpha = 0$  in the regression

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \alpha \mathbf{P}_Z \mathbf{y} + \text{residuals}, \tag{1}$$

where  $\mathbf{P}_Z \equiv \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ , so that  $\mathbf{P}_Z \mathbf{y}$  is the vector of fitted values from OLS estimation of  $H_2$ . Asymptotically, the  $J$  statistic is distributed as  $N(0, 1)$  under  $H_1$ .

In practice, the  $t(n - k - 1)$  distribution is often used for finite-sample inference, although there is, in general, no formal justification for doing so.

The  $J$  statistic for testing  $H_1$  can be written as

$$\hat{J} = \frac{\mathbf{y}^\top \mathbf{P}_Z \mathbf{M}_X \mathbf{y}}{\hat{s}^2 (\mathbf{y}^\top \mathbf{P}_Z \mathbf{M}_X \mathbf{P}_Z \mathbf{y})^{1/2}}, \quad (2)$$

where  $\hat{s}$  is the usual estimated standard error from regression (1),  $\mathbf{P}_X$  is the projection matrix  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ , and  $\mathbf{M}_X \equiv \mathbf{I} - \mathbf{P}_X$ . Since  $\hat{J}$  depends on the regressor matrices only through the projections  $\mathbf{M}_X$ ,  $\mathbf{P}_X$ , and  $\mathbf{P}_Z$ , it is invariant to any changes in  $\mathbf{X}$  and  $\mathbf{Z}$  that do not change the subspaces spanned by the columns of these matrices.

In order to bootstrap the  $J$  test, we need to generate  $B$  bootstrap samples from a DGP that approximates what  $H_1$  would be if it had actually generated the data. The natural choice for  $\beta$  is the OLS estimator  $\hat{\beta}$ . Then the  $j^{\text{th}}$  bootstrap sample will be

$$\mathbf{y}_j^* = \mathbf{X} \hat{\beta} + \mathbf{u}_j^*, \quad (3)$$

where the elements of  $\mathbf{u}_j^*$  can be simulated in various ways. One possibility would be to draw them from the  $N(0, s^2)$  distribution, where  $s$  is the standard error of the vector  $\hat{\mathbf{u}} \equiv \mathbf{M}_X \mathbf{y}$  of residuals from OLS estimation of  $H_1$ . However, this approach is based on the assumption that the error terms are normally distributed, which may be uncomfortably strong. Since we are using the bootstrap, it is natural to generate  $\mathbf{u}^*$  by resampling from the vector

$$\tilde{\mathbf{u}} \equiv \left( \frac{n}{n - k} \right)^{1/2} \hat{\mathbf{u}}. \quad (4)$$

Here we have rescaled the residuals so that the average squared residual has expectation  $\sigma_1^2$ . In many applications of the bootstrap, this step is omitted. However, as we will see in Section 4, in the case of the  $J$  test it is very important to resample from the rescaled residuals (4) rather than from the ordinary residuals  $\hat{\mathbf{u}}$ . There are other, more complicated, ways to rescale the residuals, which take into account the different leverage of each observation, but we have not observed any advantage to using them.

In writing (3), we have implicitly assumed that all the regressors are exogenous. If there are lagged dependent variables, it will be necessary to generate the elements of the vector  $\mathbf{y}_j^*$  recursively, using the value of the dependent variable in observation 0 to begin the recursion. In this case, the regressor matrix  $\mathbf{X}$  will be different for every bootstrap sample, as will the regressor matrix  $\mathbf{Z}$  if it too includes a lagged dependent variable. Of course, when it is  $H_1$  that is being tested, the lagged dependent variable that appears in  $\mathbf{Z}$  for the bootstrap samples must be generated from  $H_1$ .

Ideally,  $B$  should be a reasonably large number, and it should be chosen so that  $\alpha(B + 1)$  is an integer for all levels  $\alpha$  of interest. One common choice is  $B = 999$ , although for a test as easy to compute as the  $J$  test, it might be reasonable to use an

even larger number, such as 4999 or even 9999; see Davidson and MacKinnon (2000) for a practical way to choose  $B$  endogenously so as to minimize simulation error.

