

Davidson, Russell; MacKinnon, James G.

Working Paper

Bootstrap Inference in a Linear Equation Estimated by Instrumental Variables

Queen's Economics Department Working Paper, No. 1024

Provided in Cooperation with:

Queen's University, Department of Economics (QED)

Suggested Citation: Davidson, Russell; MacKinnon, James G. (2006) : Bootstrap Inference in a Linear Equation Estimated by Instrumental Variables, Queen's Economics Department Working Paper, No. 1024, Queen's University, Department of Economics, Kingston (Ontario)

This Version is available at:

<https://hdl.handle.net/10419/189308>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



Queen's Economics Department Working Paper No. 1024

Bootstrap Inference in a Linear Equation Estimated by Instrumental Variables

Russell Davidson
GREQAM and McGill University

James MacKinnon
Queen's University

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

2-2006

Bootstrap Inference in a Linear Equation Estimated by Instrumental Variables

Russell Davidson

GREQAM
Centre de la Vieille Charité
2 rue de la Charité
13236 Marseille cedex 02, France

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

email: Russell.Davidson@mcgill.ca

and

James G. MacKinnon

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

email: jgm@econ.queensu.ca

Abstract

We study several tests for the coefficient of the single right-hand-side endogenous variable in a linear equation estimated by instrumental variables. We show that all the test statistics — Student's t , Anderson-Rubin, Kleibergen's K , and likelihood ratio (LR) — can be written as functions of six random quantities. This leads to a number of interesting results about the properties of the tests under weak-instrument asymptotics. We then propose several new procedures for bootstrapping the three non-exact test statistics and a conditional version of the LR test. These use more efficient estimates of the parameters of the reduced-form equation than existing procedures. When the best of these new procedures is used, K and conditional LR have excellent performance under the null, and LR also performs very well. However, power considerations suggest that the conditional LR test, bootstrapped using this new procedure when the sample size is not large, is probably the method of choice.

This research was supported, in part, by grants from the Social Sciences and Humanities Research Council of Canada. We are grateful to seminar participants at the University of New South Wales, the University of Sydney, and the University of California Santa Barbara, and to three referees, for comments on earlier versions.

Revised, July, 2007

1. Introduction

This paper is concerned with tests for the value of the coefficient of the single right-hand-side endogenous variable in a linear structural equation estimated by instrumental variables. We consider the Wald (or t) test, Kleibergen’s (2002) K test, which can be thought of as an LM test, and the likelihood ratio (LR) test, as well as its conditional variant (Moreira, 2003), which we refer to as CLR. Both asymptotic and bootstrap versions of these tests are studied, and their relationships to the Anderson-Rubin (AR) test (Anderson and Rubin, 1949) are explored. The analysis allows for instruments that may be either strong or weak.

The major theoretical contributions of the paper depend on a simple way of writing all the test statistics of interest as functions of six random quantities. This makes it easy to understand the properties of all the tests under both weak and strong instruments. Our theoretical results are supported by extensive simulations.

Our results also make it inexpensive to simulate and bootstrap all the test statistics under the assumption of normally distributed disturbances. Although, in practice, it is generally preferable to use bootstrap methods based on resampling residuals, our experiments suggest that results from the parametric bootstrap under normality provide a very good guide to the performance of methods that resample residuals.

The paper’s main practical contribution is to propose new procedures for bootstrapping the test statistics we study. These use more efficient estimates of the parameters of the reduced-form equation than existing procedures, and what seems to be the best procedure also employs a form of bias correction. Using this procedure instead of more conventional ones greatly improves the performance under the null of all the tests. The improvement is generally greatest for the Wald test and least for the K test, because the latter already works very well in most cases. Using the new procedures also severely reduces the apparent power of the Wald test when the instruments are weak, making its power properties much more like those of the other tests.

The practical conclusions we come to differ only modestly from those of other recent papers, such as Andrews, Moreira, and Stock (2006, 2007) and Moreira, Porter, and Suarez (2004). The K test, when bootstrapped using one of our new methods, seems to be the most reliable procedure under the null when the instruments are weak and the sample size is small. However, the CLR test of Moreira (2003), which involves simulation, also seems to be very reliable when the instruments are weak and the sample size is reasonably large. When the CLR test is bootstrapped, it performs as well as the K test, except in very small samples. Power considerations suggest that either the conditional LR test or the ordinary LR test, when the latter is bootstrapped using the same new method, seem to be the best procedures overall.

In the next section, we discuss the four test statistics and show that they are all functions of six random quantities. Then, in Section 3, we show how all the statistics can be simulated very efficiently under the assumption of normally distributed disturbances. In Section 4, we consider the asymptotic properties of the statistics under both strong and weak instruments. In Section 5, we discuss some new and old ways of bootstrapping the statistics and show how, in some cases, the properties of bootstrap

tests differ from those of asymptotic tests. Finally, in Section 6, we present extensive simulation evidence on the performance of asymptotic and bootstrap tests based on all of the test statistics.

2. The Four Test Statistics

The model treated in this paper consists of just two equations,

$$\mathbf{y}_1 = \beta \mathbf{y}_2 + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}_1, \text{ and} \quad (1)$$

$$\mathbf{y}_2 = \mathbf{W}\boldsymbol{\pi} + \mathbf{u}_2. \quad (2)$$

Here \mathbf{y}_1 and \mathbf{y}_2 are n -vectors of observations on endogenous variables, \mathbf{Z} is an $n \times k$ matrix of observations on exogenous variables, and \mathbf{W} is an $n \times l$ matrix of instruments such that $\mathcal{S}(\mathbf{Z}) \subset \mathcal{S}(\mathbf{W})$, where the notation $\mathcal{S}(\mathbf{A})$ means the linear span of the columns of the matrix \mathbf{A} . The disturbances are assumed to be serially uncorrelated and, for many of the analytical results, normally distributed. We assume that $l > k$, so that the model is either exactly identified or, more commonly, overidentified.

The parameters of this model are the scalar β , the k -vector $\boldsymbol{\gamma}$, the l -vector $\boldsymbol{\pi}$, and the 2×2 contemporaneous covariance matrix of the disturbances \mathbf{u}_1 and \mathbf{u}_2 :

$$\boldsymbol{\Sigma} \equiv \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}. \quad (3)$$

Equation (1) is the structural equation we are interested in, and equation (2) is a reduced-form equation for the second endogenous variable \mathbf{y}_2 . We wish to test the hypothesis that $\beta = 0$. There is no loss of generality in considering only this null hypothesis, since we could test the hypothesis that $\beta = \beta_0$ for any nonzero β_0 by replacing the left-hand side of (1) by $\mathbf{y}_1 - \beta_0 \mathbf{y}_2$.

Since we are not directly interested in the parameters contained in the l -vector $\boldsymbol{\pi}$, we may without loss of generality suppose that $\mathbf{W} = [\mathbf{Z} \ \mathbf{W}_1]$, with $\mathbf{Z}^\top \mathbf{W}_1 = \mathbf{O}$. Notice that \mathbf{W}_1 can easily be constructed by projecting the columns of \mathbf{W} that do not belong to $\mathcal{S}(\mathbf{Z})$ off \mathbf{Z} .

The 2SLS (or IV) estimate $\hat{\beta}$ from (1), with instruments the columns of \mathbf{W} , satisfies the estimating equation

$$\mathbf{y}_2^\top \mathbf{P}_1 (\mathbf{y}_1 - \hat{\beta} \mathbf{y}_2) = 0, \quad (4)$$

where $\mathbf{P}_1 \equiv \mathbf{P}_{\mathbf{W}_1}$ is the matrix that projects on to $\mathcal{S}(\mathbf{W}_1)$. This follows because $\mathbf{Z}^\top \mathbf{W}_1 = \mathbf{O}$. The vector $\hat{\mathbf{u}}_1$ of 2SLS residuals from (1) is orthogonal to \mathbf{Z} , and so

$$\hat{\mathbf{u}}_1 = \mathbf{M}_{\mathbf{Z}} (\mathbf{y}_1 - \hat{\beta} \mathbf{y}_2). \quad (5)$$

We consider four test statistics: an asymptotic t statistic on which we may base a Wald test, the Anderson-Rubin (AR) statistic, Kleibergen's K statistic, and a likelihood ratio (LR) statistic.

From the estimating equations (4) and from expression (5) for the residuals $\hat{\mathbf{u}}_1$, it can be shown that the asymptotic t statistic for a test of the hypothesis that $\beta = 0$ is

$$t = \frac{n^{1/2} \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_1}{\|\mathbf{P}_1 \mathbf{y}_2\| \left\| \mathbf{M}_Z \left(\mathbf{y}_1 - \frac{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_1}{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2} \mathbf{y}_2 \right) \right\|}. \quad (6)$$

It can be seen that the right-hand side of (6) is homogeneous of degree zero with respect to \mathbf{y}_1 and also with respect to \mathbf{y}_2 . Consequently, the distribution of the statistic is invariant to the scales of each of the endogenous variables. In addition, the expression is unchanged if \mathbf{y}_1 and \mathbf{y}_2 are replaced by the projections $\mathbf{M}_Z \mathbf{y}_1$ and $\mathbf{M}_Z \mathbf{y}_2$, since $\mathbf{P}_1 \mathbf{M}_Z = \mathbf{M}_Z \mathbf{P}_1 = \mathbf{P}_1$, given the orthogonality of \mathbf{W}_1 and \mathbf{Z} . The second factor in the denominator can be rewritten as

$$\left(\mathbf{y}_1^\top \mathbf{M}_Z \mathbf{y}_1 - 2 \mathbf{y}_1^\top \mathbf{M}_Z \mathbf{y}_2 \frac{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_1}{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2} + \mathbf{y}_2^\top \mathbf{M}_Z \mathbf{y}_2 \left(\frac{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_1}{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2} \right)^2 \right)^{1/2}.$$

Thus it is not hard to see that the statistic (6) depends on the data only through the six quantities

$$\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1, \quad \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2, \quad \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2, \quad \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1, \quad \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2, \quad \text{and} \quad \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2; \quad (7)$$

notice that $\mathbf{y}_i^\top \mathbf{M}_Z \mathbf{y}_j = \mathbf{y}_i^\top (\mathbf{M}_W + \mathbf{P}_1) \mathbf{y}_j$, for $i, j = 1, 2$.

These six quantities, which we can think of as sufficient statistics, can easily be calculated by means of four OLS regressions on just two sets of regressors. By regressing \mathbf{y}_i on \mathbf{Z} and \mathbf{W} for $i = 1, 2$, we obtain four sets of residuals. Using the fact that $\mathbf{P}_1 \mathbf{y}_i = (\mathbf{M}_W - \mathbf{M}_Z) \mathbf{y}_i$, all six quantities can be obtained as sums of squared residuals, differences of sums of squared residuals, inner products of residual vectors, or inner products of differences of residual vectors.

Another way to test a hypothesis about β is to use the famous test statistic of Anderson and Rubin (1949). The Anderson-Rubin statistic for the hypothesis that $\beta = \beta_0$ can be written as

$$\text{AR}(\beta_0) = \frac{n-l}{l-k} \frac{(\mathbf{y}_1 - \beta_0 \mathbf{y}_2)^\top \mathbf{P}_1 (\mathbf{y}_1 - \beta_0 \mathbf{y}_2)}{(\mathbf{y}_1 - \beta_0 \mathbf{y}_2)^\top \mathbf{M}_W (\mathbf{y}_1 - \beta_0 \mathbf{y}_2)}. \quad (8)$$

Notice that, when $\beta_0 = \beta$, the AR statistic depends on the data only through the first and fourth of the six quantities (7). Under the normality assumption, this statistic is exactly distributed as $F(l-k, n-l)$ under the null hypothesis. However, because it has $l-k$ degrees of freedom, it has lower power than statistics with only one degree of freedom when $l-k > 1$.

Kleibergen (2002) therefore proposed a modification of the Anderson-Rubin statistic which has only one degree of freedom. His statistic for testing $\beta = 0$, which can also be interpreted as an LM statistic, is

$$K = (n-l) \frac{\mathbf{y}_1^\top \mathbf{P}_{\mathbf{M}_Z \mathbf{W}} \mathbf{y}_1}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1}, \quad (9)$$

which is asymptotically distributed as $\chi^2(1)$ under the null hypothesis that $\beta = 0$. The matrix $\mathbf{P}_{\mathbf{M}_Z \mathbf{W} \tilde{\boldsymbol{\pi}}}$ projects orthogonally on to the one-dimensional subspace generated by the vector $\mathbf{M}_Z \mathbf{W} \tilde{\boldsymbol{\pi}}$, where $\tilde{\boldsymbol{\pi}}$ is a vector of efficient estimates of the reduced-form parameters. These estimates will be discussed below in the context of bootstrapping.

Under our assumptions, the vector $\mathbf{M}_Z \mathbf{W} \tilde{\boldsymbol{\pi}}$ is equal to the vector $\mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1$, where $\tilde{\boldsymbol{\pi}}_1$ is the vector of OLS estimates from the artificial regression

$$\mathbf{M}_Z \mathbf{y}_2 = \mathbf{W}_1 \boldsymbol{\pi}_1 + \delta \mathbf{M}_Z \mathbf{y}_1 + \text{residuals}. \quad (10)$$

Noting that $\mathbf{M}_1 \mathbf{M}_Z = \mathbf{M}_W$, we see that the estimate of δ from this regression is given by the FWL regression¹

$$\mathbf{M}_W \mathbf{y}_2 = \delta \mathbf{M}_W \mathbf{y}_1 + \text{residuals},$$

so that $\hat{\delta} = \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 / \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1$. Substituting this result into (10) then shows that

$$\mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1 = \mathbf{P}_1 (\mathbf{M}_Z \mathbf{y}_2 - \hat{\delta} \mathbf{M}_Z \mathbf{y}_1) = \mathbf{P}_1 \left(\mathbf{y}_2 - \frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \mathbf{y}_1 \right). \quad (11)$$

From this, it can be seen that the K statistic, like the t statistic and the AR statistic, depends on the data only through the six quantities (7). Allowing for notational differences and our special assumptions, equation (11) is the same as the expression given on page 1785 of Kleibergen (2002).

The explicit expression of K in terms of these six quantities is

$$K = \frac{(n-l)(\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 - \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2)^2}{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 (\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1)^3 + \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 (\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2)^2 - 2 \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 (\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1)^2}. \quad (12)$$

We see directly from this expression that the statistic K , like the t and AR statistics, is invariant to the scales of \mathbf{y}_1 and \mathbf{y}_2 .

It is well known that, except for an additive constant, the concentrated loglikelihood function for the model specified by (1), (2), and (3) can be written as

$$-\frac{n}{2} \log \left(1 + \frac{l-k}{n-l} \text{AR}(\beta) \right), \quad (13)$$

where $\text{AR}(\beta)$ is the Anderson-Rubin statistic (8) evaluated at β . Using (8) and the fact that $\mathbf{M}_Z = \mathbf{M}_W + \mathbf{P}_1$, the loglikelihood function (13) becomes

$$\ell(\beta) \equiv -\frac{n}{2} \log \left(\frac{\|\mathbf{M}_Z(\mathbf{y}_1 - \mathbf{y}_2 \beta)\|^2}{\|\mathbf{M}_W(\mathbf{y}_1 - \mathbf{y}_2 \beta)\|^2} \right). \quad (14)$$

¹ For an exposition of the FWL regression, see, for instance, Davidson and MacKinnon (2004, Chapter 2).

The LR statistic for a test of $\beta = 0$ is $2(\ell(\hat{\beta}_{\text{ML}}) - \ell(0))$, where $\hat{\beta}_{\text{ML}}$ is the ML estimate. As shown by Anderson and Rubin (1949), the statistic can also be defined in terms of the value of the ratio that is the argument of the logarithm in (14), minimized with respect to β . We denote this minimized value by $\hat{\kappa}$. The LR statistic is then $n(\log \kappa_0 - \log \hat{\kappa})$, where κ_0 is the ratio evaluated at $\beta = 0$.

