

Davidson, Russell; MacKinnon, James

**Working Paper**

## Artificial regressions

Queen's Economics Department Working Paper, No. 1038

**Provided in Cooperation with:**

Queen's University, Department of Economics (QED)

Suggested Citation: Davidson, Russell; MacKinnon, James (2001) : Artificial regressions, Queen's Economics Department Working Paper, No. 1038, Queen's University, Department of Economics, Kingston (Ontario)

This Version is available at:

<http://hdl.handle.net/10419/67825>

**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. 1038

## Artificial Regressions

Russell Davidson  
McGill University

James MacKinnon  
Queen's University

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

1-2001

# Artificial Regressions

by

**Russell Davidson**

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

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

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

and

**James G. MacKinnon**

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

**jgm@econ.queensu.ca**

## Abstract

Associated with every popular nonlinear estimation method is at least one “artificial” linear regression. We define an artificial regression in terms of three conditions that it must satisfy. Then we show how artificial regressions can be useful for numerical optimization, testing hypotheses, and computing parameter estimates. Several existing artificial regressions are discussed and are shown to satisfy the defining conditions, and a new artificial regression for regression models with heteroskedasticity of unknown form is introduced.

A slightly earlier version of this paper appeared as GREQAM Document de Travail 99a04. This research was supported, in part, by the Social Sciences and Humanities Research Council of Canada.

January, 2001

## 1. Introduction

All popular nonlinear estimation methods, including nonlinear least squares (NLS), maximum likelihood (ML), and the generalized method of moments (GMM), yield estimators which are asymptotically linear. Provided the sample size is large enough, the behavior of these nonlinear estimators in the neighborhood of the true parameter values closely resembles the behavior of the ordinary least squares (OLS) estimator. A particularly illuminating way to see the relationship between any nonlinear estimation method and OLS is to formulate the **artificial regression** that corresponds to the nonlinear estimator.

An artificial regression is a linear regression in which the regressand and regressors are constructed as functions of the data and parameters of the nonlinear model that is really of interest. In addition to helping us understand the asymptotic properties of nonlinear estimators, artificial regressions are often extremely useful as calculating devices. Among other things, they can be used to estimate covariance matrices, as key ingredients of nonlinear optimization methods, to compute one-step efficient estimators, and to calculate test statistics.

In the next section, we discuss the defining properties of an artificial regression. In the subsequent section, we introduce the Gauss-Newton regression, which is probably the most popular artificial regression. Then, in Section 4, we illustrate a number of uses of artificial regressions, using the Gauss-Newton regression as an example. In Section 5, we develop the most important use of artificial regressions, namely, hypothesis testing. We go beyond the Gauss-Newton regression in Sections 6 and 7, in which we introduce two quite generally applicable artificial regressions, one for models estimated by maximum likelihood, and one for models estimated by the generalized method of moments. Section 8 shows how artificial regressions may be modified to take account of the presence of heteroskedasticity of unknown form. Then, in Sections 9 and 10, we discuss double-length regressions and artificial regressions for binary response models, respectively.

## 2. The Concept of an Artificial Regression

Consider a fully parametric nonlinear model that is characterized by a parameter vector  $\boldsymbol{\theta}$  which belongs to a parameter space  $\Theta \subseteq \mathbb{R}^k$  and which can be estimated by minimizing a criterion function  $Q(\boldsymbol{\theta})$  using  $n$  observations. In the case of a nonlinear regression model estimated by nonlinear least squares,  $Q(\boldsymbol{\theta})$  would be one half the sum of squared residuals, and in the case of a model estimated by maximum likelihood,  $Q(\boldsymbol{\theta})$  would be minus the loglikelihood function.

If an artificial regression exists for such a model, it always involves two things: a regressand,  $\boldsymbol{r}(\boldsymbol{\theta})$ , and a matrix of regressors,  $\boldsymbol{R}(\boldsymbol{\theta})$ . The number of regressors for the artificial regression is equal to  $k$ , the number of parameters. The number of “observations” for the artificial regression is often equal to  $n$ , but it may also be equal to a small integer, such as 2 or 3, times  $n$ . We can write a generic artificial

regression as

$$\mathbf{r}(\boldsymbol{\theta}) = \mathbf{R}(\boldsymbol{\theta})\mathbf{b} + \text{residuals}, \quad (1)$$

where  $\mathbf{b}$  is a  $k$ -vector of coefficients. “Residuals” is used here as a neutral term to avoid any implication that (1) is a statistical model. The regressand and regressors in (1) can be evaluated at any point  $\boldsymbol{\theta} \in \Theta$ , and the properties of the artificial regression will depend on the point at which they are evaluated. In many cases, we will want to evaluate (1) at a vector of estimates  $\hat{\boldsymbol{\theta}}$  that is root- $n$  consistent. This means that, if the true parameter vector is  $\boldsymbol{\theta}_0 \in \Theta$ , then  $\hat{\boldsymbol{\theta}}$  approaches  $\boldsymbol{\theta}_0$  at a rate proportional to  $n^{-1/2}$ . One such vector that is of particular interest is  $\hat{\boldsymbol{\theta}}$ , the vector of estimates which minimizes the criterion function  $Q(\boldsymbol{\theta})$ .

For (1) to constitute an artificial regression, the vector  $\mathbf{r}(\boldsymbol{\theta})$  and the matrix  $\mathbf{R}(\boldsymbol{\theta})$  must satisfy certain defining properties. These may be stated in a variety of ways, which depend on the class of models to which the artificial regression is intended to apply. For the purposes of this paper, we will say that (1) is an artificial regression if it satisfies the following three conditions:

- (i) The estimator  $\hat{\boldsymbol{\theta}}$  is defined, uniquely in a neighborhood in  $\Theta$ , by the  $k$  equations  $\mathbf{R}^\top(\hat{\boldsymbol{\theta}})\mathbf{r}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$ ;
- (ii) for any root- $n$  consistent  $\hat{\boldsymbol{\theta}}$ , a consistent estimate of  $\text{Var}(\text{plim } n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0))$  is given by the inverse of  $n^{-1}\mathbf{R}^\top(\hat{\boldsymbol{\theta}})\mathbf{R}(\hat{\boldsymbol{\theta}})$ . Formally,

$$\text{Var}\left(\text{plim}_{n \rightarrow \infty} n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\right) = \text{plim}_{n \rightarrow \infty} (n^{-1}\mathbf{R}^\top(\hat{\boldsymbol{\theta}})\mathbf{R}(\hat{\boldsymbol{\theta}}))^{-1};$$

- (iii) if  $\hat{\mathbf{b}}$  denotes the vector of estimates from the artificial regression (1) with regressand and regressors evaluated at  $\hat{\boldsymbol{\theta}}$ , then

$$\hat{\boldsymbol{\theta}} + \hat{\mathbf{b}} = \hat{\boldsymbol{\theta}} + o_p(n^{-1/2}).$$

Many artificial regressions actually satisfy a stronger version of condition (i):

$$(i') \quad \mathbf{g}(\boldsymbol{\theta}) = -\mathbf{R}^\top(\boldsymbol{\theta})\mathbf{r}(\boldsymbol{\theta}),$$

where  $\mathbf{g}(\boldsymbol{\theta})$  denotes the gradient of the criterion function  $Q(\boldsymbol{\theta})$ . Clearly, condition (i') implies condition (i), but not *vice versa*. The minus sign in (i') is due to the arbitrary choice that the estimator is defined by minimizing  $Q(\boldsymbol{\theta})$  rather than maximizing it.

Condition (ii) has been written in a particularly simple form, and some nonstandard artificial regressions do not actually satisfy it. However, as we will see, this does not prevent them from having essentially the same properties as artificial regressions that do satisfy it.

Condition (iii), which is perhaps the most interesting of the three conditions, will be referred to as the **one-step property**. It says that, if we take one step from

an initial consistent estimator  $\hat{\theta}$ , where the step is given by the coefficients  $\hat{b}$  from the artificial regression, we will obtain an estimator that is asymptotically equivalent to  $\hat{\theta}$ .

The implications of these three conditions will become clearer when we study specific artificial regressions in the remainder of this paper. These conditions differ substantially from the conditions used to define an artificial regression in Davidson and MacKinnon (1990), because that paper was concerned solely with artificial regressions for models estimated by maximum likelihood.

### 3. The Gauss-Newton Regression

Associated with every nonlinear regression model is a somewhat nonstandard artificial regression which is probably more widely used than any other. Consider the univariate, nonlinear regression model

$$y_t = x_t(\boldsymbol{\beta}) + u_t, \quad u_t \sim \text{IID}(0, \sigma^2), \quad t = 1, \dots, n, \quad (2)$$

where  $y_t$  is the  $t^{\text{th}}$  observation on the dependent variable, and  $\boldsymbol{\beta}$  is a  $k$ -vector of parameters to be estimated. The scalar function  $x_t(\boldsymbol{\beta})$  is a nonlinear regression function. It determines the mean value of  $y_t$  as a function of unknown parameters  $\boldsymbol{\beta}$  and, usually, of explanatory variables, which may include lagged dependent variables. The explanatory variables are not shown explicitly in (2), but the  $t$  subscript on  $x_t(\boldsymbol{\beta})$  reminds us that they are present. The model (2) may also be written as

$$\mathbf{y} = \mathbf{x}(\boldsymbol{\beta}) + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (3)$$

where  $\mathbf{y}$  is an  $n$ -vector with typical element  $y_t$ , and  $\mathbf{x}(\boldsymbol{\beta})$  is an  $n$ -vector of which the  $t^{\text{th}}$  element is  $x_t(\boldsymbol{\beta})$ .