For each bootstrap sample, a bootstrap test statistic is computed in exactly the same way as  $\hat{J}$  was computed from the original data. We denote the bootstrap statistics by  $J_j^*$ ,  $j = 1, \dots, B$ . Then we can compute a bootstrap  $P$  value by the formula

$$\hat{p}^*(\hat{J}) = \frac{1}{B} \sum_{j=1}^B I(J_j^* \geq \hat{J}), \quad (5)$$

where  $I(\cdot)$  is an indicator function, equal to 1 if its argument is true and equal to 0 otherwise. This assumes that the test is a one-tailed test which rejects in the upper tail, as is usually the case with the  $J$  test and other nonnested tests. For a two-tailed test, the indicator function in (5) would become  $I(|J_j^*| \geq |\hat{J}|)$ .

Since the statistic  $\hat{J}$  is not pivotal, a bootstrap test based upon it will not be exact. The problem is that the true  $P$  value depends on the unknown true distribution of  $\hat{J}$ , while the bootstrap  $P$  value (5) is based on the distribution of the bootstrap statistics  $J_j^*$ , which depends on the bootstrap DGP. These two distributions will differ whenever a test statistic is not pivotal and the parameter estimates used in the bootstrap DGP differ from the true values of the parameters. However, provided the test statistic is asymptotically pivotal, the bootstrap distribution will converge to the true one as the sample size increases. In consequence, as Beran (1988) showed, the bootstrap  $P$  value will converge to the true  $P$  value at a rate faster than the asymptotic  $P$  value converges to it.

In Davidson and MacKinnon (2002), we derived an expression for  $\hat{J}$  as a function of various random variables and quantities that depend on  $\mathbf{X}$ ,  $\mathbf{Z}$ , and the parameters of  $H_1$ , under the assumptions that the error terms are normally distributed and the regressor matrices are exogenous. This expression is rather complicated, but a key determinant of the finite-sample distribution of  $\hat{J}$  turned out to be the quantity

$$\|\boldsymbol{\theta}\|^2 \equiv \|\mathbf{M}_{\mathbf{X}} \mathbf{P}_{\mathbf{Z}} \mathbf{X} \boldsymbol{\beta}\|^2 / \sigma_1^2. \quad (6)$$

The numerator of (6) is the squared length of the part of  $\mathbf{X}\boldsymbol{\beta}$  that is explained by  $\mathbf{Z}$ , projected off  $\mathbf{X}$ . The denominator is just the variance of the error terms. Other things being equal, the larger is  $\|\boldsymbol{\theta}\|^2$ , the closer the finite-sample distribution of  $\hat{J}$  will be to the standard normal distribution.

### 3. The Fast Double Bootstrap

The fast double bootstrap proposed by Davidson and MacKinnon (2001) can be thought of as an approximation to the double bootstrap of Beran (1988). It involves calculating two different bootstrap statistics for each replication. These are based on two different bootstrap datasets drawn from two different bootstrap DGPs. The first bootstrap DGP is the one already described, in which, for the  $j^{\text{th}}$  replication, a bootstrap sample  $\mathbf{y}_j^*$  is drawn from (3), with the error terms obtained by resampling

from the vector  $\tilde{\mathbf{u}}$  defined in (4). This vector is used to compute a statistic  $J_j^*$ . For the FDB procedure, these data are also used to compute estimates  $\hat{\beta}_j^*$  and rescaled residuals  $\tilde{\mathbf{u}}_j^*$ , which characterize the second bootstrap DGP. A second bootstrap sample  $\mathbf{y}_j^{**}$  is then drawn from this DGP in precisely the same way as the first bootstrap sample was drawn from the DGP characterized by  $\hat{\beta}$  and  $\tilde{\mathbf{u}}$ , and this sample is used to compute a statistic  $J_j^{**}$ .

The fast double bootstrap  $P$  value is easily calculated from the actual test statistic  $\hat{J}$ , the  $B$  first-level bootstrap test statistics  $J_j^*$ , and the  $B$  second-level bootstrap test statistics  $J_j^{**}$ . We first calculate the single-bootstrap  $P$  value  $\hat{p}^*$  using expression (5). Next, we calculate the  $1 - \hat{p}^*$  quantile of the  $J_j^{**}$ , denoted by  $\hat{Q}^*(1 - \hat{p}^*)$  and defined implicitly by the equation

$$\frac{1}{B} \sum_{j=1}^B I(J_j^{**} > \hat{Q}^*(1 - \hat{p}^*)) = \hat{p}^*. \quad (7)$$

Of course, for finite  $B$ , there will be a range of values of  $Q^*(1 - \hat{p}^*)$  that satisfy (7), and we will need to choose one of them in a somewhat arbitrary manner. Then the FDB  $P$  value is

$$\hat{p}^{**} = \frac{1}{B} \sum_{j=1}^B I(J_j^* > \hat{Q}^*(1 - \hat{p}^*)). \quad (8)$$

Thus, instead of seeing how often the bootstrap test statistics are more extreme than the actual test statistic, we see how often they are more extreme than the  $1 - \hat{p}^*$  quantile of the  $J_j^{**}$ .

