

Davidson, Russell; MacKinnon, James G.

Working Paper

Testing for Consistency using Artificial Regressions

Queen's Economics Department Working Paper, No. 687

Provided in Cooperation with:

Queen's University, Department of Economics (QED)

Suggested Citation: Davidson, Russell; MacKinnon, James G. (1987) : Testing for Consistency using Artificial Regressions, Queen's Economics Department Working Paper, No. 687, Queen's University, Department of Economics, Kingston (Ontario)

This Version is available at:

<https://hdl.handle.net/10419/189088>

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. 687

Testing for Consistency Using Artificial Regressions

Russell Davidson
Queen's University

James G. MacKinnon
Queen's University

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

1987

Testing for Consistency Using Artificial Regressions

Russell Davidson

and

James G. MacKinnon

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

Abstract

We consider several issues related to Durbin-Wu-Hausman tests, that is tests based on the comparison of two sets of parameter estimates. We first review a number of results about these tests in linear regression models, discuss what determines their power, and propose a simple way to improve power in certain cases. We then show how in a general nonlinear setting they may be computed as “score” tests by means of slightly modified versions of any artificial linear regression that can be used to calculate Lagrange Multiplier tests, and explore some of the implications of this result. In particular, we show how to create a variant of the information matrix test that tests for parameter consistency. We examine the conventional information matrix test and our new version in the context of binary choice models, and provide a simple way to compute both tests using artificial regressions.

This research was supported, in part, by grants from the Social Sciences and Humanities Research Council of Canada. We are grateful to Paul Ruud, three anonymous referees, and seminar participants at GREQE (Marseille), Free University of Brussels, the University of Bristol, and Nuffield College, Oxford, for comments on earlier versions. This paper appeared in *Econometric Theory*, Volume 5, 1989, pp. 363–384. This version is similar to the published version, but it includes the Monte Carlo results from the original working paper.

April, 1987

1. Introduction

There are at least two distinct questions we may ask when we test an econometric model. The first is simply whether certain parametric restrictions hold. This question is what standard t and F tests attempt to answer in the case of regression models, and what the three classical tests, Wald, LM, and LR, attempt to answer in models estimated by maximum likelihood. The second is whether the parameters of interest have been estimated consistently. Hausman (1978), in a very influential paper, introduced a family of tests designed to answer this second question and called them “specification tests.” The basic idea of Hausman’s tests, namely that one may base a test on a “vector of contrasts” between two sets of estimates, one of which will be consistent under weaker conditions than the other, dates back to a relatively neglected paper by Durbin (1954). Wu (1974) also made use of a test for possible correlation between errors and regressors in linear regression models which was based on a vector of contrasts. We shall therefore refer to all tests of this general type as Durbin-Wu-Hausman, or DWH, tests.

There has been a good deal of work on DWH tests in recent years; see the survey paper by Ruud (1984). In this paper, we consider several issues related to tests of this type. In section 2, we review a number of results on DWH tests in linear regression models. The primary function of this section is to present results for the simplest possible case; these should then serve as an aid to intuition. We also present some new material on the distribution of DWH test statistics when the model being tested is false, and on a simple way to improve the power of the tests in certain cases.

In section 3, we provide a simple and intuitive exposition of results, originally due to Ruud (1984) and Newey (1985), on the calculation of DWH tests in nonlinear models as “score” tests by means of artificial linear regressions. We go beyond previous work by showing that any artificial regression which can be used to compute LM tests can be modified so as to compute DWH tests. An immediate implication of our argument is Holly’s (1982) result on the equivalence of DWH and classical tests in certain cases. They will be equivalent whenever the number of restrictions tested by the classical test is no greater than the number of parameters the consistency of which is being tested by the DWH test, provided that those parameters would actually be estimated inconsistently if the restrictions were incorrect. We also show that there are circumstances in which the DWH and classical tests will be equivalent (in finite samples) even when incorrect restrictions would not prevent the parameters in question from being estimated consistently. Thus rejection of the null by a DWH test does not always indicate parameter inconsistency.

In section 4, we build on results of Davidson and MacKinnon (1987a) to show how to compute a DWH version of any score-type test based on an artificial regression, even one not designed against any explicit alternative. We show how this procedure may be applied to tests such as the information matrix test (White, 1982; Chesher, 1984), and Newey’s (1985) conditional moment tests. In section 5, we discuss the power of DWH tests as compared with classical tests, in the case where the two are not identical.

It is seen that in many cases the DWH test, with fewer degrees of freedom than the corresponding classical test, will be more powerful. Finally, in section 6, we discuss the information matrix test and its DWH version in the context of binary choice models. We provide a simple way to compute both tests based on artificial regressions.

2. The Case of Linear Regression Models

Suppose the model to be tested is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (1)$$

where there are n observations and k regressors. When conducting asymptotic analysis, we shall assume that $\text{plim}(n^{-1}\mathbf{X}^\top \mathbf{u}) = \mathbf{0}$ and that $\text{plim}(n^{-1}\mathbf{X}^\top \mathbf{X})$ is a positive definite matrix. When conducting finite-sample analysis, we shall further assume that \mathbf{X} is fixed in repeated samples and that the u_t are normally distributed.

The basic idea of the DWH test is to compare the OLS estimator

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

with some other linear estimator

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A} \mathbf{y}, \quad (2)$$

where \mathbf{A} is a symmetric $n \times n$ matrix assumed for simplicity to have rank no less than k . If (1) actually generated the data, these two estimators will have the same probability limit; they will have the same expectation if \mathbf{X} is fixed in repeated samples or independent of \mathbf{u} .

The test is based on the vector of contrasts

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A} \mathbf{y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{A} \mathbf{y} - \mathbf{X}^\top \mathbf{A} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) \\ &= (\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A} \mathbf{M}_\mathbf{X} \mathbf{y}, \end{aligned} \quad (3)$$

where $\mathbf{M}_\mathbf{X} \equiv \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}$ is the orthogonal projection onto the orthogonal complement of the span of the columns of the matrix \mathbf{X} . The complementary orthogonal projection will be denoted $\mathbf{P}_\mathbf{X}$, and throughout the paper the notations \mathbf{P} and \mathbf{M} subscripted by a matrix expression will denote orthogonal projections, respectively, onto, and onto the orthogonal complement of, the span of the columns of that expression.

The first factor in (3), $(\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1}$, is simply a $k \times k$ matrix with full rank. Hence what we want to do is to test whether

$$\text{plim}(n^{-1} \mathbf{X}^\top \mathbf{A} \mathbf{M}_\mathbf{X} \mathbf{y}) = \mathbf{0}. \quad (4)$$

The vector $\mathbf{X}^\top \mathbf{A} \mathbf{M}_\mathbf{X} \mathbf{y}$ has k elements, but even if $\mathbf{A} \mathbf{X}$ has full rank, not all those elements may be random, because $\mathbf{M}_\mathbf{X}$ may annihilate some columns of $\mathbf{A} \mathbf{X}$. Suppose

that k^* is the number of linearly independent columns of $\mathbf{A}\mathbf{X}$ not annihilated by $\mathbf{M}_{\mathbf{X}}$. Then if we let these columns be denoted by \mathbf{X}^* , testing (4) is equivalent to testing

$$\text{plim}(n^{-1}\mathbf{X}^{*\top}\mathbf{A}\mathbf{M}_{\mathbf{X}}\mathbf{y}) = \mathbf{0}. \quad (5)$$

Now consider the artificial regression

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{A}\mathbf{X}^*\boldsymbol{\delta} + \text{residuals}.$$

The ordinary F statistic for $\boldsymbol{\delta} = \mathbf{0}$ in (5) is

$$\frac{\mathbf{y}^\top \mathbf{P}_{\mathbf{M}_{\mathbf{X}}\mathbf{A}\mathbf{X}^*} \mathbf{y} / k^*}{\mathbf{y}^\top \mathbf{M}_{\mathbf{X}; \mathbf{A}\mathbf{X}^*} \mathbf{y} / (n - k - k^*)}. \quad (6)$$

If (1) actually generated the data, this statistic will certainly be valid asymptotically, since the denominator will then consistently estimate σ^2 . It will be exactly distributed as $F(k^*, n - k - k^*)$ in finite samples if the u_t in (1) are normally distributed.

There are many possible choices for \mathbf{A} . In the case originally studied by Durbin (1954), $\hat{\boldsymbol{\beta}}$ is an IV estimator formed by first projecting \mathbf{X} onto the space spanned by a matrix of instruments \mathbf{W} , so that $\mathbf{A} = \mathbf{P}_{\mathbf{W}}$; see Wu (1974), Hausman (1978), Nakamura and Nakamura (1981) and Fisher and Smith (1985). The test is then often *interpreted* as a test for the exogeneity of those components of \mathbf{X} not in the space spanned by \mathbf{W} . This interpretation is misleading, since what is being tested is not the exogeneity or endogeneity of some components of \mathbf{X} , but rather the effect on the estimates of $\boldsymbol{\beta}$ of any possible endogeneity.

Alternatively, $\hat{\boldsymbol{\beta}}$ may be the OLS estimator for $\boldsymbol{\beta}$ in the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}, \quad (7)$$

where \mathbf{Z} is an $n \times l$ matrix of regressors not in the span of the columns of \mathbf{X} , so that $\mathbf{A} = \mathbf{M}_{\mathbf{Z}}$. This form of the test thus asks whether the estimates of $\boldsymbol{\beta}$ when \mathbf{Z} is excluded from the model are consistent. It is a simple example of the case examined, in a much more general context, by Holly (1982); see Section 3 below.

