



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Journal of Econometrics 122 (2004) 27–46

JOURNAL OF
Econometrics

www.elsevier.com/locate/econbase

On the harm that ignoring pretesting can cause

Dmitry Danilov^a, Jan R. Magnus^{b,*}

^a*Eurandom, Eindhoven University of Technology, P.O. Box 513,
5600 MB Eindhoven, The Netherlands*

^b*CentER, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, The Netherlands*

Accepted 23 October 2003

Abstract

In econometrics the same data set is typically used to select the model and to estimate the parameters in the selected model. In applied econometrics practice, however, one typically acts as if the model had been given a priori, thus ignoring the fact that the estimators are in fact *pretest* estimators. Hence one assumes incorrectly that the estimator is unbiased, and that the reported variance, conditional on the selected model, is equal to its unconditional variance.

In this paper, we find the *unconditional* first and second moments of the pretest estimator (in fact, of a more general estimator, the WALS estimator), taking full account of the fact that model selection and estimation are an integrated procedure, and show that the error in not reporting the correct moments can be large. We also show that this error can vary substantially between different model selection procedures. Finally, we ask how the error increases when the number of auxiliary regressors increases.

© 2003 Elsevier B.V. All rights reserved.

JEL classification: C13; C51

Keywords: Pretest estimator; Model selection; Mean squared error

1. Introduction

In econometrics, but also in many other disciplines, the same data set is commonly used for model selection and for estimation. Standard statistical theory is therefore not directly applicable, since the properties of most estimators depend not only on the stochastic nature of the selected model, but also on the way the model was selected.

* Corresponding author. Tel.: +31-13-466-3092; fax: +31-13-466-3066.

E-mail addresses: danilov@eurandom.tue.nl (D. Danilov), magnus@uvt.nl (J.R. Magnus).

The simplest example of this situation is the standard linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{z} + \boldsymbol{\varepsilon}$, where we are uncertain whether to include \mathbf{z} or not.¹ The usual procedure is to compute the t -statistic of γ , and then, depending on whether $|t|$ is ‘large’ or ‘small’, decide to use the unrestricted or the restricted model. We then estimate $\boldsymbol{\beta}$ from the selected model. This estimator is a *pretest* estimator, but we commonly report its properties as if estimation had not been preceded by model selection. Thus we report no bias and an incorrect variance.

This is clearly wrong. Our view is *not* that we should avoid pretesting, even though it is well-known that pretest estimators have poor properties, inadmissibility being only one of them. This would be near-impossible in applied work. Our view is simply that we should report estimates of the correct bias and variance (or mean squared error), taking full account of the fact that model selection and estimation are an integrated procedure. This paper attempts to do this.

The literature on pretesting starts with Bancroft’s (1944) famous article. Bancroft is mostly concerned with the bias introduced by pretests of homogeneity of variances and pretests of a regression coefficient. He considers the simplest case, in our notation $\mathbf{y} = \beta\mathbf{x} + \gamma\mathbf{z} + \boldsymbol{\varepsilon}$ (one β , one γ), where he wishes to estimate β while being uncertain about whether \mathbf{z} should be in the regression or not. He then investigates the bias of the pretest estimator of β . Mosteller (1948) considers the special case $\mathbf{x}' = (\mathbf{1}', \mathbf{1}')$, $\mathbf{z}' = (\mathbf{0}', \mathbf{1}')$, where $\mathbf{1}$ denotes the vector of ones. Thus, Mosteller considers pooling: if $\gamma = 0$ we pool, otherwise we do not pool. In this context, he calculates the mean squared error of the pretest estimator. Huntsberger (1955) extends Mosteller’s paper by explicitly writing the pretest estimator as a (continuous) weighted average of the restricted ($\gamma = 0$) and unrestricted estimator, where the weights are functions of the relevant t -statistic. The fact that the pretest estimator has many undesirable properties is highlighted by Sclove et al. (1972). The early literature is discussed in detail in Judge and Bock’s (1978) important monograph.

Lovell (1983) asks what will be the true significance level of a t -test after pretesting, and recommends a simple rule-of-thumb. Roehrig (1984) establishes the relationship between the mean squared error of the pretest estimator and the mean squared error of the estimator of the nuisance parameters, a result later generalized by Magnus and Durbin (1999). Mittelhammer (1984) compares the risk functions of several estimators (including the pretest) under model misspecification, and concludes *inter alia* that all alternatives to OLS can be inferior to OLS in terms of prediction risk. The literature of this period is well summarized in Judge and Bock (1983) and in the special issue of the *Journal of Econometrics* (1984), edited by George Judge.

Asymptotic aspects are considered in Sen (1979), Pötscher (1991), Zhang (1992), and Pötscher and Novak (1998). While most studies, including ours, are confined to the first two moments of the pretest statistics, Giles and Srivastava (1993) and Leeb and Pötscher (2003) derive the distribution of the traditional pretest estimator. Summaries of the latest developments are given in Giles and Giles (1993), Chatfield (1995), and Magnus (1999).

¹ We follow the notation proposed by Abadir and Magnus (2002).

Different model selection strategies are discussed by Hoover and Perez (1999), who favor the general-to-specific procedure. Hendry (2001) advertises computer-automated general-to-specific procedures and claims that these procedures perform well in Monte Carlo experiments.

In spite of all this literature, we are still far removed from having a fully integrated procedure of model selection and parameter estimation. The current paper attempts to narrow this gap. Our main tool is a generalization of the ‘equivalence theorem’ of Magnus and Durbin (1999). We derive the bias, variance, and mean squared error of the pretest estimator (in fact, of a more general estimator, the so-called WALS estimator), and estimate the error of not reporting the correct moments (the ‘underreporting error’). This error can be very substantial. We also show that there can be large differences in underreporting between different model selection procedures. Finally, we investigate (in a special case) how the underreporting error increases when the number of auxiliary regressors z_1, \dots, z_m increases.

The paper is organized as follows. We define the formal framework and the notation in Section 2, where we also provide a generalization of the ‘equivalence theorem’ of Magnus and Durbin (1999). This generalization forms the basis of the subsequent analysis. In Section 3 we discuss underreporting and its bounds. Section 4 discusses the simplest case, where there is only one auxiliary regressor z , and hence only one possible pretest procedure (using the t -statistic). In Sections 5 and 6 we address the more difficult case where we have two auxiliary regressors. Then, there is no unique selection procedure. We show, *inter alia*, that there can be large differences between general-to-specific and specific-to-general model selection. Section 7 concludes. The appendix contains two useful results on constrained least squares.

2. The WALS estimator and the equivalence theorem

In this paper we will consider pretest-type estimators at two levels of generality. In the current section we consider only the most general type, the so-called WALS estimator, to be defined shortly. Our framework is the standard linear regression model

$$y = X\beta + Z\gamma + \varepsilon, \quad (1)$$

where y ($n \times 1$) is the vector of observations, X ($n \times k$) and Z ($n \times m$) are matrices of nonrandom regressors, ε ($n \times 1$) is a random vector of unobservable disturbances, and β ($k \times 1$) and γ ($m \times 1$) are unknown nonrandom parameter vectors. We assume that $k \geq 1$, $m \geq 1$, $k + m \leq n - 1$, that the design matrix $(X : Z)$ has full column-rank $k + m$, and that the disturbances $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ are i.i.d. $N(0, \sigma^2)$.

The reason for distinguishing between X and Z is that X contains explanatory variables which we want in the model on theoretical or other grounds (irrespective of the found t -values of the β -parameters), while Z contains additional explanatory variables of which we are less certain. Our focus is the estimation of β . Hence the only role for Z is to improve the estimation of β , while γ is a vector of nuisance parameters. The columns of X are called ‘focus’ regressors, and the columns of Z

‘auxiliary’ regressors. We define the matrices

$$M = I_n - X(X'X)^{-1}X', \quad Q = (X'X)^{-1}X'Z(Z'MZ)^{-1/2}$$

and the scaled and normalized parameter vector $\eta = (Z'MZ)^{1/2}\gamma/\sigma$. The matrix Q can be interpreted as the (scaled) correlation between X and Z . Clearly, $Q = \mathbf{0}$ if and only if Z is orthogonal to X . The least-squares (LS) estimators of β and γ are

$$b_u = b_r - Q\hat{\theta}, \quad \hat{\gamma} = (Z'MZ)^{-1}Z'My,$$

where $b_r = (X'X)^{-1}X'y$ and $\hat{\theta} = (Z'MZ)^{1/2}\hat{\gamma}$. The subscripts ‘ u ’ and ‘ r ’ denote ‘unrestricted’ and ‘restricted’ (with $\gamma = \mathbf{0}$), respectively. Letting $\hat{\eta} = \hat{\theta}/\sigma$, we see that $\hat{\eta} \sim N(\eta, I_m)$. Notice that $\hat{\eta}$ is only observable when σ is known, while $\hat{\theta}$ is observable whether σ is known or not.

