

Bootstrap Performance with Heteroskedasticity

by

Russell Davidson

Department of Economics and CIREQ
McGill University
Montreal, Quebec, Canada
H3A 2T7

Aix-Marseille Université
CNRS, EHESS, AMSE
13205 Marseille cedex 01, France

email: russell.davidson@mcgill.ca

and

Andrea Monticini

Catholic University
Via Necchi 5
20123 Milan, Italy

email: andrea.monticini@unicatt.it

Abstract

The aim of this paper is to illustrate more than one instance of poor bootstrap performance, and to see how available diagnostic techniques can indicate reliably when and how this poor performance can arise. Two particular features that seem to be important to explain bootstrap discrepancy are illustrated by some Monte Carlo experiments.

Keywords: Bootstrap inference, fast double bootstrap, conditional fast double bootstrap, heteroskedasticity

JEL codes: C12, C22, C32

This research was supported by the Canada Research Chair program (Chair in Economics, McGill University) and by grants from the Fonds de Recherche du Québec - Société et Culture. This work was also supported by the French National Research Agency Grant ANR-17-EURE-0020

November, 2023

1. Introduction

Although the bootstrap is in many ways a polyvalent, multi-purpose, technique for obtaining reliable statistical inference, bootstrap failure can happen. A diagnostic tool for detecting bootstrap failure was proposed by Beran (1997), and other diagnostics were proposed in Davidson (2017). The latter reference suggests reasons for which the bootstrap may yield less than satisfactory results, but stops short of proposing remedies. It may also happen that the bootstrap works poorly even when the diagnostics are relatively positive.

In this paper, we carry out a rather thorough investigation of a particular case in which it can be difficult to devise a bootstrap procedure that is reliable under weak regularity conditions. The model under consideration is a linear regression model, with exogenous regressors and normal disturbances, and the null hypothesis is that there is no conditional heteroskedasticity of type ARCH or GARCH. But we wish the tests to be robust to the possible presence of *unconditional* heteroskedasticity.

Consider a linear regression model

$$y_t = \mathbf{X}_t\boldsymbol{\beta} + u_t, \quad t = 1, \dots, n, \quad (1)$$

where the regressors \mathbf{X}_t include a constant, and where the disturbances u_t may be either unconditionally or conditionally heteroskedastic, or both. The null hypothesis to be tested is that there is no conditional heteroskedasticity. A common way to perform the test is to run the regression by OLS, save the residuals \hat{u}_t , and then run the testing regression

$$\hat{u}_t^2 = a + b\hat{u}_{t-1}^2 + \text{residual}, \quad t = 2, \dots, n. \quad (2)$$

Among others, a suitable test statistic is n times the centred R^2 from this testing regression.

In many cases, reliability of the test is enhanced by use of the bootstrap. Let $R_t = \mathbf{X}_t\hat{\boldsymbol{\beta}}$ denote the fitted value for observation t from regression (1). A bootstrap DGP takes the form

$$y_t^* = R_t + u_t^* \quad t = 1, \dots, n,$$

where there are several possible ways of constructing the bootstrap disturbances u_t^* . In order to construct the bootstrap test statistic, one first regresses the y_t^* on the regressors \mathbf{X}_t , saving the residuals \hat{u}_t^* . Then the test statistic is n times the centred R^2 from the following regression, analogous to (2):

$$(\hat{u}_t^*)^2 = a + b(\hat{u}_{t-1}^*)^2 + \text{residual}, \quad t = 2, \dots, n.$$

If one does not wish to allow for unconditional heteroskedasticity of the disturbances u_t , one way to generate the u_t^* is by resampling the residuals \hat{u}_t . This makes the u_t^* an IID sequence. Alternatively, one might shuffle the residuals \hat{u}_t to get the u_t^* , in a sort of permutation test.

The possibility of unconditional heteroskedasticity can be taken account of by use of some form of wild bootstrap, which we can write as $u_t^* = s_t^* \hat{u}_t$, where the s_t^* are IID drawings from a distribution with expectation zero and variance one, independent of the observed data. Normally a good choice for the distribution of the s_t^* is the Rademacher distribution, where $s_t^* = \pm 1$, each possibility with probability one half. The variance of the Rademacher distribution is one, its third moment is zero, and its fourth moment is one. But then $(u_t^*)^2 = \hat{u}_t^2$, and so it can be expected that the $(\hat{u}_t^*)^2$ and the \hat{u}_t^2 will be highly correlated, leading to a strong correlation between the statistic computed from the original data and the bootstrap statistic. It is known that the bootstrap discrepancy can be very large in such circumstances. The correlation can be broken by shuffling or resampling the $s_t^* \hat{u}_t$, but that ignores the possibility of unconditional heteroskedasticity.

Originally, Mammen’s suggestion for the wild bootstrap was to draw the s_t^* from the two-point distribution

$$s_t^* = \begin{cases} -(\sqrt{5} - 1)/2 & \text{with probability } (\sqrt{5} + 1)/2\sqrt{5}, \\ (\sqrt{5} + 1)/2 & \text{with probability } (\sqrt{5} - 1)/2\sqrt{5}. \end{cases}$$

It can be checked that this distribution has expectation zero, variance one, third moment one, and fourth moment two. It could be hoped that use of this distribution instead of the Rademacher would break the correlation between the statistic and its bootstrap counterpart.

A different choice for the s_t^* that has found some favour is for them to be drawn from the standard normal distribution, with expectation and third moment zero, variance one, and fourth moment three. Since this is a continuous distribution rather than a two-point one, it can be hoped that it would succeed in breaking the troubling correlation. Yet another possibility is to draw the s_t^* from a continuous distribution that shares its first three moments with Mammen’s two-point distribution.

It is not clear at first glance why any of these possibilities should be better than the others and in what circumstances. It *is* clear why the wild bootstrap using the Rademacher distribution can be expected to perform very poorly, but it is not easy to provide a theoretical explanation of the performance of other possibilities. Some preliminary suggestions are made in [Section 2](#).

Whatever choice is made for generating the bootstrap disturbances, over and above the traditional single bootstrap, there is the fast double bootstrap (FDB) – see Davidson and MacKinnon (2007) – and the conditional fast double bootstrap (CFDB) – see Davidson and Monticini (2023). In favourable circumstances, the fast double bootstraps are more reliable than the single bootstrap, in the sense that the bootstrap discrepancy is smaller, at the cost of roughly doubling computing time for a fixed number of bootstrap repetitions. But by no means all circumstances are favourable, and it is often difficult to find theoretical reasons for why this should be so. Consequently, it is always advisable to conduct simulation experiments in order to see whether the fast bootstraps do indeed improve the reliability of inference. Some simulation results are presented in [Section 3](#).