There is an interesting relationship between the “exogeneity” and omitted-variables variants of the DWH test. In the former, $\mathbf{A} = \mathbf{P}_{\mathbf{W}}$, and $\mathbf{P}_{\mathbf{W}}\mathbf{X}^*$ consists of all columns of $\mathbf{P}_{\mathbf{W}}\mathbf{X}$ that do not lie in the span of \mathbf{X} , so that the test regression is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{P}_{\mathbf{W}}\mathbf{X}^*\boldsymbol{\delta} + \text{residuals}. \quad (8)$$

In the latter, provided that the matrix $[\mathbf{X} \ \mathbf{Z}]$ has full rank, $\mathbf{M}_{\mathbf{Z}}\mathbf{X}^* = \mathbf{M}_{\mathbf{Z}}\mathbf{X}$. Now suppose that we expand \mathbf{Z} so that it equals \mathbf{W} . Evidently, \mathbf{X}^* will then consist of those columns of \mathbf{X} which are not in the span of \mathbf{W} , so that the test regression is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{M}_{\mathbf{W}}\mathbf{X}^*\boldsymbol{\delta} + \text{residuals}. \quad (9)$$

But it is evident that (8) and (9) will have exactly the same explanatory power. This means that exactly the same test can be *interpreted* as a test for exogeneity or as a test for the consistency of parameter estimates when certain variables have been omitted. For more on this, see Ruud (1984).

The matrix \mathbf{A} could also be almost any sort of $n \times n$ covariance matrix, so that (2) would then be a GLS estimator. Note that in this case the DWH test is not testing for a non-scalar covariance matrix, but rather for misspecification of the regression function. A special case of this is the differencing specification test, where \mathbf{A} is a matrix such that β is a vector of estimates based on first-differenced data; see Plosser, Schwert, and White (1982) and Davidson, Godfrey, and MacKinnon (1985). In this case, there are a few minor complications caused by the fact that $\mathbf{X}^\top \mathbf{A} \mathbf{X}$ does not have full rank. A similar procedure when the null hypothesis involves estimation by GLS was proposed by Boothe and MacKinnon (1986). Breusch and Godfrey (1986) discuss a variety of tests of this sort and call them “data transformation tests.”

It is often claimed that DWH tests may fruitfully be used when the null hypothesis is *not* that the data were generated by (1), but simply that the OLS estimates $\tilde{\beta}$ from (1) are consistent. While this is true up to a point, there are two major difficulties with trying to use DWH tests in this way. The first problem is that DWH tests cannot directly test the hypothesis that parameters are estimated consistently. Suppose that the model under test is

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{u}. \quad (10)$$

If the data were generated by

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

where $\mathbf{M}_\mathbf{X} \mathbf{Z} \neq \mathbf{0}$ and $\gamma \neq \mathbf{0}$, and in the testing regression (8)

$$\mathbf{P}_\mathbf{W} \mathbf{M}_\mathbf{X} \mathbf{Z} = \mathbf{0},$$

then estimates of β from the null model (1) will be inconsistent, but the power of the test based on (8) will be equal to its size. The problem here is that, for certain choices of \mathbf{W} , $\mathbf{W} \mathbf{M}_\mathbf{X} \mathbf{Z}$ may equal zero, even though $\mathbf{M}_\mathbf{X} \mathbf{Z} \neq \mathbf{0}$.

Conversely, a rejection by a DWH test does not necessarily imply that parameter estimates are inconsistent, as may be seen from the following simple example. Let the null model be (10) as before, and let the data be generated by

$$\mathbf{y} = \mathbf{X}\beta + \gamma \mathbf{z} + \mathbf{u}, \quad (11)$$

with the $n \times k$ random matrix \mathbf{X} and the $n \times 1$ random vectors \mathbf{z} and \mathbf{u} being distributed in such a way that $\text{plim}(n^{-1} \mathbf{X}^\top \mathbf{z}) = 0$ and $\text{plim}(n^{-1} \mathbf{X}^\top \mathbf{u}) = 0$. Clearly, OLS estimation of (10) will yield consistent estimates of β . The DWH test may be based on the regression

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{z}(\mathbf{z}^\top \mathbf{z})^{-1} \mathbf{z}^\top \mathbf{y} + \mathbf{u}, \quad (12)$$

where \mathbf{x}^* may be any column of \mathbf{X} . Unless $\mathbf{z}^\top \mathbf{x}^*$ happens to be exactly equal to zero, in which case the test cannot be computed, a t test for $\delta = 0$ in (12) will be numerically identical to a t test for $\gamma = 0$ in (11). Thus, if $\gamma = 0$, and the sample is large enough, the DWH test will reject the null hypothesis with probability one, even though β is in fact consistent! The reason for this is that in a finite sample we have computed a DWH test that is meaningless asymptotically, because the regressor $\mathbf{z}(\mathbf{z}^\top \mathbf{z})^{-1} \mathbf{z} \mathbf{x}^*$ vanishes. Unfortunately, it is often possible to do this. In such circumstances, results from DWH tests may easily be misinterpreted.

The second problem with using DWH tests of (1) when neither (1) nor (5) actually generated the data is that the denominator of (6) will then provide an overestimate of the amount of noise in the actual data generating process, or DGP. Specifically, if the data are generated by the process

$$\mathbf{y} = \mathbf{X}\beta_0 + \mathbf{a}_0 + \mathbf{u}, \quad \mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}),$$

where \mathbf{a}_0 may be thought of as a linear combination of omitted variables, then the F statistic for $\delta = \mathbf{0}$ in (5) will be distributed as *doubly* non-central $F(k^*, n - k - k^*)$ with numerator and denominator non-centrality parameters (NCPs)

$$\frac{\mathbf{a}_0^\top \mathbf{P}_{\mathbf{M}_\mathbf{X}} \mathbf{A} \mathbf{X}^* \mathbf{a}_0}{\sigma_0^2} \tag{13}$$

and

$$\frac{\mathbf{a}_0^\top \mathbf{M}_{\mathbf{X}; \mathbf{A} \mathbf{X}^*} \mathbf{a}_0}{\sigma_0^2}, \tag{14}$$

respectively. If one considers the artificial linear regression

$$\mathbf{a}_0 = \mathbf{X}\alpha + \mathbf{A} \mathbf{X}^* \psi + \text{residuals}, \tag{15}$$

then (14) is the sum of squared residuals from (15), and (13) is the reduction in the sum of squared residuals due to the inclusion of $\mathbf{A} \mathbf{X}^*$ in (15). When regression (15) fits perfectly, this means that \mathbf{X} and $\mathbf{A} \mathbf{X}^*$ jointly explain all the variation in \mathbf{a}_0 . The numerator NCP (13) then simplifies to

$$\frac{\mathbf{a}_0^\top \mathbf{M}_\mathbf{X} \mathbf{a}_0}{\sigma_0^2}, \tag{16}$$

and the denominator NCP (14) is equal to zero. The test will then have as much power as any test with k^* degrees of freedom could have. However, when (15) fits less than perfectly, the numerator NCP (13) is smaller than (16), and the denominator NCP (14) is greater than zero, both of which cause the test to have less power; see Thursby and Schmidt (1977).

In certain cases, it is possible to improve the estimate of σ^2 , reducing the denominator NCP and hence increasing power. Consider again the case where $\mathbf{A} = \mathbf{M}_\mathbf{Z}$. Whenever

$\rho(\mathbf{A}\mathbf{X}) = \rho(\mathbf{M}_Z\mathbf{X}) < \rho(\mathbf{Z})$, the DWH test differs from the classical test for $\boldsymbol{\gamma} = \mathbf{0}$ in (7), and the test regression

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{M}_Z\mathbf{X}\boldsymbol{\delta} + \mathbf{u} \quad (17)$$

fits less well than regression (7), because the latter has additional regressors. Instead of using the ordinary F statistic for $\boldsymbol{\delta} = \mathbf{0}$ in (17), then, one might use the test statistic

$$\frac{\mathbf{y}^\top \mathbf{P}_{\mathbf{M}_Z\mathbf{X}} \mathbf{y} / k^*}{\mathbf{y}^\top \mathbf{M}_{\mathbf{X}; \mathbf{Z}} \mathbf{y} / (n - k - l)}. \quad (18)$$

The denominator here is the estimate of σ^2 from (7); the ordinary F statistic would use the estimate of σ^2 from (17).

The test statistic (18) will be asymptotically valid whenever (7) generated the data, and it will have the $F(k^*, n - k - l)$ distribution in finite samples when the null hypothesis that $E(\hat{\boldsymbol{\beta}} = \tilde{\boldsymbol{\beta}})$ is true, the regressors are fixed, and the errors are normal. Reducing the number of degrees of freedom in the denominator of an F test has the effect of reducing power; see Das Gupta and Perlman (1974). Thus, if the data were generated by (17), the modified F test (18) would have less power than the ordinary F test. However, in some cases where (1) is false, (7) may fit much better than (17), thus yielding a much lower estimate of σ^2 . In such cases, the modified F test (18) will be more powerful than the ordinary one.

3. General Nonlinear Models

Since the work of Hausman (1978), it has been well known that DWH tests may be used in the context of very general classes of models involving maximum likelihood estimation. There are three principal theoretical results in this literature. The first, due to Hausman, is that the (asymptotic) covariance matrix of a vector of contrasts is equal to the difference between the (asymptotic) covariance matrices of the two vectors of parameter estimates, provided that one of the latter is (asymptotically) efficient under the null hypothesis. This is essentially a corollary of the Cramér-Rao bound.

The second principal result, due to Holly (1982), is that when the two parameter vectors being contrasted correspond to restricted and unrestricted ML estimates (the vectors consisting only of those parameters which are estimated under the restrictions), the DWH test will in certain circumstances be equivalent to the three classical test statistics, Wald, LM, and LR. Whether this equivalence holds or not will depend on the numbers of parameters in the restricted and unrestricted models, and on the rank of certain matrices; as we show below, the results are completely analogous to those on whether the DWH test based on (17) is equivalent to the F test based on (7).

The third principal result, due to White (1982), Ruud (1984), and Newey (1985), is that tests asymptotically equivalent to DWH tests can be computed as score tests. As Ruud and Newey recognized, this implies that various artificial regressions can be

used to compute these tests. The only artificial regression which has been explicitly suggested for this purpose is the so-called outer-product-of-the-gradient, or OPG, regression, in which a vector of ones is regressed on the matrix of contributions from single observations to the gradient of the loglikelihood function. This regression is widely used for calculating LM tests; see Godfrey and Wickens (1981). It has more recently been suggested by Newey (1985) as an easy way to calculate “conditional moment” tests, including some which are DWH tests. Unfortunately, the OPG regression is known to have poor and often disastrous finite-sample properties. See Davidson and MacKinnon (1983, 1985a), Bera and McKenzie (1986), Chesher and Spady (1988), and Kennan and Neumann (1988); the last two provide spectacular examples. As we shall now show, any artificial regression that can be used to compute LM tests can also be used to compute DWH tests. In view of the poor properties of the OPG regression, this result may be important for applied work.