The nonlinear least squares (NLS) estimator  $\hat{\boldsymbol{\beta}}$  for model (3) minimizes the sum of squared residuals. It is convenient to use for the criterion function to be minimized this sum divided by 2. Thus we define

$$Q(\boldsymbol{\beta}) = \frac{1}{2}(\mathbf{y} - \mathbf{x}(\boldsymbol{\beta}))^\top (\mathbf{y} - \mathbf{x}(\boldsymbol{\beta})). \quad (4)$$

The Gauss-Newton regression can be derived as an approximation to **Newton's Method** for the minimization of  $Q(\boldsymbol{\beta})$ . In this case, Newton's Method consists of the following iterative procedure. One starts from some suitably chosen starting value,  $\boldsymbol{\beta}_{(0)}$ . At step  $m$  of the procedure,  $\boldsymbol{\beta}_{(m)}$  is updated by the formula

$$\boldsymbol{\beta}_{(m+1)} = \boldsymbol{\beta}_{(m)} - \mathbf{H}_{(m)}^{-1} \mathbf{g}_{(m)},$$

where the  $k \times 1$  vector  $\mathbf{g}_{(m)}$  and the  $k \times k$  matrix  $\mathbf{H}_{(m)}$  are, respectively, the gradient and the Hessian of  $Q(\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$ , evaluated at  $\boldsymbol{\beta}_{(m)}$ . For general  $\boldsymbol{\beta}$ , we have

$$\mathbf{g}(\boldsymbol{\beta}) = -\mathbf{X}^\top(\boldsymbol{\beta})(\mathbf{y} - \mathbf{x}(\boldsymbol{\beta})),$$

where the matrix  $\mathbf{X}(\boldsymbol{\beta})$  is an  $n \times k$  matrix with  $ti^{\text{th}}$  element the derivative of  $x_t(\boldsymbol{\beta})$  with respect to  $\beta_i$ , the  $i^{\text{th}}$  component of  $\boldsymbol{\beta}$ . A typical element of the Hessian  $\mathbf{H}(\boldsymbol{\beta})$  is

$$H_{ij}(\boldsymbol{\beta}) = - \sum_{t=1}^n \left( (y_t - x_t(\boldsymbol{\beta})) \frac{\partial X_{ti}(\boldsymbol{\beta})}{\partial \beta_j} - X_{ti}(\boldsymbol{\beta}) X_{tj}(\boldsymbol{\beta}) \right), \quad i, j = 1, \dots, k. \quad (5)$$

The Gauss-Newton procedure is one of the set of so-called **quasi-Newton** procedures, in which the exact Hessian is replaced by an approximation. Here, only the second term in (5) is used, so that the  $\mathbf{H}(\boldsymbol{\beta})$  of Newton's method is replaced by the matrix  $\mathbf{X}^\top(\boldsymbol{\beta})\mathbf{X}(\boldsymbol{\beta})$ . Thus the Gauss-Newton updating formula is

$$\boldsymbol{\beta}_{(m+1)} = \boldsymbol{\beta}_{(m)} + (\mathbf{X}_{(m)}^\top \mathbf{X}_{(m)})^{-1} \mathbf{X}_{(m)}^\top (\mathbf{y} - \mathbf{x}_{(m)}), \quad (6)$$

where we write  $\mathbf{X}_{(m)} = \mathbf{X}(\boldsymbol{\beta}_{(m)})$  and  $\mathbf{x}_{(m)} = \mathbf{x}(\boldsymbol{\beta}_{(m)})$ . The updating term on the right-hand side of (6) is the set of OLS parameter estimates from the **Gauss-Newton regression**, or **GNR**,

$$\mathbf{y} - \mathbf{x}(\boldsymbol{\beta}) = \mathbf{X}(\boldsymbol{\beta})\mathbf{b} + \text{residuals}, \quad (7)$$

where the variables  $\mathbf{r}(\boldsymbol{\beta}) \equiv \mathbf{y} - \mathbf{x}(\boldsymbol{\beta})$  and  $\mathbf{R}(\boldsymbol{\beta}) \equiv \mathbf{X}(\boldsymbol{\beta})$  are evaluated at  $\boldsymbol{\beta}_{(m)}$ . Notice that there is no regressor in (7) corresponding to the parameter  $\sigma^2$ , because the criterion function  $Q(\boldsymbol{\beta})$  does not depend on  $\sigma^2$ . This is one of the features of the GNR that makes it a nonstandard artificial regression.

The GNR is clearly a linearization of the nonlinear regression model (3) around the point  $\boldsymbol{\beta}$ . In the special case in which the original model is linear,  $\mathbf{x}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$ , where  $\mathbf{X}$  is the matrix of independent variables. Since  $\mathbf{X}(\boldsymbol{\beta})$  is equal to  $\mathbf{X}$  for all  $\boldsymbol{\beta}$  in this special case, the GNR will simply be a regression of the vector  $\mathbf{y} - \mathbf{X}\boldsymbol{\beta}$  on the matrix  $\mathbf{X}$ .

An example is provided by the nonlinear regression model

$$y_t = \beta_1 Z_{t1}^{\beta_2} Z_{t2}^{1-\beta_2} + u_t, \quad u_t \sim \text{IID}(0, \sigma^2), \quad (8)$$

where  $Z_{t1}$  and  $Z_{t2}$  are independent variables. The regression function here is nonlinear and has the form of a Cobb-Douglas production function. In many cases, of course, it would be reasonable to assume that the error term is multiplicative, and it would then be possible to take logarithms of both sides and use ordinary least squares. But if we wish to estimate (8) as it stands, we must use nonlinear least squares. The GNR that corresponds to (8) is

$$y_t - \beta_1 Z_{t1}^{\beta_2} Z_{t2}^{1-\beta_2} = b_1 Z_{t1}^{\beta_2} Z_{t2}^{1-\beta_2} + b_2 \beta_1 Z_{t2} \left( \frac{Z_{t1}}{Z_{t2}} \right)^{\beta_2} \log \left( \frac{Z_{t1}}{Z_{t2}} \right) + \text{residual}.$$

The regressand is  $y_t$  minus the regression function, the first regressor is the derivative of the regression function with respect to  $\beta_1$ , and the second regressor is the derivative of the regression function with respect to  $\beta_2$

Now consider the defining conditions of an artificial regression. We have

$$\mathbf{R}^\top(\boldsymbol{\theta}) \mathbf{r}(\boldsymbol{\theta}) = \mathbf{X}^\top(\boldsymbol{\beta})(\mathbf{y} - \mathbf{x}(\boldsymbol{\beta})), \quad (9)$$

which is just minus the gradient of  $Q(\boldsymbol{\beta})$ . Thus condition (i') is satisfied.

Next, consider condition (iii). Let  $\hat{\boldsymbol{\beta}}$  denote a vector of initial estimates, which are assumed to be root- $n$  consistent. The GNR (7) evaluated at these estimates is

$$\mathbf{y} - \hat{\mathbf{x}} = \hat{\mathbf{X}}\mathbf{b} + \text{residuals},$$

where  $\hat{\mathbf{x}} \equiv \mathbf{x}(\hat{\boldsymbol{\beta}})$  and  $\hat{\mathbf{X}} \equiv \mathbf{X}(\hat{\boldsymbol{\beta}})$ . The estimate of  $\mathbf{b}$  from this regression is

$$\hat{\mathbf{b}} = (\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top (\mathbf{y} - \hat{\mathbf{x}}). \quad (10)$$

The **one-step efficient estimator** is then defined to be

$$\hat{\boldsymbol{\beta}} \equiv \hat{\boldsymbol{\beta}} + \hat{\mathbf{b}}. \quad (11)$$

By Taylor expanding the expression  $n^{-1/2} \hat{\mathbf{X}}^\top (\mathbf{y} - \hat{\mathbf{x}})$  around  $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ , where  $\boldsymbol{\beta}_0$  is the true parameter vector, and using standard asymptotic arguments, it can be shown that, to leading order,

$$n^{-1/2} \hat{\mathbf{X}}^\top (\mathbf{y} - \hat{\mathbf{x}}) = n^{-1/2} \mathbf{X}_0^\top \mathbf{u} - n^{-1} \mathbf{X}_0^\top \mathbf{X}_0 n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0),$$

where  $\mathbf{X}_0 \equiv \mathbf{X}(\boldsymbol{\beta}_0)$ . This relation can be solved to yield

$$n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = (n^{-1} \mathbf{X}_0^\top \mathbf{X}_0)^{-1} (n^{-1/2} \mathbf{X}_0^\top \mathbf{u} - n^{-1/2} \hat{\mathbf{X}}^\top (\mathbf{y} - \hat{\mathbf{x}})). \quad (12)$$

Now it is a standard result that, asymptotically,

$$n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = (n^{-1} \mathbf{X}_0^\top \mathbf{X}_0)^{-1} (n^{-1/2} \mathbf{X}_0^\top \mathbf{u}); \quad (13)$$

see, for example, Davidson and MacKinnon (1993, Section 5.4). By (10), the second term on the right-hand side of (12) is asymptotically equivalent to  $-n^{1/2} \hat{\mathbf{b}}$ . Thus (12) implies that

$$n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) = n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) - n^{1/2} \hat{\mathbf{b}}.$$

Rearranging this and using the definition (11), we see that, to leading order asymptotically,

$$n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = n^{1/2} (\hat{\boldsymbol{\beta}} + \hat{\mathbf{b}} - \boldsymbol{\beta}_0) = n^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0).$$



In other words, after both are centered and multiplied by  $n^{1/2}$ , the one-step estimator  $\hat{\beta}$  and the NLS estimator  $\hat{\beta}$  tend to the same random variable asymptotically. This is just another way of writing condition (iii) for model (3).