2. Suggestions

In Davidson (2017), some diagnostic tools are proposed for seeing whether the bootstrap works badly, and, if so, why. The first tool proceeds as for the fast double bootstrap, and, for some chosen method of setting up a bootstrap DGP, generates a set of IID paired realisations of the bootstrap statistic and the second-level bootstrap statistic. Denote a typical pair by (τ_i, τ_i^1) , $i = 1, \dots, N$, where i indexes the realised pair. Then the second-level statistics τ_i^1 are regressed by OLS on a constant and the τ_i . Both the t statistic for the coefficient of the τ_i and the centred R^2 from this regression can serve as indicators of the extent to which these statistics are correlated. It is shown there, and in somewhat more detail in Davidson and Monticini (2023), that a positive (negative) correlation is associated with under-(over-)rejection of the bootstrap test (under the null hypothesis) for conventional significance levels. The other diagnostic is a straightforward comparison of the empirical distributions of the two statistics. This allows one to gauge the extent to which the bootstrap DGP mimics the true DGP used in the simulation.

Obviously these diagnostics are at best qualitative, but a more detailed, quantitative, account of the bootstrap discrepancy is very difficult except in very special, usually trivial, cases. Consequently, most of the discussion in this section makes no effort to arrive at definitive quantitative conclusions.

If it is the correlation between the \hat{u}_t^2 and the $(\hat{u}_t^*)^2$ that induces correlation between the two levels of bootstrap statistics, as seems very likely, anything that serves to lower this correlation should improve the performance of all the bootstrap tests: the simple bootstrap, the fast double bootstrap, and the conditional fast double bootstrap. As the number of regressors, k , in model (1) increases for a fixed sample size n , the OLS residuals from (1) and those from its bootstrap counterpart will surely become less correlated even if the Rademacher wild bootstrap is used. One might expect, therefore, that the reliability of bootstrap inference should improve with increasing k .

The wild bootstrap with Mammen’s two-point distribution gives rise to bootstrap residuals whose squares will automatically be less correlated with those from the original residuals than the bootstrap residuals from the Rademacher wild bootstrap, and this, too, should help bootstrap performance.

The standard normal distribution is continuous, unlike the Rademacher and Mammen two-point distributions, and use of it rather than either of the discrete distributions can be expected to reduce the correlation between the original and bootstrap squared residuals. Other continuous distributions would presumably have the same effect, in particular, any continuous distribution that shares the lower-order moments with Mammen’s skewed two-point distribution.

It was shown in Davidson and Flachaire (2008) that performance of the wild bootstrap, with either the Rademacher or Mammen distributions, is degraded by skewed regressors combined with unconditional heteroskedasticity. Moreover, in Davidson, Monticini, and Peel (2007) it was shown, using a new class of two-point distributions, that the Rademacher distribution, preserving the original skewness, ought to be pre-

ferred to Mammen’s distribution. On general grounds, therefore, it may be that, in the absence of these perturbing factors, the bootstrap will be more reliable.

The underlying theory of the FDB – see Davidson and MacKinnon (2007) – shows that its advantages are greater in circumstances in which the test statistic and the bootstrap DGP are only weakly correlated. This result is corroborated by Davidson and Monticini (2023), where it is proposed that the CFDB can improve matters when the correlation is stronger.

Finally, since all the methods discussed here have quite strong asymptotic justifications, it is of interest to see to what extent inference is more reliable with larger sample sizes.

The evidence uncovered in the [next section](#) shows that, at least for the particular setup considered here, there are two main features of the DGP and its bootstrap counterpart that contribute to the bootstrap discrepancy, and that they may act in opposing directions, so that an appearance of reliability may arise when the effect of one feature offsets that of the other. Of course, this may be simply coincidental, so that the same effect may disappear with relatively slight changes in the DGP and bootstrap DGP.

The first feature is the one already alluded to, namely a correlation between the first- and second-level bootstrap statistics, which betrays a correlation between the actual test statistic and the bootstrap DGP. It was shown in Davidson and MacKinnon (1999) that this correlation leads to slower convergence to ideal bootstrap performance as the sample size tends to infinity.

The second feature can be thought of as **bias**. The bootstrap DGP is constructed as an estimate of the true unobserved DGP, and it may suffer from bias in the sense that the expectation of the distribution of the test statistic under the bootstrap DGP is biased, in one direction or the other, away from the expectation under the true DGP. Since such a bias can readily be detected and estimated using Davidson’s (2017) diagnostic techniques, it may be that attempts to debias the bootstrap distribution can lead to more reliable inference. We do not explore this possibility in this note.

3. Simulation evidence

Most of the experiments of which the results are reported here are for a sample size of 50. The regressors are the constant and one other drawn from the standard normal distribution, but with the third observation equal to 5, so as to create a high-leverage point. The disturbances are standard normal multiplied by this non-constant regressor, with, therefore, considerable unconditional heteroskedasticity. There were 10,000 replications, with 399 bootstrap repetitions each. The P value discrepancy plots (above) show results for the single bootstrap in red, the FDB in green, and the CFDB in blue. In red in the kernel density plots (below) are estimated densities of the statistic, in green of the bootstrap statistic, and in blue of the second-level bootstrap statistic.