There are many classes of models for which artificial linear regressions other than the OPG regression are available. These include univariate and multivariate nonlinear regression models (Engle 1982, 1984), possibly with heteroskedasticity of unknown form (Davidson and MacKinnon, 1985b), probit and logit models (Davidson and MacKinnon, 1984b), and a rather general class of nonlinear models, with nonlinear transformations on the dependent variable(s), for which “double-length” artificial regressions with $2n$ “observations” are appropriate (Davidson and MacKinnon, 1984a, 1988). To the extent that evidence is available, these all appear to have better finite-sample properties than the OPG regression.

We shall deal with the following general case. There is a sample of size n which gives rise to a loglikelihood function

$$\mathcal{L}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \sum_{t=1}^n \ell(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2), \quad (19)$$

where $\boldsymbol{\theta}_1$ is a k -vector and $\boldsymbol{\theta}_2$ an l -vector of parameters, the latter equal to zero if the model is correctly specified. Maximum likelihood estimates of the vector $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top]^\top$ under the restriction $\boldsymbol{\theta}_2 = \mathbf{0}$ will be denoted $\tilde{\boldsymbol{\theta}}$, while unrestricted estimates will be denoted $\hat{\boldsymbol{\theta}}$. The scores with respect to $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ are denoted by $\mathbf{g}_1(\boldsymbol{\theta})$ and $\mathbf{g}_2(\boldsymbol{\theta})$; thus

$$\mathbf{g}_i(\boldsymbol{\theta}) = \sum_{t=1}^n \frac{\partial \ell_t(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)}{\partial \boldsymbol{\theta}_i}, \quad i = 1, 2.$$

A hat or a tilde over any quantity indicates that it is evaluated at $\hat{\boldsymbol{\theta}}$ or $\tilde{\boldsymbol{\theta}}$, respectively. The model represented by (19) is assumed to satisfy all the usual conditions for maximum likelihood estimation and inference to be asymptotically valid; see, for example, Amemiya (1985, Chapter 4). In particular, we assume that the true parameter vector $\boldsymbol{\theta}^0$ is interior to a compact parameter space, and that the information matrix

$$\mathcal{J}(\boldsymbol{\theta}) \equiv \lim (\mathbf{E}(\mathbf{g}\mathbf{g}^\top))$$

is a finite, nonsingular matrix. The submatrix of \mathcal{J} corresponding to $\boldsymbol{\theta}_i$ will be denoted \mathcal{J}_{ii} ; the corresponding submatrix of \mathcal{J}^{-1} will be denoted $(\mathcal{J}^{-1})_{ii}$.

Taking Taylor series approximations to the first-order conditions for $\tilde{\boldsymbol{\theta}}_1$ and $(\hat{\boldsymbol{\theta}}_1, \hat{\boldsymbol{\theta}}_2)$ around $\boldsymbol{\theta}^0$, and applying a suitable law of large numbers, we find that

$$n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0) \cong \mathcal{J}_{11}^{-1}[\mathbf{I}_k \ \mathbf{O}]n^{-1/2}\mathbf{g}(\boldsymbol{\theta}^0)$$

and

$$n^{1/2}(\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0) \cong [\mathbf{I}_k \ \mathbf{O}]\mathcal{J}^{-1}n^{-1/2}\mathbf{g}(\boldsymbol{\theta}^0),$$

where \mathbf{I}_k is a $k \times k$ identity matrix, and \mathbf{O} is a $k \times l$ matrix of zeros. It follows that

$$\begin{aligned} n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0) &\cong (\mathcal{J}_{11}^{-1}[\mathbf{I}_k \ \mathbf{O}] - [\mathbf{I}_k \ \mathbf{O}]\mathcal{J}^{-1})n^{-1/2}\mathbf{g}(\boldsymbol{\theta}^0) \\ &= (\mathcal{J}_{11}^{-1} - (\mathcal{J}^{-1})_{11})n^{-1/2}\mathbf{g}_1(\boldsymbol{\theta}_0) - ((\mathcal{J}^{-1})_{12})n^{-1/2}\mathbf{g}_2(\boldsymbol{\theta}_0). \end{aligned} \quad (20)$$

From (20) it is easy to show that the asymptotic covariance matrix of $n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \hat{\boldsymbol{\theta}}_1)$ is

$$\begin{aligned} &((\mathcal{J}_{11})^{-1}[\mathbf{I}_k \ \mathbf{O}] - [\mathbf{I}_k \ \mathbf{O}]\mathcal{J}^{-1})\mathcal{J}((\mathcal{J}_{11})^{-1}[\mathbf{I}_k \ \mathbf{O}] - [\mathbf{I}_k \ \mathbf{O}]\mathcal{J}^{-1})^\top \\ &= (\mathcal{J}^{-1})_{11} - (\mathcal{J}_{11})^{-1}. \end{aligned} \quad (21)$$

The first term in (21) is the asymptotic covariance matrix of $n^{1/2}(\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0)$, and the second is the asymptotic covariance matrix of $n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0)$, so that (21) is a special case of Hausman's principal result.

Standard results on partitioned matrices tell us that

$$(\mathcal{J}^{-1})_{11} = (\mathcal{J}_{11} - \mathcal{J}_{12}(\mathcal{J}_{22}^{-1})\mathcal{J}_{21})^{-1}$$

and

$$(\mathcal{J}^{-1})_{12} = -(\mathcal{J}_{11} - \mathcal{J}_{12}(\mathcal{J}_{22}^{-1})\mathcal{J}_{21})^{-1}\mathcal{J}_{12}\mathcal{J}_{22}^{-1},$$

and substituting these into (20) yields the following expression for $n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \hat{\boldsymbol{\theta}}_1)$:

$$\mathcal{J}_{11}^{-1}n^{-1/2}\mathbf{g}_1 + (\mathcal{J}_{11} - \mathcal{J}_{12}(\mathcal{J}_{22}^{-1})\mathcal{J}_{21})^{-1}(\mathcal{J}_{12}\mathcal{J}_{22}^{-1}n^{-1/2}\mathbf{g}_2 - n^{-1/2}\mathbf{g}_1).$$

This expression allows us to derive easily computed test statistics based on the general notion of an artificial regression.

In the usual case of testing restrictions in the context of maximum likelihood estimation, an artificial regression involves two things: a regressand, say $\mathbf{r}(\boldsymbol{\theta})$, and a matrix of regressors, say $\mathbf{R}(\boldsymbol{\theta})$, partitioned as $[\mathbf{R}_1 \ \mathbf{R}_2]$, which have the properties that

- (i) $\mathbf{R}^\top(\boldsymbol{\theta})\mathbf{r}(\boldsymbol{\theta})$ is the gradient of the loglikelihood function at $\boldsymbol{\theta}$, and
- (ii) $n^{-1}\mathbf{R}^\top(\ddot{\boldsymbol{\theta}})\mathbf{R}(\ddot{\boldsymbol{\theta}})$ consistently estimates the information matrix whenever $\ddot{\boldsymbol{\theta}}$ consistently estimates $\boldsymbol{\theta}$.

Replacing the gradients and information sub-matrices in (22) by their finite-sample analogues, evaluated at $\tilde{\boldsymbol{\theta}}$, and ignoring factors of n , yields the expression

$$(\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{r}} - (\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{M}}_2 \tilde{\mathbf{r}} = -(\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{M}}_2 \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}}, \quad (23)$$

where $\tilde{\mathbf{M}}_i$ denotes the matrix which projects onto the orthogonal complement of the subspace spanned by the columns of $\tilde{\mathbf{R}}_i$, for $i = 1, 2$. Notice that the left-hand side of (23) resembles the expression for a restricted OLS estimator minus an unrestricted one. Think of $\tilde{\mathbf{r}}$ as the regressand, $\tilde{\mathbf{R}}_1$ as the matrix of regressors for the null hypothesis, and $\tilde{\mathbf{M}}_2$ as the matrix which projects onto the orthogonal complement of the space spanned by the additional regressors whose coefficients are zero under the null.

Now consider the artificial regression

$$\tilde{\mathbf{r}} = \tilde{\mathbf{R}}_1 \mathbf{b}_1 + \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1^* \mathbf{b}_2 + \text{residuals}, \quad (24)$$

where the $n \times k^*$ matrix \mathbf{R}_1^* consists of as many columns of \mathbf{R}_1 as possible subject to the condition that the matrix $[\tilde{\mathbf{R}}_1 \ \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1^*]$ have full rank. The explained sum of squares from this regression is

$$\tilde{\mathbf{r}}^\top \mathbf{P}_{\tilde{\mathbf{R}}_1; \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1^*} \tilde{\mathbf{r}} = \tilde{\mathbf{r}}^\top \mathbf{P}_{\tilde{\mathbf{M}}_1 \tilde{\mathbf{M}}_2 \tilde{\mathbf{R}}_1^*} \tilde{\mathbf{r}},$$

since $\mathbf{R}_1^\top \tilde{\mathbf{r}} = \mathbf{0}$ by the first-order conditions. Under suitable regularity conditions, it is easily shown that this statistic is asymptotically distributed as $\chi^2(k^*)$ under the null hypothesis that $\boldsymbol{\theta}_2 = \mathbf{0}$. This result also extends to any situation where the data are generated by a sequence of local DGPs with $\boldsymbol{\theta}_2 \mathcal{I}_{21} = \mathbf{0}$ which tends to $\boldsymbol{\theta}_2 = \mathbf{0}$, *provided* that \mathcal{I}_{21} has full rank; we discuss this important proviso below.