Magnus and Durbin (1999) considered the estimation of β in model (1) and proposed a weighted-average least-squares (WALS) estimator of β of the form $b = \lambda b_u + (1 - \lambda)b_r$, where $\lambda = \lambda(\hat{\theta}, s_u^2)$ and s_u^2 denotes the estimator for σ^2 in the unrestricted model. This includes the usual pretest estimator as a special case, but only when one restricts the choice of model to the fully restricted and the fully unrestricted case. In the case $m = 1$, there are indeed two models: the unrestricted and the restricted (where $\gamma = 0$). But when $m = 2$, there are already four possible models: the unrestricted model, two partially restricted models (one of the two γ 's is zero), and the restricted model (both γ 's are zero). In general, there are 2^m models to consider. Thus we need a generalization of the WALS estimator, which includes as special cases not only the unrestricted estimator b_u and the restricted estimator b_r (where *all* γ 's are set equal to zero), but also many or all intermediate estimators where *some* of the γ 's are set equal to zero.

Let S_i be an $m \times r_i$ selection matrix of rank $r_i \geq 0$, so that $S_i' = (I_{r_i}; \mathbf{0})$ or a column-permutation thereof. Let \mathcal{M}_i denote the linear model (1) under the restriction $S_i'\gamma = \mathbf{0}$, and denote LS estimators of β and γ in model \mathcal{M}_i by $b_{(i)}$ and $c_{(i)}$. In Lemma A1 in the appendix we show that the partially restricted estimator $b_{(i)}$ can be written as $b_{(i)} = b_r - QW_i\hat{\theta}$, where W_i is an idempotent matrix, depending only on X , Z , and S_i . That is, $b_{(i)}$ is a linear function of the two *independent* vectors b_r and $\hat{\theta}$.² Moreover, $c_{(i)}$ is a linear function of $\hat{\theta}$ only and hence independent of b_r .

We now define the WALS estimator of β as $b = \sum_i \lambda_i b_{(i)}$, where the sum is taken over all 2^m different models obtained by setting a subset of the γ 's equal to zero, and the λ_i are weight-functions satisfying certain minimal regularity conditions, namely

$$\lambda_i \geq 0, \quad \sum_{i=1}^{2^m} \lambda_i = 1, \quad \lambda_i = \lambda_i(My).$$

The WALS estimator can then be written as $b = b_r - QW\hat{\theta}$, where $W = \sum_i \lambda_i W_i$. Notice that, while W_i is nonrandom, W is random.

² The independence of b_r and $\hat{\theta}$ follows immediately from the fact that $X'y$ and $Z'My$ are independent. In fact, even if the observations y_1, \dots, y_n are not normal and the data-generating process is unknown, b_r and $\hat{\theta}$ will still be uncorrelated, as long as the $\{y_i\}$ are uncorrelated with constant variance (Leeb and Pötscher (2003), Lemma A.1).

A few words about the regularity conditions are in order. If σ^2 is known, then most or all pretest procedures will use statistics (such as t - and F -statistics) which depend on $\hat{\theta}$ only. If σ^2 is not known and estimated by s_u^2 , then all t - and F -statistics will depend on $(\hat{\theta}, s_u^2)$. Now, it is a basic result in least-squares theory that s_u^2 is independent of $(\mathbf{b}_r, \hat{\gamma})$. It follows that \mathbf{b}_r is independent of s_u^2 . Hence, \mathbf{b}_r will be independent of $(\hat{\theta}, s_u^2)$. Finally, if σ^2 is not known and estimated by $s_{(i)}^2$ (the estimator of σ^2 in model \mathcal{M}_i), then it is no longer true that all t - and F -statistics depend only on $(\hat{\theta}, s_u^2)$. However, they still depend only on $\mathbf{M}\mathbf{y}$, because Lemma A1 implies that both $c_{(i)}$ and the residuals $e_{(i)}$ from model \mathcal{M}_i are linear functions of $\mathbf{M}\mathbf{y}$. We conclude that the regularity conditions on λ_i are reasonable and mild.³

Theorem 1 (Equivalence Theorem, generalized). *If the regularity conditions on λ_i are satisfied, then*

$$E(\mathbf{b}) = \boldsymbol{\beta} - \sigma \mathbf{Q}E(\mathbf{W}\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}), \quad \text{var}(\mathbf{b}) = \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \mathbf{Q}\text{var}(\mathbf{W}\hat{\boldsymbol{\eta}})\mathbf{Q}')$$

and hence

$$\text{MSE}(\mathbf{b}) = \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \mathbf{Q}\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})\mathbf{Q}').$$

Proof. Since \mathbf{b}_r and $\mathbf{M}\mathbf{y}$ are independent, we have

$$E(\mathbf{b}_r | \mathbf{M}\mathbf{y}) = E(\mathbf{b}_r), \quad \text{var}(\mathbf{b}_r | \mathbf{M}\mathbf{y}) = \text{var}(\mathbf{b}_r).$$

Hence,

$$\begin{aligned} E(\mathbf{b} | \mathbf{M}\mathbf{y}) &= E(\mathbf{b}_r | \mathbf{M}\mathbf{y}) - \mathbf{Q}E(\mathbf{W}\hat{\boldsymbol{\theta}} | \mathbf{M}\mathbf{y}) \\ &= E(\mathbf{b}_r) - \sigma \mathbf{Q}\mathbf{W}\hat{\boldsymbol{\eta}} = \boldsymbol{\beta} - \sigma \mathbf{Q}(\mathbf{W}\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \end{aligned}$$

and

$$\text{var}(\mathbf{b} | \mathbf{M}\mathbf{y}) = \text{var}(\mathbf{b}_r | \mathbf{M}\mathbf{y}) = \text{var}(\mathbf{b}_r) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

The unconditional mean and variance of \mathbf{b} and hence its mean squared error follow. \square

Theorem 1 provides a nontrivial generalization, using a simpler proof, of Theorem 2 in Magnus and Durbin (1999). Apparently, the properties of the complicated WALS estimator \mathbf{b} of $\boldsymbol{\beta}$ depend critically on the properties of the less complicated estimator $\mathbf{W}\hat{\boldsymbol{\eta}}$ of $\boldsymbol{\eta}$.⁴ We notice that neither the bias, nor the variance or the mean squared error of \mathbf{b} depend on $\boldsymbol{\beta}$. They do, however, depend on $\boldsymbol{\gamma}$ or, more accurately, on $\boldsymbol{\eta}$.

Theorem 1 allows us to obtain bounds for the mean squared error. If we let ξ_{\min} and ξ_{\max} denote the smallest and largest eigenvalue of $\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})$, respectively, then

$$\sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \xi_{\min}\mathbf{Q}\mathbf{Q}') \leq \text{MSE}(\mathbf{b}) \leq \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \xi_{\max}\mathbf{Q}\mathbf{Q}').$$

³ The regularity conditions allow not only all standard pretest procedures, but also inequality-constrained least squares. Thus, Theorem 1 below explains the ‘surprising symmetry’ found by Thomson and Schmidt (1982, p. 176).

⁴ For the pretest estimators studied in Giles and Srivastava (1993) and Leeb and Pötscher (2003) the MSE can also be obtained from the finite-sample densities derived in these papers.

For the traditional pretest estimator we shall see later (Fig. 1) that, for the case $m = 1$, the smallest and largest eigenvalues are 0.28 and 2.46, respectively. For the case $m \geq 2$, however, the bounds can become infinite, depending on the pretest procedure. In this paper our interest is not in the bounds for the true MSE, but in the difference between the true MSE and the reported MSE.

It is also interesting to compare the WALS estimator \mathbf{b} (model selection and pretesting properly taken into account) with the unrestricted estimator \mathbf{b}_u (no model selection, always use the unrestricted model). Since, $\text{MSE}(\mathbf{b}_u) = \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \mathbf{Q}\mathbf{Q}')$, we find that

$$\text{MSE}(\mathbf{b}_u) - \text{MSE}(\mathbf{b}) = \sigma^2(\mathbf{Q}(\mathbf{I}_m - \text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})\mathbf{Q}')).$$

Hence, in theory it is possible that $\text{MSE}(\mathbf{b})$ is uniformly (that is, for all $\boldsymbol{\eta}$) larger (or smaller) than $\text{MSE}(\mathbf{b}_u)$. This happens if all eigenvalues of $\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})$ are larger (smaller) than one. This situation is unlikely to occur in practical situations. In the traditional pretest context, it does not occur for $m = 1$ (see the discussion after Fig. 1) and $m = 2$ (checked numerically, but not reported in the paper).

3. Pretesting and underreporting

While we considered the general WALS estimator in the previous section, we now consider the less general (traditional) pretest estimator. In the idealized context of the linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \text{N}(\mathbf{0}, \sigma^2\mathbf{I}_n)$, we define a *pretest procedure* as a two-step procedure. In step 1 we select the model. In step 2 we estimate the unknown parameters $\boldsymbol{\beta}$ (and σ^2) from the selected model. This yields the pretest estimators \mathbf{b} (and s^2). In a pretest procedure thus defined, all λ_i are zero except one which is one. We require (as in Theorem 1) that the model selection criterion depends on \mathbf{y} only through $\mathbf{M}\mathbf{y}$, the residuals from the restricted model. From here on we shall also assume that σ^2 is known (but see our comments in the conclusion).