Standard arguments show that $\hat{\kappa}$ is the smaller of the two roots of the determinantal equation

$$\begin{vmatrix} \mathbf{y}_1^\top (\mathbf{M}_Z - \kappa \mathbf{M}_W) \mathbf{y}_1 & \mathbf{y}_1^\top (\mathbf{M}_Z - \kappa \mathbf{M}_W) \mathbf{y}_2 \\ \mathbf{y}_2^\top (\mathbf{M}_Z - \kappa \mathbf{M}_W) \mathbf{y}_1 & \mathbf{y}_2^\top (\mathbf{M}_Z - \kappa \mathbf{M}_W) \mathbf{y}_2 \end{vmatrix} = 0.$$

Some tedious algebra then shows that

$$\text{LR} = n \log(1 + SS/n) - n \log \left(1 + \frac{SS + TT}{2n} - \frac{1}{2n} \sqrt{(SS - TT)^2 + 4ST^2} \right), \quad (15)$$

where

$$\begin{aligned} SS &\equiv n \frac{\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1}, \\ ST &\equiv \frac{n}{\Delta^{1/2}} \left(\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 - \frac{\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \right), \\ TT &\equiv \frac{n}{\Delta} \left(\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 - 2 \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 + \frac{\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 (\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2)^2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \right), \end{aligned} \quad (16)$$

and

$$\Delta \equiv \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 - (\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2)^2. \quad (17)$$

The notation is chosen so as to be reminiscent of that used by Moreira (2003) in his discussion of a conditional LR test. His development is different from ours in that he assumes for most of his analysis that the contemporaneous disturbance correlation matrix Σ is known. We do however discuss his conditional approach to the LR test in what follows. Moreira also introduces a simplified statistic, LR_0 , which is obtained by Taylor expanding the logarithms in (15) and discarding terms of order smaller than unity as $n \rightarrow \infty$. This procedure yields

$$\text{LR}_0 = \frac{1}{2} \left(SS - TT + \sqrt{(SS - TT)^2 + 4ST^2} \right). \quad (18)$$

We see that both LR and LR_0 are invariant to the scales of \mathbf{y}_1 and \mathbf{y}_2 and depend only on the six quantities (7).

Some more tedious algebra shows that the Kleibergen statistic (12) can also be expressed in terms of the quantities ST and TT as follows.

$$K = \frac{n - l}{n} \frac{ST^2}{TT}. \quad (19)$$

Moreira (2003) demonstrates this relation (without the degrees-of-freedom adjustment) for the case in which the disturbance covariance matrix is assumed known. Finally, it is worth noting that, except for the initial deterministic factors, SS is equal to the Anderson-Rubin statistic $AR(0)$.

3. Simulating the Test Statistics

Now that we have expressions for the test statistics of interest in terms of the six quantities (7), we can explore the properties of these statistics and how to simulate them efficiently. Our results will also be used in the next two sections when we discuss asymptotic and bootstrap tests.

In view of the scale invariance that we have established for all the statistics, the contemporaneous covariance matrix of the disturbances \mathbf{u}_1 and \mathbf{u}_2 can without loss of generality be set equal to

$$\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}, \quad (20)$$

with both variances equal to unity. Thus we can represent the disturbances in terms of two independent n -vectors, say \mathbf{v}_1 and \mathbf{v}_2 , of independent standard normal elements, as follows:

$$\mathbf{u}_1 = \mathbf{v}_1, \quad \mathbf{u}_2 = \rho\mathbf{v}_1 + r\mathbf{v}_2, \quad (21)$$

where $r \equiv (1 - \rho^2)^{1/2}$. We now show that we can write all the test statistics as functions of \mathbf{v}_1 , \mathbf{v}_2 , the exogenous variables, and just three parameters.

With the specification (21), we see from (2) that

$$\begin{aligned} \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 &= (\rho\mathbf{v}_1 + r\mathbf{v}_2)^\top \mathbf{M}_W (\rho\mathbf{v}_1 + r\mathbf{v}_2) \\ &= \rho^2 \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1 + r^2 \mathbf{v}_2^\top \mathbf{M}_W \mathbf{v}_2 + 2\rho r \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2, \end{aligned} \quad (22)$$

and

$$\begin{aligned} \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 &= \boldsymbol{\pi}_1^\top \mathbf{W}_1^\top \mathbf{W}_1 \boldsymbol{\pi}_1 + 2\boldsymbol{\pi}_1^\top \mathbf{W}_1^\top (\rho\mathbf{v}_1 + r\mathbf{v}_2) \\ &\quad + \rho^2 \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + r^2 \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2 + 2\rho r \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2. \end{aligned} \quad (23)$$

Now let $\mathbf{W}_1 \boldsymbol{\pi}_1 = a\mathbf{w}_1$, with $\|\mathbf{w}_1\| = 1$. The square of the parameter a is the so-called scalar concentration parameter; see Phillips (1983, p. 470) and Stock, Wright, and Yogo (2002). Further, let $\mathbf{w}_1^\top \mathbf{v}_i = x_i$, for $i = 1, 2$. Clearly, x_1 and x_2 are independent standard normal variables. Then

$$\boldsymbol{\pi}_1^\top \mathbf{W}_1^\top \mathbf{W}_1 \boldsymbol{\pi}_1 = a^2 \text{ and } \boldsymbol{\pi}_1^\top \mathbf{W}_1^\top \mathbf{v}_i = ax_i, \quad i = 1, 2. \quad (24)$$

Thus (23) becomes

$$\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 = a^2 + 2a(\rho x_1 + r x_2) + \rho^2 \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + r^2 \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2 + 2\rho r \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2. \quad (25)$$

From (1), we find that

$$\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 = \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1 + 2\beta(\rho\mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1 + r\mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2) + \beta^2 \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2. \quad (26)$$

Similarly,

$$\begin{aligned} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 &= \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + 2\beta \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{v}_1 + \beta^2 \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 \\ &= \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + 2\beta(ax_1 + \rho \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + r \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2) + \beta^2 \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2. \end{aligned} \quad (27)$$

Further, from both (1) and (2),

$$\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 = \rho \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1 + r \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2 + \beta \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2, \quad \text{and} \quad (28)$$

$$\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 = ax_1 + \rho \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + r \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2 + \beta \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2. \quad (29)$$

The relations (22), (25), (26), (27), (28), and (29) show that the six quantities (7) can be generated in terms of eight random variables and three parameters. The eight random variables are x_1 and x_2 , along with six quadratic forms of the same sort as those in (7),

$$\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1, \quad \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2, \quad \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2, \quad \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1, \quad \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2, \quad \text{and} \quad \mathbf{v}_2^\top \mathbf{M}_W \mathbf{v}_2, \quad (30)$$

and the three parameters are a , ρ , and β . Under the null hypothesis, of course, $\beta = 0$. Since $\mathbf{P}_1 \mathbf{M}_W = \mathbf{O}$, the first three variables of (30) are independent of the last three.

If we knew the distributions of the eight random variables on which all the statistics depend, we could simulate them directly. We now characterize these distributions. The symmetric matrix

$$\begin{bmatrix} \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 & \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2 \\ \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_1 & \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2 \end{bmatrix} \quad (31)$$

follows the Wishart distribution $W(\mathbf{I}_2, l - k)$, and the matrix

$$\begin{bmatrix} \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1 & \mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2 \\ \mathbf{v}_2^\top \mathbf{M}_W \mathbf{v}_1 & \mathbf{v}_2^\top \mathbf{M}_W \mathbf{v}_2 \end{bmatrix}$$

follows the distribution $W(\mathbf{I}_2, n - l)$. It follows from the analysis of the Wishart distribution in Anderson (1984, Section 7.2) that $\mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_1$ is equal to a random variable t_{11}^M which follows the chi-squared distribution with $n - l$ degrees of freedom, $\mathbf{v}_1^\top \mathbf{M}_W \mathbf{v}_2$ is the square root of t_{11}^M multiplied by a standard normal variable z_M independent of it, and $\mathbf{v}_2^\top \mathbf{M}_W \mathbf{v}_2$ is z_M^2 plus a chi-squared variable t_{22}^M with $n - l - 1$ degrees of freedom, independent of z_M and t_{11}^M .

The elements of the matrix (31) can, of course, be characterized in the same way. However, since the elements of the matrix are not independent of x_1 and x_2 , it is preferable to define $\mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2$ as x_2^2 plus t_{22}^P , $\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2$ as $x_1 x_2$ plus the square root of t_{22}^P times z_P , and $\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1$ as $x_1^2 + z_P^2 + t_{11}^P$. Here t_{11}^P and t_{22}^P are both chi-squared, with $l - k - 2$ and $l - k - 1$ degrees of freedom, respectively, and z_P is standard normal. All these variables are mutually independent, and they are also independent of x_1 and x_2 . Of course, if $l - k \leq 2$, chi-squared variables with zero or negative degrees of freedom are to be set to zero, and if $l - k = 0$, then $z_P = 0$.

An alternative way to simulate the test statistics, which does not require the normality assumption, is to make use of a much simplified model. This model may help to provide an intuitive understanding of the results in the next two sections. The simplified model is

$$\mathbf{y}_1 = \beta \mathbf{y}_2 + \mathbf{u}_1, \quad (32)$$

$$\mathbf{y}_2 = a \mathbf{w}_1 + \mathbf{u}_2, \quad (33)$$

where the disturbances are generated according to (21). Here the n -vector $\mathbf{w}_1 \in \mathcal{S}(\mathbf{W})$ with $\|\mathbf{w}_1\| = 1$, \mathbf{W} being, as before, an $n \times l$ matrix of instruments. By normalizing \mathbf{w}_1 in this way, we are implicitly using weak-instrument asymptotics; see Staiger and Stock (1997). Clearly, we may choose $a \geq 0$. The DGPs of this simple model, which are completely characterized by the parameters β , ρ , and a , can generate the six quantities (7) so as to have the same distributions as those generated by any DGP of the more complete model specified by (1), (2), and (3).

If the disturbances are not Gaussian, the distributions of the statistics depend not only on the parameters a , ρ , and β but also on the vector \mathbf{w}_1 and the linear span of the instruments. We may suspect, however, that this dependence is weak, and limited simulation evidence (not reported) strongly suggests that this is indeed the case. The distribution of the disturbances seems to have a much greater effect on the distributions of the test statistics than the features of \mathbf{W} .

4. Asymptotic Theory

To fix ideas, we begin with a short discussion of the conventional asymptotic theory of the tests discussed in Section 2. By “conventional”, we mean that the instruments are assumed to be strong, in a sense made explicit below. Under this assumption, the tests are all classical. In particular, Kleibergen (2002) shows that the K statistic is a version of the Lagrange Multiplier test.

The reduced-form equation (33) of the simplified model of the previous section is written in terms of an instrumental variable \mathbf{w}_1 such that $\|\mathbf{w}_1\| = 1$. Conventional asymptotics would set $\|\mathbf{w}_1\|^2 = n$ and let the parameter a be independent of the sample size. Our setup is better suited to the weak-instrument asymptotics of Staiger and Stock (1997). For conventional asymptotics, we may suppose that $a = n^{1/2}\alpha$, for α constant as the sample size $n \rightarrow \infty$.

Under the null, $\beta = 0$. Under local alternatives, we let $\beta = n^{-1/2}b$, for b constant as $n \rightarrow \infty$. Conventional asymptotics applied to (22), (25), (26), (27), (28), and (29) then give

$$\begin{aligned} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 &\stackrel{a}{=} (x_1 + \alpha b)^2 + t_{11}^P, & n^{-1} \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 &\stackrel{a}{=} 1, \\ n^{-1/2} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 &\stackrel{a}{=} \alpha x_1 + \alpha^2 b, & n^{-1} \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 &\stackrel{a}{=} \rho, \\ n^{-1} \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 &\stackrel{a}{=} \alpha^2, & n^{-1} \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 &\stackrel{a}{=} 1. \end{aligned} \quad (34)$$

Using these results, it is easy to check that, for $\beta = 0$, the statistics t^2 , K , and LR, given by (6), (12), and (15), respectively, are all asymptotically equal to $(x_1 + \alpha b)^2$. They have a common asymptotic distribution of χ^2 with one degree of freedom and noncentrality parameter $\alpha^2 b^2 = a^2 \beta^2$. Incidentally, we can also see that the Anderson-Rubin statistic $\text{AR}(0)$, as given by (8), is asymptotically equal to $(x_1 + \alpha b)^2 + z_P^2 + t_{11}^P$, with $l - k$ degrees of freedom and the same noncentrality parameter. Thus $\text{AR}(0)$ is asymptotically equal to the same noncentral $\chi^2(1)$ random variable as the other three statistics, plus an independent central $\chi^2(l - k - 1)$ random variable.

We now turn to the more interesting case of weak-instrument asymptotics, for which a is kept constant as $n \rightarrow \infty$. The three right-hand results of (34) are unchanged, but the left-hand ones have to be replaced by the following equations, which involve no asymptotic approximation, but hold even in finite samples:

$$\begin{aligned} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 &= x_1^2 + z_P^2 + t_{11}^P, \\ \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 &= ax_1 + \rho(x_1^2 + z_P^2 + t_{11}^P) + r(x_1 x_2 + z_P \sqrt{t_{22}^P}), \text{ and} \\ \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 &= a^2 + 2a(\rho x_1 + r x_2) + \rho^2(x_1^2 + z_P^2 + t_{11}^P) \\ &\quad + r^2(x_2^2 + t_{22}^P) + 2\rho r(x_1 x_2 + z_P \sqrt{t_{22}^P}). \end{aligned} \tag{35}$$

Since the Anderson-Rubin statistic $\text{AR}(0)$ is exactly pivotal for the model we are studying, its distribution under the null that $\beta = 0$ depends neither on a nor on ρ . Since the quantity SS in (16) is equal to $\text{AR}(0)$ except for degrees-of-freedom factors, it too is exactly pivotal. Its asymptotic distribution under weak-instrument asymptotics is that of $\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1$. Thus, as we see from the first line of (35),

$$SS \stackrel{a}{=} x_1^2 + z_P^2 + t_{11}^P, \tag{36}$$

which follows the central $\chi^2(l - k)$ distribution.

Although Kleibergen's K statistic is not exactly pivotal, it is asymptotically pivotal under both weak-instrument and strong-instrument asymptotics. From (12), and using (34) and (35), we can see, after some algebra, that, under weak-instrument asymptotics and under the null,

$$\begin{aligned} K &\stackrel{a}{=} \frac{(\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 - \rho \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1)^2}{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 - 2\rho \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 + \rho^2 \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1} \\ &= \frac{(ax_1 + r(x_1 x_2 + z_P \sqrt{t_{22}^P}))^2}{a^2 + 2arx_2 + r^2(x_2^2 + t_{22}^P)} = \frac{(x_1(a + rx_2) + rz_P \sqrt{t_{22}^P})^2}{(a + rx_2)^2 + r^2 t_{22}^P}. \end{aligned} \tag{37}$$

Although the last expression above depends on a and ρ , it is in fact just a chi-squared variable with one degree of freedom. To see this, argue conditionally on all random variables except x_1 and z_P , recalling that all the random variables in the expression are mutually independent. The numerator is the square of a linear combination of the standard normal variables x_1 and z_P , and the denominator is the conditional variance

of this linear combination. Thus the conditional asymptotic distribution of K is χ_1^2 , and so also its unconditional distribution. As Kleibergen (2002) remarks, this implies that K is asymptotically pivotal in all configurations of the instruments, including that in which $a = 0$ and the instruments are completely invalid.

For the LR statistic, we can write down expressions asymptotically equal to the quantities ST and TT in (16). First, from (17), we have

$$\Delta/n^2 \stackrel{a}{=} 1 - \rho^2.$$

It is then straightforward to check that

$$ST \stackrel{a}{=} \frac{1}{(1 - \rho^2)^{1/2}} \left(ax_1 + r(x_1x_2 + z_p \sqrt{t_{22}^P}) \right),$$

and

$$TT \stackrel{a}{=} \frac{1}{1 - \rho^2} (a^2 + 2arx_2 + r^2(x_2^2 + t_{22}^P)). \quad (38)$$

Comparison with (37) then shows that, in accordance with (19), $K \stackrel{a}{=} ST^2/TT$.