Notice that (24) may be rewritten as

$$\tilde{\mathbf{r}} = \tilde{\mathbf{R}}_1(\mathbf{b}_1 + \mathbf{b}_2) - \tilde{\mathbf{R}}_2(\tilde{\mathbf{R}}_2^\top \tilde{\mathbf{R}}_2)^{-1} \tilde{\mathbf{R}}_2^\top \tilde{\mathbf{R}}_1 \mathbf{b}_2 + \text{residuals}.$$

Thus, as with the linear case, it makes no difference whether we use (24) or

$$\tilde{\mathbf{r}} = \tilde{\mathbf{R}}_1 \mathbf{c}_1 + \tilde{\mathbf{P}}_2 \tilde{\mathbf{R}}_1^* \mathbf{c}_2 + \text{residuals} \quad (25)$$

for the purpose of computing a test.

The classical LM test can of course be computed as the explained sum of squares from the artificial regression

$$\tilde{\mathbf{r}} = \tilde{\mathbf{R}}_1 \mathbf{b}_1 + \tilde{\mathbf{R}}_2 \mathbf{b}_2 + \text{residuals} \quad (26)$$

The equivalence result of Holly (1982) now follows immediately. Suppose that $l < k$, so that there are fewer restrictions than parameters under the null hypothesis, and that $\tilde{\mathbf{R}}_2^\top \tilde{\mathbf{R}}_1$ has rank l . Then it must be the case that $\tilde{\mathbf{R}}_1$ and

$$\tilde{\mathbf{P}}_2 \mathbf{R}_1 = \tilde{\mathbf{R}}_2(\tilde{\mathbf{R}}_2^\top \tilde{\mathbf{R}}_2)^{-1} \tilde{\mathbf{R}}_2^\top \tilde{\mathbf{R}}_1$$

span the same space as $\tilde{\mathbf{R}}_1$ and $\tilde{\mathbf{R}}_2$, so that (26) and (25) will have exactly the same explanatory power. The LM and DWH tests will then be numerically identical. Provided that $\mathcal{J}_{21} = \text{plim}(n^{-1}\tilde{\mathbf{R}}_2^\top\tilde{\mathbf{R}}_1)$ has full rank l , the asymptotic equivalence of all forms of classical and DWH tests, which is Holly's result, then follows immediately from the numerical equality of these two tests.

When \mathcal{J}_{21} does not have full rank, some elements (or linear combinations of elements) of the vector $\boldsymbol{\theta}_1$ will be estimated consistently by $\tilde{\boldsymbol{\theta}}_1$ even when the restrictions are locally false. More precisely, if we assume that $\boldsymbol{\theta}_2 = n^{-1/2}\boldsymbol{\eta}$, so that the discrepancy between $\boldsymbol{\theta}_2$ and $\mathbf{0}$ is $O(n^{-1/2})$, then when \mathcal{J}_{21} does not have full rank, certain linear combinations of the components of the vector $n^{1/2}(\tilde{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1^0)$ will have mean zero, regardless of the value of $\boldsymbol{\eta}$. Ordinarily, when \mathcal{J}_{21} has full rank, this could be true only for certain values of $\boldsymbol{\eta}$ (including $\boldsymbol{\eta} = \mathbf{0}$), and the DWH test would have power whenever $\boldsymbol{\eta}$ did not have those values. This local result is of course true globally for linear regression models.

In this situation, the results of DWH tests may easily be misinterpreted. When \mathcal{J}_{21} does not have full rank, one must drop as many columns of $\tilde{\mathbf{P}}_2\tilde{\mathbf{R}}_1$ as necessary and reduce the degrees of freedom for the test accordingly. In practice, however, $\tilde{\mathbf{R}}_2^\top\tilde{\mathbf{R}}_1$ may well have full rank even though \mathcal{J}_{21} does not, so that the investigator may not realize there is a problem. As a result, the null hypothesis of consistency may well be rejected even when $\tilde{\boldsymbol{\theta}}_1$ is in fact consistent. The key to understanding this is to recognize that, even though the null hypothesis of the DWH version of a classical test is $\boldsymbol{\theta}_2\mathcal{J}_{21} = \mathbf{0}$ rather than $\boldsymbol{\theta}_2 = \mathbf{0}$, the test is still testing a hypothesis about $\boldsymbol{\theta}_2$ and *not* a hypothesis \mathcal{J}_{21} . When the test is done by an artificial regression, \mathcal{J}_{21} is merely *estimated* by $n^{-1}\tilde{\mathbf{R}}_2^\top\tilde{\mathbf{R}}_1$, and if \mathcal{J}_{21} does not have full rank, the estimate by itself will almost never reveal that fact (although combined with an estimate of the *variance* of $\tilde{\mathbf{R}}_2^\top\tilde{\mathbf{R}}_1$, it could do so). This is precisely analogous to the case of the linear regression model (10), where incorrectly omitting the regressor \mathbf{Z} had no effect on the consistency of $\hat{\boldsymbol{\beta}}$.

Note that this is a problem for all forms of the DWH test, and not simply for the score form. In cases where the information matrix is block-diagonal between the parameters which are estimated under the null and the parameters which are restricted, the former will always be estimated consistently even when the restrictions are locally false in the sense discussed above. This implies that the asymptotic covariance matrix of the vector of contrasts, expression (21), must be a zero matrix. But the finite-sample analogue of (21) will almost never be a zero matrix, and it is usually computed in such a way as to ensure that it is positive semi-definite. As a result, it will be just as possible to compute, and misinterpret, the DWH statistic in its original form as in its score form.

4. DWH Tests in Other Directions

In Davidson and MacKinnon (1987a), we showed that the Holly result is perfectly general when the null hypothesis is estimated by maximum likelihood. The reason for this is that when one set of estimates is asymptotically efficient if the model is

correctly specified, the other set is always asymptotically equivalent (locally) to ML estimates with *some* set of restrictions removed; Holly's result then shows that, when the number of restrictions removed is no greater than the number of parameters estimated under the null, and the information matrix satisfies certain conditions, a DWH test is equivalent to a classical test of those restrictions.

As a corollary of this result, we can start with any score-type test and derive a DWH variant of it, similar to the test based on regression (24). Consider an artificial regression analogous to (26), but with $\tilde{\mathbf{R}}_2$ replaced by an $n \times m$ matrix $\tilde{\mathbf{Z}} = \mathbf{Z}(\tilde{\boldsymbol{\theta}})$:

$$\tilde{\mathbf{r}} = \tilde{\mathbf{R}}_1 \mathbf{c}_1 + \tilde{\mathbf{Z}} \mathbf{c}_2 + \text{residuals}. \quad (27)$$

The matrix \mathbf{Z} must satisfy certain conditions, which essentially give it the same properties as $\tilde{\mathbf{R}}_2$; these are discussed below. Provided it does so, and assuming that the matrix $[\tilde{\mathbf{R}}_1 \ \tilde{\mathbf{Z}}]$ has full rank, the explained sum of squares from this regression will be asymptotically distributed as $\chi^2(m)$ when the DGP is (19) with $\boldsymbol{\theta}_2 = \mathbf{0}$.

The variety of tests covered by (27) is very great. In addition to LM tests based on all known artificial regressions, tests of this form include Newey's (1985) conditional moment tests, all the score-type DWH tests discussed in sections 2 and 3 above, White's (1982) information matrix test in the OPG form suggested by Lancaster (1984), and Ramsey's (1969) RESET test.

We now briefly indicate how to prove the above proposition. The proof is similar to standard proofs for LM tests based on artificial regressions, and the details are therefore omitted. As noted above, it is necessary that $\tilde{\mathbf{Z}}$ satisfy certain conditions, so that it essentially has the same properties as $\tilde{\mathbf{R}}_2$. First, we require that $\text{plim}(n^{-1} \tilde{\mathbf{r}}^\top \tilde{\mathbf{Z}}) = \mathbf{0}$ under the null hypothesis; if this condition were not satisfied, we obviously could not expect \mathbf{c}_2 in (27) to be zero. Second, we require that

$$\text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{Z}}) = \text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}) \quad (28)$$

and

$$\text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{R}}_1) = \text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{R}}_1), \quad (29)$$

which are similar to the condition that

$$\text{plim}(n^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{R}}_1) = \text{plim}(n^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1); \quad (30)$$

(30) does not have to be assumed because it is a consequence of property (ii) and the consistency of $\tilde{\boldsymbol{\theta}}$. Finally, we require that a central limit theorem be applicable to the vector

$$n^{-1/2} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}} \quad (31)$$

and that laws of large numbers be applicable to the quantities whose probability limits appear on the right-hand sides of (28), (29), and (30).

Consider the vector (31). Asymptotically, it has mean zero under the null hypothesis, and its asymptotic covariance matrix is

$$\text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}),$$

which is equal to

$$\begin{aligned} \text{plim} \Big(n^{-1} \Big(\tilde{\mathbf{Z}}^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{Z}} - \tilde{\mathbf{Z}}^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{R}}_1 (\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{Z}} - \tilde{\mathbf{Z}}^\top \tilde{\mathbf{R}}_1 (\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{Z}} \\ + \tilde{\mathbf{Z}}^\top \tilde{\mathbf{R}}_1 (\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{R}}_1 (\tilde{\mathbf{R}}_1^\top \tilde{\mathbf{R}}_1)^{-1} \tilde{\mathbf{R}}_1^\top \tilde{\mathbf{Z}} \Big) \Big). \end{aligned} \quad (32)$$

Rewriting (32) so that each term is a product of probability limits which are $O(1)$, using (28), (29), and (30), and simplifying, we find that

$$\text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}} \tilde{\mathbf{r}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}) = \text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}).$$

This plus the asymptotic normality of (31) implies that the statistic

$$(n^{-1/2} \tilde{\mathbf{r}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}) (\text{plim}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}))^{-1} (n^{-1/2} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}}) \quad (33)$$

is asymptotically distributed as $\chi^2(m)$. But since our assumptions imply that a law of large numbers can be applied to $n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}}$, the explained sum of squares from regression (27), which is

$$\tilde{\mathbf{r}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{M}}_1 \tilde{\mathbf{r}},$$

will asymptotically be the same random variable as (33).