The mean squared error of the pretest estimator \mathbf{b} is, according to Theorem 1,

$$\text{MSE}(\mathbf{b}) = \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \mathbf{Q}\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})\mathbf{Q}').$$

In applied econometrics practice the same estimator \mathbf{b} is selected, but the effects of pretesting are ignored, the reported bias is zero, and hence the reported MSE equals the reported variance. The reported MSE equals

$$\widetilde{\text{MSE}}(\mathbf{b}) = \sigma^2((\mathbf{X}'\mathbf{X})^{-1} + \mathbf{Q}\mathbf{W}\mathbf{Q}'),$$

according to Lemma A1, since $\mathbf{W} = \mathbf{W}_i$ if the i th model is selected, and σ^2 is assumed known. Notice that $\widetilde{\text{MSE}}(\mathbf{b})$ is random since \mathbf{W} is random. Let $\boldsymbol{\omega}'\boldsymbol{\beta}$ be our focus parameter, where $\boldsymbol{\omega}$ is an arbitrary nonzero $k \times 1$ vector. In order to compare

$$\text{MSE}(\boldsymbol{\omega}'\mathbf{b}) = \sigma^2(\boldsymbol{\omega}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\omega} + \boldsymbol{\omega}'\mathbf{Q}\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}})\mathbf{Q}'\boldsymbol{\omega}) \tag{2}$$

with

$$\widetilde{\text{MSE}}(\boldsymbol{\omega}'\mathbf{b}) = \sigma^2(\boldsymbol{\omega}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\omega} + \boldsymbol{\omega}'\mathbf{Q}\mathbf{W}\mathbf{Q}'\boldsymbol{\omega}), \tag{3}$$

we define the *underreporting ratio* UR as one minus the ratio of (3) and (2). Thus,

$$UR = 1 - \frac{\widetilde{MSE}(\omega' \mathbf{b})}{MSE(\omega' \mathbf{b})} = \frac{\mathbf{q}'(\mathbf{R} - \mathbf{W})\mathbf{q}}{\mathbf{q}'\mathbf{R}\mathbf{q} + (1/q_0^2)}, \tag{4}$$

where

$$\mathbf{R} = \mathbf{R}(\boldsymbol{\eta}) = MSE(\mathbf{W}\hat{\boldsymbol{\eta}}), \quad \mathbf{q} = \frac{\mathbf{Q}'\boldsymbol{\omega}}{\sqrt{\boldsymbol{\omega}'\mathbf{Q}\mathbf{Q}'\boldsymbol{\omega}}}, \quad q_0^2 = \frac{\boldsymbol{\omega}'\mathbf{Q}\mathbf{Q}'\boldsymbol{\omega}}{\boldsymbol{\omega}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\omega}}.$$

Notice that $\mathbf{q}'\mathbf{q} = 1$. The UR is a random variable, since it depends on \mathbf{W} , which depends on $\hat{\boldsymbol{\eta}}$. Both the UR and its expectation are unobservable, since they depend on $\boldsymbol{\eta}$ via $\mathbf{R}(\boldsymbol{\eta})$.

One would expect that the matrix $MSE(\mathbf{b})$ is at least as large as the matrix $E(\widetilde{MSE}(\mathbf{b}))$ (in the sense that their difference is positive semidefinite), because pretesting introduces additional noise which is ignored in the reported MSE. Since $E(\mathbf{W}) = \sum_i (E\lambda_i)\mathbf{W}_i$ and

$$MSE(\mathbf{W}\hat{\boldsymbol{\eta}}) = \sum_{i=1}^{2^m} E(\lambda_i(\mathbf{W}_i\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\mathbf{W}_i\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})'),$$

this is guaranteed if the matrix

$$\sum_{i=1}^{2^m} E(\lambda_i((\mathbf{W}_i\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\mathbf{W}_i\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})' - \mathbf{W}_i)) \tag{5}$$

is positive semidefinite. We will see in Section 4 that it is possible to devise pretest procedures which do not satisfy this requirement. Such procedures, however, tend to be rather silly. We will say that a pretest procedure is *viable* if the matrix in (5) is positive semidefinite over the whole parameter space.⁵ For any viable pretest procedure, $E(UR)$ is a number between zero and one. When q_0^2 (known to the investigator) is zero, then there is no underreporting: $E(UR) = 0$.⁶ But when q_0^2 is large, $E(UR)$ can be close to one.

The $m \times m$ matrix $E(\mathbf{W})$ is a weighted average of idempotent matrices, and hence is bounded: all its elements are ≤ 1 in absolute value, and all its diagonal elements (and all its eigenvalues) lie in the interval $[0, 1]$. In fact,

$$0 \leq \pi_u \leq \xi_j(E\mathbf{W}) \leq 1 - \pi_r \leq 1 \quad (j = 1, \dots, m),$$

where $\xi_j(\mathbf{A})$ denotes the j th eigenvalue of \mathbf{A} , π_u is the probability of choosing the unrestricted model ($\mathbf{P}_i = \mathbf{0}$), and π_r the probability of choosing the restricted model ($\mathbf{P}_i = \mathbf{I}_m$).

The $E(UR)$ is a function of \mathbf{q} (normalized by $\mathbf{q}'\mathbf{q} = 1$), q_0^2 , $\boldsymbol{\eta}$, and $\mathbf{Z}'\mathbf{M}\mathbf{Z}$ (and m). Maximizing over \mathbf{q} gives the inequality

$$E(UR) \leq q_0^2 \max_{1 \leq j \leq m} \xi_j((\mathbf{I}_m + q_0^2\mathbf{R})^{-1/2}(\mathbf{R} - E\mathbf{W})(\mathbf{I}_m + q_0^2\mathbf{R})^{-1/2}).$$

⁵ The term ‘underreporting ratio’ might therefore be criticized since it can become negative. In the class of viable procedures which we study, the expected ratio is always positive.

⁶ This happens when $\mathbf{X}'\mathbf{Z} = \mathbf{0}$, but also (more generally and less trivially) when $\mathbf{Q}'\boldsymbol{\omega} = \mathbf{0}$. In either case $\mathbf{b} = \mathbf{b}_r$ whatever pretesting we do.

Then, letting $E^*(UR) := \max_{q, q_0^2} E(UR)$, we find, as $q_0^2 \rightarrow \infty$,

$$E^*(UR) = 1 - \min_{1 \leq j \leq m} \xi_j(\mathbf{R}^{-1/2}(\mathbf{E}\mathbf{W})\mathbf{R}^{-1/2}) \leq 1 - \frac{\pi_u}{\max_j \xi_j(\mathbf{R})}, \tag{6}$$

which depends on $\boldsymbol{\eta}$ and $\mathbf{Z}'\mathbf{M}\mathbf{Z}$ (and m). We see from (6) that the expected UR can be arbitrarily close to 1 if the mean squared error \mathbf{R} fails to be bounded in $\boldsymbol{\eta}$. This cannot happen when $m = 1$ (unless we choose the restricted model with probability $\pi_r \geq \bar{\pi}_r > 0$, where $\bar{\pi}_r$ is constant, and thus does not depend on the observed t -value), but it can happen when $m \geq 2$, as we will see in Section 6.

Finally, since $E(UR)$ depends on $\mathbf{Z}'\mathbf{M}\mathbf{Z}$, we briefly consider the role of this matrix. Without loss of generality, we may scale all \mathbf{z} -variables so that $\mathbf{z}'_j\mathbf{M}\mathbf{z}_j = 1$ for all $j = 1, \dots, m$. In the special case where we can choose the \mathbf{z} -variables to be ‘orthogonal’ (in the sense that $\mathbf{M}\mathbf{z}_i$ and $\mathbf{M}\mathbf{z}_j$ are orthogonal for every $i \neq j$), we have $\mathbf{Z}'\mathbf{M}\mathbf{Z} = \mathbf{I}_m$, and major simplifications occur.