The intuition behind this procedure is as follows. Suppose, for concreteness, that the  $J_j^{**}$  tend to be less extreme than the  $J_j^*$ . This suggests that the  $J_j^*$  will tend to be less extreme than they would be if they were drawn from  $\mu_0$  instead of from  $\mu^*$ . Therefore, the ordinary bootstrap  $P$  value will be too small, and the bootstrap test will overreject. In this situation,  $\hat{Q}^*(1 - \hat{p}^*)$  will be less extreme than  $\hat{J}$  itself, and  $\hat{p}^{**}$  will consequently be larger than  $\hat{p}^*$ . Thus it appears that using  $\hat{p}^{**}$  instead of  $\hat{p}^*$  will be a step in the right direction.

The properties of  $\hat{p}^{**}$ , not specialized to the  $J$  test, were studied in Davidson and MacKinnon (2001). It is valid under quite weak conditions, but it can be expected to be more accurate than  $\hat{p}^*$  only when the bootstrap DGP is asymptotically independent of the test statistic. Since the quantities on which the bootstrap DGP depends (namely,  $\hat{\beta}$  and  $\tilde{\mathbf{u}}$ ) are either efficient estimates under the null hypothesis that  $\alpha = 0$  or functions of those estimates, the  $J$  test statistic must be asymptotically independent of them; see Davidson and MacKinnon (1999). Therefore, the theory suggests that, for the  $J$  test,  $\hat{p}^{**}$  will be more accurate than  $\hat{p}^*$ . Simulation results to be presented in the next section strongly support this hypothesis.

#### 4. Evidence from Monte Carlo Experiments

In this section, we present some simulation results for models constructed so that the bootstrap  $J$  test works relatively poorly. Our results should not be considered at all

typical for  $J$  tests computed using real data. We consider a pair of linear regression models of the form

$$H_1: y_t = \mathbf{X}_t\boldsymbol{\beta} + \delta_1 y_{t-1} + u_t \quad (9)$$

$$H_2: y_t = \mathbf{Z}_t\boldsymbol{\gamma} + \delta_2 y_{t-1} + v_t, \quad (10)$$

where the error terms for  $H_1$ , which is assumed to have generated the data, are  $t(5)$  rescaled to have variance  $\sigma_1^2$ . The first elements of  $\mathbf{X}_t$  and  $\mathbf{Z}_t$  are unity, and their dimensions are  $k - 1$  and  $l - 1$ , respectively. In all the experiments, the components of  $\mathbf{X}_t$ , except for the constant term, were distributed as  $N(0, 1)$ . Each component of  $\mathbf{Z}_t$  was also normally distributed and correlated with the corresponding component of  $\mathbf{X}_t$ , with correlation  $\rho$ . When  $k > l$  or  $l > k$ , any extra components of  $\mathbf{X}_t$  or  $\mathbf{Z}_t$  were uncorrelated with everything else. We experimented with the choice of  $\delta_1$ ,  $\sigma_1$ ,  $\rho$ ,  $k$ , and  $l$ . We found that the values of  $\delta_1$  and  $\rho$  had relatively little effect on the performance of the bootstrap  $J$  test, and we settled on base-case values of  $\delta_1 = 0.8$  and  $\rho = 0.5$ .

In the main set of experiments, the results of which we report here,  $k = l = 7$ , all the  $\beta_i$  are equal to 1, and  $\sigma_1$  takes on the values 1, 2, 4, and 8. Because there are five variables in each model that are not in the other, and the sample sizes we investigate are small, the asymptotic  $J$  test does not work particularly well. As  $\sigma_1$  increases,  $\|\boldsymbol{\theta}\|$  becomes smaller, and the performance of both the asymptotic and bootstrap  $J$  tests deteriorates.

In order to limit experimental randomness, which would make it hard to detect small departures from the nominal level of the tests, each of the experiments used 100,000 replications. This choice implies that, if the true rejection frequency is .05, the standard error of the estimated rejection frequency will be .00069. The number of bootstrap samples was always 999.

