

# PROPERTIES OF PREDICTORS IN OVERDIFFERENCED NEARLY NONSTATIONARY AUTOREGRESSION

BY ISMAEL SÁNCHEZ AND DANIEL PEÑA

*Universidad de Alicante and Universidad Carlos III de Madrid*

*First version received June 1998*

**Abstract.** We analyze the effect of overdifferencing a stationary  $AR(p+1)$  process whose largest root is near unity. It is found that, if the process is nearly nonstationary, the estimators of the overdifferenced model  $ARIMA(p, 1, 0)$  are root- $T$  consistent. It is also found that this misspecified  $ARIMA(p, 1, 0)$  has lower predictive mean squared error, to terms of small order, than the properly specified  $AR(p+1)$  model due to its parsimony. The advantage of the overdifferenced predictor depends on the remaining roots, the prediction horizon and the mean of the process.

**Keywords.** Autoregressive processes; near nonstationarity; overdifferencing; parsimony; predictive mean squared error; unit roots.

## 1. INTRODUCTION

In this paper, we investigate the consequences in estimation and prediction of overdifferencing a stationary  $AR(p+1)$  with a root close to unity. Differencing is normally used to transform a homogeneous linear nonstationary time series into a stationary process that is often modeled as an  $ARMA(p, q)$  process. It is said, then, that the original series follows an  $ARIMA(p, d, q)$  process, where  $d$  is the number of differences required to obtain stationarity. We assume that the process is not a long-memory process (see, for instance, Granger and Joyeux, 1980) and thus  $d$  is an integer equal to the number of unit roots in the autoregressive characteristic equation. When a stationary process has an autoregressive characteristic equation with a root close to unity it is said to be nearly nonstationary. Given a small or moderate sample of this process, it is very likely to be concluded, due to the low power of unit root tests in this case, that a difference should be applied. The differenced series will be noninvertible and the process is called overdifferenced.

Since the work of Fuller (1976) and Dickey and Fuller (1970), there has been a vast literature concerning the detection of unit roots in autoregressive polynomials. This literature notes the difficulty of a correct detection in near nonstationary processes. In spite of this, relatively little has been written on the consequences of a wrong detection. Previous work on the effect of overdifferencing can be found in Plosser and Schwert (1977, 1978), Harvey (1981), Campbell and Perron (1991) and Stock (1996). Plosser and Schwert (1977) examine, using Monte Carlo techniques, the

effect of overdifferencing in two cases: processes with a deterministic linear trend and stochastic regression models. They conclude that, in these situations, the loss in efficiency in both parameter estimation and prediction is not substantial, provided a moving-average parameter is included. Harvey (1981), assuming known parameters, also concludes that overdifferencing does not need to have serious implications for prediction, provided a finite-sample prediction procedure is used and a moving-average parameter is included. Campbell and Perron (1991) and Stock (1996) compare, using Monte Carlo simulations, the prediction accuracy of an AR(1) and a random walk. The empirical results of these authors show that the random walk can produce forecasts with lower prediction mean squared error (PMSE) than the AR(1) if the root is close to unity.

In this paper, we justify theoretically the advantages of the overdifferenced predictor, found empirically by Campbell and Perron (1991) and Stock (1996), in a general autoregression, and analyze the effect of other factors like the remaining roots, sample size ( $T$ ) and horizon ( $H$ ). We will assume that a root of the AR( $p + 1$ ) is close to unity, and thus we will adopt as a more plausible overdifferenced predictor the ARIMA( $p, 1, 0$ ) model, where no moving-average component is involved.