Theorem 2. *Assume that $\mathbf{Z}'\mathbf{M}\mathbf{Z} = \mathbf{I}_m$ and that σ^2 is known.⁷ Then, (a) all models which include \mathbf{z}_j as a regressor yield the same t -statistic of γ_j ; and (b) suppose we select \mathbf{z}_j if and only if the t -statistic $\hat{\eta}_j$ is significant in the sense that $|\hat{\eta}_j| > c$ for some given $c > 0$ (such as $c = 1.96$). Then, letting $\lambda(x) = 1$ if $|x| > c$ and 0 otherwise, \mathbf{W} is a diagonal matrix with $w_{jj} = \lambda(\hat{\eta}_j)$, and $\text{MSE}(\mathbf{W}\hat{\boldsymbol{\eta}}) = \mathbf{V} + \mathbf{d}\mathbf{d}'$, where \mathbf{V} is a diagonal $m \times m$ matrix and \mathbf{d} an $m \times 1$ vector with typical elements*

$$v_{jj} = \text{var}(\lambda(\hat{\eta}_j)\hat{\eta}_j), \quad d_j = E(\lambda(\hat{\eta}_j)\hat{\eta}_j - \eta_j).$$

Proof. Lemma A2 directly implies (a). To prove (b) we note that \mathbf{W} is diagonal, because all \mathbf{W}_i are diagonal. Its j th diagonal element w_{jj} is either 0 (if \mathbf{z}_j is excluded from the model) or 1 (if \mathbf{z}_j is included), that is, $w_{jj} = \lambda(\hat{\eta}_j)$. Also, the components of $\mathbf{W}\hat{\boldsymbol{\eta}}$ are independent of each other, and hence the variance matrix is diagonal. \square

Since we will see that the choice of model selection procedure may matter a lot for the properties of the estimated focus parameters, it is advisable—if at all possible—to choose the auxiliary regressors such that $\mathbf{Z}'\mathbf{M}\mathbf{Z} = \mathbf{I}_m$. In the most common cases, this will make the pretest estimator independent of the chosen model selection procedure. It also allows us to obtain explicit analytical expressions for the moments of the estimator, and it guarantees bounded risk for any value of m . This is important, because, in general, risk is not necessarily bounded when $m \geq 2$; see Section 6.

4. Underreporting: one nuisance parameter

In the case of X nuisance parameter, the model becomes $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{z} + \boldsymbol{\varepsilon}$, where the nuisance parameter γ is a scalar. We have only two models to compare: the unrestricted

⁷ This is too restrictive. If σ^2 is always estimated by s_u^2 , the LS estimator in the unrestricted model, then the theorem still holds. In more general cases, an approximate version of the theorem holds.

($W_1 = 1$, $b_{(1)} = b_u$, $\lambda_1 = \lambda$) and the restricted ($W_2 = 0$, $b_{(2)} = b_r$, $\lambda_2 = 1 - \lambda$). As a result we find

$$b = \lambda b_u + (1 - \lambda) b_r, \quad W = \lambda,$$

and

$$\text{MSE}(W\hat{\eta}) = \text{MSE}(\lambda\hat{\eta}) = E(\lambda\hat{\eta} - \eta)^2, \quad E(W) = E(\lambda).$$

The underreporting ratio is thus

$$\text{UR}(\hat{\eta}, \eta) = \frac{R(\eta) - \lambda(\hat{\eta})}{R(\eta) + (1/q_0^2)},$$

where $\lambda(\hat{\eta}) = 1$ if $|\hat{\eta}| > c$ for some $c > 0$, and 0 otherwise, and

$$R(\eta) = E(\lambda\hat{\eta} - \eta)^2, \quad q_0^2 = \frac{(z'X(X'X)^{-1}\omega)^2}{(z'Mz)(\omega'(X'X)^{-1}\omega)}.$$

Assuming again that σ^2 is known and that c is given (say, $c = 1.96$), the λ -function depends only on $\hat{\eta}$, R depends only on η , and hence the UR depends on q_0^2 and $\hat{\eta}$ (both known to the investigator), and η (unknown).

It is easy to see that the larger is $R(\eta)$, the larger is UR. The random variable $\lambda\hat{\eta}$, considered as an estimator of η , thus plays a crucial role in determining the amount of underreporting. We consider its squared bias, variance and MSE in Fig. 1.⁸

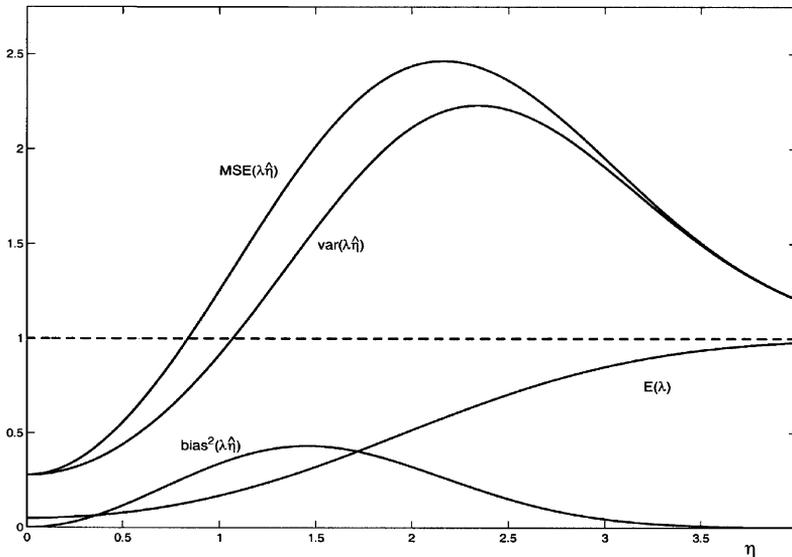


Fig. 1. Moments of $\lambda\hat{\eta}$ and λ compared ($m = 1$, $c = 1.96$).

⁸ All results reported in Figs. 1 and 2 are based on exact analytical formulas, with the exception of the locus curve in Fig. 2 which was obtained numerically (in Matlab 5.3) using the Nelder-Mead simplex method.

The bias of $\lambda\hat{\eta}$ is negative for $\eta > 0$ and reaches its minimum -0.66 at $\eta = 1.46$. The variance reaches its minimum 0.28 at $\eta = 0$ and its maximum 2.23 at $\eta = 2.34$. The MSE $R(\eta)$ is shaped similarly to the variance. It reaches its minimum at $\eta = 0$ and its maximum 2.46 at $\eta = 2.16$. The variance of $\lambda\hat{\eta}$ is large relative to its bias, suggesting that variance-reduction is more important than bias-reduction. We see that $\text{MSE}(\lambda\hat{\eta}) < 1$ if and only if $|\eta| < 0.84$. This means that the unrestricted estimator \mathbf{b}_u is better (has lower mean squared error) than the pretest estimator \mathbf{b} as soon as $|\eta| > 0.84$.

We also graph the expectation of the reported MSE of $\lambda\hat{\eta}$, that is, $E(\lambda)$, as a function of η for $c = 1.96$, and the MSE of the unrestricted estimator of η , that is $\text{MSE}(\hat{\eta})$ (the dashed line, constant at 1). Since λ only takes the values 0 and 1, $E(\lambda)$ denotes the probability of choosing the unrestricted model ($\lambda = 1$). But λ also denotes the reported variance (MSE). We see that $E(\lambda) \equiv \text{Pr}(|\hat{\eta}| > c)$ increases monotonically between 0.05 at $\eta = 0$ and 1 at $\eta = \infty$. Since $\text{MSE}(\lambda\hat{\eta}) > E(\lambda)$, the pretest procedure is viable.⁹

In Fig. 2 we graph the expected underreporting ratio $E(\text{UR})$ for five different values of q_0^2 : 0, 0.1, 1, 10, and ∞ . At $q_0^2 = 0$ there is no underreporting and $E(\text{UR}) = 0$. At $q_0^2 = 1$, we see that $E(\text{UR})$ is 0.18 at $\eta = 0$, reaches a maximum 0.57 at $\eta = 1.73$, and varies substantially with η (from 0 to 0.57), implying that on average the pretest MSE can be 2.3 times the reported MSE ($1/(1 - 0.57) = 2.3$). At $q_0^2 = \infty$, $E(\text{UR})$ is large; the maximum occurs at $\eta = 0.82$ where $E(\text{UR}) = 0.87$, so that the reported variance should be multiplied by about 7.5 in order to obtain the true MSE of the pretest estimator.

Clearly, the maximum of $E(\text{UR})$ depends on q_0 . This dependence is graphed as the dashed line in Fig. 2. We conclude that the effect of not reporting the true bias and variance of the pretest estimator can lead to serious misrepresentation of the results, *even in the case* $m = 1$. The larger is q_0^2 (known to the investigator), the larger will be the expected UR.

5. Model selection: general-to-specific and specific-to-general

When $m = 1$ pretesting is simple: look at the t -statistic for γ in the unrestricted model. If $|t| > c$, choose the unrestricted model (leading to \mathbf{b}_u); otherwise choose the restricted model (leading to \mathbf{b}_r). When $m > 1$ there are many ways to pretest. We consider the case $m = 2$ under the following conditions: model selection is based on t -statistics only, in the selected model all t -statistics are ‘significant’ (as defined in Theorem 2), and σ^2 is known.

Without loss of generality we normalize \mathbf{z}_1 and \mathbf{z}_2 , the regressors associated with the nuisance parameters γ_1 and γ_2 , by setting $\mathbf{z}_i' \mathbf{M} \mathbf{z}_i = 1$ for $i = 1, 2$. Then,

$$\mathbf{Z}' \mathbf{M} \mathbf{Z} = \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix},$$

⁹ However, not all λ -functions lead to a viable procedure. For example, the—admittedly silly—procedure defined by $\lambda = 1$ if $|\hat{\eta}| \leq c$ and 0 otherwise is not viable, since $\text{MSE}(\lambda\hat{\eta}) < E(\lambda)$ at $\eta = 0$ for any $c > 0$.