Rejection frequencies for the asymptotic  $J$  test (based on the standard normal distribution) at the nominal .05 level for ten sample sizes ( $n = 10, 15, \dots, 45, 50$ ) are shown in Figure 1. It is evident that the asymptotic  $J$  test overrejects very severely indeed, especially for the larger values of  $\sigma_1$ . The larger is  $\sigma_1$ , the less rapidly does the performance of the test improve as the sample size increases. The figure suggests that, for the two largest values of  $\sigma_1$ , the sample size would have to be very large indeed for the asymptotic  $J$  test to perform well.

When  $\sigma_1 = 1$ , the ordinary (or single) bootstrap test worked almost perfectly, and the FDB test worked even better. These results are therefore not reported. Results for  $\sigma_1 = 2$  and  $\sigma_1 = 8$ , which are much more interesting, are shown in Figures 2 and 3, respectively. When  $\sigma_1 = 2$ , the single bootstrap test overrejects quite noticeably for small sample sizes, but its performance improves rapidly as  $n$  increases. In contrast, the FDB test performs about as well as any test could be expected to perform, given some experimental error, for all sample sizes. When  $\sigma_1 = 8$ , the single bootstrap test overrejects for all sample sizes. The FDB test also overrejects, but very much less severely. The performance of both tests, mirroring that of the asymptotic test, improves only very slowly as  $n$  increases.

We remarked in the previous section that it is very important to rescale the residuals prior to resampling from them. The consequences of not doing so are shown in Figure 4, which deals with the same cases, and uses the same random numbers, as Figures 2 and 3. We see that the single bootstrap performs dramatically worse, especially for smaller sample sizes, when the residuals are not rescaled. In contrast, the FDB procedure performs only a little worse for the smaller sample sizes, and its performance is almost unchanged for the larger ones. Thus it appears that the correction implicit in the FDB procedure can compensate for flaws in the underlying method of bootstrapping.

In one final set of experiments, we took the experimental design to an extreme, making the situation very unfavorable for the finite-sample performance of the  $J$  test. We set  $k = 8$  and  $l = 9$  (the largest values that allow calculation of the  $J$  test for  $n = 10$ ), and we set  $\sigma_1 = 16$ . In these experiments, the values of  $\|\boldsymbol{\theta}\|^2$  ranged from 0.0043 (for  $n = 10$ ) to 0.2152 (for  $n = 50$ ), and the rejection frequencies for the asymptotic tests at the .05 level ranged from 0.81 to 0.75. Results for the single bootstrap and FDB tests, with and without rescaled residuals, are shown in Figure 5. As before, the FDB procedures always work very much better than the corresponding single bootstrap procedures. Curiously, the FDB procedure that uses ordinary residuals works slightly better, except when  $n = 10$ , than the FDB procedure that uses rescaled residuals. This is not true for the single bootstrap.

## 5. Conclusion

In this paper, we have proposed a simple bootstrap procedure for the  $J$  test that works extraordinarily well, even in extreme cases where the usual single bootstrap procedure overrejects quite noticeably. Our simulation experiments deliberately focused on extreme cases in which the asymptotic test rejects more than half the time and the single bootstrap test does not work well. In practice, we would expect the single bootstrap to work extremely well, and our FDB procedure to work nearly perfectly, in virtually every case that an applied worker would be likely to encounter.

## References

- Beran, R. (1988). “Prepivoting test statistics: a bootstrap view of asymptotic refinements,” *Journal of the American Statistical Association*, 83, 687–697.
- Davidson, R. and J. G. MacKinnon (1981). “Several tests for model specification in the presence of alternative hypotheses,” *Econometrica*, 49, 781–793.
- Davidson, R. and J. G. MacKinnon (1999). “The size distortion of bootstrap tests,” *Econometric Theory*, 15, 361–376
- Davidson, R. and J. G. MacKinnon (2000). “Bootstrap tests: How many bootstraps?,” *Econometric Reviews*, 19, 55–68.
- Davidson, R. and J. G. MacKinnon (2001). “Improving the reliability of bootstrap tests,” Queen’s University Institute for Economic Research Discussion Paper No. 995, revised.
- Davidson, R. and J. G. MacKinnon (2002). “Bootstrap  $J$  tests of nonnested linear regression models,” *Journal of Econometrics*, forthcoming.
- Fan, Y., and Q. Li (1995). “Bootstrapping  $J$ -type tests for non-nested regression models,” *Economics Letters*, 48, 107–112.
- Godfrey, L. G. (1998). “Tests of non-nested regression models: Some results on small sample behaviour and the bootstrap,” *Journal of Econometrics*, 84, 59–74.
- McAleer, M. (1995). “The significance of testing empirical non-nested models,” *Journal of Econometrics*, 67, 149–171.