We will prove that the PMSE of the overdifferenced model ARIMA( $p, 1, 0$ ) is lower, to terms of small order, than the PMSE of the correct model AR( $p + 1$ ) if the root that is closer to unity,  $\rho^{-1}$ , follows  $\rho = \exp(-c/T^\beta)$ ,  $\beta > 1$ . The advantage of the overdifferenced predictor is due to its parsimony. Therefore, it is larger if the AR( $p + 1$ ) process has a non-zero mean, since this will vanish in the overdifferenced model. The remaining roots also affect the advantage of the overdifferenced predictor. Positive roots increase the advantage of the overdifferenced model, whereas negative roots have the opposite effect. The advantage of the overdifferenced model is small in the short term, but can increase with the horizon.

An important consequence of these results is that, for forecasting purposes, it is better to overdifference than to underdifference. Therefore, the possible low power of unit root tests in autoregression is not as important in forecasting as in model identification, since we can still obtain an efficient predictor.

The paper is organized as follows. In Section 2 we introduce the model and notation. In Section 3 we define nearly nonstationary processes. The consequences of overdifferencing in estimation are analyzed in Section 4, and the effect on the PMSE for each predictor is analyzed in Section 5. In Section 6 we compare the PMSE of the competing models and extract further results from the AR(1) case. A simulation study is presented in Section 7 to illustrate the results.

## 2. THE MODEL AND NOTATION

Let  $\{y_t\}$  be the following stationary AR( $p + 1$ ) process:

$$\varphi(B)y_t = \phi(B)(1 - \rho B)y_t = \alpha + a_t \quad (2.1)$$

where  $B$  is the backshift operator, and  $\varphi(B) = (1 - \sum_{i=1}^{p+1} \varphi_i B^i)$  is a polynomial operator on  $B$  such that  $\varphi(B) = 0$  has all its roots outside the unit circle, with  $\rho^{-1}$  being the root closer to unity. Let  $a_t$  be a sequence of independent identically distributed random variables with zero mean and variance  $\sigma^2$ . Let  $\mu = E(y_t)$ ; then  $\alpha = \mu\varphi(1)$ . We make the following assumption.

A1. For some  $s_0 > 2$ ,  $E\{|a_t|^{s_0}\} < \infty$ .

It is well known that this model can be represented in first-order vector autoregressive form as follows:

$$Y_t = A_\alpha Y_{t-1} + U_{t,p+2} \tag{2.2}$$

with  $Y_t = (y_t, \dots, y_{t-p}, 1)'$ ,  $U_{t,p+2} = (a_t, 0, \dots, 0)'$ , where the subindex  $p + 2$  indicates the dimension of the vector and

$$A_\alpha = \begin{pmatrix} \varphi_1 & \varphi_2 & \dots & \varphi_p & \varphi_{p+1} & \alpha \\ 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 & 1 \end{pmatrix}.$$

Then  $y_t = e'_{p+2} Y_t$ , with  $e_{p+2} = (1, 0, \dots, 0)'$ . Let  $\Gamma_y = E(Y_t Y_t')$  and  $\gamma_y = E(Y_t y_{t+1})$ . If we represent the process in deviations from the mean, we obtain  $\tilde{Y}_t = A_0 \tilde{Y}_{t-1} + U_{t,p+1}$ , where  $\tilde{Y}_t = (\tilde{y}_t, \tilde{y}_{t-1}, \dots, \tilde{y}_{t-p})'$ ,  $\tilde{y}_t = y_t - \mu$ , and  $A_0$  is the first  $(p + 1) \times (p + 1)$  submatrix of  $A_\alpha$ . We will also denote  $\Gamma_{\tilde{y}} = E(\tilde{Y}_t \tilde{Y}_t')$ . If a difference is applied to  $y_t$ , the series obtained,  $w_t = (1 - B)y_t$ , can be represented as

$$\phi(B)(1 - \rho B)w_t = (1 - B)a_t \tag{2.3}$$

which is noninvertible. The process  $w_t$  has the following vector representation (Lütkepohl, 1991, p. 223):

$$Z_t = A_1 Z_{t-1} + U_{t,p+2}^* \tag{2.4}$$

with  $Z_t = (W_t', a_t)'$ ,  $W_t = (w_t, \dots, w_{t-p})'$ ,  $U_{t,p+2}^* = (a_t, 0, \dots, 0, a_t)'$  and