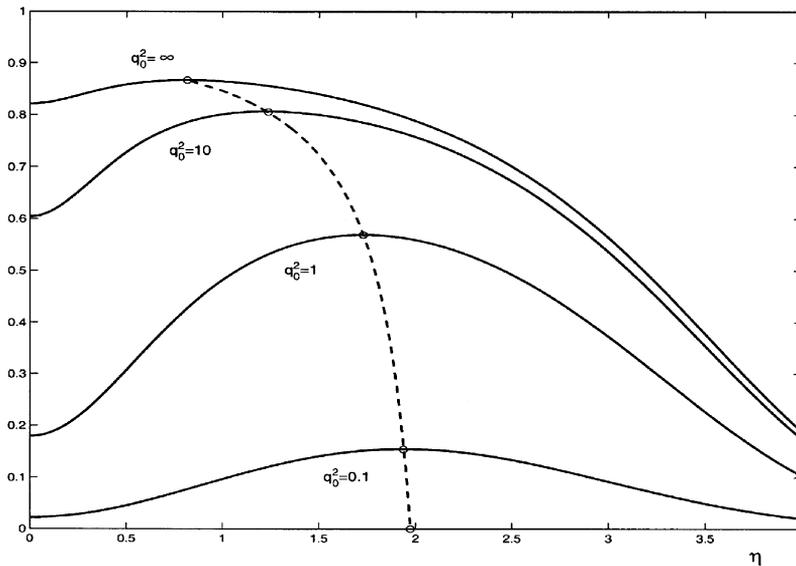


Fig. 2. E(UR) and locus of max(E(UR)) ($m = 1, c = 1.96$).

where $|r| < 1$, and

$$(\mathbf{Z}'\mathbf{M}\mathbf{Z})^{-1/2} = \frac{1}{\sqrt{1-r^2}} \begin{pmatrix} \alpha & -\rho \\ -\rho & \alpha \end{pmatrix},$$

with

$$\alpha = \frac{\sqrt{1+r} + \sqrt{1-r}}{2}, \quad \rho = \frac{\sqrt{1+r} - \sqrt{1-r}}{2}.$$

There are four t -statistics to consider: two in the unrestricted model (denoted t_1 and t_2), one in the model where $\gamma_2 = 0$ (denoted $t_{(1)}$), and one in the model where $\gamma_1 = 0$ (denoted $t_{(2)}$). Let $\hat{\eta}_1$ and $\hat{\eta}_2$ denote the components of $\hat{\boldsymbol{\eta}}$. Then, each of the four t -statistics is a linear function of $\hat{\eta}_1$ and $\hat{\eta}_2$ in accordance with Lemma A1:

$$t_1 = \alpha\hat{\eta}_1 - \rho\hat{\eta}_2, \quad t_2 = -\rho\hat{\eta}_1 + \alpha\hat{\eta}_2$$

and

$$t_{(1)} = \alpha\hat{\eta}_1 + \rho\hat{\eta}_2, \quad t_{(2)} = \rho\hat{\eta}_1 + \alpha\hat{\eta}_2.$$

Of course, since $\alpha^2 + \rho^2 = 1$, all four t -statistics are normally distributed with unit variance and, under the appropriate null hypothesis, mean zero. Also, $t_{(1)}$ is independent of t_2 and $t_{(2)}$ is independent of t_1 , for the same reason that \mathbf{b}_r and $\hat{\boldsymbol{\eta}}$ are independent. Further,

$$\text{corr}(t_1, t_{(1)}) = \text{corr}(t_2, t_{(2)}) = \sqrt{1-r^2} > 0$$

and

$$\text{corr}(t_1, t_2) = -r, \quad \text{corr}(t_{(1)}, t_{(2)}) = r.$$

Finally,

$$|t_1| > |t_2| \Leftrightarrow |t_{(1)}| > |t_{(2)}| \Leftrightarrow |\hat{\eta}_1| > |\hat{\eta}_2|.$$

We will investigate two pretest procedures that are in common use: ‘general-to-specific’ and ‘specific-to-general’. Let \mathcal{M}_0 denote the restricted model, \mathcal{M}_1 the model with only z_1 ($\gamma_2=0$), \mathcal{M}_2 the model with only z_2 ($\gamma_1=0$), and \mathcal{M}_{12} the unrestricted model. Then we define the general-to-specific (or ‘backward’ or ‘top–down’) procedure as follows: (a) Estimate the unrestricted model \mathcal{M}_{12} . This yields t -statistics t_1 and t_2 , (b) choose \mathcal{M}_{12} if both t_1 and t_2 are significant, (c) otherwise, (i) if $|t_1| > |t_2|$ estimate \mathcal{M}_1 , yielding $t_{(1)}$. If $t_{(1)}$ is significant choose \mathcal{M}_1 , otherwise choose \mathcal{M}_0 ; (ii) if $|t_1| \leq |t_2|$ estimate \mathcal{M}_2 , yielding $t_{(2)}$. If $t_{(2)}$ is significant choose \mathcal{M}_2 , otherwise choose \mathcal{M}_0 .

Similarly, we define the specific-to-general (or ‘forward’ or ‘bottom–up’) procedure as follows: (a) Estimate both partially restricted models \mathcal{M}_1 and \mathcal{M}_2 . This yields t -statistics $t_{(1)}$ and $t_{(2)}$; (b) choose \mathcal{M}_0 if neither $t_{(1)}$ nor $t_{(2)}$ is significant; (c) otherwise, estimate the unrestricted model yielding t_1 and t_2 , and choose \mathcal{M}_{12} if t_1 and t_2 are both significant; (d) in all other cases choose \mathcal{M}_1 (if $|t_{(1)}| > |t_{(2)}|$) or \mathcal{M}_2 (if $|t_{(1)}| \leq |t_{(2)}|$).

Since the two cases ($|t_{(1)}| \leq c < |t_1|$, $|t_2| \leq c < |t_{(2)}|$) and ($|t_{(2)}| \leq c < |t_2|$, $|t_1| \leq c < |t_{(1)}|$) cannot occur, we see that both procedures are identical, except for the case where t_1 and t_2 are both significant, while $t_{(1)}$ and $t_{(2)}$ are both not significant. In that case, the general-to-specific procedure chooses the unrestricted model and the specific-to-general procedure chooses the restricted model. In the special case $r = 0$, we find $t_1 = t_{(1)} = \hat{\eta}_1$ and $t_2 = t_{(2)} = \hat{\eta}_2$, and the two pretest procedures coincide. When $|r| \rightarrow 1$, the difference between the two procedures is at its largest. In spite of the seemingly small difference between the two pretest procedures, the effect of pretesting on underreporting will be surprisingly different for the two procedures.

6. Underreporting: two nuisance parameters

In the case $m = 1$ the expected underreporting ratio $E(\text{UR})$ depends (for fixed c) on two parameters: q_0^2 (known to the investigator) and η (unknown). In the case $m = 2$, $E(\text{UR})$ depends, after normalization, on five parameters: q_0^2 , q_1 (the first component of q) and r (known), and η_1 and η_2 (unknown). In addition, $E(\text{UR})$ depends on the procedure.

We have four models to compare: the unrestricted \mathcal{M}_{12} , the partially restricted \mathcal{M}_1 ($\gamma_2=0$) and \mathcal{M}_2 ($\gamma_1=0$), and the restricted \mathcal{M}_0 ($\gamma_1=\gamma_2=0$). This implies selection matrices $\mathbf{S}_0 = \mathbf{I}_2$, $\mathbf{S}_1 = (0, 1)'$, and $\mathbf{S}_2 = (1, 0)'$ (the matrix \mathbf{S}_{12} has no columns), and hence $\mathbf{W}_0 = \mathbf{0}$, $\mathbf{W}_{12} = \mathbf{I}_2$,

$$\mathbf{W}_1 = \frac{1}{2} \begin{pmatrix} 1 + \sqrt{1 - r^2} & r \\ r & 1 - \sqrt{1 - r^2} \end{pmatrix}$$

and

$$\mathbf{W}_2 = \frac{1}{2} \begin{pmatrix} 1 - \sqrt{1 - r^2} & r \\ r & 1 + \sqrt{1 - r^2} \end{pmatrix}.$$