Ordinary resampling and the permutation test

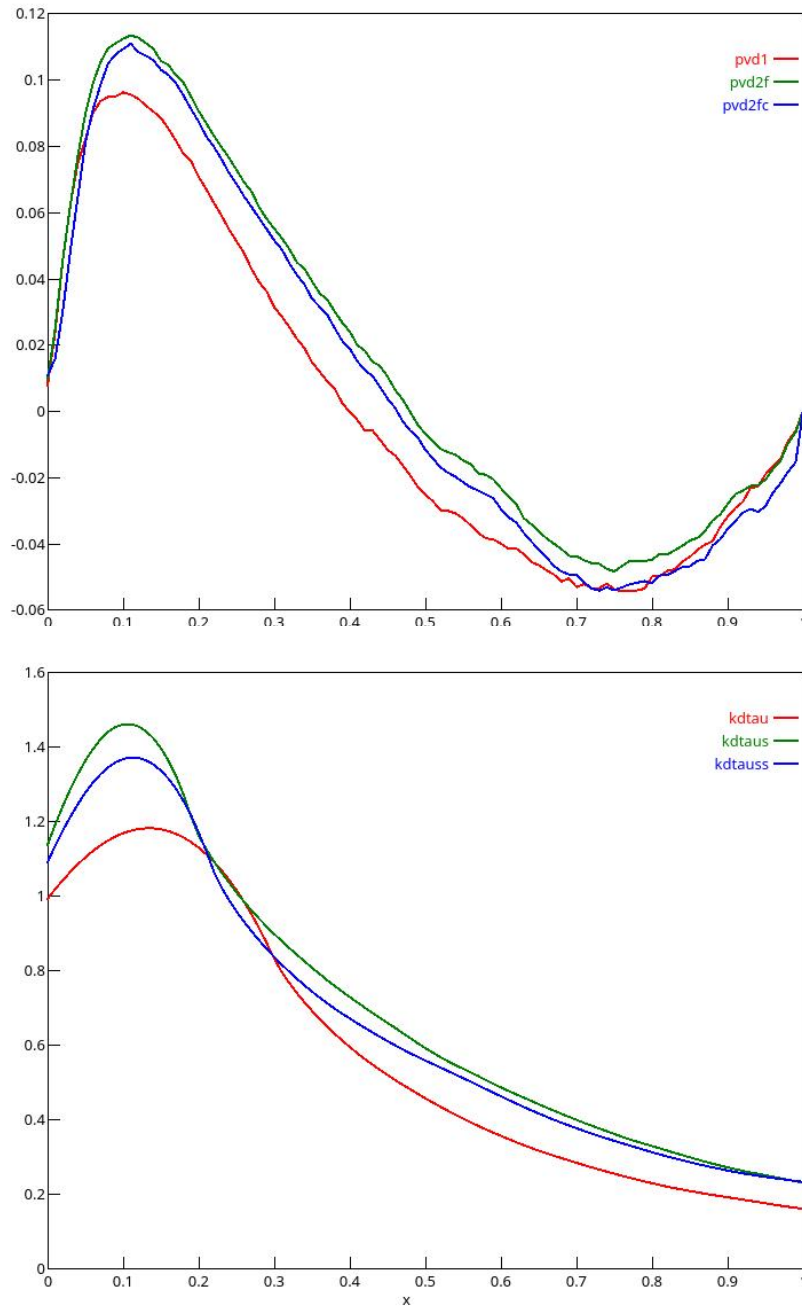


Figure 1: Resampling residuals

Figure 1 shows results where the bootstrap disturbances are simply resampled from the residuals from (1). With resampling of residuals in the presence of heteroskedasticity, there is a lot of size distortion, with considerable over-rejection in the region of interest with conventional significance levels. The fast bootstraps are even worse than the single bootstrap. The densities of both levels of bootstrap statistic seem to be somewhat less

skewed and heavier-tailed than that of the statistic itself.

The OLS diagnostic shows that there is very little correlation between the statistic and its bootstrap counterpart, with a centred R^2 of just 0.000189. However, the estimated constant in the regression is 0.736, and is highly significant. This indicates that the bootstrap discrepancy is mainly due to a negative bias in the distribution of the bootstrap statistic thought of as an estimate of the distribution of the statistic itself, rather than any correlation between the statistic and the bootstrap statistic. This is borne out by the estimated means of the three distributions: 1.279, 0.726, and 0.762, for the statistic, the bootstrap statistic, and the second-level bootstrap statistic respectively. Critical values for the bootstrap distribution are smaller than those for the true distribution; hence the observed over-rejection.

We do not report detailed results for the permutation test, in which the residuals are simply shuffled, because it is apparently still worse than the resampling bootstrap test. All three bootstraps perform similarly, and equally badly.

The Rademacher wild bootstrap

We expect under-rejection when the Rademacher wild bootstrap is used, on account of the strong positive correlation mentioned [above](#). This expectation is confirmed by what is seen in the P value discrepancy plots in [Figure 2](#). Although the FDB is more distorted than the single bootstrap, the CFDB, as expected, provides a small correction to both of these, although by no means enough for reliable inference.

The kernel density plots, on the other hand, indicate that all three statistics have very similar distributions. In other words, the bootstrap mimics the distribution of the statistic well. The distortion is mainly due to the correlation. The OLS diagnostic regression does show some bias, with an estimated constant of 0.329, but the salient feature is the R^2 of 0.469. The means of the three distributions are 1.222, 1.109, and 0.946, in the same order as [before](#).