Finally, consider condition (ii). Since  $\mathbf{X}(\beta)$  plays the role of  $\mathbf{R}(\theta)$ , we see that

$$\frac{1}{n}\mathbf{R}^\top(\theta)\mathbf{R}(\theta) = \frac{1}{n}\mathbf{X}^\top(\beta)\mathbf{X}(\beta). \quad (14)$$

If the right-hand side of (14) is evaluated at any root- $n$  consistent estimator  $\hat{\beta}$ , it must tend to the same probability limit as  $n^{-1}\mathbf{X}_0^\top\mathbf{X}_0$ . It is a standard result, following straightforwardly from (13), that, if  $\hat{\beta}$  denotes the NLS estimator for the model (3), then

$$\lim_{n \rightarrow \infty} \text{Var}(n^{1/2}(\hat{\beta} - \beta_0)) = \sigma_0^2 \text{plim}_{n \rightarrow \infty} (n^{-1}\mathbf{X}_0^\top\mathbf{X}_0)^{-1}, \quad (15)$$

where  $\sigma_0^2$  is the true variance of the error terms; see, for example, Davidson and MacKinnon (1993, Chapter 5). Thus the GNR would satisfy condition (ii) except that there is a factor of  $\sigma_0^2$  missing. However, this factor is automatically supplied by the regression package. The estimated covariance matrix will be

$$\widehat{\text{Var}}(\hat{\mathbf{b}}) = \hat{s}^2(\hat{\mathbf{X}}^\top\hat{\mathbf{X}})^{-1}, \quad (16)$$

where  $\hat{s}^2 = \text{SSR}/(n - k)$  is the estimate of  $\sigma^2$  from the artificial regression. It is not hard to show that  $\hat{s}^2$  estimates  $\sigma_0^2$  consistently, and so it is clear from (15) that (16) provides a reasonable way to estimate the covariance matrix of  $\hat{\beta}$ .

It is easy to modify the GNR so that it actually satisfies condition (ii). We just need to divide both the regressand and the regressors by  $s$ , the standard error from the original, nonlinear regression. When this is done, (14) becomes

$$\frac{1}{n}\mathbf{R}^\top(\theta)\mathbf{R}(\theta) = \frac{1}{ns^2}\mathbf{X}^\top(\beta)\mathbf{X}(\beta),$$

and condition (ii) is seen to be satisfied. However, there is rarely any reason to do this in practice.

Although the GNR is the most commonly encountered artificial regression, it differs from most artificial regressions in one key respect: There is one parameter,  $\sigma^2$ , for which there is no regressor. This happens because the criterion function,  $Q(\beta)$ , depends only on  $\beta$ . The GNR therefore has only as many regressors as  $\beta$  has components. This feature of the GNR is responsible for the fact that it does not quite satisfy condition (ii). The fact that  $Q(\beta)$  does not depend on  $\sigma^2$  also causes the asymptotic covariance matrix to be block diagonal between the  $k \times k$  block that corresponds to  $\beta$  and the  $1 \times 1$  block that corresponds to  $\sigma^2$ .

## 4. Uses of the GNR

The GNR, like other artificial regressions, has several uses, depending on the parameter values at which the regressand and regressors are evaluated. If we evaluate them at  $\hat{\beta}$ , the vector of NLS parameter estimates, regression (7) becomes

$$\mathbf{y} - \hat{\mathbf{x}} = \hat{\mathbf{X}}\mathbf{b} + \text{residuals}, \quad (17)$$

where  $\hat{\mathbf{x}} \equiv \mathbf{x}(\hat{\beta})$  and  $\hat{\mathbf{X}} \equiv \mathbf{X}(\hat{\beta})$ . By condition (i), which follows from the first-order conditions for NLS estimation, the OLS estimate  $\hat{\mathbf{b}}$  from this regression is a zero vector. In consequence, the explained sum of squares, or ESS, from regression (17) will be 0, and the SSR will be equal to

$$\|\mathbf{y} - \hat{\mathbf{x}}\|^2 = (\mathbf{y} - \hat{\mathbf{x}})^\top (\mathbf{y} - \hat{\mathbf{x}}),$$

which is the SSR from the original nonlinear regression.

Although it may seem curious to run an artificial regression all the coefficients of which are known in advance to be zero, there can be two very good reasons for doing so. The first reason is to check that the vector  $\hat{\beta}$  reported by a program for NLS estimation really does satisfy the first-order conditions. Computer programs for calculating NLS estimates do not yield reliable answers in every case; see McCullough (1999). The GNR provides an easy way to see whether the first-order conditions are actually satisfied. If all the  $t$  statistics for the GNR are not less than about  $10^{-4}$ , and the  $R^2$  is not less than about  $10^{-8}$ , then the value of  $\hat{\beta}$  reported by the program should be regarded with suspicion.

The second reason to run the GNR (17) is to calculate an estimate of  $\text{Var}(\hat{\beta})$ , the covariance matrix of the NLS estimates. The usual OLS covariance matrix from regression (17) is

$$\widehat{\text{Var}}(\hat{\mathbf{b}}) = s^2(\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1}, \quad (18)$$

which is similar to (16) except that everything is now evaluated at  $\hat{\beta}$ . Thus running the GNR (17) provides an easy way to calculate what is arguably the best estimate of  $\text{Var}(\hat{\beta})$ . Of course, for (18) to provide an asymptotically valid covariance matrix estimate, it is essential that the error terms in (2) be independent and identically distributed, as we have assumed so far. We will discuss ways to drop this assumption in Section 7.

Since the GNR satisfies the one-step property, it and other artificial regressions can evidently be used to obtain one-step efficient estimates. However, although one-step estimation is of considerable theoretical interest, it is generally of modest practical interest, for two reasons. Firstly, we often do not have a root- $n$  consistent estimator to start from and, secondly, modern computers are so fast that the savings from stopping after just one step are rarely substantial.

What is often of great practical interest is the use of the GNR as part of a numerical minimization algorithm to find the NLS estimates  $\hat{\beta}$  themselves. In

practice, the classical Gauss-Newton updating procedure (6) should generally be replaced by

$$\boldsymbol{\beta}_{(m)} = \boldsymbol{\beta}_{(m-1)} + \alpha_{(m)} \mathbf{b}_{(m)},$$

where  $\alpha_{(m)}$  is a scalar that is chosen in various ways by different algorithms, but always in such a way that  $Q(\boldsymbol{\beta}_{(m+1)}) < Q(\boldsymbol{\beta}_{(m)})$ . Numerical optimization methods are discussed by Press et al. (1992), among many others. Artificial regressions other than the GNR allow these methods to be used more widely than just in the least squares context.

## 5. Hypothesis Testing with Artificial Regressions

Artificial regressions like the GNR are probably employed most frequently for hypothesis testing. Suppose we wish to test a set of  $r$  equality restrictions on  $\boldsymbol{\theta}$ . Without loss of generality, we can assume that these are zero restrictions. This allows us to partition  $\boldsymbol{\theta}$  into two subvectors,  $\boldsymbol{\theta}_1$  of length  $k-r$ , and  $\boldsymbol{\theta}_2$  of length  $r$ , the restrictions being that  $\boldsymbol{\theta}_2 = \mathbf{0}$ . If the estimator  $\hat{\boldsymbol{\theta}}$  is not only root- $n$  consistent but also asymptotically normal, an appropriate statistic for testing these restrictions is

$$\hat{\boldsymbol{\theta}}_2^\top (\widehat{\text{Var}}(\hat{\boldsymbol{\theta}}_2))^{-1} \hat{\boldsymbol{\theta}}_2, \quad (19)$$

which will be asymptotically distributed as  $\chi^2(r)$  under the null if  $\widehat{\text{Var}}(\hat{\boldsymbol{\theta}}_2)$  is a suitable estimate of the covariance matrix of  $\hat{\boldsymbol{\theta}}_2$ .

Suppose that  $\mathbf{r}(\boldsymbol{\theta})$  and  $\mathbf{R}(\boldsymbol{\theta})$  define an artificial regression for the estimator  $\hat{\boldsymbol{\theta}}$ . Let  $\hat{\boldsymbol{\theta}} \equiv [\hat{\boldsymbol{\theta}}_1 \ ; \ \mathbf{0}]$  be a vector of root- $n$  consistent estimates under the null. Then, if the variables of the artificial regression are evaluated at  $\hat{\boldsymbol{\theta}}$ , the regression can be expressed as

$$\mathbf{r}(\hat{\boldsymbol{\theta}}_1, \mathbf{0}) = \mathbf{R}_1(\hat{\boldsymbol{\theta}}_1, \mathbf{0})\mathbf{b}_1 + \mathbf{R}_2(\hat{\boldsymbol{\theta}}_2, \mathbf{0})\mathbf{b}_2 + \text{residuals}, \quad (20)$$

where the partitioning of  $\mathbf{R} = [\mathbf{R}_1 \ \mathbf{R}_2]$  corresponds to the partitioning of  $\boldsymbol{\theta}$  as  $[\boldsymbol{\theta}_1 \ ; \ \boldsymbol{\theta}_2]$ . Regression (20) will usually be written simply as

$$\hat{\mathbf{r}} = \hat{\mathbf{R}}_1 \mathbf{b}_1 + \hat{\mathbf{R}}_2 \mathbf{b}_2 + \text{residuals},$$

although this notation hides the fact that  $\hat{\boldsymbol{\theta}}$  satisfies the null hypothesis.