It is obvious how to construct a DWH version of this test, and it is now obvious that such a test will be asymptotically valid. We obtain the DWH version simply by replacing $\tilde{\mathbf{Z}}$ in (27) with $\tilde{\mathbf{M}}_Z \tilde{\mathbf{R}}_1$ or $\tilde{\mathbf{P}}_Z \tilde{\mathbf{R}}_1$. It is evident that if $\tilde{\mathbf{Z}}$ satisfies the conditions imposed on it above, then so will $\tilde{\mathbf{P}}_Z \tilde{\mathbf{R}}_1$, because it is simply the projection of $\tilde{\mathbf{R}}_1$ onto the space spanned by $\tilde{\mathbf{Z}}$. As usual, the number of degrees of freedom of the test will in regular cases be m if $n \leq k$, in which case the DWH and ordinary score test statistics will be numerically identical. When $m > k$, however, the DWH test will have fewer degrees of freedom than the ordinary score test (i.e., at most k).

The DWH versions of score tests may be particularly useful when m is large. Consider White's (1982) information matrix (IM) test. As Lancaster (1984) has shown, this can easily be computed via the OPG regression, which is a special case of regression (27). In this case, $\tilde{\mathbf{r}}$ is an n -vector of ones, $\tilde{\mathbf{R}}_1$ is the matrix $\tilde{\mathbf{G}}_1$, the ti^{th} element of which is $\partial \ell(\theta) / \partial \theta_i$, evaluated at $\tilde{\theta}$, and $\tilde{\mathbf{Z}}$ is a matrix of which a typical element is

$$\frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} + \frac{\partial \ell(\theta)}{\partial \theta_i} \frac{\partial \ell(\theta)}{\partial \theta_j}, \quad i = 1, \dots, m, \quad j = 1, \dots, i, \quad (34)$$

evaluated at $\tilde{\theta}$. The number of columns in $\tilde{\mathbf{Z}}$ is $\frac{1}{2}(k^2 + k)$, although in practice some columns often have to be dropped if $[\tilde{\mathbf{G}}_1 \ \tilde{\mathbf{Z}}]$ has less-than-full rank.

Except when k is very small, the IM test is likely to involve a very large number of degrees of freedom. Various ways to reduce this have been suggested; one could, for example, simply restrict attention to the diagonal elements of the information matrix, setting $j = i$ in (34). But this seems arbitrary. Moreover, as Chesher (1984) has shown, the implicit alternative of the IM test is a form of random parameter variation which will not necessarily be of much economic interest. People frequently employ the test not to check for this type of parameter variation, but because it is thought to have power against a wide range of types of model misspecification. Model misspecification is often of little concern if it does not affect parameter estimates. A possible way to reduce the number of degrees of freedom of the IM test, then, is to use a DWH version of it. This can easily be accomplished by replacing $\tilde{\mathbf{Z}}$ in the artificial regression (27) by $\tilde{\mathbf{P}}_{\mathbf{Z}}\tilde{\mathbf{G}}_1$.

In many circumstances, we believe, the DWH version of the IM test will be more useful than the original. Instead of asking whether there is evidence that the OPG and Hessian estimates of the information matrix differ, the test asks whether there is evidence that they differ for a reason which affects the consistency of the parameter estimates. One could reasonably expect the DWH version of the test to have more power in many cases, since it will have at most k degrees of freedom, instead of $\frac{1}{2}(k^2 + k)$ for the usual IM test. Note, however, that it will still be impossible to compute the test when $n < \frac{1}{2}(k^2 + k)$, since $\tilde{\mathbf{P}}_{\mathbf{Z}}\tilde{\mathbf{G}}_1$ would in that case equal $\tilde{\mathbf{G}}_1$. Even in its DWH version, then, the IM test remains a procedure to be used only when the sample size is reasonably large.

Of course, it makes sense to do a DWH version of the IM test only when the full test is testing in directions which affect parameter consistency. This is by no means always so, since the directions in which the IM test tests are those which affect the consistency of the estimate of the *covariance matrix* of the parameter estimates. (And this implies that even the full IM test will have no power against misspecifications that affect the consistency of the estimates of neither the parameters nor their covariance matrix.) Consider, for instance, the case of linear regression models, where the IM test is implicitly testing for certain forms of heteroskedasticity, skewness, and kurtosis; see Hall (1987). For a linear regression model with normal errors, the contribution to the loglikelihood function from the t^{th} observation is

$$\ell_t = \frac{1}{2} \log(2\pi) - \log(\sigma) - (y_t - \mathbf{X}_t\boldsymbol{\beta})^2 / (2\sigma^2) \quad (35)$$

where $\boldsymbol{\beta}$ is a p -vector so that $k = p + 1$. The contributions to the gradient for $\boldsymbol{\beta}$ and σ , respectively, are

$$G_{ti}(y_t - \mathbf{X}_t\boldsymbol{\beta})X_{ti}/\sigma^2, \quad (36)$$

and

$$G_{t\sigma} = -1/\sigma + (y_t - \mathbf{X}_t\boldsymbol{\beta})^3 / \sigma^2. \quad (37)$$

The second derivatives of (35) are

$$\frac{\partial^2 \ell_t}{\partial \sigma \partial \sigma} = -1/\sigma^2 - 3(y_t - \mathbf{X}_t \boldsymbol{\beta})^2/\sigma^4,$$

$$\frac{\partial^2 \ell_t}{\partial \sigma \partial \beta_i} = -2(y_t - \mathbf{X}_t \boldsymbol{\beta})X_{ti}/\sigma^3,$$

and

$$\frac{\partial^2 \ell_t}{\partial \beta_i \partial \beta_j} = -X_{ti}X_{tj}/\sigma^2.$$

The right-hand side of the OPG regression consists of the test regressors \tilde{Z}_{tj} plus p regressors \tilde{G}_{ti} , which correspond to the β_i , and one regressor $\tilde{G}_{t\sigma}$, which corresponds to σ (\tilde{G}_{ti} and $\tilde{G}_{t\sigma}$ are evaluated at OLS estimates $\tilde{\boldsymbol{\beta}}$ and $\tilde{\sigma}$, the latter using n rather than $n-p$ in the denominator). The test regressor corresponding to any pair of parameters is the sum of the second derivative of ℓ_t with respect to those parameters and the product of the corresponding first derivatives, again evaluated at $\tilde{\boldsymbol{\beta}}$ and $\tilde{\sigma}$.

We simplify all these expressions by using the fact that, since the test statistic is an explained sum of squares, multiplying any regressor by a constant will have no effect on it, and by defining e_t as $\tilde{u}_t/\tilde{\sigma}$. The regressors for the OPG version of the IM test are thus seen to be:

$$\text{for } \beta_i : e_i X_{ti}, \tag{38}$$

$$\text{for } \sigma : e_t^2 - 1, \tag{39}$$

$$\text{for } \beta_i, \beta_j : (e_t^2 - 1)X_{ti}X_{tj} \tag{40}$$

$$\text{for } \beta_i, \sigma : (e_t^3 - 3e_t)X_{ti} \tag{41}$$

$$\text{for } \sigma, \sigma : e_t^4 - 5e_t^2 + 2. \tag{42}$$

When the original regression contains a constant term, (40) will be perfectly collinear with (39) when i and j both refer to the constant term, so that one of them will have to be dropped, and the degrees of freedom for the test reduced by one to $\frac{1}{2}(p^2 + 3p)$.

It is evident that the (β_i, β_j) regressors are testing in directions which correspond to heteroskedasticity of the type that White's (1980) test is designed to detect (namely, heteroskedasticity that affects the consistency of the OLS covariance matrix estimator) and that the (β_i, σ) regressors are testing in directions that correspond to skewness interacting with the X_{ti} . If we subtract (39) from (42), the result is $e_t^4 - 6e_t^2 + 3$, from which we see that the linearly independent part of the (σ, σ) regressor is testing in the kurtosis direction. The IM test is thus seen to be testing for heteroskedasticity, skewness, and kurtosis, none of which prevent $\boldsymbol{\beta}$ from being consistent. Hence it would make no sense to compute a DWH variant of the IM test in this case, and indeed it would be impossible to do so asymptotically. If one did do such a test in practice, one would run into precisely the problem discussed in the previous section: The test might well reject if the model suffered from heteroskedasticity, skewness, or kurtosis, but the rejection would not say anything about the consistency of $\tilde{\boldsymbol{\beta}}$ or of $\tilde{\sigma}$.

5. The Power of DWH and Classical Tests

When the DWH version of a classical test differs from the original, the former may or may not be more powerful than the latter. Although this fact and the reasons for it are reasonably well-known, it seems worthwhile to include a brief discussion which, we hope, makes the issues clear. We shall deal with the general case of section 3, and we will rely heavily on results in Davidson and MacKinnon (1987a).