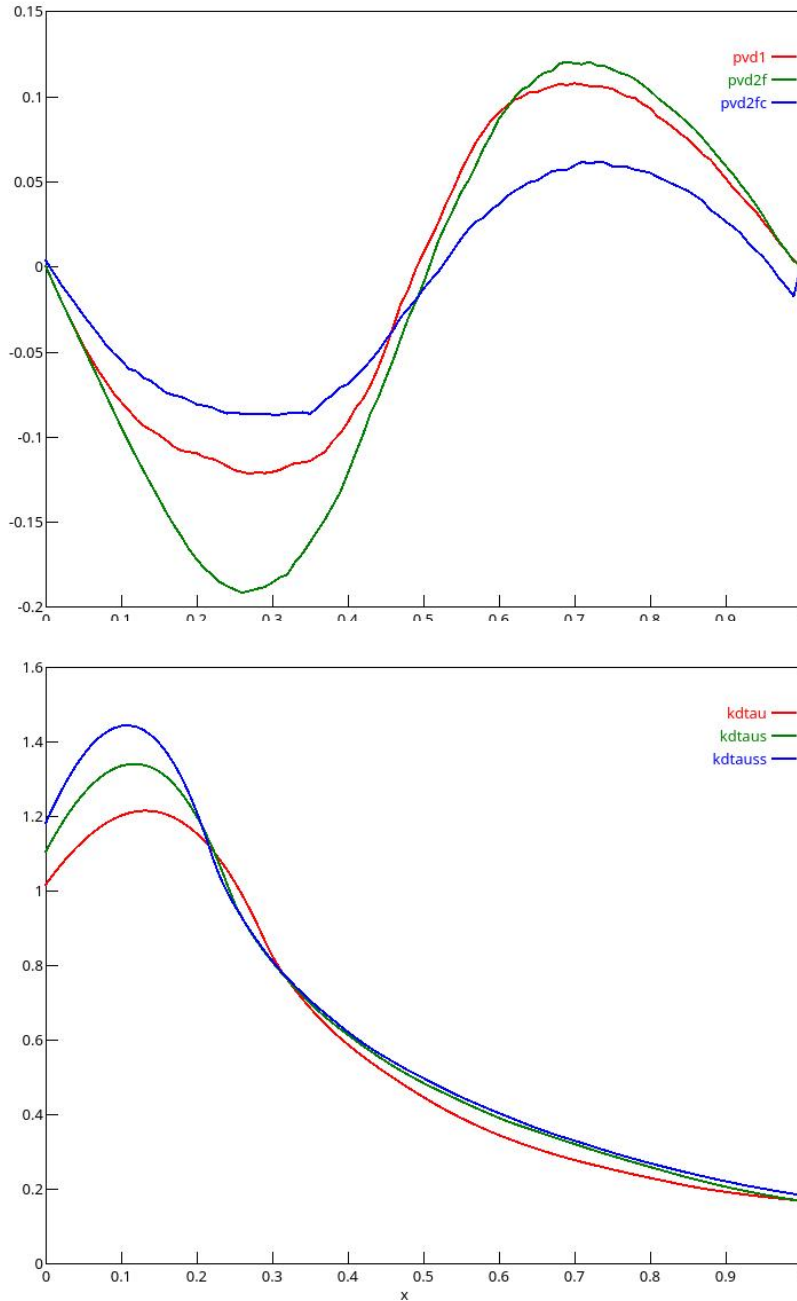


Figure 2: Rademacher wild bootstrap

The Mammen wild bootstrap

The behaviour of the wild bootstrap with the Mammen distribution is quite different, as seen in [Figure 3](#). The kernel density plots show very considerable differences in the distributions of the three statistics. By itself, this would suggest that bootstrap performance would be poor, without necessarily indicating just how. The diagnostic regression shows both a significant constant of 0.550, which leads to bias, and an R^2

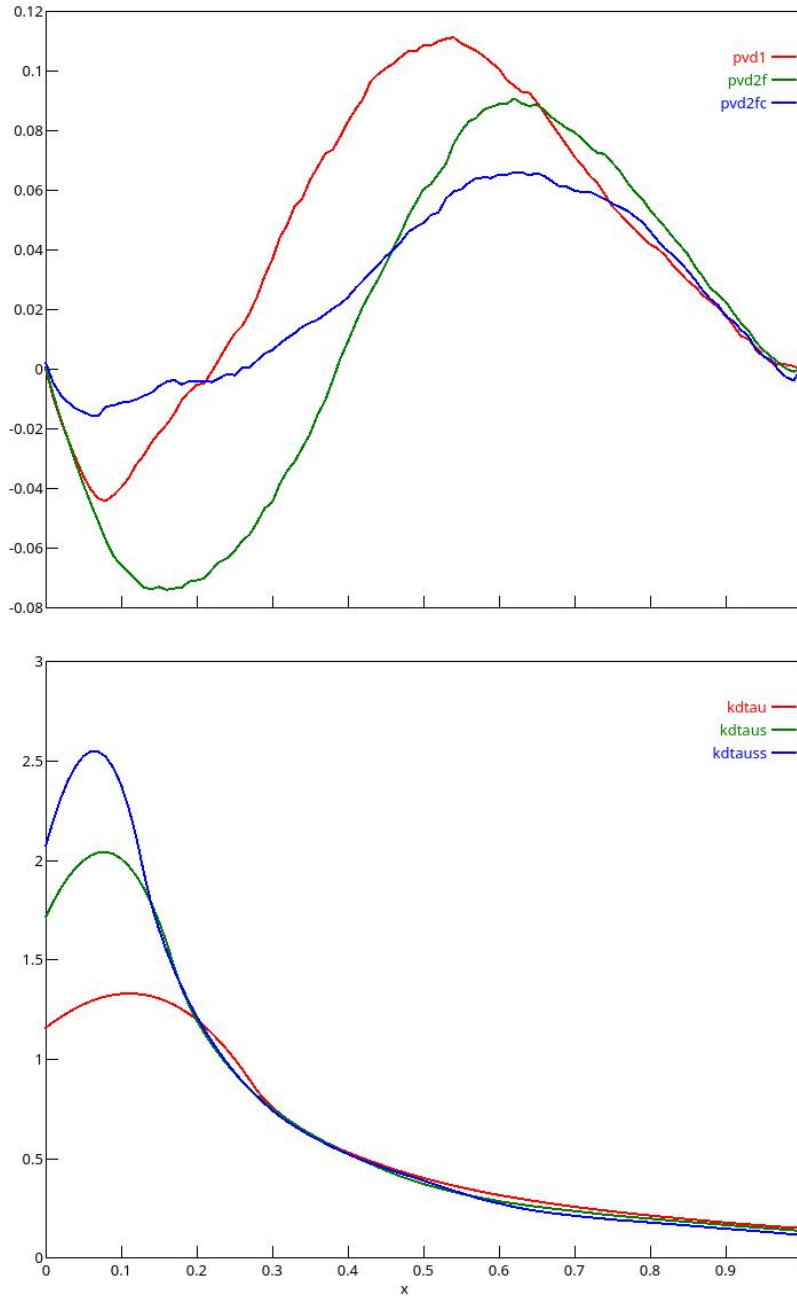


Figure 3: Mammen wild bootstrap

of 0.101 significantly different from zero, but nowhere near what was seen with the Rademacher wild bootstrap. Since the means of the three statistics, in the same order as *before*, are 1.279, 0.882, and 0.726, there are competing forces at work, the bias leading to the over-rejection in the middle of the distribution of the P value, and the correlation leading to the under-rejection in the region of interest. The CFDB test leads to considerable improvement over the other two tests, having less of both over- and under-rejection.

Wild bootstrap with standard normal distribution

If the s_t^* for the wild bootstrap are drawn from the standard normal distribution instead of from either of the two-point distributions, then the fact that this is a continuous distribution may help to reduce the correlation that gave rise to the severe distortion observed with the Radamacher distribution. That it does so to a certain extent can be seen from the graphs in [Figure 4](#).