By the one-step property,  $\hat{\mathbf{b}}_2$  from (20) is asymptotically equivalent under the null to the estimator  $\hat{\boldsymbol{\theta}}_2$ , since under the null the true value of  $\boldsymbol{\theta}_2$  is zero. This suggests that we may replace  $\hat{\boldsymbol{\theta}}_2$  in (19) by  $\hat{\mathbf{b}}_2$ . By property (ii), the asymptotic covariance matrix of  $n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$  is estimated by  $(n^{-1} \hat{\mathbf{R}}^\top \hat{\mathbf{R}})^{-1}$ . A suitable estimate of the covariance matrix of  $\hat{\boldsymbol{\theta}}_2$  can be obtained from this by use of the Frisch-Waugh-Lovell (FWL) theorem: See Davidson and MacKinnon (1993, Chapter 1) for a full treatment of the FWL Theorem. The estimate is  $(\hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2)^{-1}$ , where the orthogonal projection matrix  $\hat{\mathbf{M}}_1$  is defined by

$$\hat{\mathbf{M}}_1 = \mathbf{I} - \hat{\mathbf{R}}_1 (\hat{\mathbf{R}}_1^\top \hat{\mathbf{R}}_1)^{-1} \hat{\mathbf{R}}_1^\top. \quad (21)$$

By the same theorem, we have that

$$\hat{\mathbf{b}}_2 = (\hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2)^{-1} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}}. \quad (22)$$

Thus the artificial regression version of the test statistic (19) is

$$\hat{\mathbf{b}}_2^\top \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2 \hat{\mathbf{b}}_2 = \hat{\mathbf{r}}^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2 (\hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2)^{-1} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}}. \quad (23)$$

The following theorem demonstrates the asymptotic validity of (23).

**Theorem 1:** If the regressand  $\mathbf{r}(\boldsymbol{\theta})$  and the regressor matrix  $\mathbf{R}(\boldsymbol{\theta})$  define an artificial regression for the root- $n$  consistent, asymptotically normal, estimator  $\hat{\boldsymbol{\theta}}$ , and if the partition  $\mathbf{R} = [\mathbf{R}_1 \ \mathbf{R}_2]$  corresponds to the partition  $\boldsymbol{\theta} = [\boldsymbol{\theta}_1 \ \vdots \ \boldsymbol{\theta}_2]$ , then the statistic (23), computed at any root- $n$  consistent  $\hat{\boldsymbol{\theta}} = [\hat{\boldsymbol{\theta}}_1 \ \vdots \ \mathbf{0}]$ , is asymptotically distributed as  $\chi^2(r)$  under the null hypothesis that  $\boldsymbol{\theta}_2 = \mathbf{0}$ , and is asymptotically equivalent to the generic statistic (19).

**Proof:** To prove this theorem, we need to show two things. The first is that

$$n^{-1} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2 = n^{-1} \mathbf{R}_2(\boldsymbol{\theta}_0) \mathbf{M}_1(\boldsymbol{\theta}_0) \mathbf{R}_2(\boldsymbol{\theta}_0) + o_p(1),$$

where  $\boldsymbol{\theta}_0$  is the true parameter vector, and  $\mathbf{M}_1(\boldsymbol{\theta}_0)$  is defined analogously to (21). This result follows by standard asymptotic arguments based on the one-step property. The second is that the vector

$$n^{-1/2} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}} = n^{-1/2} \mathbf{R}_2^\top(\boldsymbol{\theta}_0) \mathbf{M}_1(\boldsymbol{\theta}_0) \mathbf{r}(\boldsymbol{\theta}_0) + o_p(1)$$

is asymptotically normally distributed. The equality here also follows by standard asymptotic arguments. The asymptotic normality of  $\hat{\boldsymbol{\theta}}$  implies that  $\hat{\mathbf{b}}$  is asymptotically normally distributed. Therefore, by (22),  $n^{-1/2} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}}$  must also be asymptotically normally distributed. These two results imply that, asymptotically under the null hypothesis, the test statistic (23) is a quadratic form in a normally distributed  $r$ -vector, the mean of which is zero, and the inverse of its covariance matrix. Such a quadratic form follows the  $\chi^2(r)$  distribution. ■

**Remarks:** The statistic (23) can be computed as the difference between the sums of squared residuals (SSR) from the regressions

$$\hat{\mathbf{r}} = \hat{\mathbf{R}}_1 \mathbf{b}_1 + \text{residuals, and} \quad (24)$$

$$\hat{\mathbf{r}} = \hat{\mathbf{R}}_1 \mathbf{b}_1 + \hat{\mathbf{R}}_2 \mathbf{b}_2 + \text{residuals.} \quad (25)$$

Equivalently, it can be computed as the difference between the explained sums of squares (ESS), with the opposite sign, or as the ESS from the FWL regression corresponding to (25):

$$\hat{\mathbf{M}}_1 \hat{\mathbf{r}} = \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2 \mathbf{b}_2 + \text{residuals.}$$

If  $\text{plim } n^{-1} \hat{\mathbf{r}}^\top \hat{\mathbf{r}} = 1$  for all root- $n$  consistent  $\hat{\boldsymbol{\theta}}$ , there are other convenient ways of computing (23), or statistics asymptotically equivalent to it. One is the ordinary  $F$  statistic for  $\mathbf{b}_2 = \mathbf{0}$  in regression (25):

$$F = \frac{\hat{\mathbf{r}}^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2 (\hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{R}}_2)^{-1} \hat{\mathbf{R}}_2^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}} / r}{\hat{\mathbf{r}}^\top \hat{\mathbf{M}}_1 \hat{\mathbf{r}} / (n - k)}, \quad (26)$$

which works because the denominator tends to a probability limit of 1 as  $n \rightarrow \infty$ . This statistic is, of course, in  $F$  rather than  $\chi^2$  form.

Another frequently used test statistic is available if  $\hat{\boldsymbol{\theta}}$  is actually the vector of restricted estimates, that is, the estimator that minimizes the criterion function when the restriction that  $\boldsymbol{\theta}_2 = \mathbf{0}$  is imposed. In this case,  $n$  times the uncentered  $R^2$  from (25) is a valid test statistic. With this choice of  $\hat{\boldsymbol{\theta}}$ , the ESS from (24) is zero, by property (i). Thus (23) is just the ESS from (25). Since  $nR^2 = \text{ESS}/(\text{TSS}/n)$ , where TSS denotes the total sum of squares, and since  $\text{TSS}/n \rightarrow 1$  as  $n \rightarrow \infty$ , it follows that this statistic is asymptotically equivalent to (23).

Even though the GNR does not satisfy condition (ii) when it is expressed in its usual form with all variables not divided by the standard error  $s$ , the  $F$  statistic (26) and the  $nR^2$  statistic are still valid test statistics, because they are both ratios. In fact, variants of the GNR are routinely used to perform many types of specification tests. These include tests for serial correlation similar to the ones proposed by Godfrey (1978), nonnested hypothesis tests where both models are parametric (Davidson and MacKinnon, 1981), and nonnested hypothesis tests where the alternative model is nonparametric (Delgado and Stengos, 1994). They also include several Durbin-Wu-Hausman, or DWH, tests, in which an efficient estimator is compared with an inefficient estimator that is consistent under weaker conditions; see Sections 7.9 and 11.4 of Davidson and MacKinnon (1993).

## 6. The OPG Regression

By no means all interesting econometric models are regression models. It is therefore useful to see if artificial regressions other than the GNR exist for wide classes of models. One of these is the **outer-product-of-the-gradient regression**, or **OPG regression**, a particularly simple artificial regression that can be used with most models that are estimated by maximum likelihood. Suppose we are interested in a model of which the loglikelihood function can be written as

$$\ell(\boldsymbol{\theta}) = \sum_{t=1}^n \ell_t(\boldsymbol{\theta}), \quad (27)$$

where  $\ell_t(\cdot)$  denotes the contribution to the loglikelihood function associated with observation  $t$ . This is the log of the density of the dependent variable(s) for

observation  $t$ , conditional on observations  $1, \dots, t-1$ . Thus lags of the dependent variable(s) are allowed. The key feature of (27) is that  $\ell(\boldsymbol{\theta})$  is a sum of contributions from each of the  $n$  observations.

Now let  $\mathbf{G}(\boldsymbol{\theta})$  be the matrix with typical element

$$G_{ti}(\boldsymbol{\theta}) \equiv \frac{\partial \ell_t(\boldsymbol{\theta})}{\partial \theta_i}; \quad t = 1, \dots, n, \quad i = 1, \dots, k.$$

The matrix  $\mathbf{G}(\boldsymbol{\theta})$  is called the **matrix of contributions to the gradient**, or the **CG matrix**, because the derivative of the sample loglikelihood (27) with respect to  $\theta_i$ , the  $i^{\text{th}}$  component of  $\boldsymbol{\theta}$ , is the sum of the elements of column  $i$  of  $\mathbf{G}(\boldsymbol{\theta})$ . The OPG regression associated with (27) can be written as

$$\boldsymbol{\nu} = \mathbf{G}(\boldsymbol{\theta})\mathbf{b} + \text{residuals}, \quad (28)$$

where  $\boldsymbol{\nu}$  denotes an  $n$ -vector of 1s.

It is easy to see that the OPG regression (28) satisfies the conditions for it to be an artificial regression. Condition (i') is evidently satisfied, since  $\mathbf{R}^\top(\boldsymbol{\theta})\mathbf{r}(\boldsymbol{\theta}) = \mathbf{G}^\top(\boldsymbol{\theta})\boldsymbol{\nu}$ , the components of which are the derivatives of  $\ell(\boldsymbol{\theta})$  with respect to each of the  $\theta_i$ . Condition (ii) is also satisfied, because, under standard regularity conditions, if  $\boldsymbol{\theta}$  is the true parameter vector,

$$\text{plim}_{n \rightarrow \infty} (n^{-1} \mathbf{R}^\top(\boldsymbol{\theta})\mathbf{R}(\boldsymbol{\theta})) = \text{plim}_{n \rightarrow \infty} (n^{-1} \mathbf{G}^\top(\boldsymbol{\theta})\mathbf{G}(\boldsymbol{\theta})) = \mathcal{J}(\boldsymbol{\theta}).$$

Here  $\mathcal{J}(\boldsymbol{\theta})$  denotes the information matrix, defined as

$$\mathcal{J}(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n E(\mathbf{G}_t^\top(\boldsymbol{\theta})\mathbf{G}_t(\boldsymbol{\theta})),$$

where  $\mathbf{G}_t(\cdot)$  is the  $t^{\text{th}}$  row of  $\mathbf{G}(\cdot)$ . Since, as is well-known, the asymptotic covariance matrix of  $n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$  is given by the inverse of the information matrix, condition (ii) is satisfied under the further weak regularity condition that  $\mathcal{J}(\boldsymbol{\theta})$  should be continuous in  $\boldsymbol{\theta}$ . Condition (iii) is also satisfied, since it can be shown that one-step estimates from the OPG regression are asymptotically equivalent to maximum likelihood estimates. The proof is quite similar to the one for the GNR given in Section 3.