Since $\mathbf{W} = \lambda_0 \mathbf{W}_0 + \lambda_1 \mathbf{W}_1 + \lambda_2 \mathbf{W}_2 + \lambda_{12} \mathbf{W}_{12}$, we thus find

$$\mathbf{W} = \frac{1}{2} \begin{pmatrix} \text{tr}(\mathbf{W}) + \sqrt{1-r^2}(\lambda_1 - \lambda_2) & r(\lambda_1 + \lambda_2) \\ r(\lambda_1 + \lambda_2) & \text{tr}(\mathbf{W}) - \sqrt{1-r^2}(\lambda_1 - \lambda_2) \end{pmatrix},$$

where $\text{tr}(\mathbf{W}) = \lambda_1 + \lambda_2 + 2\lambda_{12}$. As before, let $\lambda(x) = 1$ if $|x| > c$ and 0 otherwise. Then,

$$\begin{aligned} \lambda_0 &= (1 - \lambda(t_{(1)}))(1 - \lambda(t_{(2)})) - \delta B_1, & \lambda_1 &= \lambda(t_{(1)})(1 - \lambda(t_{(2)})) - (1 - \mu)B_2, \\ \lambda_2 &= \lambda(t_{(2)})(1 - \lambda(t_{(1)})) - \mu B_2, & \lambda_{12} &= \lambda(t_{(1)})\lambda(t_{(2)}) - (1 - \delta)B_1, \end{aligned}$$

with

$$\begin{aligned} B_1 &= \lambda(t_1)\lambda(t_2)(1 - \lambda(t_{(1)}))(1 - \lambda(t_{(2)})), \\ B_2 &= \lambda(t_{(1)})\lambda(t_{(2)})(1 - \lambda(t_1))(1 - \lambda(t_2)). \end{aligned}$$

Here, $\mu = 1$ if $|\hat{\eta}_1| > |\hat{\eta}_2|$ and 0 otherwise, and $\delta = 1$ if the pretest procedure is general-to-specific and 0 if the procedure is specific-to-general.

Because $E(\text{UR})$ depends on 5 parameters, only a six-dimensional plot would do full justice to its behavior. This task being beyond us, let us first consider the mean squared error $\mathbf{R} = \text{MSE}(\mathbf{W}\hat{\eta})$ and the expected reported variance $E(\mathbf{W})$ for the two procedures. Both functions depend on η_1, η_2 , and r . The $E(\mathbf{W})$ is always bounded, as noted in Section 3. The matrix \mathbf{R} is also bounded in the general-to-specific procedure, but \mathbf{R} can be unbounded in the specific-to-general procedure. More specifically,

$$\max_{\eta_1, \eta_2} \mathbf{R}(\eta_1, \eta_2, r) \rightarrow \infty \quad \text{as } r \rightarrow 1,$$

when the procedure is specific-to-general. This very different behavior of \mathbf{R} in the two procedures is reflected in Fig. 3, where we consider

$$E^{**}(\text{UR}) := \max_{\boldsymbol{\eta}} E^*(\text{UR}) = 1 - \min_{\boldsymbol{\eta}} \min_{1 \leq j \leq m} \xi_j(\mathbf{R}^{-1/2}(\mathbf{E}\mathbf{W})\mathbf{R}^{-1/2}),$$

as a function of r .¹⁰

For both procedures the function $E^{**}(\text{UR})$ is symmetric around $r = 0$. For $r = 0$ the two procedures are the same and the function value is almost 0.90. In the specific-to-general procedure, $E^{**}(\text{UR})$ increases monotonically to 1 as r increases from 0 to 1. The general-to-specific procedure has a uniformly lower $E^{**}(\text{UR})$, its behavior is non-monotonic, and it converges to 0.87 as $r \rightarrow 1$, the same maximum value as in the case $m = 1$ (depicted as a horizontal line in the figure). The difference between the two procedures is especially large when r is close to 1, that is when $\mathbf{M}\mathbf{z}_1$ and $\mathbf{M}\mathbf{z}_2$ are strongly correlated. This can be understood as follows. Let $r = 1$ and let $\eta_1 = -\eta_2 = \bar{\eta}$, say. Then, for large $\bar{\eta}$, the probability of choosing one of the partially restricted models \mathcal{M}_1 or \mathcal{M}_2 approaches 0. In the specific-to-general case, we will choose the restricted model \mathcal{M}_0 with probability approaching 0.95 and model \mathcal{M}_{12}

¹⁰ The $E(\text{UR})$ in Figs. 3, 5 and 6 was computed using routines D01FCF and D01GBF of the Fortran NAG Library, Mark 18. The (very computer-intensive) maximization in Fig. 3 is achieved by a grid method. Fig. 4 was obtained numerically in Matlab 5.3 using the Nelder-Mead simplex method.

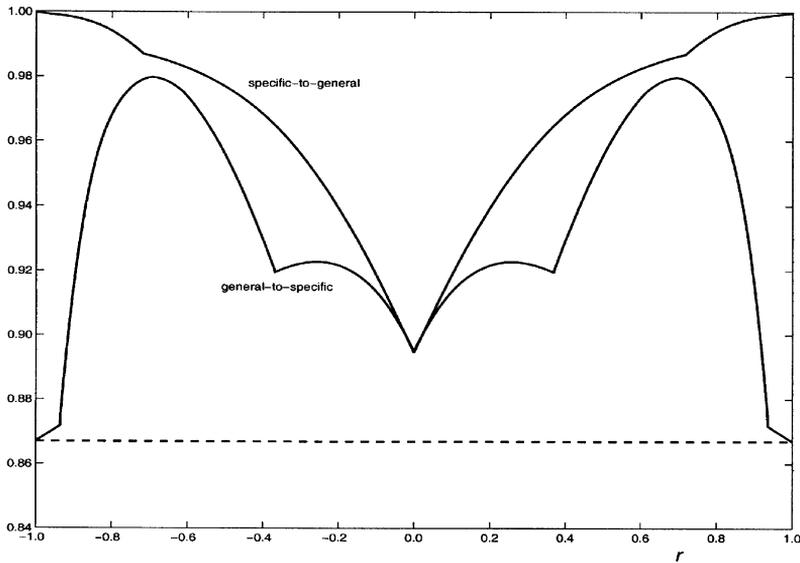


Fig. 3. $\max(E(UR))$ as a function of r ($m = 2$).

with probability approaching 0.05. Hence, for $r = 1$ and $\bar{\eta} \rightarrow \infty$, we find that $E(UR)$ approaches 1 for any q_0^2 . (In fact, the MSE of the pretest estimator is unbounded and proportional to $\bar{\eta}^2$ when $\bar{\eta}$ approaches ∞ .) But in the general-to-specific case, the MSE is always bounded and hence $E^*(UR) < 1$, using (6).

Although the functions are continuous, there are various kinks. This is the result of the fact that there exist various local maxima. At a kink we move from one local maximum to another local maximum. Clearly, underreporting can be a very serious problem and, for $m \geq 2$, can be essentially unbounded, depending on the chosen pretest procedure.

For $r = 0$ the worst case gives $E^{**}(UR) = 0.8669$ for $m = 1$ and 0.8953 for $m = 2$. We now ask how underreporting depends on m . There are 2^m models to consider and one may think therefore that ‘badness’ increases by a factor of 2^m . On the other hand, all t -statistics are functions of only m random variables $\hat{\eta}_1, \dots, \hat{\eta}_m$, so that ‘badness’ increases possibly only by a factor of m . We consider only the special case where $Z'MZ = I_m$. Then all vectors Mz_i are orthogonal, and the m -dimensional problem collapses to m one-dimensional problems (Theorem 2). The maximum $E^{**}(UR)$ is plotted in Fig. 4 as a function of m .

The figure reveals that $E^{**}(UR)$ increases with m but less than linearly. Although this result is valid only when $Z'MZ = I_m$, it nevertheless suggests that the increase in ‘badness’ is not as fast as one might have feared.

In a practical situation, we know q_0^2 , q , and r , but not η_1 and η_2 . Let us analyze a typical and representative situation where $q_0^2 = 2$, $q = (1/3, (2/3)\sqrt{2})'$ (so that $q'q = 1$), and $r = 0.8$.

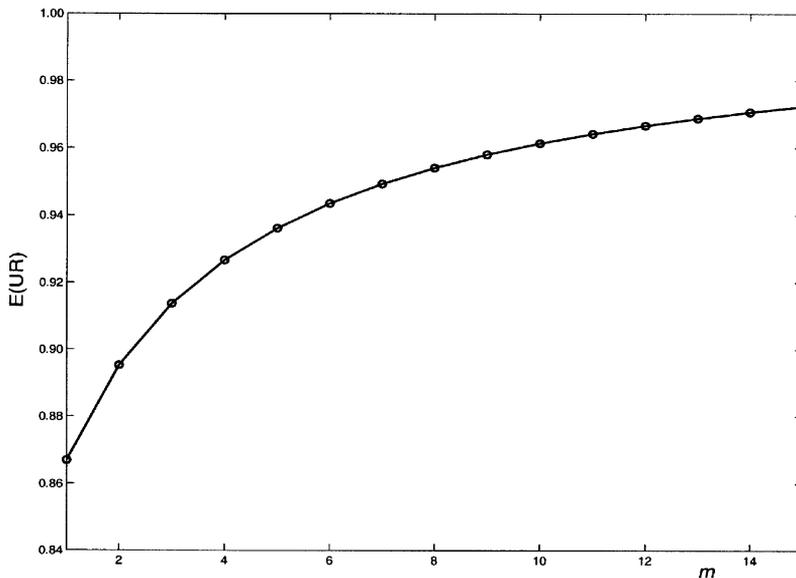


Fig. 4. $\max(E(\text{UR}))$ as a function of m ($Z'MZ = I_m$).