It is clear from (38) and (36) that SS and TT are asymptotically independent, since the former depends only on the random variables x_1 , z_P , and t_{11}^P , while the latter depends only on x_2 and t_{22}^P . The discussion based on (37) shows that, conditional on TT , ST is distributed as \sqrt{TT} times a standard normal variable and that K is asymptotically distributed as χ_1^2 .

Even though SS and ST are not conditionally independent, the variables ST and $SS - ST^2/TT$ are so asymptotically. This follows because, conditionally on x_2 and t_{22}^P , the normally distributed variable $x_1(a + rx_2) + rz_P \sqrt{t_{22}^P}$ is a linear combination of the standard normal variables x_1 and z_P that partially constitute the asymptotically chi-squared variable SS . These properties led Moreira (2003) to suggest that the distribution of the statistics LR and LR₀, which are deterministic functions of SS , ST , and TT , conditional on a given value of TT , say tt , can be estimated by a simulation experiment in which ST and $SS - ST^2/TT$ are generated as independent variables distributed respectively as $N(0, tt)$ and χ_{l-k-1}^2 . The variable SS is then generated by combining these two variables, replacing TT by tt . Such an experiment has a bootstrap interpretation that we develop in the next section.

It may be helpful to make explicit the link between the quantities SS , ST , and TT , defined in (16), and the vectors \mathbf{S} and \mathbf{T} used in Andrews, Moreira, and Stock (2006). For the simplified model given by (32) and (33), these vectors can be expressed as

$$\mathbf{S} = \mathbf{W}^\top \mathbf{v}_1 \text{ and } \mathbf{T} = \frac{1}{r} \mathbf{W}^\top (\mathbf{y}_2 - \rho \mathbf{y}_1).$$

It is straightforward to check that, with these definitions, $\mathbf{S}^\top \mathbf{S}$, $\mathbf{S}^\top \mathbf{T}$, and $\mathbf{T}^\top \mathbf{T}$ are what the expressions SS , ST , and TT would become if the three quadratic forms in the right panel of (34) were replaced by their asymptotic limits.

The results (35) are independent of b . This shows that, for local alternatives of the sort used in conventional asymptotic theory, no test statistic that depends only on the six quantities (7) can have asymptotic power greater than asymptotic size under weak-instrument asymptotics. However, if instead we consider fixed alternatives, with parameter β independent of n , then the expressions do depend on β .

For notational ease, denote the three right-hand sides in (35) by Y_{11} , Y_{12} , and Y_{22} , respectively. Then it can be seen that the weak-instrument results become

$$\begin{aligned} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 &= Y_{11} + 2\beta Y_{12} + \beta^2 Y_{22}, & n^{-1} \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 &\stackrel{a}{=} 1 + 2\beta\rho + \beta^2, \\ \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 &= Y_{12} + \beta Y_{22}, & n^{-1} \mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 &\stackrel{a}{=} \rho + \beta, \\ \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 &= Y_{22}, & n^{-1} \mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 &\stackrel{a}{=} 1. \end{aligned} \quad (39)$$

Notice that, if we specialize the above results, letting β be $O(n^{-1/2})$ and letting a be $O(n^{1/2})$, then we obtain the conventional strong-instrument results (34).

We have not written down the weak-instrument asymptotic expression for the Wald t statistic given in (6), because it is complicated and not very illuminating. Suffice it to say that it depends nontrivially on the parameters a and ρ , as does its distribution. Consequently, the statistic t is not asymptotically pivotal. Indeed, in the terminology of Dufour (1997), it is not even boundedly pivotal, by which we mean that rejection probabilities of tests based on it cannot be bounded away from one. We will see this explicitly in a moment.

The estimating equations (4) imply that the IV estimate of β is $\hat{\beta} = \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 / \mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2$. Under weak-instrument asymptotics, we see from (39) that

$$\hat{\beta} \stackrel{a}{=} \beta + Y_{12}/Y_{22}. \quad (40)$$

Since $E(Y_{12}/Y_{22}) \neq 0$, it follows that $\hat{\beta}$ is biased and inconsistent. The square of the t statistic (6) can be seen to be asymptotically equal to

$$\frac{Y_{22}(Y_{12} + \beta Y_{22})^2}{Y_{22}^2 - 2\rho Y_{12} Y_{22} + Y_{12}^2} \quad (41)$$

under weak-instrument asymptotics. Observe that this expression is of the order of β^2 as $\beta \rightarrow \infty$. Thus, for fixed a and ρ , the distributions of t^2 for $\beta = 0$ and $\beta \neq 0$ can be arbitrarily far apart.

For $\beta = 0$, however, the distribution of t^2 , for a and ρ sufficiently close to 0 and 1, respectively, is also arbitrarily far from that with fixed $a \neq 0$ and $\rho \neq 1$. It is this fact that leads to the failure of t^2 to be boundedly pivotal. Let a and $r = (1 - \rho^2)^{1/2}$ be treated as small quantities, and then expand the denominator of expression (41) through the second order in small quantities. Note that, to this order, $\rho = 1 - r^2/2$. Then we see that, to the desired order,

$$\begin{aligned} Y_{12} &= \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + (ax_1 + r\mathbf{v}_1^\top \mathbf{P}_2 \mathbf{v}_2) - \frac{1}{2}r^2 \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1, \text{ and} \\ Y_{22} &= \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 + 2(ax_1 + r\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2) + (a^2 + 2arx_2 + r^2 \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2 - \frac{1}{2}r^2 \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1). \end{aligned}$$

Consequently, to the order at which we are working,

$$Y_{22}^2 - 2\rho Y_{12}Y_{22} + Y_{12}^2 = r^2(\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1)^2.$$

To leading order, the numerator of (41) for $\beta = 0$ is just $(\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1)^3$, and so, in the neighborhood of $a = 0$ and $r = 0$, we have from (41) that

$$t^2 \stackrel{a}{=} \mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_1 / r^2. \quad (42)$$

The numerator here is just a chi-squared variable, but the denominator can be arbitrarily close to zero. Thus the distribution of t^2 can be moved arbitrarily far away from any finite distribution by letting a tend to zero and ρ tend to one.

The points in the parameter space at which $a = 0$ and $\rho = \pm 1$, which implies that $r = 0$, are points at which β is completely unidentified. To see this, consider the DGP from model (32) and (33) that corresponds to these parameter values. The DGP can be written as

$$\mathbf{y}_2 = \mathbf{v}_1, \quad \mathbf{y}_1 = (1 \pm \beta)\mathbf{y}_2. \quad (43)$$

It follows from (40) that $\hat{\beta} = 1 \pm \beta$. This is not surprising, since the second equation in (43) fits perfectly. This fact then accounts for the t statistic tending to infinity.

All the other tests have power that does not tend to 1 when $\beta \rightarrow \infty$ under weak-instrument asymptotics. For K , some algebra shows that

$$K \stackrel{a}{=} \frac{(Y_{12} - \rho Y_{11} + \beta(Y_{22} - Y_{11}) + \beta^2(\rho Y_{22} - Y_{12}))^2}{D(1 + 2\beta\rho + \beta^2)}, \quad (44)$$

where

$$D \equiv Y_{22} - 2\rho Y_{12} + \rho^2 Y_{11} + 2\beta(\rho Y_{22} - (1 + \rho^2)Y_{12} + \rho Y_{11}) + \beta^2(Y_{11} - 2\rho Y_{12} + \rho^2 Y_{22}). \quad (45)$$

As $\beta \rightarrow \infty$, the complicated expression (44) tends to the much simpler limit of

$$\frac{(\rho Y_{22} - Y_{12})^2}{Y_{11} - 2\rho Y_{12} + \rho^2 Y_{22}}. \quad (46)$$

Thus, unlike t^2 , the K statistic does not become unbounded as $\beta \rightarrow \infty$. Consequently, under weak-instrument asymptotics, the test based on K is inconsistent for any nonzero β , in the sense that the rejection probability does not tend to 1 however large the sample size.

A similar result holds for the AR test. It is easy to see that

$$SS \stackrel{a}{=} \frac{Y_{11} + 2\beta Y_{12} + \beta^2 Y_{22}}{1 + 2\beta\rho + \beta^2} \xrightarrow{\beta \rightarrow \infty} Y_{22}, \quad (47)$$

which does not depend on β . Thus, since AR is proportional to SS, we see that the asymptotic distribution of AR does not depend on β under weak-instrument asymptotics.

Similar results also hold for the LR and LR₀ statistics. By an analysis like the one that produced (47), we have that

$$ST \stackrel{a}{\underset{\beta \rightarrow \infty}{\longrightarrow}} \frac{Y_{12} - \rho Y_{11} + \beta(Y_{22} - Y_{11}) + \beta^2(\rho Y_{22} - Y_{12})}{r(1 + 2\beta\rho + \beta^2)} \underset{\beta \rightarrow \infty}{\longrightarrow} \frac{\rho Y_{22} - Y_{12}}{r}, \text{ and}$$

$$TT \stackrel{a}{\underset{\beta \rightarrow \infty}{\longrightarrow}} \frac{D}{r^2(1 + 2\beta\rho + \beta^2)} \underset{\beta \rightarrow \infty}{\longrightarrow} \frac{Y_{11} - 2\rho Y_{12} + \rho^2 Y_{22}}{r^2},$$

where D is given by (45). From (18), therefore,

$$\text{LR}_0 \underset{\beta \rightarrow \infty}{\longrightarrow} \frac{1}{2r^2} \left(Y_{22}(1 - 2\rho^2) + 2\rho Y_{12} - Y_{11} + \right. \\ \left. (Y_{11}^2 - 2Y_{11}Y_{22}(1 - 2\rho^2) - 4\rho Y_{11}Y_{12} + Y_{22}^2 + 4Y_{12}^2 - 4\rho Y_{12}Y_{22})^{1/2} \right).$$

The inconsistency of the LR₀ test follows from the fact that this random variable has a bounded distribution. This is true for the LR test as well, but we will spare readers the details.

We saw above that, when $a = 0$ and $\rho = 1$, the parameter β is unidentified. We expect, therefore, that a test statistic for the hypothesis that $\beta = 0$ would have the same distribution whatever the value of β . This turns out to be the case for the K statistic. If one computes the limit of expression (44) for $a = 0$, $r \rightarrow 0$, the limiting expression is just $(\mathbf{v}_1^\top \mathbf{P}_1 \mathbf{v}_2)^2 / \mathbf{v}_2^\top \mathbf{P}_1 \mathbf{v}_2$, independently of the value of β . Presumably a more complicated calculation would show that the same is true for LR and LR₀.

The result that the AR, K , and LR tests are inconsistent under weak-instrument asymptotics appears to contradict some of the principal results of Andrews, Moreira, and Stock (2006). The reason for this apparent contradiction is that we have made a different, and in our view more reasonable, assumption about the covariance matrix of the disturbances. We assume that the matrix Σ , defined in (3) as the covariance matrix of the disturbances in the structural equation (1) and the reduced form equation (2), remains constant as β varies. In contrast, Andrews, Moreira, and Stock (2006) assumes that the covariance matrix of the reduced form disturbances does so.

In terms of our parametrization, the covariance matrix of the disturbances in the two reduced form equations is

$$\begin{bmatrix} \sigma_1^2 + 2\rho\beta\sigma_1\sigma_2 + \beta^2\sigma_2^2 & \rho\sigma_1\sigma_2 + \beta\sigma_2^2 \\ \rho\sigma_1\sigma_2 + \beta\sigma_2^2 & \sigma_2^2 \end{bmatrix}. \quad (48)$$

This expression depends on β in a nontrivial way. In order for it to remain constant as β changes, both ρ and σ_1 must be allowed to vary. Thus the assumption that (48) is fixed, which was made by Andrews, Moreira, and Stock (2006), implies that (3) cannot remain constant as $|\beta| \rightarrow \infty$.

A little algebra shows that, as $\beta \rightarrow \pm\infty$ with the covariance matrix (48) held fixed, the parameters β and ρ of the observationally equivalent DGP of the model given by

(32) and (33) along with (21) tend to ± 1 and ∓ 1 respectively. Thus the omnipresent denominator $1 + 2\beta\rho + \beta^2$ tends to zero in either of these limits. But it is clear from (39) that this means that the estimate of σ_1^2 from the full model (1) and (2) tends to zero.

Based on the above remark, our view is that it is more reasonable that (3) should remain constant than that (48) should. The parameter ρ is a much more interesting parameter than the correlation in (48). Even when $\rho = 0$, in which case the OLS estimator of β is consistent, the correlation between the two reduced form disturbances tends to ± 1 as $\beta \rightarrow \pm\infty$. Thus we believe that the latter correlation is not a sensible quantity to hold fixed. In any case, it is rather disturbing that something as seemingly innocuous as the parametrization of the covariance matrix of the disturbances can have profound consequences for the analysis of power when the instruments are weak.

5. Bootstrapping the Test Statistics

There are several ways to bootstrap the non-exact test statistics that we have been discussing (Wald, K , and LR). In this section, we discuss five different parametric bootstrap procedures, three of which are new. We obtain a number of interesting theoretical results. We also briefly discuss how to convert the new procedures into semiparametric bootstraps. In the next section, we will see that two of the new procedures, and one of them in particular, perform extremely well when used with all three of the test statistics.

For all the statistics, we perform B bootstrap simulations and calculate the bootstrap P value as

$$\hat{p}^*(\hat{\tau}) = \frac{1}{B} \sum_{j=1}^B I(\tau_j^* > \hat{\tau}), \quad (49)$$

where $\hat{\tau}$ denotes the actual test statistic and τ_j^* denotes the statistic calculated using the j^{th} bootstrap sample.

Dufour (1997) makes it clear that bootstrapping is not in general a cure for the difficulties associated with the Wald statistic t . However, since the Wald statistic is still frequently used in practice when there is no danger of weak instruments, it is interesting to look at the performance of the bootstrapped Wald test when instruments are strong. When they are weak, we confirm Dufour's result about the ineffectiveness of bootstrapping. In the context of our new bootstrap methods, this manifests itself in an almost complete loss of power, for reasons that we analyze.

Since the K statistic is asymptotically pivotal under weak-instrument asymptotics, it should respond well to bootstrapping, at least under the null. The LR statistic is not asymptotically pivotal, but, as shown by Moreira (2003), a conditional LR test gives the asymptotically pivotal statistic we call CLR. As we explain below, the implementation of this conditional likelihood ratio test is in fact a form of bootstrapping, so that bootstrapping the conditional test is a sort of double bootstrap, and thus much more computationally intensive than the other bootstrap tests we consider.

For any bootstrapping procedure, the first task, and usually the most important one, is to choose a suitable bootstrap DGP; see Davidson and MacKinnon (2006a). An obvious but important point is that the bootstrap DGP must be able to handle both of the endogenous variables, that is, \mathbf{y}_1 and \mathbf{y}_2 . A straightforward, conventional approach is to estimate the parameters β , γ , $\boldsymbol{\pi}$, σ_1 , σ_2 , and ρ of the model specified by (1), (2), and (3), and then to generate simulated data using these equations with the estimated parameters.

However, the conventional approach estimates more parameters than it needs to. The bootstrap DGP should take advantage of the fact that the simple model specified by (32) and (33) can generate statistics with the same distributions as those generated by the full model. Equation (32) becomes especially simple when the null hypothesis is imposed: It says simply that $\mathbf{y}_1 = \mathbf{u}_1$. If this approach is used, then only the parameters a and ρ need to be estimated. In order to estimate a , we may substitute an estimate of $\boldsymbol{\pi}$ into the definition (24) with an appropriate scaling factor to take account of the fact that a is defined for DGPs with unit disturbance variances.