It is particularly easy to compute an LM test by using the OPG regression. Let  $\tilde{\boldsymbol{\theta}}$  denote the constrained ML estimates obtained by imposing  $r$  restrictions when maximizing the loglikelihood. Then the ESS from the OPG regression

$$\boldsymbol{\nu} = \mathbf{G}(\tilde{\boldsymbol{\theta}})\mathbf{b} + \text{residuals}, \quad (29)$$

which is equal to  $n$  times the uncentered  $R^2$ , is the OPG form of the LM statistic.

Like the GNR, the OPG regression can be used for many purposes. The use of what is essentially the OPG regression for obtaining maximum likelihood estimates and computing covariance matrices was advocated by Berndt, Hall, Hall, and Hausman (1974). Using it to compute Lagrange Multiplier, or LM, tests was suggested by Godfrey and Wickens (1981), and using it to compute information matrix tests was proposed by Chesher (1983) and Lancaster (1984). The OPG regression is appealing for all these uses because it applies to a very wide variety of models and requires only first derivatives. In general, however, both estimated covariance matrices and test statistics based on the OPG regression are not very reliable in finite samples. In particular, a large number of papers, including Chesher and Spady (1991), Davidson and MacKinnon (1985a, 1992), and Godfrey, McAleer, and McKenzie (1988), have shown that, in finite samples, LM tests based on the OPG regression tend to overreject, often very severely.

Despite this drawback, the OPG regression provides a particularly convenient way to obtain various theoretical results. For example, suppose that we are interested in the variance of  $\hat{\theta}_2$ , the last element of  $\hat{\boldsymbol{\theta}}$ . If  $\boldsymbol{\theta}_1$  denotes a vector of the remaining  $k - 1$  elements, and  $\mathbf{G}(\boldsymbol{\theta})$  and  $\mathbf{b}$  are partitioned in the same way as  $\boldsymbol{\theta}$ , the OPG regression becomes

$$\boldsymbol{\iota} = \mathbf{G}_1(\boldsymbol{\theta})\mathbf{b}_1 + \mathbf{G}_2(\boldsymbol{\theta})\mathbf{b}_2 + \text{residuals},$$

and the FWL regression derived from this by retaining only the last regressor is

$$\mathbf{M}_1\boldsymbol{\iota} = \mathbf{M}_1\mathbf{G}_2\mathbf{b}_2 + \text{residuals},$$

where  $\mathbf{M}_1 \equiv \mathbf{I} - \mathbf{G}_1(\mathbf{G}_1^\top\mathbf{G}_1)^{-1}\mathbf{G}_1^\top$ , and the dependence on  $\boldsymbol{\theta}$  has been suppressed for notational convenience. The covariance matrix estimate from this is just

$$(\mathbf{G}_2^\top\mathbf{M}_1\mathbf{G}_2)^{-1} = (\mathbf{G}_2^\top\mathbf{G}_2 - \mathbf{G}_2^\top\mathbf{G}_1(\mathbf{G}_1^\top\mathbf{G}_1)^{-1}\mathbf{G}_1^\top\mathbf{G}_2)^{-1}. \quad (30)$$

If we divide each of the components of (30) by  $n$  and take their probability limits, we find that

$$\lim_{n \rightarrow \infty} \text{Var}(n^{1/2}(\hat{\theta}_2 - \theta_{20})) = (\mathcal{J}_{22} - \mathcal{J}_{21}\mathcal{J}_{11}^{-1}\mathcal{J}_{12})^{-1},$$

where  $\theta_{20}$  is the true value of  $\theta_2$ . This is a very well-known result, but, since its relation to the FWL theorem is not obvious without appeal to the OPG regression, it is not usually obtained in such a convenient or illuminating way.

## 7. An Artificial Regression for GMM Estimation

Another useful artificial regression, much less well known than the OPG regression, is available for a class of models estimated by the generalized method of moments (GMM). Many such models can be formulated in terms of functions  $f_t(\boldsymbol{\theta})$  of the model parameters and the data, such that, when they are evaluated at the true  $\boldsymbol{\theta}$ , their expectations conditional on corresponding information sets,  $\Omega_t$ , vanish. The  $\Omega_t$  usually contain all information available prior to the time of observation  $t$ , and so, as with the GNR and the OPG regression, lags of dependent variables are allowed.

Let the  $n \times l$  matrix  $\mathbf{W}$  denote the instruments used to obtain the GMM estimates. The  $t^{\text{th}}$  row of  $\mathbf{W}$ , denoted  $\mathbf{W}_t$ , must contain variables in  $\Omega_t$  only. The dimension of  $\boldsymbol{\theta}$  is  $k$ , as before, and, for  $\boldsymbol{\theta}$  to be identified, we need  $l \geq k$ . The GMM estimates with  $l \times l$  weighting matrix  $\mathbf{A}$  are obtained by minimizing the criterion function

$$Q(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{f}^\top(\boldsymbol{\theta}) \mathbf{W} \mathbf{A} \mathbf{W}^\top \mathbf{f}(\boldsymbol{\theta}) \quad (31)$$

with respect to  $\boldsymbol{\theta}$ . Here  $\mathbf{f}(\boldsymbol{\theta})$  is the  $n$ -vector with typical element  $f_t(\boldsymbol{\theta})$ . For the procedure known as efficient GMM, the weighting matrix  $\mathbf{A}$  is chosen so as to be proportional, asymptotically at least, to the inverse of the covariance matrix of  $\mathbf{W}^\top \mathbf{f}(\boldsymbol{\theta})$ . In the simplest case, the  $f_t(\boldsymbol{\theta})$  are serially uncorrelated and homoskedastic with variance 1, and so an appropriate choice is  $\mathbf{A} = (\mathbf{W}^\top \mathbf{W})^{-1}$ . With this choice, the criterion function (31) becomes

$$Q(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{f}^\top(\boldsymbol{\theta}) \mathbf{P}_\mathbf{W} \mathbf{f}(\boldsymbol{\theta}), \quad (32)$$

where  $\mathbf{P}_\mathbf{W}$  is the orthogonal projection on to the columns of  $\mathbf{W}$ .

Let  $\mathbf{J}(\boldsymbol{\theta})$  be the negative of the  $n \times k$  Jacobian matrix of  $\mathbf{f}(\boldsymbol{\theta})$ , so that the  $ti^{\text{th}}$  element of  $\mathbf{J}(\boldsymbol{\theta})$  is  $-\partial f_t / \partial \theta_i(\boldsymbol{\theta})$ . The first-order conditions for minimizing (32) are

$$\mathbf{J}^\top(\boldsymbol{\theta}) \mathbf{P}_\mathbf{W} \mathbf{f}(\boldsymbol{\theta}) = \mathbf{0}. \quad (33)$$

By standard arguments, it can be seen that the vector  $\hat{\boldsymbol{\theta}}$  that solves (33) is asymptotically normal and asymptotically satisfies the equation

$$n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = (n^{-1} \mathbf{J}_0^\top \mathbf{P}_\mathbf{W} \mathbf{J}_0)^{-1} n^{-1/2} \mathbf{J}_0^\top \mathbf{P}_\mathbf{W} \mathbf{f}_0, \quad (34)$$

with  $\mathbf{J}_0 = \mathbf{J}(\boldsymbol{\theta}_0)$  and  $\mathbf{f}_0 = \mathbf{f}(\boldsymbol{\theta}_0)$ . See Davidson and MacKinnon (1993, Chapter 17), for a full discussion of GMM estimation.

Now consider the artificial regression

$$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{P}_\mathbf{W} \mathbf{J}(\boldsymbol{\theta}) \mathbf{b} + \text{residuals}. \quad (35)$$

By the first-order conditions (33) for  $\boldsymbol{\theta}$ , this equation clearly satisfies condition (i), and in fact it also satisfies condition (i') for the criterion function  $Q(\boldsymbol{\theta})$



of (32). Since the covariance matrix of  $\mathbf{f}(\boldsymbol{\theta}_0)$  is just the identity matrix, it follows from (34) that condition (ii) is also satisfied. Arguments just like those presented in Section 3 for the GNR can be used to show that condition (iii), the one-step property, is also satisfied by (35).

If the  $f_t(\boldsymbol{\theta}_0)$  are homoskedastic but with unknown variance  $\sigma^2$ , regression (35) can be used in exactly the same way as the GNR. Either the regressand and regressors can be divided by a suitable consistent estimate of  $\sigma$ , or else all test statistics can be computed as ratios, in  $F$  or  $nR^2$  form, as appropriate.