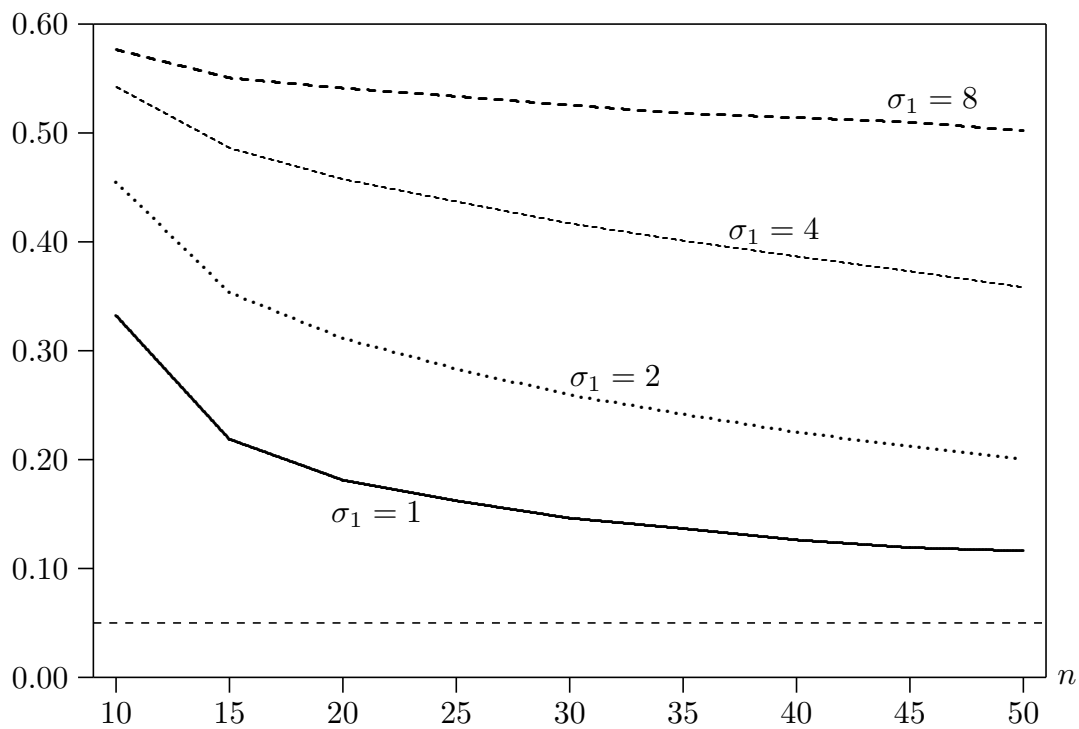


Figure 1. Rejection Frequencies for Asymptotic Tests

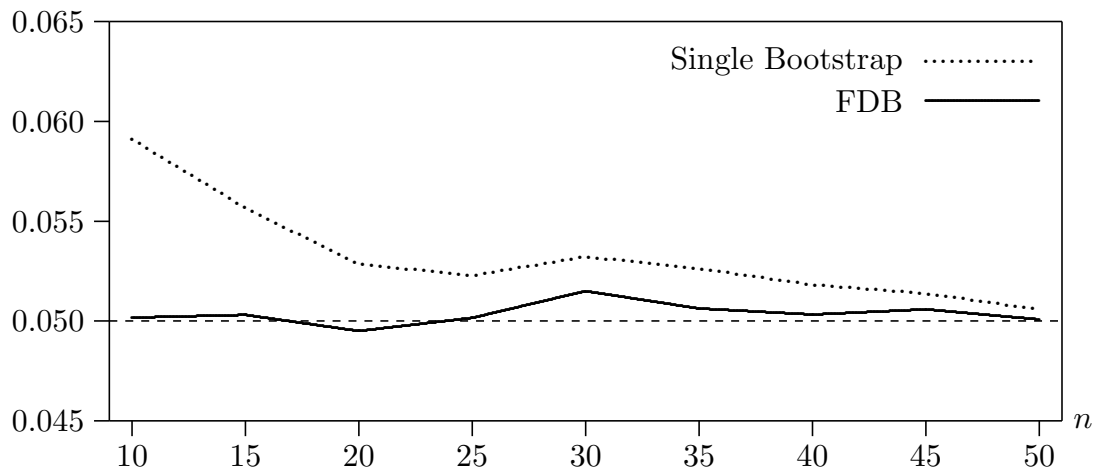


Figure 2. Rejection Frequencies for