Since we have assumed up to now that the disturbances of our model are Gaussian, it is appropriate to use a parametric bootstrap in which the disturbances are normally distributed. In practice, however, investigators will often be reluctant to make this assumption. At the end of this section, we therefore discuss semiparametric bootstrap techniques that resample the residuals. In the next section, we will see that the performance of the parametric and semiparametric bootstraps appears to be very similar when the disturbances are in fact normally distributed.

We investigate five different ways of estimating the parameters ρ and a . To estimate ρ , we just need residuals from equations (32) and (33), or, in the general case, (1) and (2). To estimate a , we need estimates of the vector $\boldsymbol{\pi}_1$ from the reduced-form equation (33), or from (2), along with residuals from that equation. If $\ddot{\mathbf{u}}_1$ and $\ddot{\mathbf{u}}_2$ denote the two residual vectors, and $\hat{\boldsymbol{\pi}}_1$ denotes the estimate of $\boldsymbol{\pi}_1$, then our estimates are

$$\hat{\rho} = \frac{\ddot{\mathbf{u}}_1^\top \ddot{\mathbf{u}}_2}{(\ddot{\mathbf{u}}_1^\top \ddot{\mathbf{u}}_1 \ddot{\mathbf{u}}_2^\top \ddot{\mathbf{u}}_2)^{1/2}}, \text{ and} \quad (50)$$

$$\hat{a} = \sqrt{n \hat{\boldsymbol{\pi}}_1^\top \mathbf{W}_1^\top \mathbf{W}_1 \hat{\boldsymbol{\pi}}_1 / \ddot{\mathbf{u}}_2^\top \ddot{\mathbf{u}}_2}. \quad (51)$$

Existing methods, and the new ones that we propose, use various estimates of $\boldsymbol{\pi}_1$ and various residual vectors.

The simplest way to estimate ρ and a is probably to use the restricted residuals

$$\tilde{\mathbf{u}}_1 = \mathbf{M}_Z \mathbf{y}_1 = \mathbf{M}_W \mathbf{y}_1 + \mathbf{P}_1 \mathbf{y}_1,$$

which, in the case of the simple model, are just equal to \mathbf{y}_1 , along with the OLS estimates $\hat{\boldsymbol{\pi}}_1$ and OLS residuals $\hat{\mathbf{u}}_2$ from the FWL regression

$$\mathbf{M}_Z \mathbf{y}_2 = \mathbf{W}_1 \boldsymbol{\pi}_1 + \mathbf{u}_2. \quad (52)$$

We call this widely-used method the RI bootstrap, for “Restricted, Inefficient”. It can be expected to work much better than the pairs bootstrap that was studied by Freedman (1983) and better than other parametric procedures that do not impose the null hypothesis.

As the name implies, the problem with the RI bootstrap is that $\hat{\pi}_1$ is not an efficient estimator. That is why Kleibergen (2002) did not use $\hat{\pi}_1$ in constructing the K statistic. Instead, he used the estimates $\tilde{\pi}_1$ from equation (10). It can be shown that these estimates are asymptotically equivalent to the ones that would be obtained by using 3SLS or FIML on the system consisting of equations (1) and (2). The estimated vector of disturbances from equation (10) is not the vector of OLS residuals but rather the vector $\tilde{\mathbf{u}}_2 = \mathbf{M}_Z \mathbf{y}_2 - \mathbf{W}_1 \tilde{\pi}_1$.

Instead of equation (10), it may be more convenient to run the regression

$$\mathbf{y}_2 = \mathbf{W}_1 \pi_1 + \mathbf{Z} \pi_2 + \delta \mathbf{M}_Z \mathbf{y}_1 + \text{residuals}. \quad (53)$$

This is just the reduced form equation augmented by the residuals from restricted estimation of the structural equation. The vector $\tilde{\mathbf{u}}_2$ is equal to the vector of OLS residuals from regression (53) plus $\hat{\delta} \mathbf{M}_Z \mathbf{y}_1$. We call the bootstrap that uses $\tilde{\mathbf{u}}_1$, $\tilde{\pi}_1$, and $\tilde{\mathbf{u}}_2$ the RE bootstrap, for “Restricted, Efficient”.

Two other bootstrap methods do not impose the restriction that $\beta = 0$ when estimating ρ and a . For the purposes of testing, it is a bad idea not to impose this restriction, as we argued in Davidson and MacKinnon (1999). However, it is quite inconvenient to impose restrictions when constructing bootstrap confidence intervals, and since confidence intervals are implicitly obtained by inverting tests, it is of interest to see how much harm is done by not imposing the restriction.

The UI bootstrap, for “Unrestricted, Inefficient”, uses the unrestricted residuals $\hat{\mathbf{u}}_1$ from IV estimation of (1), along with the estimates $\hat{\pi}_1$ and residuals $\hat{\mathbf{u}}_2$ from OLS estimation of (2). The UE bootstrap, for “Unrestricted, Efficient”, also uses $\hat{\mathbf{u}}_1$, but the other quantities come from the artificial regression

$$\mathbf{M}_Z \mathbf{y}_2 = \mathbf{W}_1 \pi_1 + \delta \hat{\mathbf{u}}_1 + \text{residuals}, \quad (54)$$

which is similar to regression (10). Of course, a regression analogous to (53) could be used instead of (54). A fifth bootstrap method will be proposed after we have obtained some results on which it depends.

It is possible to write the estimates of a and ρ used by all four of these bootstrap schemes as functions solely of the six quantities (7). This makes it possible to program the bootstrap very efficiently. Because many of the functions are quite complicated, we will spare readers most of the details. However, we need the following results for the RE bootstrap:

$$\tilde{\rho} = \frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2 + \frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1}{\left((\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1 + \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1) \left(\mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 + \left(\frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \right)^2 \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \right) \right)^{1/2}}, \quad (55)$$

and

$$\tilde{a}^2 = \frac{\mathbf{y}_2^\top \mathbf{P}_1 \mathbf{y}_2 - 2\mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_2 \frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} + \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \left(\frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \right)^2}{\mathbf{y}_2^\top \mathbf{M}_W \mathbf{y}_2 + \mathbf{y}_1^\top \mathbf{P}_1 \mathbf{y}_1 \left(\frac{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_2}{\mathbf{y}_1^\top \mathbf{M}_W \mathbf{y}_1} \right)^2}. \quad (56)$$

Davidson and MacKinnon (1999) show that the size distortion of bootstrap tests may be reduced by use of a bootstrap DGP that is asymptotically independent of the statistic that is bootstrapped. In general, this is true only for bootstrap DGPs that are based on efficient estimators. Thus it makes sense to use the efficient estimator $\tilde{\boldsymbol{\pi}}_1$ rather than the inefficient estimator $\hat{\boldsymbol{\pi}}_1$ in order to estimate a , and, via the reduced-form residuals, ρ . Either restricted or unrestricted residuals from (1) can be used as the extra regressor in estimating $\boldsymbol{\pi}_1$ without interfering with the desired asymptotic independence, but general considerations of efficiency suggest that restricted residuals are the better choice. Thus we would expect that, when conventional asymptotics yield a good approximation, the best choice for bootstrap DGP is RE.

Under weak-instrument asymptotics, things are rather different. We use the results of (35) and the right-hand results of (34) to see that, with data generated by the model (32) and (33) under the null hypothesis,

$$\tilde{\sigma}_1^2 \stackrel{a}{=} 1, \quad \tilde{\sigma}_2^2 \stackrel{a}{=} 1, \quad \text{and} \quad \tilde{\rho} \tilde{\sigma}_1 \tilde{\sigma}_2 \stackrel{a}{=} \rho.$$

Thus the RE bootstrap estimator $\tilde{\rho}$, as defined by (55), is a consistent estimator, as is the estimator used by the RI bootstrap. It can be checked that this result does *not* hold for any of the estimators that use unrestricted residuals from equation (32), since they depend on the inconsistent IV estimate of β ; recall (40).

The weak-instrument asymptotic version of (56) under the null can be seen to be

$$\tilde{a}^2 \stackrel{a}{=} a^2 + 2arx_2 + r^2(x_2^2 + t_{22}^P). \quad (57)$$

Unless $r = 0$, then, \tilde{a}^2 is inconsistent. It is also biased, the bias being equal to $r^2(l-k)$. It seems plausible, therefore, that the bias-corrected estimator

$$\tilde{a}_{\text{BC}}^2 \equiv \max(0, \tilde{a}^2 - (l-k)(1 - \tilde{\rho}^2)) \quad (58)$$

may be better for the purposes of defining the bootstrap DGP. Thus we consider a fifth bootstrap method, REC, for ‘‘Restricted, Efficient, Corrected.’’ It differs from RE in that it uses \tilde{a}_{BC} instead of \tilde{a} . This has the effect of reducing the R^2 of the reduced-form equation in the bootstrap DGP.

For the purposes of an analysis of power, it is necessary to look at the properties of the estimates $\tilde{\rho}$ and \tilde{a}^2 under the alternative, that is, for nonzero β . From (39), we see that

$$\tilde{\sigma}_1^2 \stackrel{a}{=} 1 + 2\beta\rho + \beta^2, \quad \tilde{\sigma}_2^2 \stackrel{a}{=} 1, \quad \text{and} \quad \tilde{\rho} \tilde{\sigma}_1 \tilde{\sigma}_2 \stackrel{a}{=} \rho + \beta,$$

from which we find that

$$\tilde{\rho} \stackrel{a}{=} \frac{\rho + \beta}{(1 + 2\beta\rho + \beta^2)^{1/2}}.$$

As $\beta \rightarrow \infty$, then, we see that $\tilde{\rho} \rightarrow 1$, for all values of a and ρ . For the rate of convergence, it is better to reason in terms of the parameter $r \equiv (1 - \rho^2)^{1/2}$. We have

$$\tilde{r}^2 = 1 - \tilde{\rho}^2 \stackrel{a}{=} 1 - \frac{(\rho + \beta)^2}{1 + 2\beta\rho + \beta^2} = \frac{r^2}{1 + 2\beta\rho + \beta^2}.$$

Thus $\tilde{r} = O_p(\beta^{-1})$ as $\beta \rightarrow \infty$.

The calculation for \tilde{a}^2 is a little more involved. From (56) and (39), we find that

$$\begin{aligned} \tilde{a}^2 &\stackrel{a}{=} Y_{22} - \frac{2(\rho + \beta)}{1 + 2\beta\rho + \beta^2}(Y_{12} + \beta Y_{22}) + \left(\frac{\rho + \beta}{1 + 2\beta\rho + \beta^2} \right)^2 (Y_{11} + 2\beta Y_{12} + \beta^2 Y_{22}) \\ &\stackrel{a}{=} \frac{1}{(1 + 2\beta\rho + \beta^2)^2} ((1 + \beta\rho)Y_{22} - 2(\rho + \beta)(1 + \beta\rho)Y_{12} + (\rho + \beta)^2 Y_{11}). \end{aligned}$$

Clearly, this expression is of the order of β^{-2} in probability as $\beta \rightarrow \infty$, so that $\tilde{a} \rightarrow 0$, again for all a and ρ . In fact, it is clear that $\tilde{a} = O_p(\beta^{-1})$ as $\beta \rightarrow \infty$, from which we conclude that \tilde{a} and \tilde{r} tend to zero at the same rate as $\beta \rightarrow \infty$, as in the calculation that led to (42).

These results can be understood intuitively by considering (32) and (33). Estimation of ρ uses residuals which for that model are just the vector $\mathbf{y}_1 = \beta\mathbf{y}_2 + \mathbf{v}_1$. For large β , this residual vector is almost collinear with \mathbf{y}_2 , and so also with the residual vector $\tilde{\mathbf{u}}_2$. The estimated correlation coefficient therefore tends to 1. Similarly, when \mathbf{y}_1 is introduced as an extra regressor for the estimation of a , it is highly collinear with the dependent variable and explains almost all of it, leaving no apparent explanatory power for the weak instruments.

For large β then, the RE bootstrap DGP is characterized by parameters a and ρ close to 0 and 1, respectively. As we saw [near the end](#) of the last section, at this point in the parameter space β is unidentified, and the Wald statistic has an unbounded distribution. These facts need not be worrisome for the bootstrapping of statistics that are asymptotically pivotal with weak instruments, but they mean that the bootstrap version of the Wald test, like the Kleibergen and LR tests, is inconsistent, having a probability of rejecting the null hypothesis that does not tend to one as $\beta \rightarrow \infty$.

To see this, we make use of expression (42) to see that the distribution of the Wald statistic t^2 , for a and r small and of the same order and $\beta = 0$, is of order r^{-2} . For large β , therefore, the distribution of the bootstrap Wald statistic, under the null, is of order \tilde{r}^{-2} , which we have just seen is the same order as β^2 . But the distribution of the Wald statistic itself for large β is also of order β^2 , unlike the K and LR statistics. Although the distribution of the actual statistic t^2 for large β and that of the bootstrap statistic $(t^*)^2$ are not the same, and are unbounded, the distributions of t^2/β^2 and $(t^*)^2/\beta^2$ are of order unity in probability, and, having support on the whole real line,

they overlap. Thus the probability of rejection of the null by the bootstrap test does not tend to 1 however large β may be.

This conclusion, which is borne out by the simulation experiments of the [next section](#), bears some discussion. In Horowitz and Savin (2000), it is pointed out that, unless one is working with pivotal statistics, it is not in general possible to define an empirically relevant definition of the power of a test that does not have true level equal to its nominal level. They conclude that the best measure in practice is the rejection probability of a well-constructed bootstrap test.

In Davidson and MacKinnon (2006b), we point out that, even for well-constructed bootstrap tests, ambiguity remains in general. Only when the bootstrap DGP is asymptotically independent of the asymptotically pivotal statistic being bootstrapped can level adjustment be performed unambiguously on the basis of the DGP in the null hypothesis of which the parameters are the probability limits of the estimators used to define the bootstrap DGP. This result, as proved, applies only to the parametric bootstrap, and, more importantly here, to cases in which these estimators have non-random probability limits. But, as we have seen, that is not the case here. It seems therefore that there is no theoretically satisfying measure of the power of tests for which the bootstrap DGP is asymptotically nonrandom. It is therefore pointless to try to refine our earlier result for the Wald test, whereby we learn merely that its bootstrap version is inconsistent.

Because the Wald statistic is not boundedly pivotal, a test based on it has size equal to one, as shown by Dufour (1997). Dufour also draws the conclusion that no Wald-type confidence set based on a statistic that is not at least boundedly pivotal can be valid, whether or not the confidence set is constructed by bootstrapping. If, however, instead of using a conventional bootstrap confidence set, we invert the bootstrap Wald test to obtain a confidence set that contains all parameter values which are not rejected by a bootstrap Wald test, we may well obtain confidence sets with a level of less than one, since, on account of the inconsistency of the bootstrap test, unbounded confidence sets can arise with positive probability.

We mentioned [earlier](#) that Moreira’s conditional LR test has a bootstrap interpretation. We may consider the variable TT defined in (16) as a random variable on which a bootstrap distribution is conditioned. In fact, as can be seen from (38) and (57), TT is equivalent under weak-instrument asymptotics to \tilde{a}^2/r^2 . Given a value for TT , Moreira shows that the statistics LR and LR_0 are asymptotically pivotal. Thus, rather than estimating both a and ρ and using the estimates to generate bootstrap versions of the six sufficient statistics (7), we can evaluate TT and then generate “bootstrap” versions of the two (conditionally) sufficient statistics SS and ST based on their asymptotic conditional distributions, as discussed [earlier](#). From this we obtain a conditional empirical distribution which is used to compute a P value for the CLR test in the usual way.

This procedure is not quite a real bootstrap, although it is almost as computationally intensive as a fully parametric bootstrap based on simulating the six quantities, because the “bootstrap” conditional distributions of SS and ST are not known ex-

actly. Instead, they are approximated on the basis of the distributions when the contemporaneous covariance matrix is known.