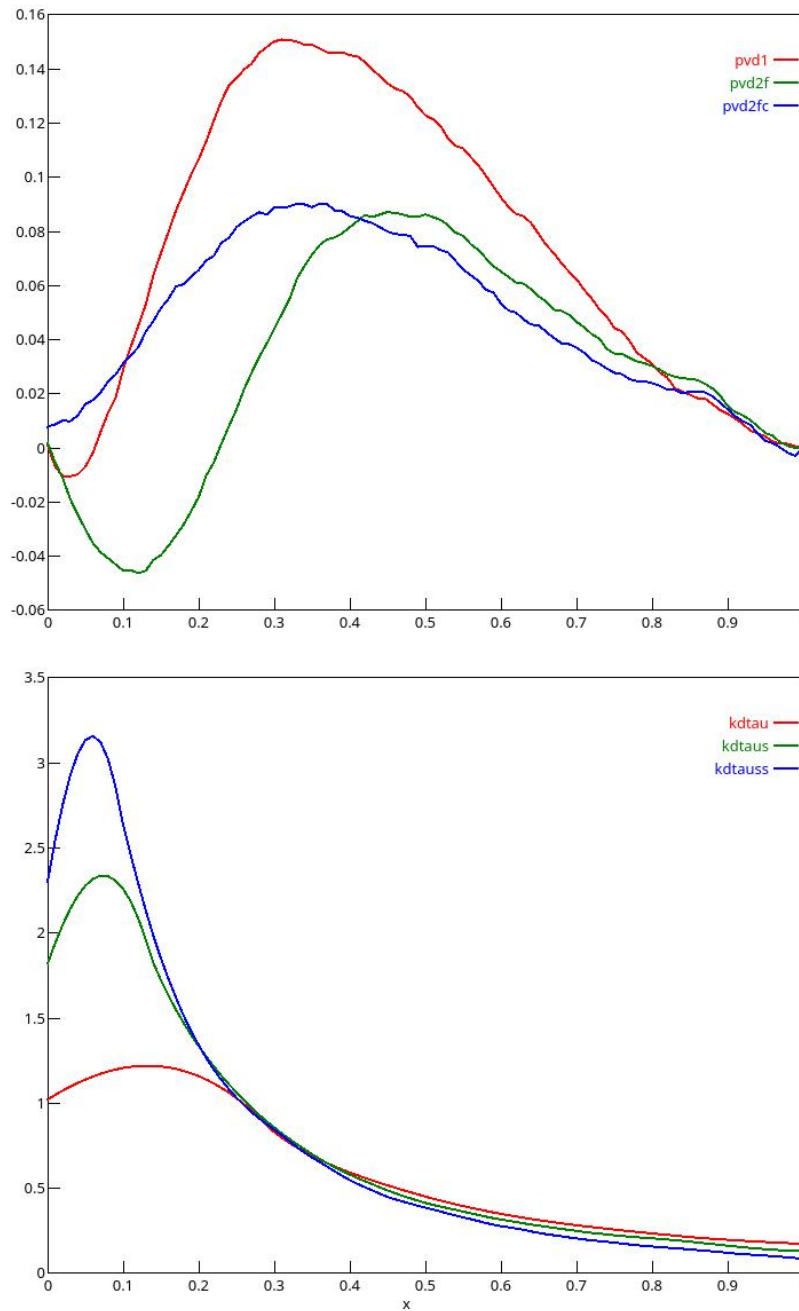


Figure 4: Wild bootstrap; standard normal

Once again, there are plainly two things at work: the correlation leads to under-rejection for small significance levels, and a bias leads to over-rejection over the rest of the distribution. This is borne out by the diagnostic regression, with two highly significant coefficients, the constant equal to 0.528, and a centred R^2 of 0.045, which, although seemingly small, is nonetheless significant. The bias is evident on looking at the means of the three statistics: 1.222, 0.718, and 0.565. The CFDB has a salutary effect on the distortion due to the correlation.

Wild bootstrap with skewed continuous distribution

So far, we have seen that use of Mammen’s skewed distribution and use of the continuous symmetrical standard normal distribution give improvements over both resampling and the Rademacher wild bootstrap. It is tempting to think that a skewed continuous distribution might improve things still further. A suitable distribution might share some moments with Mammen’s distribution.

This can be achieved by use of a method based on the Cornish-Fisher expansion proposed by Maillard (2018). Starting from a standard normal variable Z , the transformation needed to generate a variable with the desired moments is

$$X = Z + S(Z^2 - 1)/6 + K(Z^3 - 3Z)/24 - S^2(2Z^3 - 5Z)/36,$$

where the constants S and K are adjusted so as to yield the desired moments, namely 0, 1, 1, and 5, where the excess kurtosis, namely $5 - 3 = 2$, is the smallest obtainable by this method that is compatible with the first three moments. Maillard provides a table the entries of which give the values of S and K needed for the desired skewness and kurtosis, and from this we see that the appropriate choices are $S = 0.866$ and $K = 1.618$.

Graphical results can be seen in Figure 5. They are qualitatively similar to those obtained with other wild bootstraps, with under-rejection for conventional levels, and over-rejection elsewhere. The CFDB appears to be pretty reliable for levels up to around 10%.

The diagnostic regression again yields two highly significant coefficients, and the means of the three statistics are 1.222, 0.688, and 0.545. There seems to be little difference between the symmetric standard normal distribution and this skewed distribution.