An important special case of (35) is provided by the class of regression models, linear or nonlinear, estimated with instrumental variables (IV). Such a model can be written in the form (3), but it will be estimated by minimizing, not the criterion function (4) related to the sum of squared residuals, but rather

$$Q(\boldsymbol{\beta}) \equiv \frac{1}{2}(\mathbf{y} - \mathbf{x}(\boldsymbol{\beta}))^\top \mathbf{P}_W (\mathbf{y} - \mathbf{x}(\boldsymbol{\beta})),$$

where  $\mathbf{W}$  is an  $n \times l$  matrix of instrumental variables. This criterion function has exactly the same form as (32), with  $\boldsymbol{\beta}$  instead of  $\boldsymbol{\theta}$ , and with  $\mathbf{f}(\boldsymbol{\beta}) = \mathbf{y} - \mathbf{x}(\boldsymbol{\beta})$ . In addition,  $\mathbf{J}(\boldsymbol{\beta}) = \mathbf{X}(\boldsymbol{\beta})$ , where  $\mathbf{X}(\boldsymbol{\beta})$  is defined, exactly as for the GNR, to have  $t^{\text{th}}$  element  $\partial x_t / \partial \beta_i(\boldsymbol{\beta})$ . The resulting artificial regression for the IV model, which takes the form

$$\mathbf{y} - \mathbf{x}(\boldsymbol{\beta}) = \mathbf{P}_W \mathbf{X}(\boldsymbol{\beta}) \mathbf{b} + \text{residuals}, \quad (36)$$

is often referred to as a GNR, because, except for the projection matrix  $\mathbf{P}_W$ , it is identical to (7): See Davidson and MacKinnon (1993, Chapter 7).

## 8. Artificial Regressions and Heteroskedasticity

Covariance matrices and test statistics calculated via the GNR (7), or via artificial regressions such as (35) and (36), are not asymptotically valid when the assumption that the error terms are IID is violated. Consider a modified version of the nonlinear regression model (3), in which  $E(\mathbf{u}\mathbf{u}^\top) = \boldsymbol{\Omega}$ , where  $\boldsymbol{\Omega}$  is an  $n \times n$  diagonal matrix with  $t^{\text{th}}$  diagonal element  $\omega_t^2$ . Let  $\hat{\boldsymbol{\Omega}}$  denote an  $n \times n$  diagonal matrix with the squared residual  $\hat{u}_t^2$  as the  $t^{\text{th}}$  diagonal element. It has been known since the work of White (1980) that the matrix

$$(\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \hat{\boldsymbol{\Omega}} \hat{\mathbf{X}} (\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \quad (37)$$

provides an estimator of  $\text{Var}(\hat{\boldsymbol{\beta}})$ , which can be used in place of the usual estimator,  $s^2(\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1}$ . Like the latter, this **heteroskedasticity-consistent covariance matrix estimator**, or **HCCME**, can be computed by means of an artificial regression. We will refer to this regression as the **heteroskedasticity-robust Gauss-Newton Regression**, or **HRGNR**.

In order to derive the HRGNR, it is convenient to begin with a linear regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ , and to consider the criterion function

$$Q(\boldsymbol{\beta}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

The negative of the gradient of this function with respect to  $\boldsymbol{\beta}$  is

$$\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \quad (38)$$

and its Hessian is the matrix

$$\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}, \quad (39)$$

of which the inverse is the HCCME if we replace  $\boldsymbol{\Omega}$  by  $\hat{\boldsymbol{\Omega}}$ . Equating the gradient to zero just yields the OLS estimator, since  $\mathbf{X}^\top \mathbf{X}$  and  $\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}$  are  $k \times k$  nonsingular matrices.

Let  $\mathbf{V}$  be an  $n \times n$  diagonal matrix with  $t^{\text{th}}$  diagonal element equal to  $\omega_t$ ; thus  $\mathbf{V}^2 = \boldsymbol{\Omega}$ . Consider the  $n \times k$  regressor matrix  $\mathbf{R}$  defined by

$$\mathbf{R} = \mathbf{V} \mathbf{X}(\mathbf{X}^\top \mathbf{V}^2 \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{P}_{\mathbf{V}\mathbf{X}} \mathbf{V}^{-1} \mathbf{X}, \quad (40)$$

where  $\mathbf{P}_{\mathbf{V}\mathbf{X}}$  projects orthogonally on to the columns of  $\mathbf{V}\mathbf{X}$ . We have

$$\mathbf{R}^\top \mathbf{R} = \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}, \quad (41)$$

which is just the Hessian (39). Let  $\mathbf{U}(\boldsymbol{\beta})$  be a diagonal matrix with  $t^{\text{th}}$  diagonal element equal to  $y_t - \mathbf{X}_t \boldsymbol{\beta}$ . Then, if we define  $\mathbf{R}(\boldsymbol{\beta})$  as in (40) but with  $\mathbf{V}$  replaced by  $\mathbf{U}(\boldsymbol{\beta})$ , we find that  $\hat{\mathbf{R}}^\top \hat{\mathbf{R}}$  is the HCCME (37).

In order to derive the regressand  $\mathbf{r}(\boldsymbol{\beta})$ , note that, for condition (i') to be satisfied, we require

$$\mathbf{R}^\top(\boldsymbol{\beta}) \mathbf{r}(\boldsymbol{\beta}) = \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{U}^2(\boldsymbol{\beta}) \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta});$$

recall (38). Since the  $t^{\text{th}}$  element of  $\mathbf{U}(\boldsymbol{\beta})$  is  $y_t - \mathbf{X}_t \boldsymbol{\beta}$ , this implies that

$$\mathbf{r}(\boldsymbol{\beta}) = \mathbf{U}^{-1}(\boldsymbol{\beta})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \boldsymbol{\nu}.$$

In the general nonlinear case,  $\mathbf{X}$  becomes  $\mathbf{X}(\boldsymbol{\beta})$ , and the HRGNR has the form

$$\boldsymbol{\nu} = \mathbf{P}_{\mathbf{U}(\boldsymbol{\beta})\mathbf{X}(\boldsymbol{\beta})} \mathbf{U}^{-1}(\boldsymbol{\beta}) \mathbf{X}(\boldsymbol{\beta}) \mathbf{b} + \text{residuals}, \quad (42)$$

where now the  $t^{\text{th}}$  diagonal element of  $\mathbf{U}(\boldsymbol{\beta})$  is  $y_t - x_t(\boldsymbol{\beta})$ . When  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$ , the vector of NLS estimates,

$$\begin{aligned} \hat{\mathbf{r}}^\top \hat{\mathbf{R}} &= \boldsymbol{\nu}^\top \mathbf{P}_{\hat{\mathbf{U}}\hat{\mathbf{X}}} \hat{\mathbf{U}}^{-1} \hat{\mathbf{X}} \\ &= \boldsymbol{\nu}^\top \hat{\mathbf{U}} \hat{\mathbf{X}} (\hat{\mathbf{X}}^\top \hat{\mathbf{U}} \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \hat{\mathbf{U}} \hat{\mathbf{U}}^{-1} \hat{\mathbf{X}} \\ &= \hat{\mathbf{u}}^\top \hat{\mathbf{X}} (\hat{\mathbf{X}}^\top \hat{\boldsymbol{\Omega}} \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \hat{\mathbf{X}} = \mathbf{0}, \end{aligned} \quad (43)$$

because the NLS first-order conditions give  $\hat{\mathbf{X}}^\top \hat{\mathbf{u}} = \mathbf{0}$ . Thus condition (i) is satisfied for the nonlinear case. Condition (ii) is satisfied by construction, as can be seen by putting hats on everything in (41).

For condition (iii) to hold, regression (42) must satisfy the one-step property. We will only show that this property holds for linear models. Extending the argument to nonlinear models would be tedious but not difficult. In the linear case, evaluating (42) at an arbitrary  $\hat{\boldsymbol{\beta}}$  gives

$$\hat{\mathbf{b}} = (\mathbf{X}^\top \hat{\mathbf{U}}^{-1} \mathbf{P}_{\hat{\mathbf{U}}\mathbf{X}} \hat{\mathbf{U}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{U}}^{-1} \mathbf{P}_{\hat{\mathbf{U}}\mathbf{X}} \boldsymbol{\nu}.$$

With a little algebra, it can be shown that this reduces to

$$\hat{\mathbf{b}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{u}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}, \quad (44)$$

where  $\hat{\boldsymbol{\beta}}$  is the OLS estimator. It follows that the one-step estimator  $\hat{\boldsymbol{\beta}} + \hat{\mathbf{b}}$  is equal to  $\hat{\boldsymbol{\beta}}$ , as we wished to show. In the nonlinear case, of course, we obtain an asymptotic equality rather than an exact equality.

As with the ordinary GNR, the HRGNR is particularly useful for hypothesis testing. If we partition  $\boldsymbol{\beta}$  as  $[\boldsymbol{\beta}_1 \ ; \ \boldsymbol{\beta}_2]$  and wish to test the  $r$  zero restrictions  $\boldsymbol{\beta}_2 = \mathbf{0}$ , we need to run two versions of the regression and compute the difference between the two SSRs or SSEs. The two regressions are:

$$\boldsymbol{\nu} = \mathbf{P}_{\tilde{\mathbf{U}}\tilde{\mathbf{X}}} \tilde{\mathbf{U}}^{-1} \tilde{\mathbf{X}}_1 \mathbf{b}_1 + \text{residuals, and} \quad (45)$$

$$\boldsymbol{\nu} = \mathbf{P}_{\tilde{\mathbf{U}}\tilde{\mathbf{X}}} \tilde{\mathbf{U}}^{-1} \tilde{\mathbf{X}}_1 \mathbf{b}_1 + \mathbf{P}_{\tilde{\mathbf{U}}\tilde{\mathbf{X}}} \tilde{\mathbf{U}}^{-1} \tilde{\mathbf{X}}_2 \mathbf{b}_2 + \text{residuals.} \quad (46)$$

It is important to note that the first regression is *not* the HRGNR for the restricted model, because it uses the matrix  $\mathbf{P}_{\tilde{\mathbf{U}}\tilde{\mathbf{X}}}$  rather than the matrix  $\mathbf{P}_{\tilde{\mathbf{U}}\tilde{\mathbf{X}}_1}$ . In consequence, the regressand in (45) will not be orthogonal to the regressors. This is why we need to run two artificial regressions. We could compute an ordinary  $F$  statistic instead of the difference between the SSRs from (45) and (46), but there would be no advantage to doing so, since the  $F$  form of the test merely divides by a stochastic quantity that tends to 1 asymptotically.