Of course, it is possible to bootstrap the CLR test, but this amounts to a form of double bootstrap, which is expensive and requires some care to implement correctly. Suppose we perform s simulations to calculate a CLR test P value \hat{p} . Since this P value is $1/s$ times the number of simulations for which the simulated LR statistic, conditional on TT , exceeds the actual LR statistic, it can take on only $s + 1$ different values. Now imagine bootstrapping this entire procedure. If we also use s simulations to compute each bootstrap statistic p_j^* , then p_j^* can only take on the same $s + 1$ values as \hat{p} . Thus the actual and bootstrap statistics will be equal with probability approximately $1/(s + 1)$. To avoid this problem, we either need to make s very large or use a different number of simulations, say s^* , for the bootstrap statistics. The bootstrap P value is computed by a modified version of (49), with $I(p_j^* < \hat{p})$ replacing $I(\tau_j^* > \hat{\tau})$.

Some remarks concerning bootstrap validity are in order at this point. If the statistic that is bootstrapped is asymptotically pivotal, then the bootstrap is valid asymptotically in the sense that the difference between the bootstrap distribution and the distribution under the true DGP, provided the latter satisfies the null hypothesis, converges to zero as the sample size tends to infinity; see, among many others, Davidson and MacKinnon (2006a). The bootstrap provides higher-order refinements if, in addition, the bootstrap DGP consistently estimates the true DGP under the null; see Beran (1988). Yet another level of refinement can be attained if the statistic bootstrapped is asymptotically independent of the bootstrap DGP; see Davidson and MacKinnon (1999). All of these requirements are satisfied by any of the statistics considered here under strong-instrument asymptotics if the RE or REC bootstrap is used.

Besides the AR statistic, only the K statistic and CLR test P value are asymptotically pivotal with weak instruments, and so it is only for them that we can conclude without further ado that the bootstrap is valid with weak-instrument asymptotics. However, even if the statistic bootstrapped is not asymptotically pivotal, the bootstrap may still be valid if the bootstrap DGP consistently estimates the true DGP under the null. But with weak instruments this is true of no conceivable bootstrap method, since there is no consistent estimate of the parameter a . Consequently, however large the sample size may be, the bootstrap distributions of the Wald statistic and the LR statistic are different from their true distributions.

Although our discussion has focused on parametric bootstrap procedures, we do not necessarily recommend that they should be used in practice, since the assumption of Gaussian disturbances may often be uncomfortably strong. All of the parametric bootstrap procedures that we have discussed have semiparametric analogs which do not require that the disturbances should be normally distributed. We will discuss only the RE and REC bootstraps, partly because they are new, partly because they will be seen in the next section to work very well, and partly because it will be obvious how to construct semiparametric analogs of the other procedures.

For the semiparametric RE bootstrap, we first estimate equation (1) under the null hypothesis to obtain restricted residuals $\mathbf{M}_Z \mathbf{y}_1$. We then run regression (10) or, equivalently, regression (53). The residual vector $\tilde{\mathbf{u}}_2$ that we want to resample from is the vector of residuals from the regression plus $\hat{\delta} \mathbf{M}_Z \mathbf{y}_1$. The bootstrap DGP is then

$$\begin{aligned} \mathbf{y}_1^* &= \mathbf{u}_1^* \\ \mathbf{y}_2^* &= \mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1 + \mathbf{u}_2^*, \quad [\mathbf{u}_1^* \ \mathbf{u}_2^*] \sim \text{EDF}[\tilde{\mathbf{u}}_1 \ \tilde{\mathbf{u}}_2]. \end{aligned} \quad (59)$$

Thus the bootstrap disturbances are resampled from the joint EDF of the two residual vectors. This preserves the sample correlation between them.

For the REC bootstrap, we need to use a different set of fitted values in the reduced-form equation. To do so, we first compute \tilde{a}^2 using either the formula (56) or, more conveniently,

$$\tilde{a}^2 = \frac{\tilde{\boldsymbol{\pi}}_1^\top \mathbf{W}_1^\top \mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1}{\tilde{\mathbf{u}}_2^\top \tilde{\mathbf{u}}_2 / n},$$

which is the square of (51) evaluated at the appropriate values of $\boldsymbol{\pi}_1$ and \mathbf{u}_2 . Then we calculate \tilde{a}_{BC}^2 from (58). The bootstrap DGP is almost the same as (59), except that the fitted values $\mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1$ are replaced by $\tilde{a}_{\text{BC}} / \tilde{a}$ times $\mathbf{W}_1 \tilde{\boldsymbol{\pi}}_1$. This reduces the length of the vector of fitted values somewhat. The fitted values actually shrink to zero in the extreme case in which $\tilde{a}_{\text{BC}} = 0$.

The quantity \tilde{a}^2 is very closely related to the “test statistic” for weak instruments recently proposed by Stock and Yogo (2005). When it is large, the instruments are almost certainly not weak, and even asymptotic inference should be reasonably reliable. When it is very small, however, many tests are likely to overreject severely, and those that do not are likely to be seriously lacking in power. There is evidence on these points in the next section.

6. Simulation Evidence

In this section, we report the results of a large number of simulation experiments with data generated by the simplified model (32) and (33). All experiments had 100,000 replications for each set of parameter values. In most cases, we used fully parametric bootstrap DGPs, but we also investigated the semiparametric ones discussed at the end of the preceding section.

In most of the experiments, $n = 50$. A different choice for n , but for the same value(s) of a , would have had only a modest effect on the results for most of the bootstrap tests. Many of the asymptotic tests are quite sensitive to n , however, and we later provide some evidence on this point. For the base cases, we consider every combination of $a = 2$ and $a = 8$ with $\rho = 0.1$ and $\rho = 0.9$. The limiting R^2 of the reduced-form regression (33) is $a^2 / (n + a^2)$. Thus, when $a = 2$, the instruments are very weak, and when $a = 8$, they are moderately strong. When $\rho = 0.1$, there is not much correlation between the structural and reduced-form disturbances, and when $\rho = 0.9$, there is a great deal of correlation.

All our results are presented graphically. Within each figure, the horizontal scale is always the same, but the vertical scale often changes, because otherwise it would be impossible to see many important differences between alternative tests and bootstrap methods. Readers should check the vertical scales carefully when comparing results in different panels of the same figure.

Figures 1 through 4 concern the properties of asymptotic tests. [Figure 1](#) shows the rejection frequencies for the Wald, K , LR, and CLR tests (with $s = 199$) for the four base cases as functions of $l - k$. These are all increasing functions, so test performance generally deteriorates as the number of over-identifying restrictions, $l - k - 1$, increases. Of particular note are the extremely poor performance of the Wald test when $\rho = 0.9$ and the surprisingly poor performance of the LR test when $\rho = 0.1$. It is also of interest that the Wald test underrejects severely when ρ , a , and $l - k$ are all small. This is a case that is rarely investigated in simulation experiments.

[Figure 2](#) shows rejection frequencies as functions of ρ for four values of a . In this and all subsequent figures, $l - k = 7$, so that there are six overidentifying restrictions. As expected, performance improves dramatically as a increases. The Wald test is extremely sensitive to ρ . The others, especially K and CLR, are much less so. Only when ρ is small and a is large does the Wald test perform at all well.

In [Figure 3](#), we consider values of n between 20 and 1280 that increase by factors of approximately $\sqrt{2}$. This figure, which appears to be novel, makes it clear that the mediocre performance of K and CLR evident in the first two figures is a consequence of using $n = 50$. The performance of both these tests always improves dramatically as n increases. Recall that K and CLR are asymptotically pivotal under both weak-instrument and conventional asymptotics. Thus it is not surprising that they can safely be used as asymptotic tests when the sample size is large but the instruments are weak. The performance of LR also improves as n increases, but it continues to overreject, sometimes very severely, even for large values of n . The Wald test is the least sensitive to n , but its performance often deteriorates as n increases.

[Figure 4](#) shows what happens as a varies. We consider values from $a = 1$ to $a = 90.51$ that increase by factors of $\sqrt{2}$. As expected, the performance of the Wald and LR tests improves dramatically as a increases. There is little effect on K and only a modest effect on CLR. The latter tests would instead benefit from a larger sample size, holding a constant.

Figures 5, 6, and 7 concern the properties of parametric bootstrap tests under the null hypothesis. In all cases, $B = 199$. [Figure 5](#) shows rejection frequencies for the Wald, K , and LR tests as a function of a for two values of ρ , and [Figure 6](#) shows rejection frequencies for the same tests as a function of ρ for two values of a . Results are presented for the five different bootstrap DGPs that were discussed in Section 5. Procedures with an “R” use restricted estimates of the structural equation, while procedures with a “U” use unrestricted estimates. Procedures with an “E” use efficient estimates of the reduced-form equation, while procedures with an “I” use inefficient ones. The REC procedure bias-corrects the estimate of a .

For the Wald test, our new RE and REC bootstraps perform reasonably well, although

they do lead to severe underrejection in some cases. The other three methods lead to very severe overrejection when ρ is not small. This is consistent with what other authors have observed when investigating conventional bootstrap methods applied to IV estimation; see, for example, Moreira, Porter, and Suarez (2004).

For the K test, all methods work very well, although REC is arguably the best of a remarkably good bunch. For the LR test, REC is unquestionably the best method in three out of four of the experiments, and it is arguably the best in the fourth, where everything except UI performs quite well. Using REC leads to only modest overrejection in the worst cases.

Figure 7 deals with the CLR test, bootstrapping which is about as computationally intensive (and expensive) as a double bootstrap.² We used $s = 199$ and $s^* = 299$ in these experiments; these numbers were defined above in our discussion of the bootstrap interpretation of the CLR test. Using different values would have produced somewhat different results. Results for the CLR test are quite similar to those for the LR test. Bootstrapping generally improves its performance, and the REC bootstrap is clearly the procedure of choice. Its only defect seems to be a very modest tendency to overreject when a is small.

Figures 8 and 9 concern power. Because there is no point comparing the powers of tests that do not perform reliably under the null, we consider only the REC, RE, and RI bootstraps. The last of these is included only for purposes of comparison, because it is impossible to justify for the Wald test. In these figures, we present results for the AR test (not bootstrapped, since it is exact) as well as for the three bootstrap tests. To avoid the power loss associated with small values of B , we set $B = 999$ in these experiments. Since we would also have had to use large values of s and s^* , for the same reason, including the CLR test in these experiments would have been prohibitively expensive, and we did not do so.

Figure 8 deals with the weak-instrument case in which $a = 2$. In the left-hand panels, $\rho = 0.1$, and in the right-hand panels, $\rho = 0.9$. No test has good power properties. When $\rho = 0.1$, AR is probably the best test to use. It is often more powerful than K , the only other test that is reliable under the null, and it is rarely much less powerful. The apparently much greater power of Wald and LR with the RI bootstrap is spurious, because these tests are not reliable under the null. Although it is true that the Wald test actually underrejects here, we would not know in practice that ρ is very small, and for larger values of ρ the Wald test with the RI bootstrap overrejects severely.

When $\rho = 0.9$, all the tests have very strange-looking power functions. They have far less power against positive values of β than against negative ones. The power function for K has a curious dip for some negative values of β , and in this region K can be much less powerful than AR; see Poskitt and Skeels (2007), as well as Stock,

² A recent paper by Hillier (2006) derives the exact conditional distribution of the LR statistic under the assumption that the disturbance covariance matrix Σ is known. Critical values for the test can be obtained from this conditional distribution by numerical methods. Use of these critical values might considerably reduce the computational burden of bootstrapping the CLR test.

Wright, and Yogo (2002), whose simulation results are not directly comparable to ours, since they use the same parametrization as Andrews, Moreira, and Stock (2006) and assume that Σ is known. The minimum of the power function for the Wald test also occurs well to the left of $\beta = 0$. The REC version of the LR test appears to be the best of a bad lot. It seems to be reliable under the null, and it has noticeably more power than AR against some alternatives.

Figure 9 deals with the case in which $a = 8$, so that the instruments are moderately strong. All the tests now perform very much better, and both K and LR often outperform AR. However, when $\rho = 0.9$, the power function for K once again has a curious dip for some negative values of β , and the power function for Wald has its minimum noticeably to the left of $\beta = 0$. The power function for K when $\rho = 0.1$ is even stranger; power actually declines as the absolute value of β increases beyond a certain point. The LR test does not have any of these problems, and it is quite reliable, at least when bootstrapped using REC. Thus there appears to be a fairly strong case for using the REC bootstrap version of LR. When a is reasonably large, we would expect the CLR test, bootstrapped if n is small, to have power similar to that of the REC bootstrap version of LR. See Andrews, Moreira, and Stock (2006).

As we mentioned at the end of the last section, it is probably better in practice to use a semiparametric rather than a fully parametric bootstrap, because the normality assumption is likely to be false. We therefore performed a number of experiments using semiparametric versions of the REC and RE bootstraps. In every case, the semiparametric and fully parametric bootstraps yielded very similar results. Figures 10 and 11 show how the parametric and semiparametric bootstrap tests perform as functions of the sample size. They show rejection frequencies under the null as a function of the sample size for the REC and RE bootstraps, respectively, both parametric and semiparametric, for four sets of parameter values. The similarity of the two sets of results is striking. Note that the same random numbers were used to generate the underlying data, but different ones were used for bootstrapping.

Figure 10, which shows the results for the REC bootstrap, makes it clear that the decision to focus on the case $n = 50$ in most of our experiments is not entirely inconsequential. The K test works very well indeed for all but the smallest sample sizes, as does the CLR test, although the latter generally does not perform as well when n is small. The Wald test underrejects for the smaller sample sizes in three cases out of four, but its performance improves as n increases. When $a = 2$, the LR test overrejects moderately for small sample sizes and underrejects moderately for large ones. It just happens to perform very well for $n = 50$, so that our previous results may paint a somewhat overoptimistic picture.

Figure 11 is similar to Figure 10, but it shows results for the RE bootstrap. Once again, we see that the fully parametric and semiparametric bootstraps produce almost identical results with normally distributed disturbances. At least in some cases, the RE procedure is substantially inferior to the REC one. In particular, the LR test overrejects quite severely when $a = 2$ and $\rho = 0.1$, and the Wald test now performs less well for cases where a is small and n is large. However, for the K and CLR tests, performance is once again excellent for large n , although n needs to be somewhat

larger than for the REC bootstrap.

It emerges clearly from these last two figures that the rejection frequencies of the RE bootstrap Wald and LR tests do not seem to converge to the nominal level as $n \rightarrow \infty$ when $a = 2$, whereas those of the K and CLR tests do so. This is in accord with our discussion of the [previous section](#). We echo Moreira, Porter, and Suarez (2004), however, in noting that it is remarkable that the bootstrap Wald and LR tests perform as well as they do.

7. Concluding Remarks

We have provided a detailed analysis of the properties of several tests for the coefficient of a single right-hand-side endogenous variable in a linear structural equation estimated by instrumental variables. First, we showed that the Student's t (or Wald) statistic, Kleibergen's K statistic, and the LR statistic can be written as functions of six random quantities. The Anderson-Rubin test is also a function of two of these six quantities. Using these results, we obtained explicit expressions for the asymptotic distributions of all the test statistics under both conventional and weak-instrument asymptotics.

Under weak-instrument asymptotics, we found that none of the test statistics can have any real asymptotic power against local alternatives. Even when the alternative is fixed, AR, K , and LR are not consistent tests under weak-instrument asymptotics. The t test has very different properties, however. It is unbounded as $\beta \rightarrow \infty$, so that it appears to be consistent. But it is also unbounded as certain parameters of the DGP tend to limiting values, so that it is not asymptotically pivotal, or even boundedly pivotal. Note that these results depend in an essential way on how the DGP is specified, in particular the disturbance covariance matrix.

We then proposed some new procedures for bootstrapping the three test statistics. Our RE and REC procedures use more efficient estimates of the coefficients of the reduced-form equation than existing procedures and impose the restriction of the null hypothesis. In addition, the REC procedure corrects for the tendency of the reduced-form equation to fit too well. A semiparametric version of this procedure is quite easy to implement. In most cases, the REC bootstrap outperforms the RE bootstrap, which in turn outperforms previously proposed methods. The improvement can be quite dramatic.