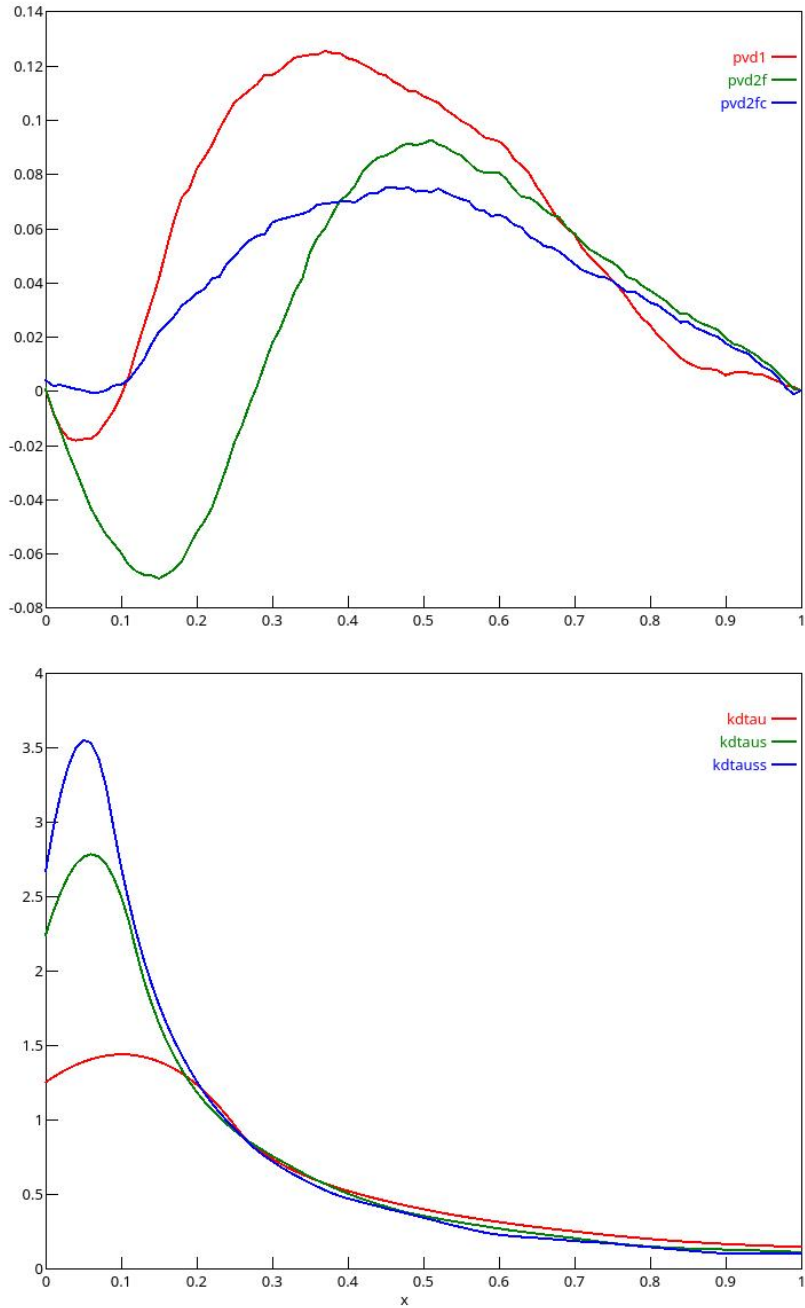


Figure 5: Wild bootstrap; skewed continuous distribution

Homoskedasticity

In the simulation results reported so far, there has been a considerable measure of unconditional heteroskedasticity. If instead the disturbances in regression (1) are homoskedastic, the wild bootstrap should still suffer from the harmful correlation between the τ_i and the τ_i^1 , but ordinary resampling should successfully break this correlation.

In Figure 6 it is seen that the performance of all three bootstrap methods is excellent. There is no distortion that can be distinguished from simulation noise. Thus we can conclude that the distortion seen in Figure 1 is entirely due to the heteroskedasticity of the disturbances combined with a skewed regressor.

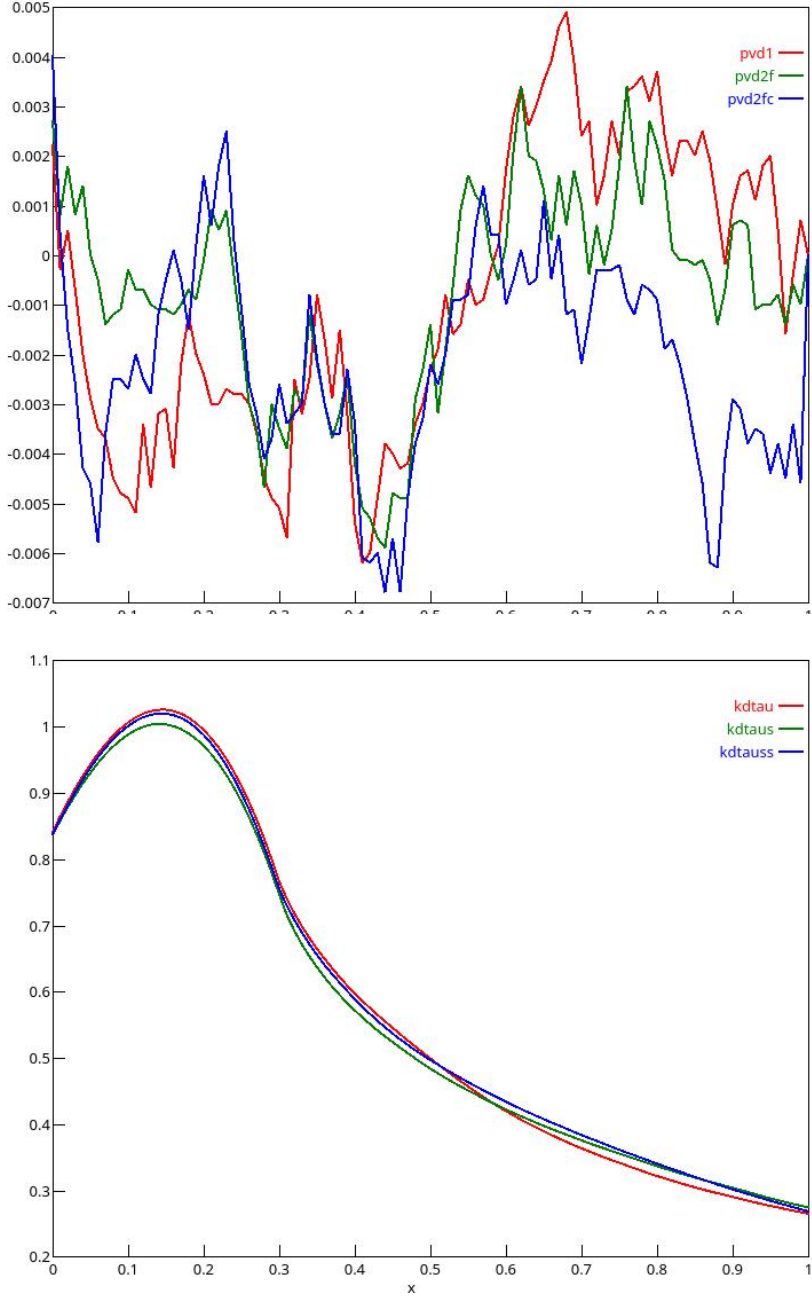


Figure 6: Resampling; homoskedasticity

The wild bootstrap with the standard normal distribution was the best wild bootstrap, albeit of a pretty bad lot. There is little reason to suppose that going from heteroskedasticity to homoskedasticity will change things by very much, since the wild bootstrap is supposedly robust to heteroskedasticity.