Bootstrap Tests,  $\sigma_1 = 2$

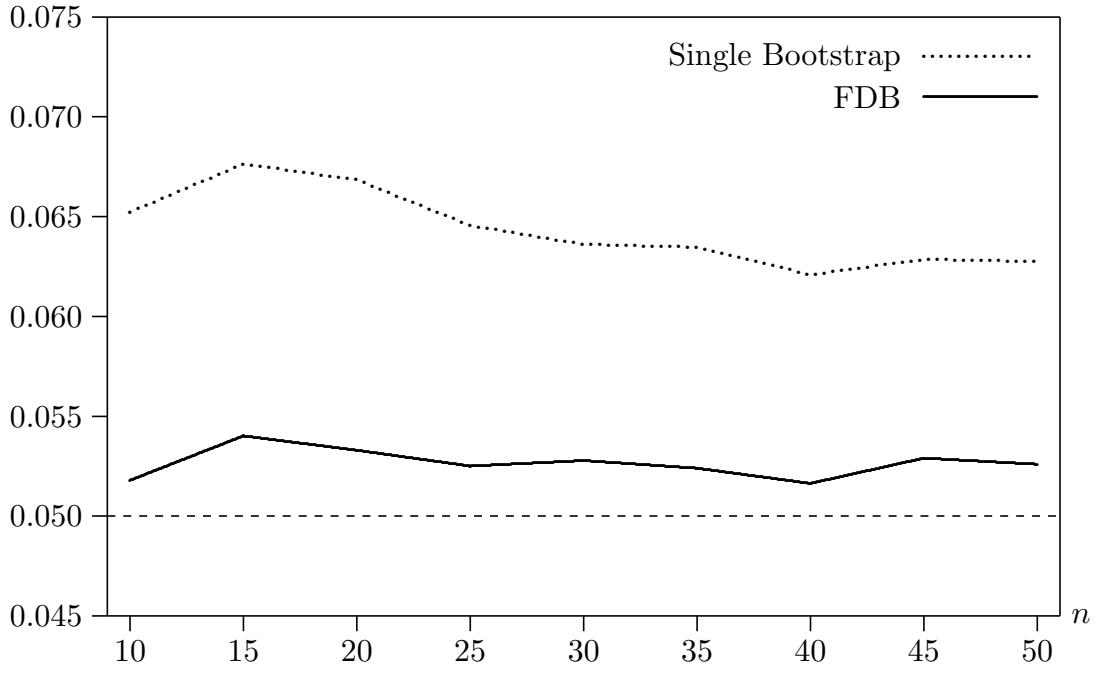


Figure 3. Rejection Frequencies for Bootstrap Tests,  $\sigma_1 = 8$

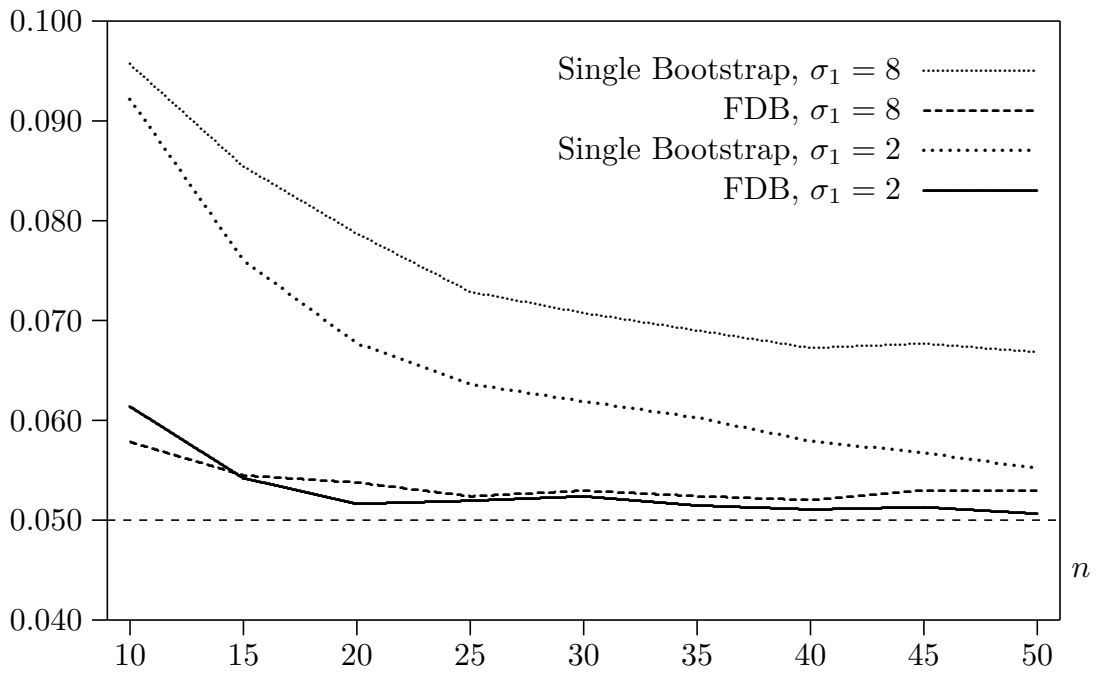