Figs. 5 and 6 give the $E(\text{UR})$ as a function of η_1 and η_2 , first for the general-to-specific procedure, then for the specific-to-general procedure. The $E(\text{UR})$ lies always between 0 and 1, and is symmetric around the point $(\eta_1, \eta_2) = (0, 0)$. The functional dependence on (η_1, η_2) is quite complicated, and also quite different for the two procedures. In the general-to-specific procedure (Fig. 5), $E(\text{UR})$ is very close to 0 at $(\eta_1, \eta_2) = (4, -4)$, but can be as large as 0.6551 at $(0.4, 1.6)$. In the specific-to-general procedure (Fig. 6), $E(\text{UR})$ varies from around 0 at $(4, 4)$ to 0.8798 around the point $(4, -4)$. In this case (and in general), the specific-to-general is more sensitive to underreporting than the general-to-specific procedure. The contours in the (η_1, η_2) plane are iso-value curves.

7. Conclusions

In this paper, we have analyzed the effect of ignoring the model selection procedure in reporting the bias, variance and mean squared error of the commonly used least-squares estimator. We conclude that underreporting can be a very serious problem. The pretest bias appears to be less of a problem than the pretest variance.

When we have m auxiliary regressors z_1, \dots, z_m , there are 2^m models to choose between. There are many different possible (viable) procedures to select the model. We find that the choice of model selection procedure matters a lot. Our results for $m = 2$ suggest that the general-to-specific procedure has more desirable properties than the specific-to-general procedure, in the sense that $\max(E(\text{UR}))$ is smaller for general-to-specific than for the specific-to-general. The influence of the selection

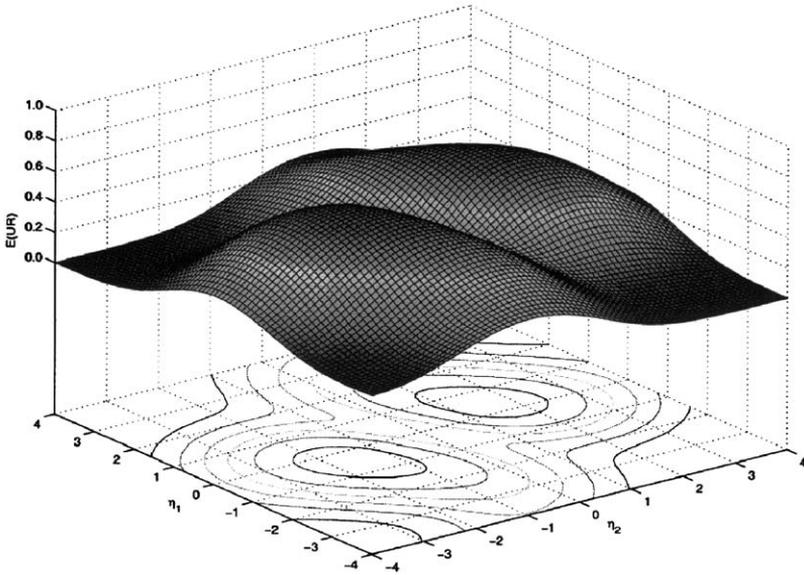


Fig. 5. $E(UR)$ as a function of η_1 and η_2 : general-to-specific.

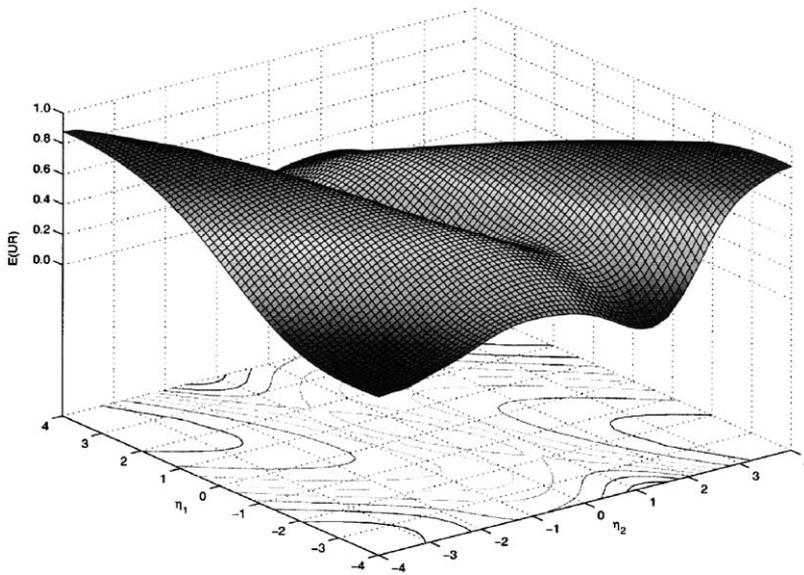


Fig. 6. $E(UR)$ as a function of η_1 and η_2 : specific-to-general.

Table 1
E(UR) as a function of the d.f. $n - k - m$ (σ^2 unknown)

$n - k - m$	η			
	0	1	2	4
10	0.76	0.83	0.77	0.26
30	0.80	0.85	0.78	0.22
50	0.81	0.86	0.79	0.21
∞	0.82	0.86	0.79	0.19

procedure is higher when the correlation between the \mathbf{z} -variables (measured by $\mathbf{Z}'\mathbf{M}\mathbf{Z}$) is high, than when it is low. If we can choose the auxiliary regressors such that they are ‘orthogonal’ (that is, $\mathbf{Z}'\mathbf{M}\mathbf{Z} = \mathbf{I}_m$), then the general-to-specific and specific-to-general pretest procedures are the same (among others), and hence the sampling properties of the estimators do not depend on the model selection procedure.

As the number of auxiliary regressors m grows, the dangers of underreporting grow as well. In the special ‘orthogonal’ case ($\mathbf{Z}'\mathbf{M}\mathbf{Z} = \mathbf{I}_m$) this growth is somewhat less than linear.

Although Theorem 1 is valid whether or not σ^2 is known, Sections 3–6 assume that σ^2 is known. This is of course unrealistic and we need to address the question how the results are affected when σ^2 is unknown. As an example, let us consider the case of Fig. 2 where $m = 1$, $q_0^2 = \infty$ and $c = 1.96$. When σ^2 is known, the E(UR) takes the values 0.82, 0.86, 0.79, and 0.19 for η equal to 0, 1, 2 and 4, respectively. When σ^2 is not known the calculations are more involved and depend on the degrees of freedom $n - k - m$. The results are summarized in Table 1.¹¹

We see that the effects of estimating σ^2 are relatively small, especially in the region of interest where $|\eta|$ is around 1 or 2. Although this example is typical for the behavior of the E(UR), more work is needed in this direction, especially for $m \geq 2$.

Future work will also have to clarify various other issues not covered in this paper. For example, there are selection procedures other than general-to-specific and specific-to-general. How does the E(UR) depend on these? All these pretest procedures are discontinuous; they choose one of 2^m models based on the values of t - and F -statistics. A continuous procedure can also be defined. Theorem 1 allows for this situation, and it will probably lead to better sampling properties; see Magnus (2002) for an analysis of the case $m = 1$.

Also, the E(UR) depends on η which is unknown. What is the best way to estimate the expected underreporting ratio? Simplest is to replace η by $\hat{\eta}$. The properties of this estimator will have to be analyzed. Another possibility is to report the worst possible case, given the observed \mathbf{y} , \mathbf{X} and \mathbf{Z} . That is, to calculate $\max_{\eta} \text{E(UR)}$. In Figs. 5 and 6, for example, the maxima are 0.6551 (general-to-specific) and 0.8798

¹¹ Table 1 was obtained by numerical integration (in Matlab 5.3, using NAG routine D01AMF) of $R(\eta/s_u)$, where s_u^2 is the LS estimator of σ^2 in the unrestricted model.

(specific-to-general), and they give an indication of how bad underreporting can be in a specific situation.

Acknowledgements

We are grateful to seminar participants at the University of Michigan, Michigan State University, the University of Wisconsin at Madison, Tilburg University, Eindhoven University of Technology, the University of Amsterdam, the University of Vienna, ESEM 2001 in Lausanne, and to two referees for constructive and useful comments.