The K and CLR tests have excellent performance under the null when bootstrapped using either the RE or REC procedures, even when the sample size is small and the instruments are weak. The LR test also performs very well in most cases, especially when REC is used. Even the Wald test performs quite well when bootstrapped using these procedures, although it sometimes underrejects fairly severely. Interestingly, however, as we show analytically, the RE and REC bootstrap versions of the Wald test are not consistent against fixed alternatives under weak-instrument asymptotics, and they seem to be less powerful than the other tests when the instruments are weak.

All of our theoretical analysis is conducted under the assumption that the disturbances are Gaussian, although some results do not in fact depend on this assumption. To

our knowledge, little work has been done on the properties of tests in the presence of weak instruments and non-Gaussian disturbances. We conjecture that the qualitative features of the tests considered in this paper, both asymptotic and bootstrap, do not greatly depend on the assumption of Gaussianity.

In the light of our results, it is tempting to conclude that, when the number of overidentifying restrictions is large, so that the AR test may suffer significant power loss, the best methods to use are variants of the LR test. When the sample size is large, it appears to be safe to use the conditional LR test without bootstrapping. When the instruments are reasonably strong, it appears to be safe to use the LR test bootstrapped using the REC procedure. However, when the sample size is small and the instruments are weak, it is best to use the bootstrap CLR test, preferably bootstrapped using the REC procedure.

References

- Anderson, T. W. (1984). *An Introduction to Multivariate Statistical Analysis*, 2nd edition, Wiley, New York.
- Anderson, T. W., and H. Rubin (1949). “Estimation of the Parameters of a Single Equation in a Complete Set of Stochastic Equations”, *The Annals of Mathematical Statistics*, 20, 46–63.
- Andrews, D. W. K., M. J. Moreira, and J. H. Stock (2006). “Optimal two-sided invariant similar tests for instrumental variables regression”, *Econometrica*, 74, 715–752.
- Andrews, D. W. K., M. J. Moreira, and J. H. Stock (2007). “Performance of conditional Wald tests in IV regression with weak instruments,” *Journal of Econometrics*, 139, 116–132.
- 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 (1999). “The Size Distortion of Bootstrap Tests”, *Econometric Theory*, 15, 361–376.
- Davidson, R. and J. G. MacKinnon (2004). *Econometric Theory and Methods*, Oxford University Press, New York.
- Davidson, R. and J. G. MacKinnon (2006a). “Bootstrap Methods in Econometrics”, Chp. 23 in *Palgrave Handbook of Econometrics*, Volume 1, eds T. C. Mills and K. Patterson, Palgrave-Macmillan, Basingstoke, 812–838.
- Davidson, R. and J. G. MacKinnon (2006b). “The Power of Bootstrap and Asymptotic Tests”, *Journal of Econometrics*, to appear.
- Dufour, J.-M., (1997). “Some Impossibility Theorems in Econometrics with Applications to Structural and Dynamic Models”, *Econometrica*, 65, 1365–1387.

- Freedman, D. A. (1984). “On bootstrapping two-stage least-squares estimates in stationary linear models”, *Annals of Statistics*, 12, 827–842.
- Hillier, G. (2006). “Exact Critical Value and Power Functions for the Conditional Likelihood Ratio and Related Tests in the IV Regression Model with Known Covariance”, working paper.
- Horowitz, J. L., and N. E. Savin (2000). “Empirically relevant critical values for hypothesis tests,” *Journal of Econometrics*, 95, 375–389.
- Kleibergen, F. (2002). “Pivotal Statistics for Testing Structural Parameters in Instrumental Variables Regression”, *Econometrica*, 70, 1781–1803.
- Moreira, M. J. (2003). “A Conditional Likelihood Ratio Test for Structural Models”, *Econometrica*, 71, 1027–1048.
- Moreira, M. J., J. R. Porter, and G. A. Suarez (2004). “Bootstrap and Higher-Order Expansion Validity when Instruments may be Weak”, NBER Working Paper No. 302.
- Phillips, P. C. B. (1983). “Exact Small Sample Theory in the Simultaneous Equations Model”, in *Handbook of Econometrics*, Vol I, eds Z. Griliches and M. D. Intriligator, North Holland.
- Poskitt, S. S., and C. L. Skeels (2007). “Approximating the distribution of the two-stage least squares estimator when the concentration parameter is small,” *Journal of Econometrics*, 139, 217–236.
- Staiger, D., and J. H. Stock (1997). “Instrumental Variables Regression with Weak Instruments”, *Econometrica*, 65, 557–586.
- Stock, J. H., J. H. Wright, and M. Yogo (2002). “A Survey of Weak Instruments and Weak Identification in Generalized Method of Moments”, *Journal of Business and Economic Statistics*, 20, 518–529.
- Stock, J. H., and M. Yogo (2005). “Testing for Weak Instruments in Linear IV Regression”, in *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, eds D. W. K. Andrews and J. H. Stock, Cambridge University Press, Cambridge, 80–108.

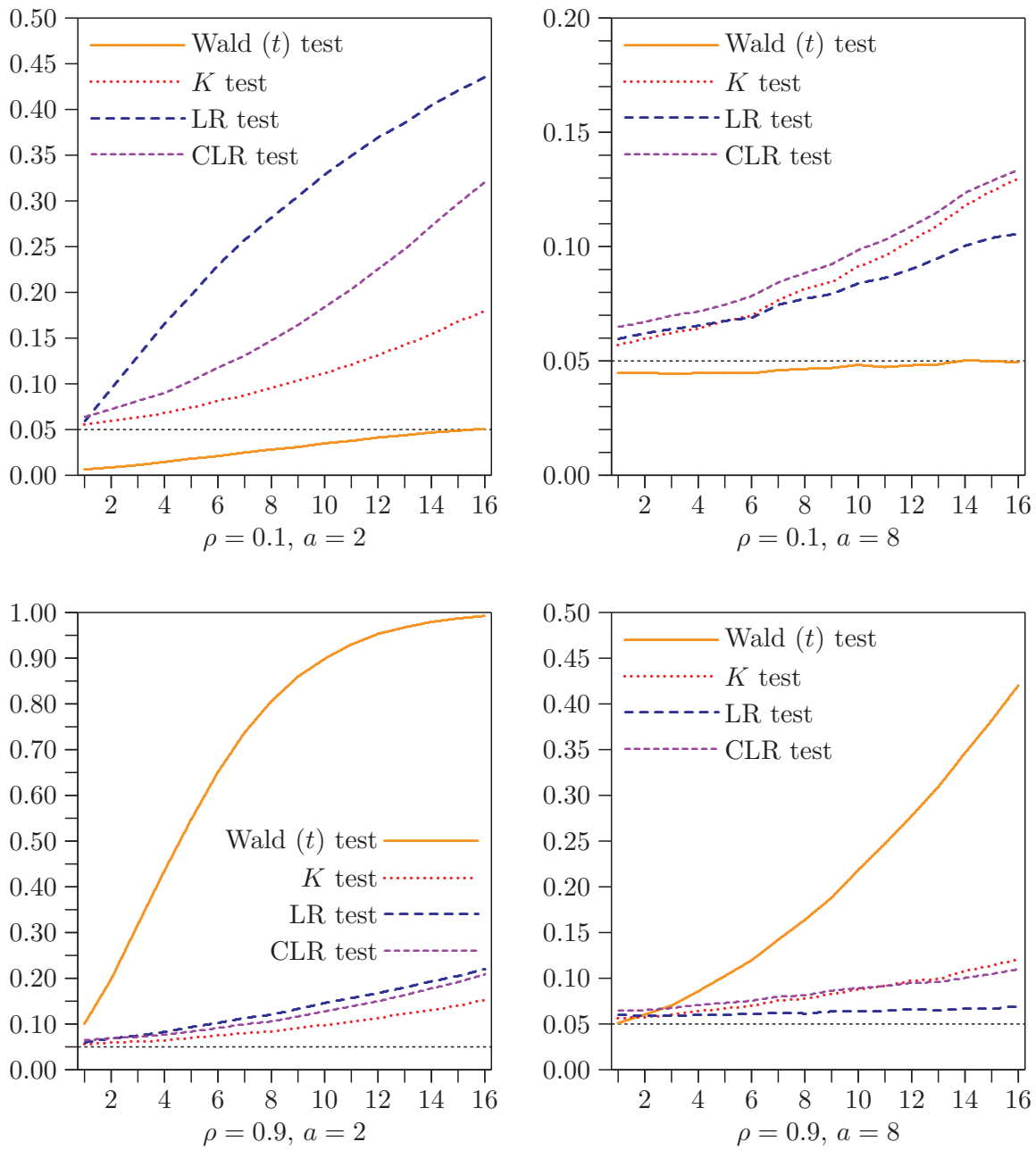


Figure 1. Rejection frequencies for asymptotic tests as functions of $l - k$, $n = 50$

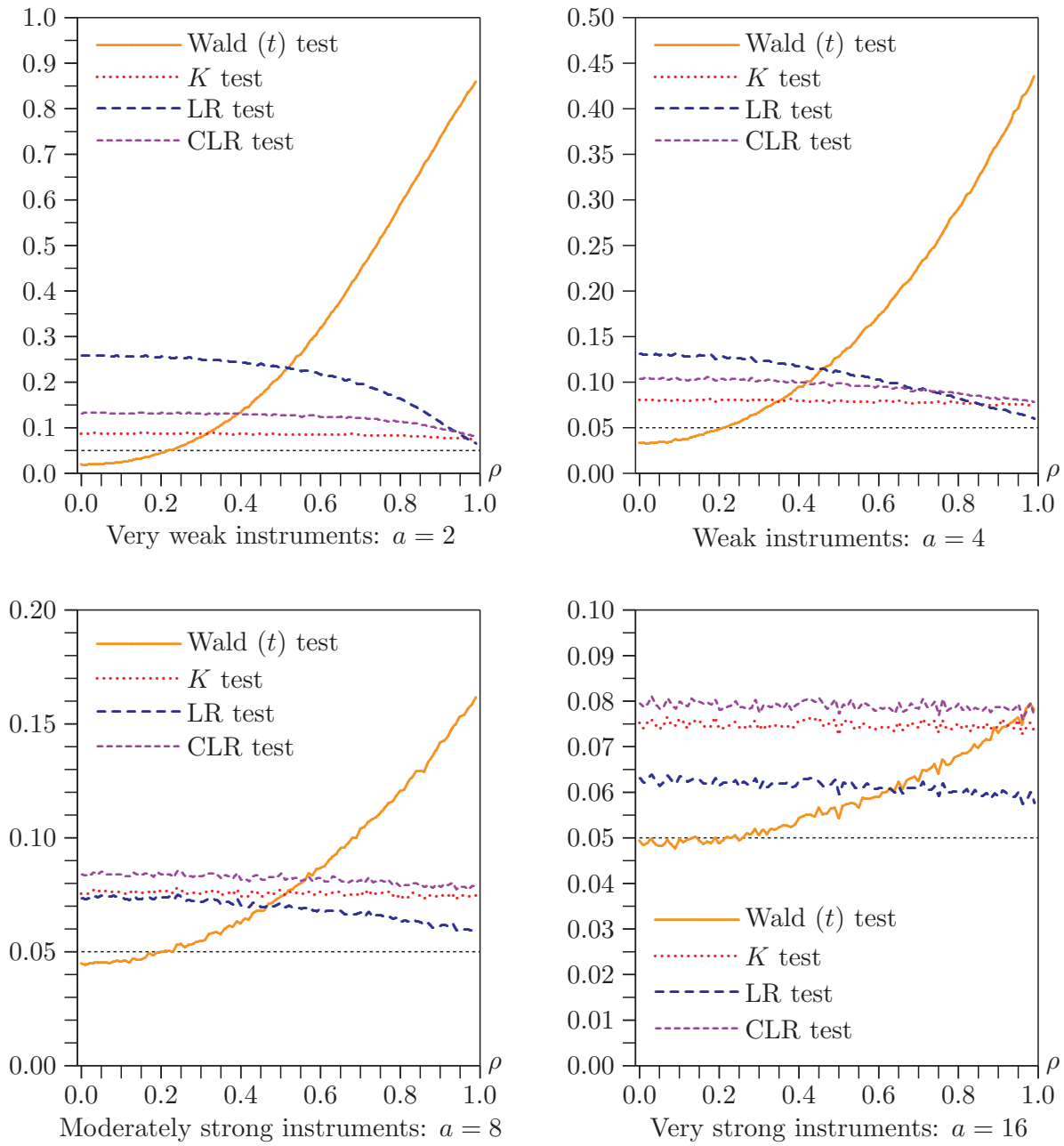


Figure 2. Rejection frequencies for asymptotic tests as functions of ρ for $l - k = 7$, $n = 50$

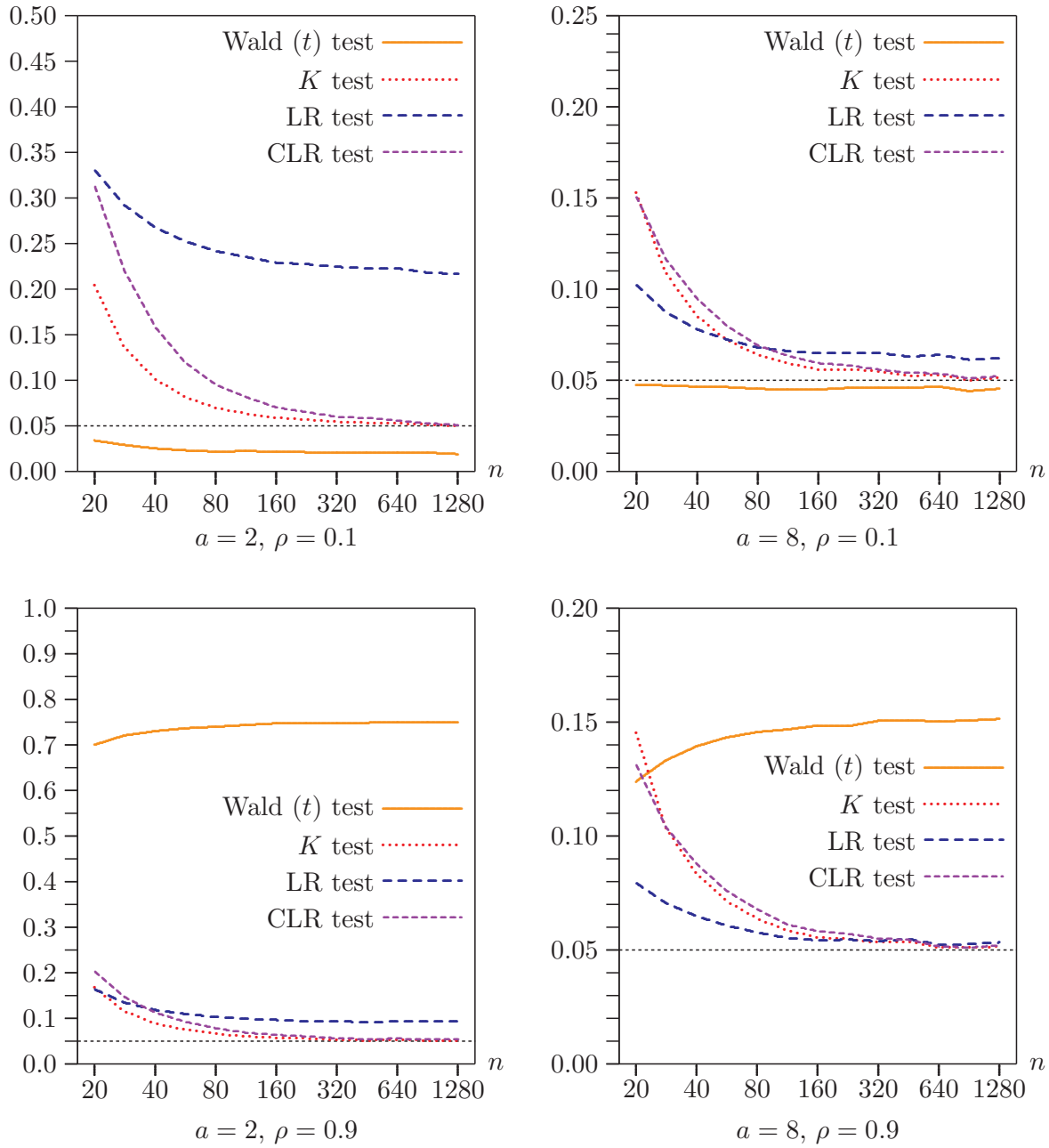


Figure 3. Rejection frequencies for weak-instrument asymptotic tests as a function of n for $l - k = 7$