Figure 4. Bootstrap Rejection Frequencies, Ordinary Residuals

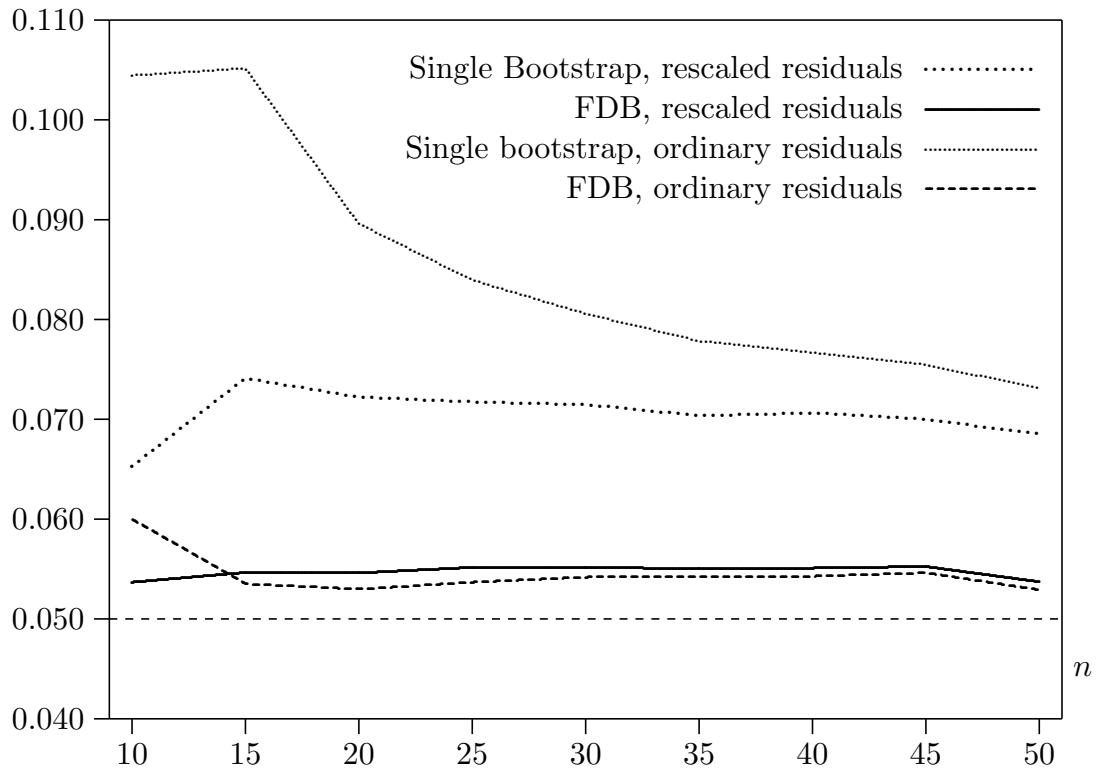


Figure 5. Bootstrap Rejection Frequencies, Extreme Case