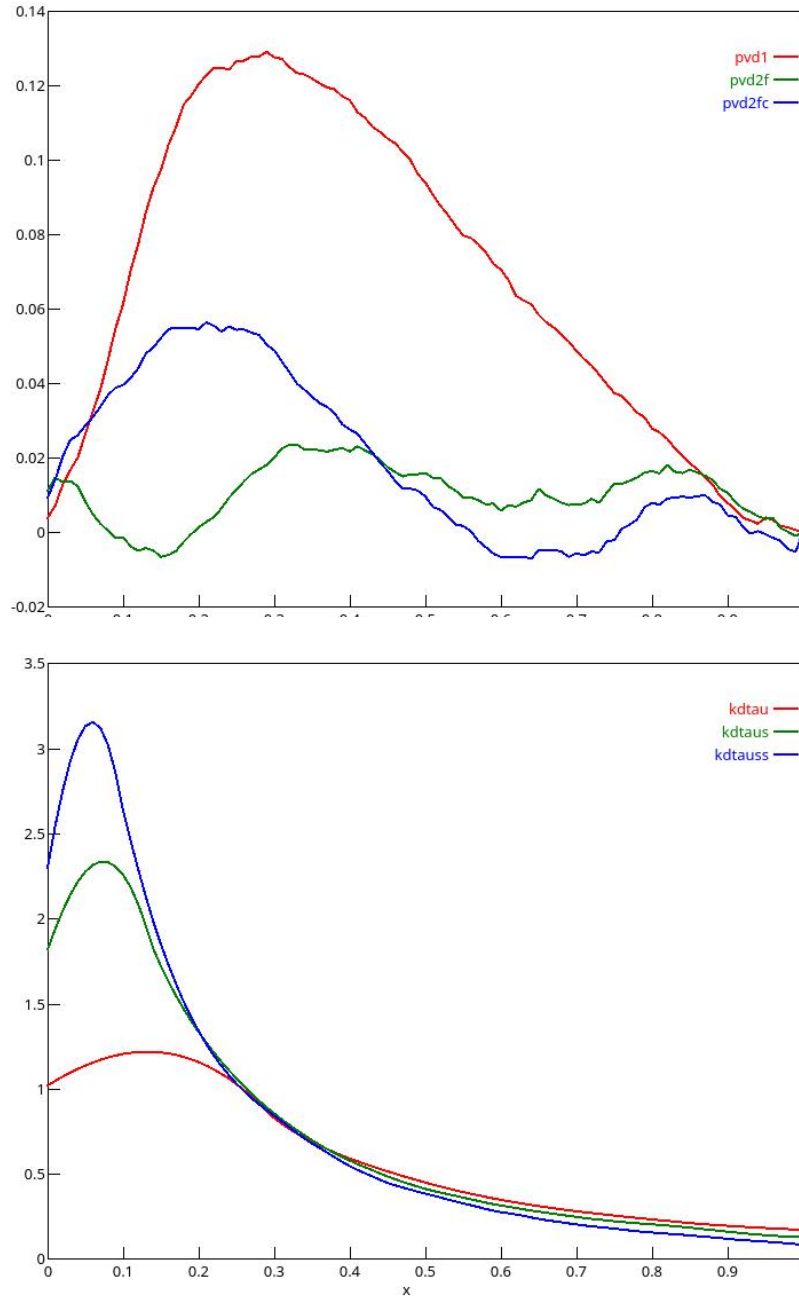


Figure 7: Standard normal wild bootstrap; homoskedasticity

It appears, from Figure 7 that this expectation is borne out by the simulation results. Comparison of Figure 7 and Figure 4 gives some evidence of this, but the CFDB is here

worse than the FDB, which suggests that the harmful correlation is not so important with homoskedasticity.

Other effects

It was suggested [earlier](#) that by increasing the number of regressors in (1) the harmful correlation might be mitigated. We undertook two experiments with the wild bootstrap to investigate this suggestion, one with the Rademacher distribution, where the distortion due to the correlation is most visible, the other with the standard normal distribution, where other effects are apparent.

Results obtained with the Rademacher distribution are shown in [Figure 8](#), to be compared with those in [Figure 2](#). The most striking thing to be seen is that the *scale* of the vertical axis is much compressed. Thus all three bootstraps are less distorted with a greater number of regressors. The overall shapes of the three P value discrepancy plots are similar, with the CFDB seemingly the least distorted, as before. The kernel density plots reveal few differences compared with the case with many fewer regressors, and all three are very similar to the others.

The diagnostic regression reveals qualitative similarity, but the indicators of distortion are different. There is a larger bias, with an estimated constant of 0.454, but a much smaller (centred) R^2 of 0.191. The means of the three statistics are 0.926, 0.735, and 0.695.

In [Figure 9](#) the results for the standard normal distribution are displayed. These should be compared with [Figure 4](#), with only two regressors, the constant and one other. In both cases, there is substantial heteroskedasticity. It is immediately clear that the number of regressors has very little impact on the performance of any of the three bootstrap procedures. In addition, results from the diagnostic regression are very similar in the two experiments.