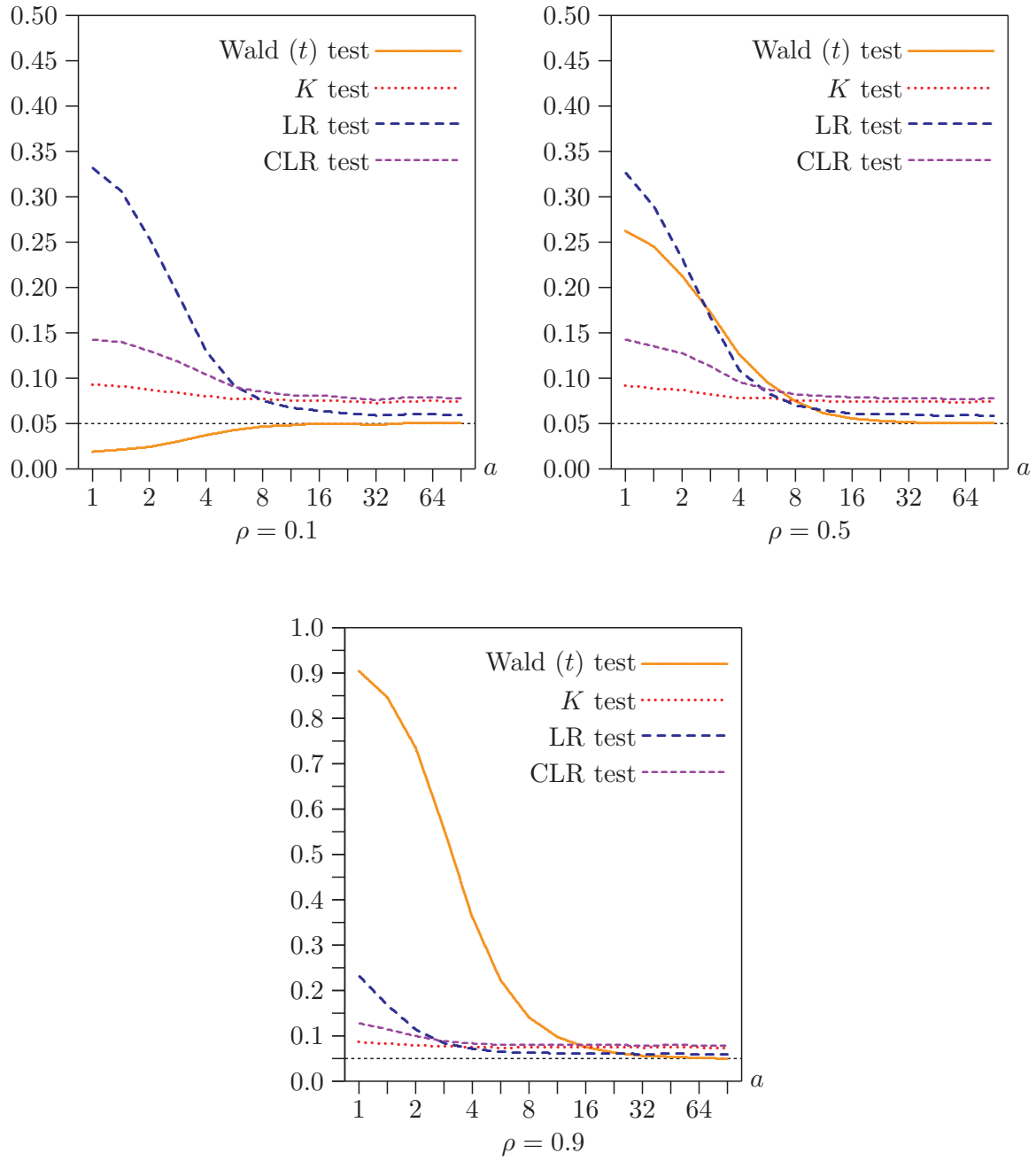


Figure 4. Rejection frequencies for asymptotic tests as functions of a , for $l - k = 7$, $n = 50$

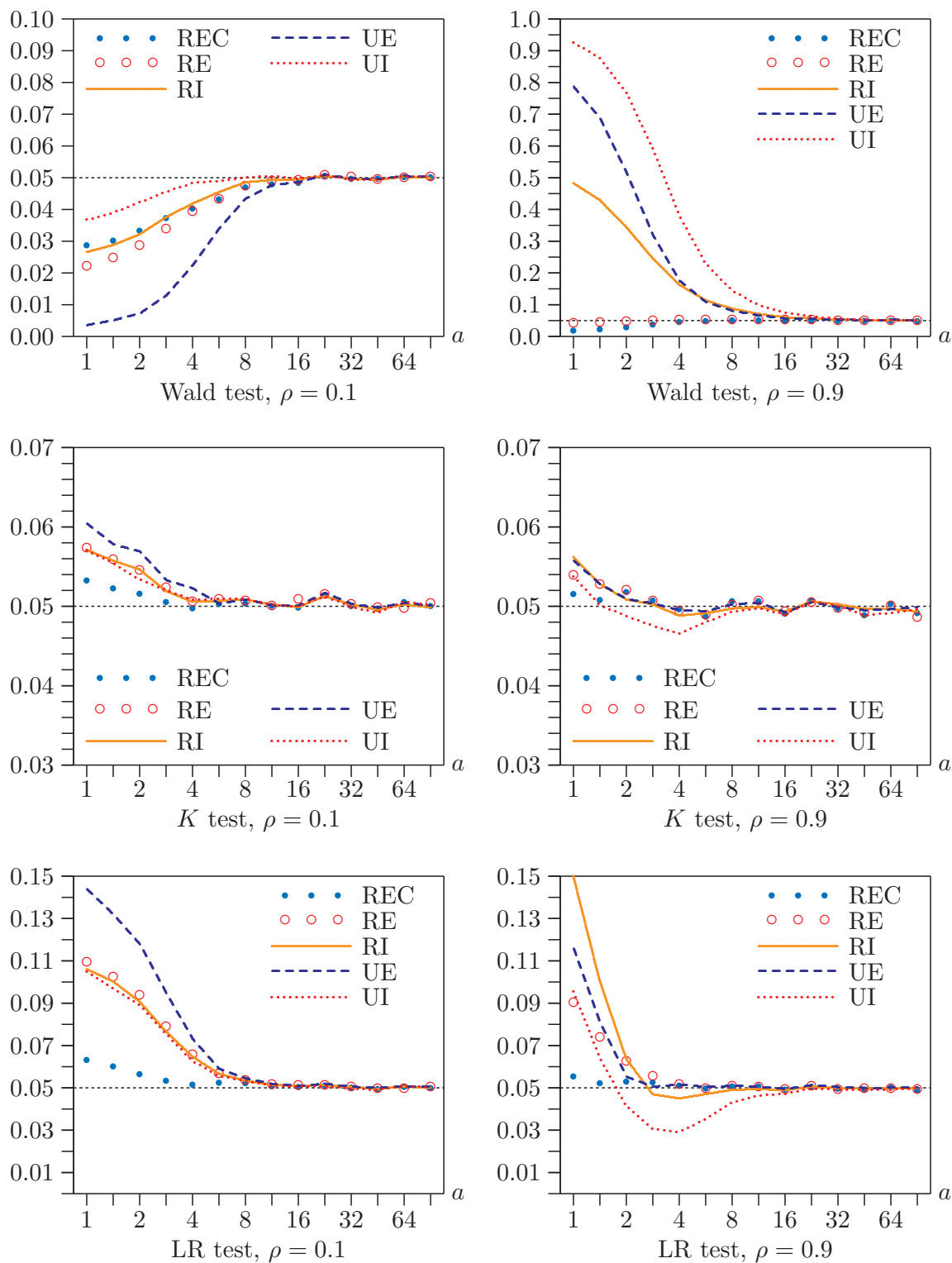


Figure 5. Rejection frequencies for bootstrap tests as functions of a , for $l - k = 7$, $n = 50$

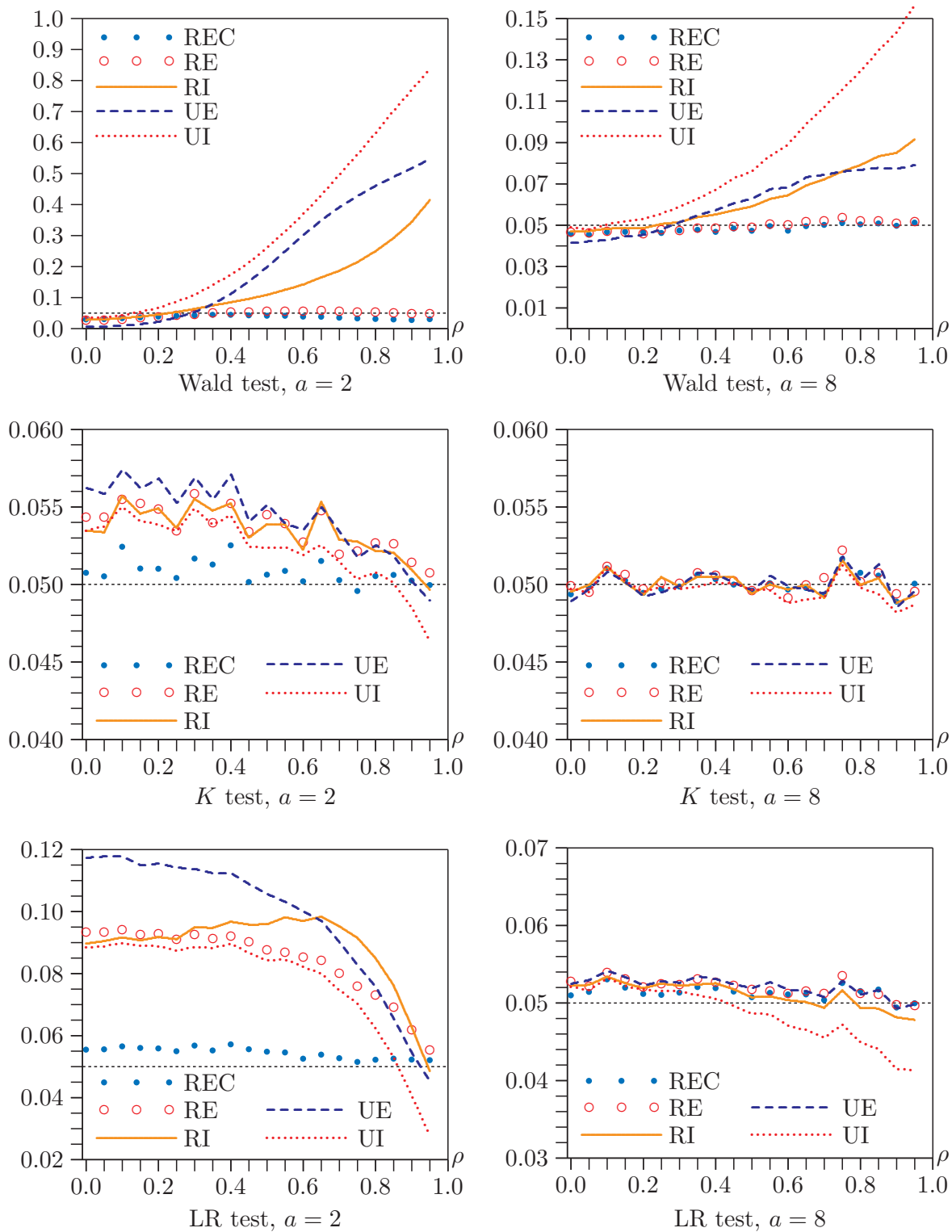


Figure 6. Rejection frequencies for bootstrap tests as functions of ρ , for $l - k = 7$, $n = 50$

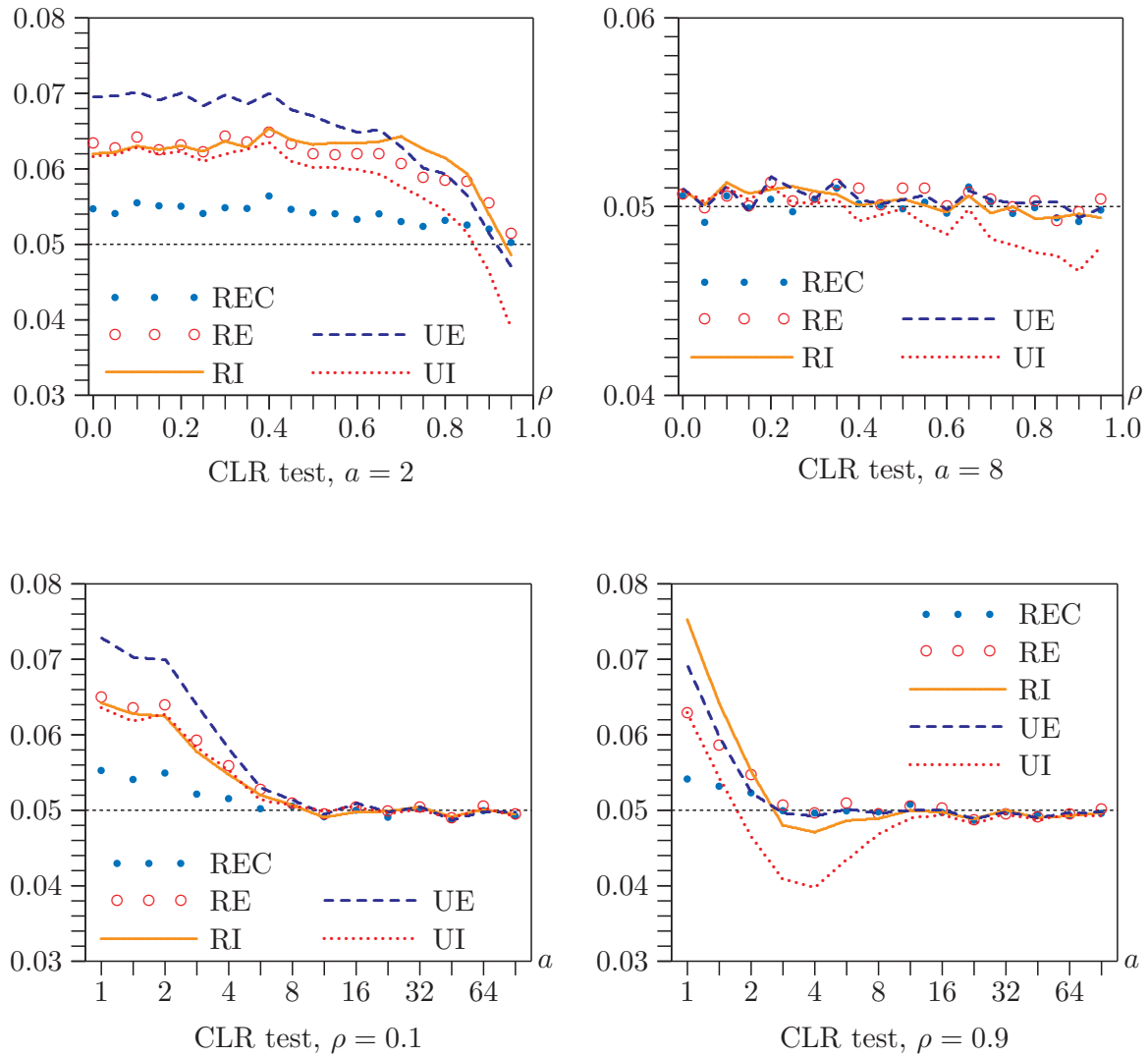


Figure 7. Rejection frequencies for bootstrap CLR tests for $l - k = 7$, $n = 50$