The HRGNR appears to be new. The trick of multiplying  $\mathbf{X}(\boldsymbol{\beta})$  by  $\mathbf{U}^{-1}(\boldsymbol{\beta})$  in order to obtain an HCCME by means of an OLS regression was used, in a different context, by Messer and White (1984). This trick does cause a problem in some cases. If any element on the diagonal of the matrix  $\mathbf{U}(\boldsymbol{\beta})$  is equal to 0, the inverse of that element cannot be computed. Therefore, it is necessary to replace any such element by a small, positive number before computing  $\mathbf{U}^{-1}(\boldsymbol{\beta})$ .

A different, and considerably more limited, type of heteroskedasticity-robust GNR, which is applicable only to hypothesis testing, was first proposed by Davidson and MacKinnon (1985b). It was later rediscovered by Wooldridge (1990, 1991) and extended to handle other cases, including regression models with error terms that have autocorrelation as well as heteroskedasticity of unknown form.

It is possible to construct a variety of artificial regressions that provide different covariance matrix estimators for regression models. From (43) and (44), it follows that any artificial regression with regressand

$$\mathbf{r}(\boldsymbol{\beta}) = \mathbf{U}^{-1}(\boldsymbol{\beta})(\mathbf{y} - \mathbf{x}(\boldsymbol{\beta}))$$

and regressors

$$\mathbf{R}(\boldsymbol{\beta}) = \mathbf{P}_{\mathbf{U}(\boldsymbol{\beta})\mathbf{X}(\boldsymbol{\beta})}\mathbf{U}^{-1}(\boldsymbol{\beta})\mathbf{X}(\boldsymbol{\beta})$$

satisfies properties (i) and (iii) for the least-squares estimator, for any nonsingular matrix  $\mathbf{U}(\boldsymbol{\beta})$ . Thus any sandwich covariance matrix estimator can be computed by choosing  $\mathbf{U}(\boldsymbol{\beta})$  appropriately; the estimator (37) is just one example. In fact, it is possible to develop artificial regressions that allow testing not only with a variety of different HCCMEs, but also with some sorts of heteroskedasticity and autocorrelation consistent (HAC) covariance matrix estimators. It is also a simple matter to use such estimators with modified versions of the artificial regression (35) used with models estimated by GMM.

## 9. Double-Length Regressions

Up to this point, the number of observations for all the artificial regressions we have studied has been equal to  $n$ , the number of observations in the data. In some cases, however, artificial regressions may have  $2n$  or even  $3n$  observations. This can happen whenever each observation makes two or more contributions to the criterion function.

The first **double-length artificial regression**, or **DLR**, was proposed by Davidson and MacKinnon (1984a). We will refer to it as *the* DLR, even though it is no longer the only artificial regression with  $2n$  observations. The class of models to which the DLR applies is a subclass of the one used for GMM estimation. Such models may be written as

$$f_t(y_t, \boldsymbol{\theta}) = \varepsilon_t, \quad t = 1, \dots, n, \quad \varepsilon_t \sim \text{NID}(0, 1), \quad (47)$$

where, as before, each  $f_t(\cdot)$  is a smooth function that depends on the data and on a  $k$ -vector of parameters  $\boldsymbol{\theta}$ . Here, however, the  $f_t$  are assumed to be normally distributed conditional on the information sets  $\Omega_t$ , as well as being of mean zero, serially uncorrelated, and homoskedastic with variance 1. Further,  $f_t$  may depend only on a scalar dependent variable  $y_t$ , although lagged dependent variables are allowed as explanatory variables.

The class of models (47) is much less restrictive than it may at first appear to be. In particular, it is not essential that the error terms follow the normal distribution, although it is essential that they follow some specified, continuous distribution, which can be transformed into the standard normal distribution, so as to allow the model to be written in the form of (47). A great many models

that involve transformations of the dependent variable can be put into the form of (47). For example, consider the Box-Cox regression model

$$\tau(y_t, \lambda) = \sum_{i=1}^k \beta_i \tau(X_{ti}, \lambda) + \sum_{j=1}^l \gamma_j Z_{tj} + u_t, \quad u_t \sim N(0, \sigma^2), \quad (48)$$

where  $\tau(x, \lambda) = (x^\lambda - 1)/\lambda$  is the Box-Cox transformation (Box and Cox, 1964),  $y_t$  is the dependent variable, the  $X_{ti}$  are independent variables that are always positive, and the  $Z_{tj}$  are additional independent variables. We can rewrite (48) in the form of (47) by making the definition

$$f_t(y_t, \boldsymbol{\theta}) = \frac{1}{\sigma} \left( \tau(y_t, \lambda) - \sum_{i=1}^k \beta_i \tau(X_{ti}, \lambda) - \sum_{j=1}^l \gamma_j Z_{tj} \right).$$

For the model (47), the contribution of the  $t^{\text{th}}$  observation to the loglikelihood function  $\ell(\mathbf{y}, \boldsymbol{\theta})$  is

$$\ell_t(y_t, \boldsymbol{\theta}) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} f_t^2(y_t, \boldsymbol{\theta}) + k_t(y_t, \boldsymbol{\theta}),$$

where

$$k_t(y_t, \boldsymbol{\theta}) \equiv \log \left| \frac{\partial f_t(y_t, \boldsymbol{\theta})}{\partial y_t} \right|$$

is a Jacobian term. Now let us make the definitions

$$F_{ti}(y_t, \boldsymbol{\theta}) \equiv \frac{\partial f_t(y_t, \boldsymbol{\theta})}{\partial \theta_i} \quad \text{and} \quad K_{ti}(y_t, \boldsymbol{\theta}) \equiv \frac{\partial k_t(y_t, \boldsymbol{\theta})}{\partial \theta_i}$$

and define  $\mathbf{F}(\mathbf{y}, \boldsymbol{\theta})$  and  $\mathbf{K}(\mathbf{y}, \boldsymbol{\theta})$  as the  $n \times k$  matrices with typical elements  $F_{ti}(y_t, \boldsymbol{\theta})$  and  $K_{ti}(y_t, \boldsymbol{\theta})$  and typical rows  $\mathbf{F}_t(\mathbf{y}, \boldsymbol{\theta})$  and  $\mathbf{K}_t(\mathbf{y}, \boldsymbol{\theta})$ . Similarly, let  $\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})$  be the  $n$ -vector with typical element  $f_t(y_t, \boldsymbol{\theta})$ .

The DLR, which has  $2n$  artificial observations, may be written as

$$\begin{bmatrix} \mathbf{f}(\mathbf{y}, \boldsymbol{\theta}) \\ \boldsymbol{\iota} \end{bmatrix} = \begin{bmatrix} -\mathbf{F}(\mathbf{y}, \boldsymbol{\theta}) \\ \mathbf{K}(\mathbf{y}, \boldsymbol{\theta}) \end{bmatrix} \mathbf{b} + \text{residuals}. \quad (49)$$

Since the gradient of  $\ell(\mathbf{y}, \boldsymbol{\theta})$  is

$$\mathbf{g}(\mathbf{y}, \boldsymbol{\theta}) = -\mathbf{F}^\top(\mathbf{y}, \boldsymbol{\theta}) \mathbf{f}(\mathbf{y}, \boldsymbol{\theta}) + \mathbf{K}^\top(\mathbf{y}, \boldsymbol{\theta}) \boldsymbol{\iota}, \quad (50)$$

we see that regression (49) satisfies condition (i'). It can also be shown that it satisfies conditions (ii) and (iii), and thus it has all the properties of an artificial regression.

The DLR can be used for many purposes, including nonnested hypothesis tests of models with different functional forms (Davidson and MacKinnon, 1984a), tests of functional form (MacKinnon and Magee, 1990), and tests of linear and loglinear regressions against Box-Cox alternatives like (48) (Davidson and MacKinnon, 1985a). The latter application has recently been extended to models with AR(1) errors by Baltagi (1999). An accessible discussion of the DLR may be found in Davidson and MacKinnon (1988). When both the OPG regression and the DLR are available, the finite-sample performance of the latter always seems to be very much better than that of the former.

As we remarked earlier, the DLR is not the only artificial regression with  $2n$  artificial observations. In particular, Orme (1995) showed how to construct such a regression for the widely-used tobit model, and Davidson and MacKinnon (1999) provided evidence that Orme's regression generally works very well. It makes sense that a double-length regression should be needed in this case, because the tobit loglikelihood is the sum of two summations, which are quite different in form. One summation involves all the observations for which the dependent variable is equal to zero, and the other involves all the observations for which it takes on a positive value.

## 10. An Artificial Regression for Binary Response Models

For binary response models such as the logit and probit models, there exists a very simple artificial regression that can be derived as an extension of the Gauss-Newton regression. It was independently suggested by Engle (1984) and Davidson and MacKinnon (1984b).