Appendix. Constrained least squares

Consider the linear regression model $y = X\beta + Z\gamma + \varepsilon$, where y is the vector of observations, X ($n \times k$) and Z ($n \times m$) are matrices of nonrandom regressors, ε is a random vector of unobservable disturbances, and β and γ are unknown nonrandom parameter vectors. Assume that $k \geq 1$, $m \geq 1$, $k + m \leq n - 1$, that $(X : Z)$ has full column-rank, and that the disturbances $\{\varepsilon_i\}$ are i.i.d. $N(0, \sigma^2)$. Define $M = I_n - X(X'X)^{-1}X'$, $Q = (X'X)^{-1}X'Z(Z'MZ)^{-1/2}$, and $\eta = (Z'MZ)^{1/2}\gamma/\sigma$. Let S_i be an $m \times r_i$ selection matrix of rank $r_i \geq 0$, so that $S_i' = (I_{r_i} : \mathbf{0})$ or a column-permutation thereof. We are interested in the constrained LS estimators of β and γ , the constraint being $S_i'\gamma = \mathbf{0}$.

Lemma A1. *Let $b_r = (X'X)^{-1}X'y$ and $\hat{\theta} = (Z'MZ)^{-1/2}Z'My$. Then the constrained LS estimators of β and γ are given by*

$$b_{(i)} = b_r - QW_i\hat{\theta}, \quad c_{(i)} = (Z'MZ)^{-1/2}W_i\hat{\theta},$$

where

$$P_i := (Z'MZ)^{-1/2}S_i(S_i'(Z'MZ)^{-1}S_i)^{-1}S_i'(Z'MZ)^{-1/2}$$

is a symmetric idempotent $m \times m$ matrices of rank r_i , and $W_i := I_m - P_i$. (If $r_i = 0$ then $P_i = \mathbf{0}$.) The distribution of $b_{(i)}$ is given by

$$b_{(i)} \sim N(\beta + \sigma QP_i\eta, \sigma^2((X'X)^{-1} + QW_iQ')),$$

the distribution of $c_{(i)}$ by

$$c_{(i)} \sim N(\sigma(Z'MZ)^{-1/2}W_i\eta, \sigma^2((Z'MZ)^{-1/2}W_i(Z'MZ)^{-1/2})),$$

and the covariance of $b_{(i)}$ and $c_{(i)}$ is $\text{cov}(b_{(i)}, c_{(i)}) = -\sigma^2 QW_i(Z'MZ)^{-1/2}$. The residual vector is $e_{(i)} = y - Xb_{(i)} - Zc_{(i)} = D_iy$, where

$$D_i = M - MZ(Z'MZ)^{-1/2}W_i(Z'MZ)^{-1/2}Z'M$$

is a symmetric idempotent matrix of rank $n - k - m + r_i$, and the distribution of $s_{(i)}^2 = e_{(i)}'e_{(i)}/(n - k - m + r_i)$ is given by

$$\frac{(n - k - m + r_i)s_{(i)}^2}{\sigma^2} \sim \chi^2(n - k - m + r_i, \eta'P_i\eta).$$

Proof. Let $X_* = (X : Z)$, $\beta_* = (\beta', \gamma')$, and $R' = (\mathbf{0} : S'_i)$. The LS estimator of β_* in the model $y = X_*\beta_* + \varepsilon$ under the restriction $R'\beta_* = \mathbf{0}$ is then given by

$$b_* = (X_*'X_*)^{-1}X_*'y - (X_*'X_*)^{-1}R'(R'(X_*'X_*)^{-1}R)^{-1}R'(X_*'X_*)^{-1}X_*'y.$$

Noting that

$$(X_*'X_*)^{-1} = \begin{pmatrix} X'X & X'Z \\ Z'X & Z'Z \end{pmatrix}^{-1} = \begin{pmatrix} (X'X)^{-1} + QQ' & -Q(Z'MZ)^{-1/2} \\ -(Z'MZ)^{-1/2}Q' & (Z'MZ)^{-1} \end{pmatrix},$$

and simplifying, the results follow. \square

Lemma A2. Under the same conditions as in Lemma A1, assume, in addition, that $Z'MZ = I_m$. The constrained LS estimators of β and γ are then given by $b_{(i)} = b_r - QW_i\hat{\theta}$ and $c_{(i)} = W_i\hat{\theta}$, where $W_i := I_m - S_iS'_i$ is a diagonal $m \times m$ matrix with $m - r_i$ ones and r_i zeros on the diagonal, such that the j th diagonal element of W_i is 0 if γ_j is constrained to be zero, and 1 otherwise. The distribution of $b_{(i)}$ is given by

$$b_{(i)} \sim N(\beta + \sigma QS_iS'_i\eta, \sigma^2((X'X)^{-1} + QW_iQ')),$$

the distribution of $c_{(i)}$ by $c_{(i)} \sim N(\sigma W_i\eta, \sigma^2 W_i)$, and the covariance of $b_{(i)}$ and $c_{(i)}$ is $\text{cov}(b_{(i)}, c_{(i)}) = -\sigma^2 QW_i$. Hence all models which include z_j as a regressor will have the same estimator of γ_j , namely $\hat{\theta}_j$, irrespective which other γ 's are estimated. Moreover, the estimators $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_m$ are independent.

Proof. We have $P_i = S_i(S'_iS_i)^{-1}S'_i$, and, since S'_i is a selection matrix of the form $(I_{r_i} : \mathbf{0})$ or a column-permutation thereof, it follows that $S'_iS_i = I_{r_i}$ and hence that P_i is a diagonal matrix with r_i ones and $m - r_i$ zeros on the diagonal, and that W_i is a diagonal matrix with $m - r_i$ ones and r_i zeros on the diagonal. The results now follow from Lemma A1. \square

References

- Abadir, K.M., Magnus, J.R., 2002. Notation in econometrics: a proposal for a standard. *The Econometrics Journal* 5, 76–90.
- Bancroft, T.A., 1944. On biases in estimation due to the use of preliminary tests of significance. *Annals of Mathematical Statistics* 15, 190–204.
- Chatfield, C., 1995. Model uncertainty, data mining and statistical inference (with discussion). *Journal of the Royal Statistical Society, Series A* 158, 419–466.
- Giles, J.A., Giles, D.E.A., 1993. Pre-test estimation and testing in econometrics: recent developments. *Journal of Economic Surveys* 7, 145–197.
- Giles, J.A., Srivastava, V., 1993. The exact distribution of a least-squares regression coefficient after a preliminary t -test. *Statistics and Probability Letters* 16, 59–64.
- Hendry, D.F., 2001. Achievements and challenges in econometric methodology. *Journal of Econometrics* 100, 7–10.
- Hoover, K.D., Perez, S.J., 1999. Data mining reconsidered; encompassing and the general-to-specific approach to specification search. *The Econometrics Journal* 2, 1–25.
- Huntsberger, D.V., 1955. A generalization of a preliminary testing procedure for pooling data. *Annals of Mathematical Statistics* 26, 734–743.

- Judge, G.G., Bock, M.E., 1978. *The Statistical Implications of Pre-Test and Stein-Rule Estimators in Econometrics*. North-Holland, Amsterdam.
- Judge, G.G., Bock, M.E., 1983. Biased estimation. In: Griliches, Z., Intriligator, M.D. (Eds.), *Handbook of Econometrics*, Vol. 1. North-Holland, Amsterdam (Chapter 10).
- Leeb, H., Pötscher, B.M., 2003. The finite-sample distribution of post-model-selection estimators, and uniform versus non-uniform approximations. *Econometric Theory* 19, 100–142.
- Lovell, M.C., 1983. Data mining. *The Review of Economics and Statistics* 65, 1–12.
- Magnus, J.R., 1999. The traditional pretest estimator. *Theory of Probability and its Applications* 44, 293–308.
- Magnus, J.R., 2002. Estimation of the mean of a univariate normal distribution with known variance. *The Econometrics Journal* 5, 225–236.
- Magnus, J.R., Durbin, J., 1999. Estimation of regression coefficients of interest when other regression coefficients are of no interest. *Econometrica* 67, 639–643.
- Mittelhammer, R.C., 1984. Restricted least squares, pre-test, OLS and Stein-rule estimators: Risk comparisons under model misspecification. *Journal of Econometrics* 25, 151–164.
- Mosteller, F., 1948. On pooling data. *Journal of the American Statistical Association* 43, 231–242.
- Pötscher, B.M., 1991. Effects of model selection on inference. *Econometric Theory* 7, 163–185.
- Pötscher, B.M., Novak, A.J., 1998. The distribution of estimators after model selection: large and small sample results. *Journal of Statistical Computation and Simulation* 60, 19–56.
- Roehrig, C.S., 1984. Optimal critical regions for pre-test estimators using a Bayes risk criterion. *Journal of Econometrics* 25, 3–14.
- Sclove, S.L., Morris, C., Radhakrishnan, R., 1972. Non-optimality of preliminary-test estimators for the mean of a multivariate normal distribution. *The Annals of Mathematical Statistics* 43, 1481–1490.
- Sen, P.K., 1979. Asymptotic properties of maximum likelihood estimators based on conditional specification. *Annals of Statistics* 7, 1019–1033.
- Thomson, M., Schmidt, P., 1982. A note on the comparison of the mean square error of inequality constrained least squares and other related estimators. *The Review of Economics and Statistics* 64, 174–176.
- Zhang, P., 1992. On the distributional properties of model selection criteria. *Journal of the American Statistical Association* 87, 732–737.