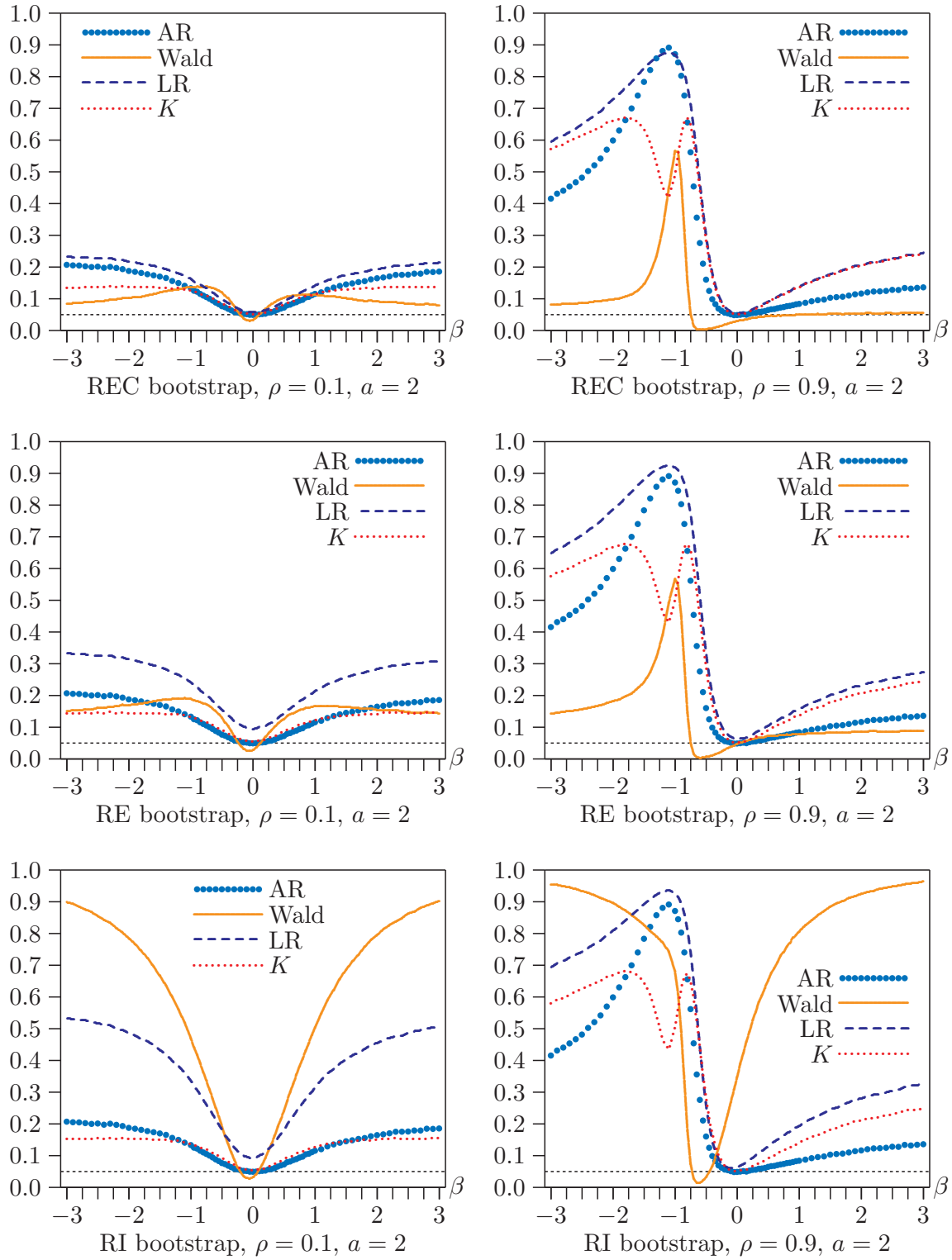


Figure 8. Power of bootstrap tests when instruments are weak, for $l - k = 7, n = 50$

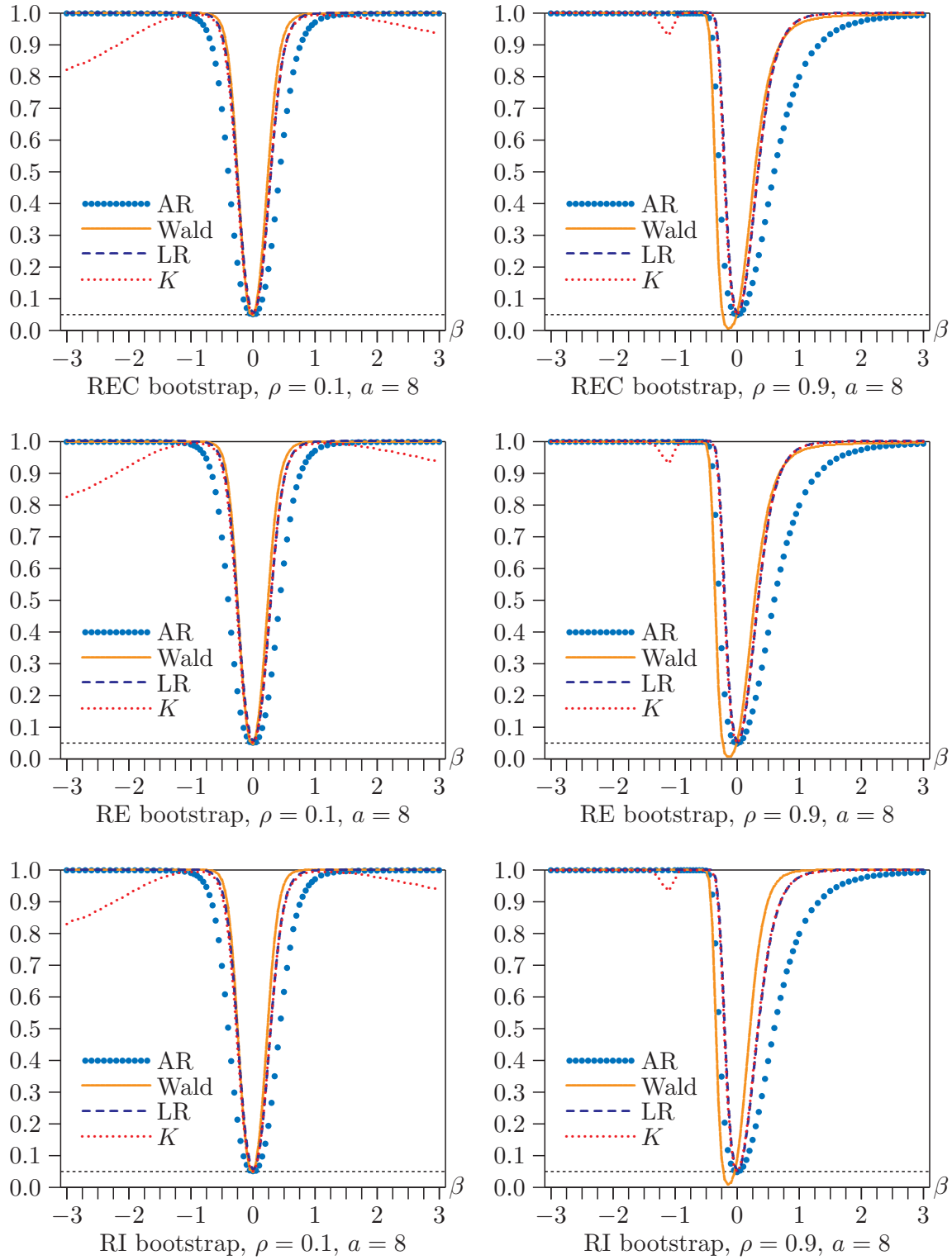


Figure 9. Power of bootstrap tests when instruments are strong, for $l - k = 7$, $n = 50$

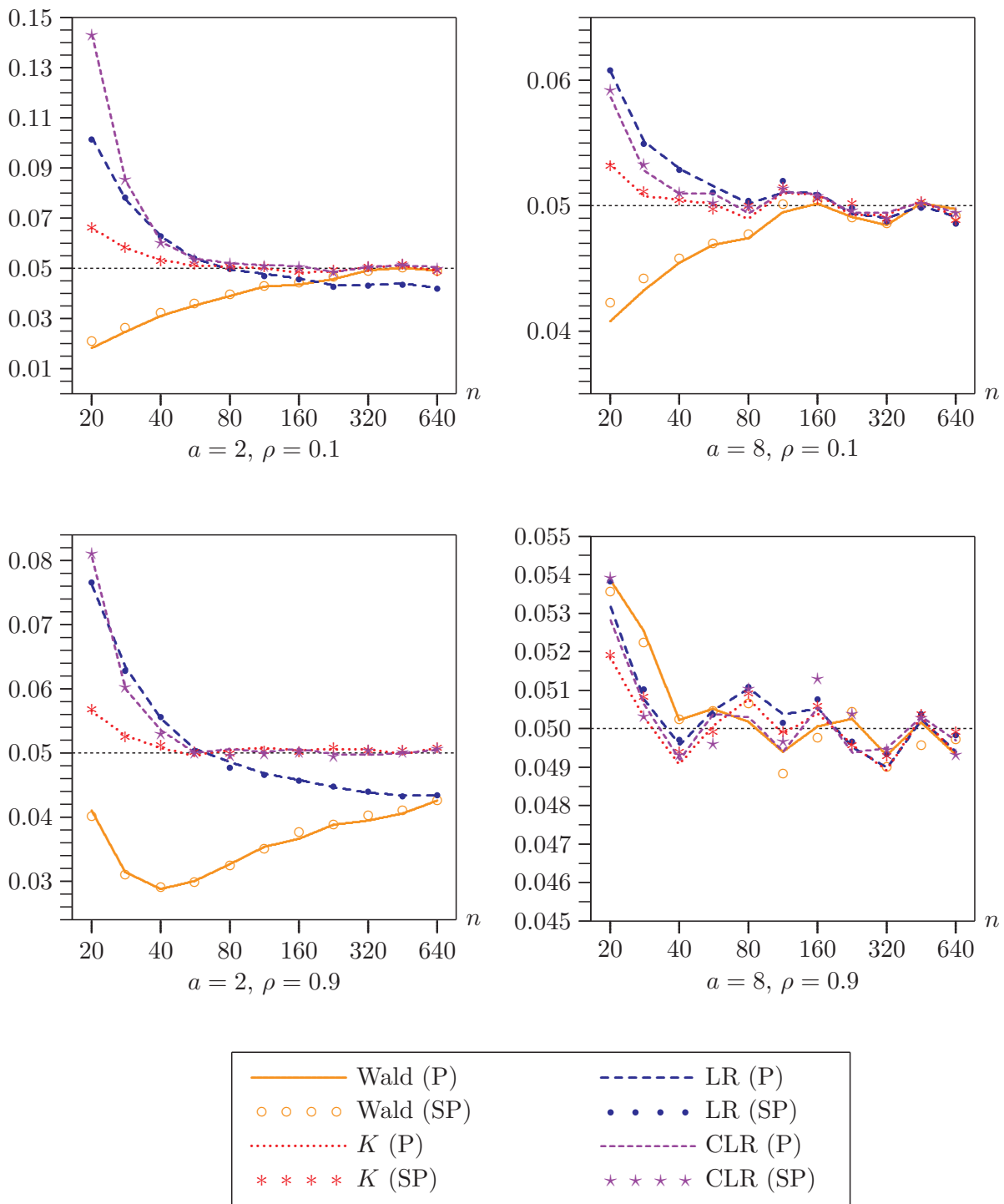


Figure 10. Rejection frequencies for REC bootstrap tests as a function of n for $l - k = 7$

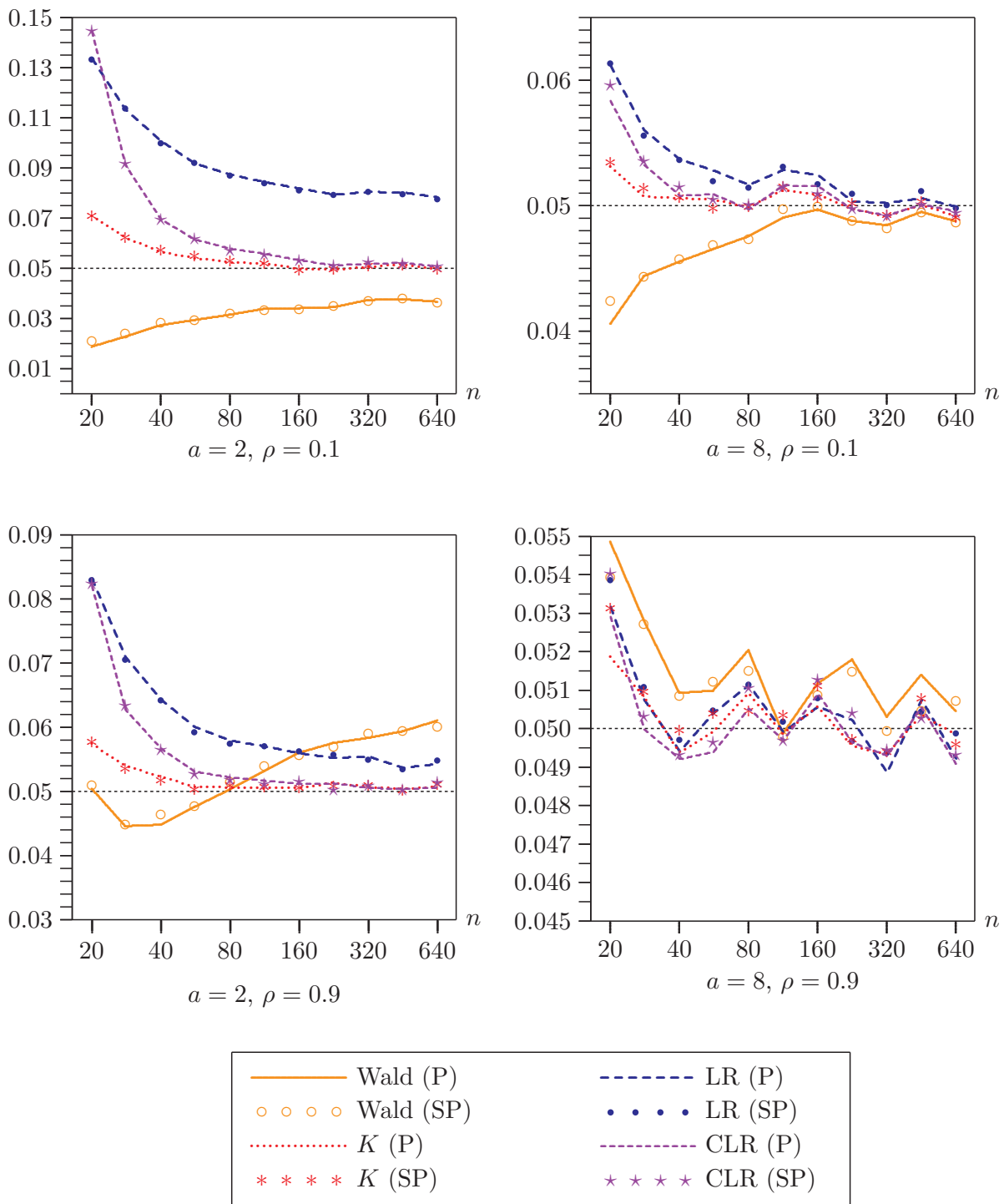


Figure 11. Rejection frequencies for RE bootstrap tests as a function of n for $l - k = 7$