Suppose the data are generated by a sequence of local DGPs which tends to the point $\boldsymbol{\theta}^0 \equiv (\boldsymbol{\theta}_1^0, \mathbf{0})$. The direction in which the null is incorrect can always be represented by a vector

$$\mathbf{M}_1(\mathbf{R}_2\mathbf{w}_2 + \mathbf{R}_3\mathbf{w}_3),$$

where $\mathbf{M}_1 \equiv \mathbf{M}_1(\boldsymbol{\theta}^0)$, $\mathbf{R}_2 \equiv \mathbf{R}_2(\boldsymbol{\theta}^0)$, and \mathbf{R}_3 is a matrix with the same properties as \mathbf{R}_1 and \mathbf{R}_2 , which represents directions other than those contained in the alternative hypothesis. The vectors \mathbf{w}_2 and \mathbf{w}_3 indicate the weights to be given to the various directions; one can think of \mathbf{w}_2 as being proportional to $\boldsymbol{\theta}_2$. Following Davidson and MacKinnon (1987a), it is possible to show that under such a sequence any of the classical test statistics for the hypothesis $\boldsymbol{\theta}_2 = \mathbf{0}$ will be asymptotically distributed as noncentral $\chi^2(l)$ with noncentrality parameter (or NCP)

$$\begin{aligned} & \text{plim} \left(\frac{1}{n} (\mathbf{w}_2^\top \mathbf{R}_2^\top + \mathbf{w}_3^\top \mathbf{R}_3^\top) \mathbf{M}_1 \mathbf{R}_2 \right) \left(\text{plim} \left(\frac{1}{n} \mathbf{R}_2^\top \mathbf{M}_1 \mathbf{R}_2 \right)^{-1} \right) \\ & \times \text{plim} \left(\frac{1}{n} \mathbf{R}_2^\top \mathbf{M}_1 (\mathbf{R}_2 \mathbf{w}_2 + \mathbf{R}_3 \mathbf{w}_3) \right). \end{aligned} \quad (43)$$

This NCP is the probability limit of n^{-1} times the explained sum of squares from the artificial regression

$$\mathbf{M}_1(\mathbf{R}_2\mathbf{w}_2 + \mathbf{R}_3\mathbf{w}_3) = \mathbf{M}_1\mathbf{R}_2\mathbf{b} + \text{residuals}. \quad (44)$$

When the DGP belongs to the alternative hypothesis, so that $\mathbf{w}_3 = \mathbf{0}$, this regression fits perfectly, and (43) simplifies to

$$\text{plim} \left(\frac{1}{n} \mathbf{w}_2^\top \mathbf{R}_2^\top \mathbf{M}_1 \mathbf{R}_2 \mathbf{w}_2 \right),$$

which is equivalent to expressions for noncentrality parameters found in standard references such as Engle (1984).

Similarly, the noncentrality parameter for the DWH variant of the classical test against $\boldsymbol{\theta}_2 = \mathbf{0}$ will be the probability limit of n^{-1} times the explained sum of squares from the artificial regression

$$\mathbf{M}_1(\mathbf{R}_2\mathbf{w}_2 + \mathbf{R}_3\mathbf{w}_3) = \mathbf{M}_1\mathbf{P}_2\mathbf{R}_2\mathbf{b}^* + \text{residuals}. \quad (45)$$

If we make the definition

$$\mathbf{C} \equiv (\mathbf{R}_2^\top \mathbf{R}_2)^{-1} \mathbf{R}_2^\top \mathbf{R}_1,$$

regression (45) can be rewritten as

$$\mathbf{M}_1(\mathbf{R}_2\mathbf{w}_2 + \mathbf{R}_3\mathbf{w}_3) = \mathbf{M}_1\mathbf{R}_2\mathbf{C}\mathbf{b}^* + \text{residuals}. \quad (46)$$

From (44) and (46), it is clear that the DWH and classical tests will have the same NCP in two circumstances. The first of these is when $l = k$ and the matrix \mathbf{C} has full rank, which is the familiar case where the classical and DWH tests are equivalent. The second is when

$$\mathbf{R}_2\mathbf{w}_2 = \mathbf{R}_2\mathbf{C}\mathbf{w}^*, \quad (47)$$

where \mathbf{w}^* is a k -vector. In both these cases, regressions (44) and (46) will have the same explained sum of squares.

When the DWH test is not equivalent to the classical tests and condition (47) does not hold, it must have a smaller NCP than the classical tests. This will be true whether or not $\mathbf{w}_3 = \mathbf{0}$, since $\mathbf{R}_2\mathbf{C}$ can never have more explanatory power than \mathbf{R}_2 . Whether the DWH test will have more or less power than the classical test then depends on whether its reduced number of degrees of freedom more than offsets its smaller NCP.

6. Binary Choice Models: An Example

In this section, we consider a simple example where a DWH variant of the IM test does make sense. Failures of distributional assumptions, of the sort which do not affect the consistency of least squares estimates, do render ML estimates of binary choice models inconsistent. It is therefore both important to test for these and interesting to see if they are affecting the parameter estimates.

We shall be concerned with the simplest type of binary choice model, in which the dependent variable y_t may be either zero or one and

$$\Pr(y_t = 1) = F(\mathbf{X}_t\boldsymbol{\beta}), \quad (48)$$

where $F(x)$ is a thrice continuously differentiable function which maps from the real line to the 0–1 interval, is weakly increasing in x , and has the properties

$$F(x) \geq 0; \quad F(-\infty) = 0; \quad F(\infty) = 1; \quad F(-x) = 1 - F(x). \quad (49)$$

Two examples are the probit model, where $F(x)$ is the cumulative standard normal distribution function, and the logit model, where $F(x)$ is the logistic function. The contribution to the loglikelihood of the t^{th} observation is

$$\ell_t(\boldsymbol{\beta}) = y_t \log(F(\mathbf{X}_t\boldsymbol{\beta})) + (1 - y_t) \log(F(-\mathbf{X}_t\boldsymbol{\beta})).$$

The contributions to the gradient for $y_t = 1$ and $y_t = 0$ are, respectively,

$$f(\mathbf{X}_t\boldsymbol{\beta})\mathbf{X}_{ti}/F(\mathbf{X}_t\boldsymbol{\beta})$$

and

$$-f(-\mathbf{X}_t\boldsymbol{\beta})\mathbf{X}_{ti}/F(-\mathbf{X}_t\boldsymbol{\beta}),$$

where $f(x)$ is the first derivative of $F(x)$. Thus the corresponding elements of the matrix $\mathbf{G}^\top\mathbf{G}$ are

$$\left(\frac{f(\mathbf{X}_t\boldsymbol{\beta})}{F(\mathbf{X}_t\boldsymbol{\beta})}\right)^2 X_{ti}X_{tj} \quad (50)$$

and

$$\left(\frac{f(-\mathbf{X}_t\boldsymbol{\beta})}{F(-\mathbf{X}_t\boldsymbol{\beta})}\right)^2 X_{ti}X_{tj}. \quad (51)$$

The second derivatives of $\ell_t(\boldsymbol{\beta})$ for $y_t = 1$ and $y_t = 0$ are, respectively,

$$\frac{f'(\mathbf{X}_t\boldsymbol{\beta})F(\mathbf{X}_t\boldsymbol{\beta}) - f^2(\mathbf{X}_t\boldsymbol{\beta})}{F^2(\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj} \quad (52)$$

and

$$\frac{-f'(\mathbf{X}_t\boldsymbol{\beta})F(-\mathbf{X}_t\boldsymbol{\beta}) - f^2(-\mathbf{X}_t\boldsymbol{\beta})}{F^2(-\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj}, \quad (53)$$

where $f'(x)$ denotes the derivative of $f(x)$, and we have used the symmetry property of (49) which implies that $f'(x) = -f'(-x)$. The sum of (50) and (52) is

$$\frac{f'(\mathbf{X}_t\boldsymbol{\beta})}{F(\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj}, \quad (54)$$

and the sum of (51) and (53) is

$$\frac{-f'(\mathbf{X}_t\boldsymbol{\beta})}{F(-\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj}. \quad (55)$$

The random variable whose two possible realizations are (54) and (55) is the difference between the OPG and minus the Hessian. If the model is correctly specified, the expectation of this random variable is

$$\begin{aligned} & F(\mathbf{X}_t\boldsymbol{\beta}) \left(\frac{f'(\mathbf{X}_t\boldsymbol{\beta})}{F(\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj} \right) + F(-\mathbf{X}_t\boldsymbol{\beta}) \left(\frac{-f'(\mathbf{X}_t\boldsymbol{\beta})}{F(-\mathbf{X}_t\boldsymbol{\beta})} X_{ti}X_{tj} \right) \\ &= f'(\mathbf{X}_t\boldsymbol{\beta})X_{ti}X_{tj} - f'(\mathbf{X}_t\boldsymbol{\beta})X_{ti}X_{tj} = 0. \end{aligned}$$

The IM test asks whether it is in fact equal to zero.

The IM test may be based on the OPG regression, as usual, or it may be based on the artificial regression proposed by Engle (1984) and Davidson and MacKinnon (1984b) specifically for binary choice models, which we shall refer to as the PL (for probit and logit) regression. Computing the IM test by means of an artificial regression other than the OPG regression may be attractive because of the very poor finite-sample

properties of the latter; see Chesher and Spady (1988), Davidson and MacKinnon (1985a), Kennan and Neumann (1988), and Orme (1987b).

The regressand for the PL artificial regression is

$$\tilde{r}_t = \frac{y_t - F(\mathbf{X}_t\tilde{\boldsymbol{\beta}})}{(F(-\mathbf{X}_t\tilde{\boldsymbol{\beta}})F(\mathbf{X}_t\tilde{\boldsymbol{\beta}}))^{1/2}}, \quad (56)$$

and the regressors corresponding to the β_i are

$$\frac{f'(\mathbf{X}_t\tilde{\boldsymbol{\beta}})X_{ti}X_{tj}}{(F(-\mathbf{X}_t\tilde{\boldsymbol{\beta}})F(\mathbf{X}_t\tilde{\boldsymbol{\beta}}))^{1/2}}. \quad (57)$$

We want to construct the test regressors so that the ij^{th} test regressor times (56) yields (54) when $y_t = 1$ and (55) when $y_t = 0$. It is thus easily seen that the ij^{th} test regressor must be

$$\tilde{Z}_{t,ij} = \frac{f'(\mathbf{X}_t\tilde{\boldsymbol{\beta}})X_{ti}X_{tj}}{(F(-\mathbf{X}_t\tilde{\boldsymbol{\beta}})F(\mathbf{X}_t\tilde{\boldsymbol{\beta}}))^{1/2}}. \quad (58)$$

This artificial regression was also derived by Orme (1987a).