The object of a binary response model is to predict the probability that the binary dependent variable,  $y_t$ , is equal to 1 conditional on some information set  $\Omega_t$ . A useful class of binary response models can be written as

$$E(y_t | \Omega_t) = \Pr(y_t = 1) = F(\mathbf{Z}_t\boldsymbol{\beta}). \quad (51)$$

Here  $\mathbf{Z}_t$  is a row vector of explanatory variables that belong to  $\Omega_t$ ,  $\boldsymbol{\beta}$  is the vector of parameters to be estimated, and  $F(x)$  is the differentiable cumulative distribution function (CDF) of some scalar probability distribution. For the probit model,  $F(x)$  is the standard normal CDF. For the logit model,  $F(x)$  is the logistic function

$$\frac{\exp(x)}{1 + \exp(x)} = (1 + \exp(-x))^{-1}.$$

The loglikelihood function for this class of binary response models is

$$\ell(\boldsymbol{\beta}) = \sum_{t=1}^n \left( (1 - y_t) \log(1 - F(\mathbf{Z}_t\boldsymbol{\beta})) + y_t \log(F(\mathbf{Z}_t\boldsymbol{\beta})) \right), \quad (52)$$

If  $f(x) = F'(x)$  is the density corresponding for the CDF  $F(x)$ , the first-order conditions for maximizing (52) are

$$\sum_{t=1}^n \frac{(y_t - \hat{F}_t) \hat{f}_t Z_{ti}}{\hat{F}_t(1 - \hat{F}_t)} = 0, \quad i = 1, \dots, k, \quad (53)$$

where  $Z_{ti}$  is the  $ti^{\text{th}}$  component of  $\mathbf{Z}_t$ ,  $\hat{f}_t \equiv f(\mathbf{Z}_t \hat{\boldsymbol{\beta}})$  and  $\hat{F}_t \equiv F(\mathbf{Z}_t \hat{\boldsymbol{\beta}})$ .

There is more than one way to derive the artificial regression that corresponds to the model (51). The easiest is to rewrite it in the form of the nonlinear regression model

$$y_t = F(\mathbf{Z}_t \boldsymbol{\beta}) + u_t. \quad (54)$$

The error term  $u_t$  here is evidently non-normal and heteroskedastic. Because  $y_t$  is like a Bernoulli trial with probability  $p$  given by  $F(\mathbf{Z}_t \boldsymbol{\beta})$ , and the variance of a Bernoulli trial is  $p(1 - p)$ , the variance of  $u_t$  is

$$v_t(\boldsymbol{\beta}) \equiv F(\mathbf{Z}_t \boldsymbol{\beta})(1 - F(\mathbf{Z}_t \boldsymbol{\beta})). \quad (55)$$

The ordinary GNR for (54) would be

$$y_t - F(\mathbf{Z}_t \boldsymbol{\beta}) = f(\mathbf{Z}_t \boldsymbol{\beta}) \mathbf{Z}_t \mathbf{b} + \text{residual},$$

but the ordinary GNR is not appropriate because of the heteroskedasticity of the  $u_t$ . Multiplying both sides by the square root of the inverse of (55) yields the artificial regression

$$v_t^{-1/2}(\boldsymbol{\beta})(y_t - F(\mathbf{Z}_t \boldsymbol{\beta})) = v_t^{-1/2}(\boldsymbol{\beta})f(\mathbf{Z}_t \boldsymbol{\beta}) \mathbf{Z}_t \mathbf{b} + \text{residual}. \quad (56)$$

This regression has all the usual properties of artificial regressions. It can be seen from (53) that it satisfies condition (i'). Because a typical element of the information matrix corresponding to (52) is

$$J_{ij}(\boldsymbol{\beta}) = \text{plim}_{n \rightarrow \infty} \left( \frac{1}{n} \sum_{t=1}^n Z_{ti} Z_{tj} \frac{f(\mathbf{Z}_t \boldsymbol{\beta})^2}{F(\mathbf{Z}_t \boldsymbol{\beta})(1 - F(\mathbf{Z}_t \boldsymbol{\beta}))} \right),$$

it is not difficult to show that regression (56) satisfies condition (ii). Finally, since (56) has the structure of a GNR, the arguments used in Section 3 show that it also satisfies condition (iii), the one-step property.

As an artificial regression, (56) can be used for all the things that other artificial regressions can be used for. In particular, when it is evaluated at restricted estimates  $\tilde{\boldsymbol{\beta}}$ , the explained sum of squares is an LM test statistic for testing the restrictions. The normalization of the regressand by its standard error means that other test statistics, such as  $nR^2$  and the ordinary  $F$  statistic for the coefficients on the regressors that correspond to the restricted parameters to be zero, are also asymptotically valid. However, they seem to have slightly poorer finite-sample properties than the ESS (Davidson and MacKinnon, 1984b). It is, of course, possible to extend regression (56) in various ways. For example, it has been extended to tests of the functional form of  $F(x)$  by Thomas (1993) and to tests of ordered logit models by Murphy (1996).

## 11. Conclusion

In this paper, we have introduced the concept of an artificial regression and discussed several examples. We have seen that artificial regressions can be useful for minimizing criterion functions, computing one-step estimates, calculating covariance matrix estimates, and computing test statistics. The last of these is probably the most common application. There is a close connection between the artificial regression for a given model and the asymptotic theory for that model. Therefore, as we saw in Section 6, artificial regressions can also be very useful for obtaining theoretical results.

Most of the artificial regressions we have discussed are quite well-known. This is true of the Gauss-Newton regression discussed in Sections 3 and 4, the OPG regression discussed in Section 6, the double-length regression discussed in Section 9, and the regression for binary response models discussed in Section 10. However, the artificial regression for GMM estimation discussed in Section 7 does not appear to have been treated previously in published work, and we believe that the heteroskedasticity-robust GNR discussed in Section 8 is new.



## References

- Baltagi, B. (1999). "Double length regressions for linear and log-linear regressions with AR(1) disturbances," *Statistical Papers*, **40**, 199–209.
- Berndt, E. R., B. H. Hall, R. E. Hall, and J. A. Hausman (1974). "Estimation and inference in nonlinear structural models," *Annals of Economic and Social Measurement*, **3**, 653–665.
- Box, G. E. P., and D. R. Cox (1964). "An analysis of transformations," *Journal of the Royal Statistical Society, Series B*, **26**, 211–252.
- Chesher, A. (1983). "The information matrix test: simplified calculation via a score test interpretation," *Economics Letters*, **13**, 45–48.
- Chesher, A., and R. Spady (1991). "Asymptotic expansions of the information matrix test statistic," *Econometrica*, **59**, 787–815.
- Davidson, R., and J. G. MacKinnon (1981). "Several tests for model specification in the presence of alternative hypotheses," *Econometrica*, **49**, 781–793.
- Davidson, R., and J. G. MacKinnon (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*, **18**, 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 (1988). "Double-length artificial regressions," *Oxford Bulletin of Economics and Statistics*, **50**, 203–217.
- Davidson, R., and J. G. MacKinnon (1990). "Specification tests based on artificial regressions," *Journal of the American Statistical Association*, **85**, 220–227.
- Davidson, R., and J. G. MacKinnon (1992). "A new form of the information matrix test," *Econometrica*, **60**, 145–157.
- Davidson, R., and J. G. MacKinnon (1993). *Estimation and Inference in Econometrics*, New York: Oxford University Press.
- Davidson, R., and J. G. MacKinnon (1999). "Bootstrap testing in nonlinear models," *International Economic Review*, **40**, 487–508.
- Delgado, M. A., and T. Stengos (1994). "Semiparametric specification testing of non-nested econometric models," *Review of Economic Studies*, **61**, 291–303.
- Engle, R. F. (1984). "Wald, Likelihood Ratio and Lagrange Multiplier Tests in Econometrics," Ch. 13 in *Handbook of Econometrics*, Vol. II. Eds.: Zvi Griliches and Michael D. Intriligator. Amsterdam: North-Holland.

- Godfrey, L. G. (1978). “Testing against general autoregressive and moving average error models when the regressors include lagged dependent variables,” *Econometrica*, **46**, 1293–1301.
- Godfrey, L. G., M. McAleer, and C. R. McKenzie (1988). “Variable addition and Lagrange Multiplier tests for linear and logarithmic regression models,” *Review of Economics and Statistics*, **70**, 492–503.
- Godfrey, L. G., M. McAleer, and C. R. McKenzie (1988). “Variable addition and Lagrange Multiplier tests for linear and logarithmic regression models,” *Review of Economics and Statistics*, **70**, 492–503.
- Godfrey, L. G., and M. R. Wickens (1981). “Testing linear and log-linear regressions for functional form,” *Review of Economic Studies*, **48**, 487–496.
- Lancaster, T. (1984). “The covariance matrix of the information matrix test,” *Econometrica*, **52**, 1051–1053.
- MacKinnon, J. G., and L. Magee (1990). “Transforming the dependent variable in regression models,” *International Economic Review*, **31**, 315–339.
- McCullough, B. D. (1999). “Econometric software reliability: EViews, LIMDEP, SHAZAM, and TSP,” *Journal of Applied Econometrics*, **14**, 191–202.
- Messer, K., and H. White (1984). “A note on computing the heteroskedasticity consistent covariance matrix using instrumental variable techniques,” *Oxford Bulletin of Economics and Statistics*, **46**, 181–184.
- Murphy, A. (1996). “Simple LM tests of mis-specification for ordered logit models,” *Economics Letters*, **52**, 137–141.
- Orme, C. (1995). “On the use of artificial regressions in certain microeconomic models,” *Econometric Theory*, **11**, 290–305.
- Press, W. H., S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, (1992). *Numerical Recipes in C*, Second edition, Cambridge, Cambridge University Press.
- Thomas, J. (1993). “On testing the logistic assumption in binary dependent variable models,” *Empirical Economics*, **18**, 381–392.
- White, H. (1980). “A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity,” *Econometrica*, **48**, 817–38.
- Wooldridge, J. M. (1990). “A unified approach to robust, regression-based specification tests,” *Econometric Theory*, **6**, 17–43.
- Wooldridge, J. M. (1991). “On the application of robust, regression-based diagnostics to models of conditional means and conditional variances,” *Journal of Econometrics*, **47**, 5–46.