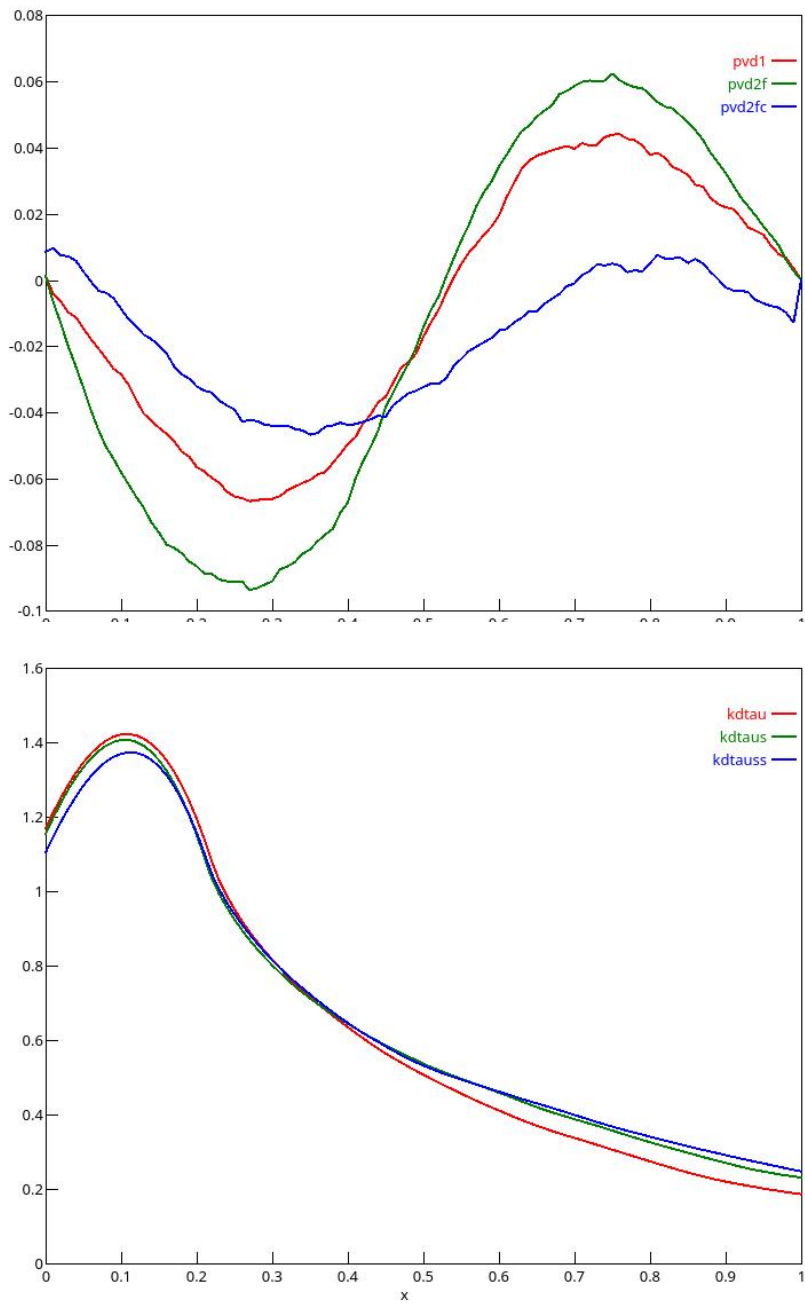


Figure 8: Rademacher wild bootstrap with 10 regressors

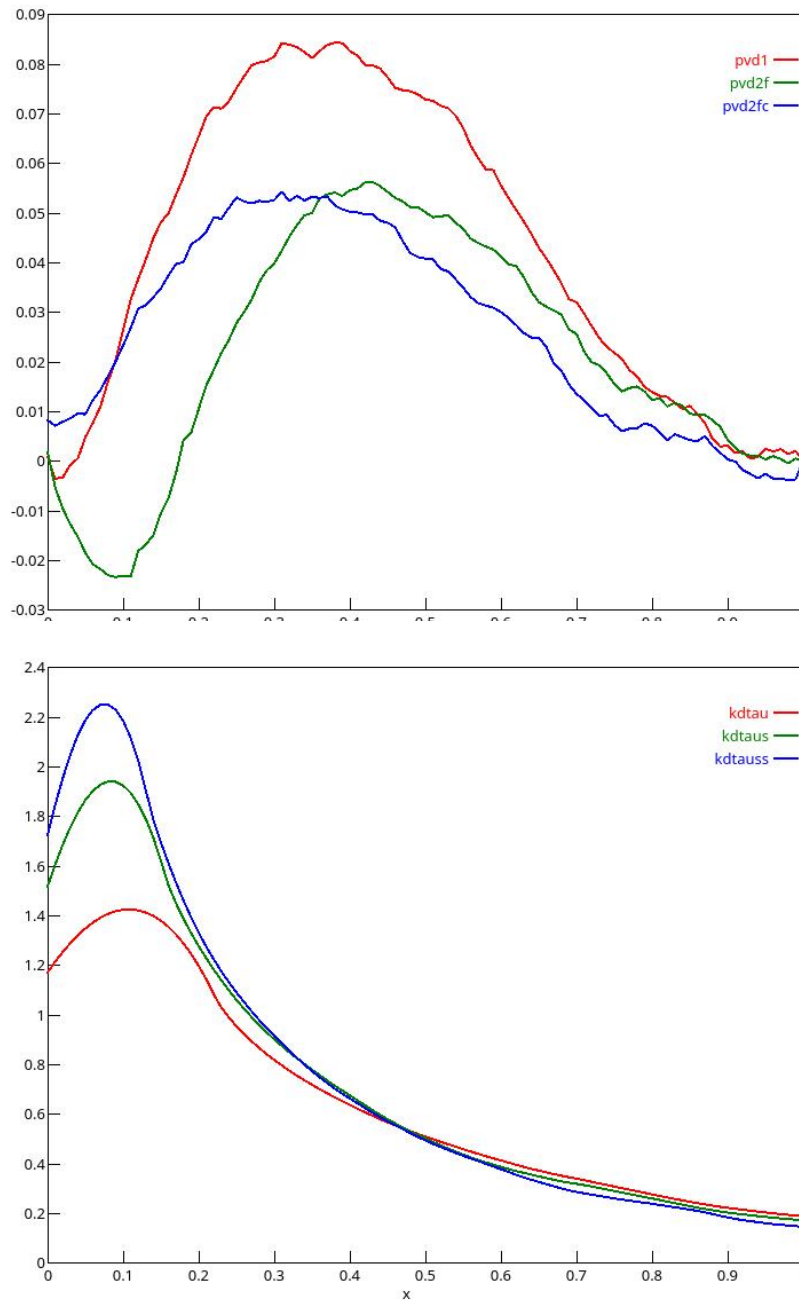


Figure 9: Standard normal wild bootstrap with 10 regressors

4. Concluding Remarks

The modest aim of this paper has been to illustrate more than one instance of poor bootstrap performance, and to see how available diagnostic techniques can indicate reliably when and how this poor performance can arise. We encountered two different

features that seem to be important in giving rise to significant bootstrap discrepancy. One, which has been documented in many places for at least two decades, is correlation between the test statistic and the bootstrap DGP, the two random elements generated from the original dataset. We see that the unfortunate effects of the correlation can be mitigated by use of the *conditional* fast double bootstrap.

The other feature that contributes to the bootstrap discrepancy is bias. The mean of the distribution of the bootstrap statistic can be seriously biased relative to the mean of the statistic generated by the true DGP. This bias can be estimated using the diagnostic techniques we have considered here, and it remains for future work to see to what extent the bootstrap discrepancy can be reduced by use of this information.

References

- Beran, R. (1997) “Diagnosing Bootstrap Success”, *Annals of the Institute of Statistical Mathematics*, **49**, 1–24.
- Davidson, J., A. Monticini, and D. Peel, (2007) “Implementing the Wild Bootstrap using a Two-points distribution”, *Economics Letters* **96,3**, 309–315
- Davidson, R. and J. G. MacKinnon (1999). “The Size Distortion of Bootstrap Tests”, *Econometric Theory*, **15**, 361-376.
- Davidson, R. (2017). “Diagnostics for the Bootstrap and Fast Double Bootstrap”, *Econometric Reviews*, **36**, 1021–1038, doi:10.1080/07474938.2017.130791
- Davidson, R. and E. Flachaire (2008). “The Wild Bootstrap, Tamed at Last”, *Journal of Econometrics*, **146**, 162–169.
- Davidson, R. and J. G. MacKinnon (2007). “Improving the Reliability of Bootstrap Tests with the Fast Double Bootstrap,” *Computational Statistics and Data Analysis*, **51**, 3259–3281.
- Davidson, R. and A. Monticini (2023) “An Improved Fast Double Bootstrap”, working paper.
- Maillard, Didier, “A User’s Guide to the Cornish Fisher Expansion (May 1, 2018)”. Available at SSRN: <https://ssrn.com/abstract=1997178> or <http://dx.doi.org/10.2139/ssrn.1997178>