In the probit case, this artificial regression has a very interesting interpretation. Since $f(\mathbf{X}_t\boldsymbol{\beta})$ is the standard normal density,

$$f'(\mathbf{X}_t\boldsymbol{\beta}) = -(2\pi)^{-1/2} \exp\left(\frac{1}{2}(\mathbf{X}_t\boldsymbol{\beta})^2\right) \mathbf{X}_t\boldsymbol{\beta} = -\mathbf{X}_t\boldsymbol{\beta} f(\mathbf{X}_t\boldsymbol{\beta}),$$

so that (58) becomes

$$\frac{-f(\mathbf{X}_t\tilde{\boldsymbol{\beta}})\mathbf{X}_t\tilde{\boldsymbol{\beta}}X_{ti}X_{tj}}{(F(-\mathbf{X}_t\tilde{\boldsymbol{\beta}})F(\mathbf{X}_t\tilde{\boldsymbol{\beta}}))^{1/2}}. \quad (59)$$

This is identical to the test regressor one would get if one did an LM test of the model (48) against the alternative

$$\Pr(y_t = 1) = F\left(\mathbf{X}_t\boldsymbol{\beta} / \exp\left(\sum_{i=1}^k \sum_{j=1}^i X_{ti}X_{tj}\gamma_{ij}\right)\right), \quad (60)$$

which can be derived from the latent variable model

$$y_t^* = \mathbf{X}_t\boldsymbol{\beta} + u_t, \quad u_t \sim N\left(0, \exp\left(2 \sum_{i=1}^k \sum_{j=1}^i X_{ti}X_{tj}\gamma_{ij}\right)\right) \quad (61)$$

$$y_t = 1 \text{ if } y_t^* > 0, \quad y_t = 0 \text{ otherwise.}$$

The model (61) is thus a special case of a model which incorporates a natural form of heteroskedasticity. The general model was considered by Davidson and MacKinnon (1984b), who derived the appropriate LM test. This model is special because the

variance of u_t depends exclusively on the cross-products of the X_{ti} . It is clear that the implicit alternative of the IM test is precisely this heteroskedastic model. Moreover, just as for ordinary regression models it is only heteroskedasticity related to the cross-products of the regressors which affects the consistency of the covariance matrix estimates, so for probit models it is only heteroskedasticity of this type which (locally) prevents the information matrix equality from holding and which thus renders ML probit estimates inconsistent. This is purely a local result, of course; if a DGP involving any form of heteroskedasticity were some fixed distance from the probit model, one could not expect ML estimates based on homoskedasticity to be consistent.

Notice that if one of the X_{ti} , say X_{tj} , is a constant term, the test regressor (59) which corresponds to X_{tj} is

$$\frac{-f(\mathbf{X}_t\tilde{\boldsymbol{\beta}})\mathbf{X}_t\tilde{\boldsymbol{\beta}}}{(F(\mathbf{X}_t\tilde{\boldsymbol{\beta}})F(-\mathbf{X}_t\tilde{\boldsymbol{\beta}}))^{1/2}},$$

which is a linear combination of the regressors (57) that correspond to the $\tilde{\beta}_i$. This test regressor must therefore be dropped, and the degrees of freedom of the test reduced to $\frac{1}{2}k(k+1) - 1$.

Newey (1985) recognized that the IM test implicitly tests against heteroskedasticity in the case of probit models, and he suggested that this test may be particularly attractive for such models. He proposed to use the OPG form of the test. The PL version discussed here is no more difficult to compute than the OPG form, however, and it seems to have much better finite-sample properties. In the Appendix, we present detailed results of a small Monte Carlo experiment designed to investigate the performance of the OPG and PL tests, in regular and DWH versions.

The Monte Carlo experiments yielded two main results. First, we found that the PL form of the IM test proposed above performed reasonably well for moderately large samples, but that it generated far too many realizations in the right-hand tail, and that its rate of convergence to its asymptotic distribution was disappointingly slow. (Since binary choice models are often estimated with large samples, this will not necessarily be a fatal drawback.) In contrast, the OPG form of the IM test performed dismally in samples of all sizes, sometimes rejecting the null more than 90% of the time at the nominal 5% level. These dreadful results are consistent with those of Chesher and Spady (1988), Davidson and MacKinnon (1985a), Kennan and Neumann (1988), and Orme (1987) for other applications of the OPG form of the IM test; the problem is primarily that $n^{-1}\tilde{\mathbf{G}}^T\tilde{\mathbf{G}}$ tends to estimate $\mathcal{J}(\boldsymbol{\theta})$ very poorly. Like Kennan and Neumann, we found that the performance of the OPG form deteriorated markedly as the number of degrees of freedom for the test was increased.

Second, we found that in some realistic cases, but not all cases, the DWH version of the IM test can have significantly more power than the ordinary version. This was most likely to be the case when the number of parameters was large, so that the DWH variant of the IM test would have many fewer degrees of freedom than the ordinary IM test; unfortunately, this was also the circumstance in which the finite-sample performance of all the tests was worst.

Based on these results, we find it difficult to endorse Newey's (1985) recommendation of the IM test for probit models. The conventional OPG form of the test should clearly not be used. Among the tests we studied, the DWH version computed via the PL regression generally performs the best, both under the null and under the alternatives we studied, but its finite-sample performance is far from ideal. It might well be more productive to test for particular, relatively simple forms of heteroskedasticity which do not involve many degrees of freedom, especially those which seem plausible for the model at hand, rather than to calculate any form of the IM test.

7. Conclusion

This paper has dealt with several aspects of Durbin-Wu-Hausman tests. The main contribution of the paper has been to show that DWH tests may be based on any artificial regression that can be used to compute score-type tests, and that any test based on such a regression can be converted into a DWH test. In particular, we have shown that this is true for the information matrix test, and we have demonstrated how to compute a DWH version of the IM test for the case of binary choice models.

References

- Amemiya, T. (1985). *Advanced Econometrics*, Cambridge, MA, Harvard University Press.
- Bera, A. K., and C. R. McKenzie (1986). "Alternative forms and properties of the score test," *Journal of Applied Statistics*, 13, 13–25.
- Boothe, P., and J. G. MacKinnon (1986). "A specification test for models estimated by GLS," *Review of Economics and Statistics*, 68, 711–714.
- Breusch, T. S., and L. G. Godfrey (1986). "Data transformation tests," *Economic Journal*, 96, 47–58.
- Chesher, A. (1984). "Testing for neglected heterogeneity," *Econometrica*, 52, 865–872.
- Chesher, A., and R. Spady (1988). "Asymptotic expansions of the information matrix test statistic," presented at the Econometric Study Group meeting, Bristol.
- Das Gupta, S., and M. D. Perlman (1974). "Power of the noncentral F test: Effect of additional variates on Hotelling's τ^2 test," *Journal of the American Statistical Association*, 69, 174–180.
- Davidson, R., L. G. Godfrey, and J. G. MacKinnon (1985). "A simplified version of the differencing test," *International Economic Review*, 26, 639–647.
- Davidson, R., and J. G. MacKinnon (1983). "Small sample properties of alternative forms of the Lagrange Multiplier test," *Economics Letters* 12, 269–275.

- Davidson, R., and J. G. MacKinnon (1984a). "Model specification tests based on artificial linear regressions, *International Economic Review*, 25, 485–502.
- Davidson, R., and J. G. MacKinnon (1984b). "Convenient specification tests for logit and probit models," *Journal of Econometrics*, 25, 241–262.
- Davidson, R., and J. G. MacKinnon (1985a). "Testing linear and loglinear regressions against Box-Cox alternatives," *Canadian Journal of Economics*, 25, 499–517.
- Davidson, R., and J. G. MacKinnon (1985b). "Heteroskedasticity-robust tests in regression directions. *Annales de l'INSEE*, 59/60, 183–218.
- Davidson, R., and J. G. MacKinnon (1987a). "Implicit alternatives and the local power of test statistics. *Econometrica*, 55, 1305–1329.
- Davidson, R., and J. G. MacKinnon (1988). "Double-length artificial regressions. *Oxford Bulletin of Economics and Statistics*, 50, 203–217.
- Durbin, J. (1954). "Errors in variables," *Review of the International Statistical Institute*, 22, 23–32.
- Engle, R. F. (1982). "A general approach to Lagrange multiplier model diagnostics, *Journal of Econometrics*, 20, 83–104.
- Engle, R. F. (1984). "Wald, likelihood ratio and Lagrange multiplier tests in econometrics," in Z. Griliches and M. Intriligator (ed.), *Handbook of Econometrics*, Amsterdam, North Holland.
- Fisher, G. R., and R. J. Smith (1985). "Least squares theory and the Hausman specification test," Queen's University Economics Working Paper No. 641.
- Godfrey, L. G., and M. R. Wickens (1981). "Testing linear and log-linear regressions for functional form," *Review of Economic Studies*, 48, 487–496.
- Hall, A. (1987). "The information matrix test for the linear model," *Review of Economic Studies*, 54, 257–263.
- Hausman, J. A. (1978). "Specification tests in econometrics," *Econometrica*, 46, 1251–1272.
- Holly, A. (1982). "A remark on Hausman's specification test. *Econometrica*, 50, 749–759.
- Kennan, J., and G. R. Neumann (1988). "Why does the information matrix test reject too often? A diagnosis of some Monte Carlo symptoms," Hoover Institution, Stanford University, Working Papers in Economics E-88-10.
- Lancaster, T. (1984). "The covariance matrix of the information matrix test," *Econometrica*, 52, 1051–1053.

- Nakamura, A., and M. Nakamura (1981). "On the relationships among several specification error tests presented by Durbin, Wu and Hausman." *Econometrica*, 49, 1583–1588.
- Newey, W. K. (1985). "Maximum likelihood specification testing and conditional moment tests," *Econometrica*, 53, 1047–1070.
- Orme, C. (1987a). "The calculation of the information matrix test for binary data models," University of York, mimeo.
- Orme, C. (1987b). "The small sample performance of the information matrix test," University of York, mimeo.
- Plosser, C. I., G. W. Schwert, and H. White (1982). "Differencing as a test of specification," *International Economic Review*, 23, 535–552.
- Ramsey, J. B. (1969). "Tests for specification errors in classical linear least-squares regression analysis," *Journal of the Royal Statistical Society, Series B*, 31, 350–371.
- Ruud, P. A. (1984). "Tests of specification in econometrics." *Econometric Reviews*, 3, 211–242.
- Thursby, J. G., and P. Schmidt (1977). "Some properties of tests for specification error in a linear regression model," *Journal of the American Statistical Association* 72, 635–641.
- White, H. (1980). "A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity," *Econometrica*, 48, 817–838.
- White, H. (1982). "Maximum likelihood estimation of misspecified models," *Econometrica*, 50, 1-25.
- Wu, D. (1973). "Alternative tests of independence between stochastic regressors and disturbances," *Econometrica*, 41, 733–750.

Appendix

In this appendix, we report the results of a small Monte Carlo experiment designed to shed light on Newey’s (1985) conjecture that the IM test may be particularly attractive for probit models. There are two main results. First, we find that the OPG form of the IM test for probit models rejects the null far too often in samples of moderate or even rather large size, while the PL form of the IM test proposed in this paper performs much better. Second, we find that in some realistic cases the DH version of the IM test may have significantly more power than the ordinary version.

In all our experiments, the matrix \mathbf{X} consisted of a constant term and one or more other regressors, which were normally distributed and equi-correlated with correlation one half. Only one set of realizations of these variables was generated, and only for 100 observations. For larger sample sizes, this set of observations was replicated as many times as necessary. This scheme reduced the costs of the simulation, made it easy to calculate NCPs (which for a given test depend only on \mathbf{X} and on the parameters of the DGP), and ensured that any changes as n was increased were not due to changes in the pattern of the exogenous variables.

We first investigated the performance under the null of the ordinary IM test and its DH version, calculated by both the OPG and PL regressions, for samples of size 100, 200, 400, 800, and 1600. We let k , the number of parameters under the null hypothesis, vary from 2 to 4, so that the number of degrees of freedom for the ordinary IM test was 2, 5, or 9, and for the DH version 2, 3, or 4. The DH and ordinary IM tests are thus identical when $k = 2$.

Results for samples of sizes 100, 400, and 1600 are shown in Table 1. The most striking result is the extreme tendency to over-reject of the OPG tests, which worsens rapidly as k increases and diminishes only slowly as the sample size increases. For $k = 4$, the OPG IM test rejects over 98% of the time at the nominal 5% level when $n = 100$, and over 50% of the time even when $n = 1600$. It is clear that the sample would have to be enormous for this test’s true size to be anywhere close to its nominal one. The DH version of the OPG test is slightly better behaved than the original, but the improvement is marginal. Previous results on the finite-sample performance of the OPG test have generally not been favorable to it, but the present results are far worse than those reported previously. Since most applications are likely to involve many more than four parameters, it seems doubtful that the OPG form of the IM test for probit models can ever yield even approximately reliable results in samples of the size that are typically used by econometricians.

The tests based on the PL regression are far better behaved than the OPG tests, but are still a long way from their asymptotic distribution even in samples of 1600. They have roughly the right mean, but their standard deviations are too high because very large values occur much more often than they should by chance. As a result, they tend to under-reject at the 10% level and over-reject at the 1% level, while being fairly close to their nominal size at 5%. Curiously, the problem of too many outliers appears initially to get worse as n increases; for $k = 4$ (the worst case), the standard deviation

for both the ordinary and DH versions is largest for $n = 400$, as is the rejection frequency at the nominal 1% level.

Since the OPG test rejects so often as to be completely useless, there is apparently no choice but to use the PL version; however, these results suggest that even it should be regarded with considerable suspicion, especially if there are more than a very few parameters and the sample size is not very large indeed.

Our second set of experiments was designed to investigate power when the data were generated by (61). Calculation of NCPs showed that, for a wide range of γ_{ij} chosen so that all cross-products contributed very roughly the same amount to the variance, the NCP for the DH version was only slightly smaller than the NCP for the ordinary IM test. In more extreme cases, such as when only one γ_{ij} was non-zero, the NCP for the DH version could be less than half as large. In the former case, the DH version should be more powerful asymptotically, since a slight reduction in the NCP is more than offset by what can be a substantial reduction in degrees of freedom, but in the latter the ordinary IM test would be more powerful.

The object of the Monte Carlo experiments was to see how accurately the asymptotic analysis of Section 5 predicted finite-sample power. We considered a single “plausible” pattern for the γ_{ij} and then scaled the latter to the sample size so that the tests would have power somewhere around 50% at the nominal 5% level. The resulting NCPs, which are of course invariant to the sample size, were 5.15 for $k = 2$, 6.18 and 5.97 (DH version) for $k = 3$, and 9.16 and 8.57 (DH version) for $k = 4$.

Results for the PL tests only are shown in Table 2; results for the OPG tests are not shown because, as one would expect from the results in Table 1, they always rejected far more often than asymptotic theory predicted. The table also shows, in rows labelled “Asymp”, the values that would be expected if the test statistics actually had their asymptotic non-central chi-squared distributions.

The behavior of the PL tests when the null is false is broadly consistent with their behavior when it is true. In particular, they reject much too often at the 1% level, and they have means which are often far too large, because there are many more extremely large values than asymptotic theory predicts. However, they do not consistently under-reject at the 10% level, and the pattern as n increases is not always monotonic. For the case considered here, asymptotic analysis predicts that the DH version will have a modest power advantage. This is usually the case in the experimental results as well, although the ordinary IM test is sometimes more powerful when n is small, especially at the 1% level.

Table 1. Performance of Alternative Tests Under the Null

k	Obs.	Test	Mean	Std. Dev.	Rej. 10%	Rej. 5%	Rej. 1%
2	100	OPG	6.94**	7.36**	46.5**	40.1**	28.2**
		PL	1.81*	2.67**	7.8*	5.2	2.3**
2	400	OPG	4.11**	5.45**	28.0**	21.4**	12.4**
		PL	1.91	2.04	8.2*	4.1	1.2
2	1600	OPG	2.74**	3.60**	17.3**	11.0**	4.9**
		PL	1.96	1.95	9.3	4.3	1.0
3	100	OPG	21.71**	11.27**	84.5**	79.5**	67.6**
		PL	4.14**	5.87**	6.6**	4.9	2.7**
		OPG-DH	15.79**	11.07**	76.7**	70.1**	57.5**
		PL-DH	2.44**	4.69**	6.0**	4.4	2.3**
3	400	OPG	13.40**	10.82**	55.7**	47.5**	32.2**
		PL	4.79	4.87**	9.8	6.4*	3.1**
		OPG-DH	9.31**	9.09**	51.2**	44.0**	30.3**
		PL-DH	2.81*	3.51**	8.8	5.7	2.4**
3	1600	OPG	8.53**	6.94**	33.6**	25.3**	14.2**
		PL	4.91	4.04**	9.7	6.2*	2.9**
		OPG-DH	5.55**	5.81**	29.3**	22.9**	12.7**
		PL-DH	2.95	3.21**	9.6	5.8	2.4**
4	100	OPG	35.37**	8.70**	99.4**	98.3**	94.2**
		PL	6.64**	5.72**	6.2**	4.5	2.6**
		OPG-DH	22.09**	10.01**	92.0**	88.7**	79.4**
		PL-DH	2.60**	3.19**	4.8**	3.0**	1.6*
4	400	OPG	37.48**	21.15**	88.9**	84.4**	75.2**
		PL	8.22**	7.92**	9.1	6.6*	3.7**
		OPG-DH	24.72**	19.82**	82.0**	75.8**	65.4**
		PL-DH	3.57**	5.17**	8.0*	5.6	2.8**
4	1600	OPG	21.94**	15.81**	62.5**	53.1**	38.5**
		PL	8.78	5.51**	9.9	6.0*	2.6**
		OPG-DH	12.92**	13.35**	55.7**	47.3**	33.3**
		PL-DH	3.76*	3.63**	8.9	5.8	2.2**

All results are based on 2000 replications.

* and ** indicate that a quantity differs from what it should be asymptotically at the .05 and .001 levels, respectively.

Degrees of freedom for the ordinary IM tests are 2 for $k = 2$, 5 for $k = 3$, and 9 for $k = 4$.

The standard deviations of χ^2 random variables with 2, 3, 4, 5, and 9 degrees of freedom are, respectively, 2, 2.45, 2.83, 3.16, and 4.24.

Table 2. Power of Alternative Tests

k	Obs.	Test	Mean	Rej. 10%	Rej. 5%	Rej. 1%
2	Asymp.	PL	7.15	64.0	51.6	28.3
	100	PL	7.40	47.4**	40.3**	27.1*
	200	PL	7.49*	53.6**	44.0**	27.7
	400	PL	7.77**	59.5**	48.9*	30.0
	800	PL	7.90**	62.1	51.0	30.1
	1600	PL	7.76**	65.4	53.2	31.4*
3	Asymp.	PL	11.18	57.4	44.5	22.6
		PL-DH	8.97	64.0	51.5	28.5
	100	PL	35.19**	55.1*	52.1**	46.8**
		PL-DH	31.43**	57.2**	51.5	44.9**
	200	PL	37.29**	63.2**	58.4**	51.6**
		PL-DH	33.02**	65.5	58.8**	49.0**
	400	PL	30.76	61.3**	56.4**	48.9**
		PL-DH	26.62**	64.9	57.5**	47.2**
	800	PL	23.22**	65.2**	58.5**	47.8**
		PL-DH	19.27**	67.2*	60.4**	46.3**
4	Asymp.	PL	18.16	64.5	51.9	28.6
		PL-DH	12.57	75.0	63.9	40.1
	100	PL	114.08**	50.8**	48.9*	45.9**
		PL-DH	114.18**	51.7**	50.1**	47.3**
	200	PL	159.94**	67.0*	64.3**	59.6**
		PL-DH	141.66**	69.8**	66.2*	61.0**
	400	PL	120.85**	68.7**	64.1**	57.8**
		PL-DH	107.29**	73.0*	68.8**	61.1**
	800	PL	69.29**	66.0	61.2**	51.4**
		PL-DH	60.12**	72.6*	67.2*	56.9**
1600		PL	33.78**	60.0**	53.4	40.5**
		PL-DH	27.15**	68.5**	60.6*	48.6**

All results are based on 2000 replications.

* and ** indicate that a quantity differs from what it should be asymptotically at the .05 and .001 levels, respectively.

Degrees of freedom for the ordinary IM tests are 2 for $k = 2$, 5 for $k = 3$, and 9 for $k = 4$.