$$A_1 = \begin{pmatrix} A_0 & -e_{p+1} \\ 0 \dots 0 & 0 \end{pmatrix}$$

with  $w_t = e'_{p+1} Z_t$ . Let  $\Gamma_w = E(W_t W_t')$  and  $\gamma_w = E(W_t w_{t+1})$ . In what follows, we will use a circumflex  $\hat{\phantom{x}}$  to denote estimates from a sample of the overdifferenced process  $\{w_t\}$  and the check symbol  $\check{\phantom{x}}$  for estimates from a sample of the original process  $\{y_t\}$ . The least squares estimator of the AR( $p + 1$ ) parameter vector  $\varphi = (\varphi_1, \dots, \varphi_{p+1}, \alpha)'$ , fitted to a sample of size  $T$  of the original process, is  $\check{\varphi} = \check{\Gamma}_y^{-1} \check{\gamma}_y$ , where  $\check{\Gamma}_y = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} Y_j Y_j'$  and  $\check{\gamma}_y = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} Y_j y_{j+1}$ . Similarly, the least squares estimator of

the parameter vector  $\phi = (\phi_1, \dots, \phi_p)'$  of a misspecified  $AR(p)$ , fitted to a sample of size  $T - 1$  of the overdifferenced process (2.3), is  $\hat{\phi} = \hat{\Gamma}_w^{-1} \hat{\gamma}_w$ , where  $\hat{\Gamma}_w = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} W_j W_j'$  and  $\hat{\gamma}_w = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} W_j w_{j+1}$ . We also make the following assumptions, where  $\|\cdot\|$  denotes the Euclidean norm.

A2.  $E(\|\check{T}_y^{-1}\|^{2k})$  ( $k = 1, 2, \dots, k_0$ ) is bounded for all finite and sufficiently large  $T$  and some  $k_0$ .

A3.  $E(\|\hat{T}_w^{-1}\|^{2k})$  ( $k = 1, 2, \dots, k_0$ ) is bounded for all finite and sufficiently large  $T$  and some  $k_0$ .

Assumptions A2 and A3 are similar to Assumption A3 of Kunitomo and Yamamoto (1985). They are also equivalent to Assumption A3 of Bhansali (1981). It should be noted that they are satisfied if the distribution is normal (see Fuller and Hasza, 1981). These assumptions are needed in several parts of this work, especially in application to the results of Kunitomo and Yamamoto (1985) and Bhansali (1981). They imply that, for a large enough sample size, the estimations of the covariance matrices are sufficiently near the true values (Bhansali, 1981, p. 590).

### 3. NEARLY NONSTATIONARY AUTOREGRESSIONS

A process is said to be nearly nonstationary (near integrated) if its autoregressive characteristic equation has a root,  $\rho^{-1}$ , very close to unity. If  $\rho$  is close enough to unity, the term  $1 - \rho B$  in (2.3) will be similar to  $1 - B$ . Therefore, although the overdifferenced process  $w_t$  is strictly a noninvertible  $ARMA(p + 1, 1)$ , an average correlogram of  $w_t$  will suggest estimating by an  $AR(p)$  instead.

The similarity between  $w_t$  and a true  $AR(p)$  process does not only depend on  $\rho$  but is influenced by the remaining roots. In order to see this point, let  $\pi_j$  be the coefficients of the polynomial  $\pi(B) = 1 - \pi_1 B - \pi_2 B^2 - \dots$ , where  $\varphi(B) = \pi(B)(1 - B)$ . These coefficients follow

$$\pi_j = \begin{cases} \phi_j + (\rho - 1)(1 - \sum_{k=1}^{j-1} \phi_k) & \text{if } j \leq p \\ (\rho - 1)(1 - \sum_{k=1}^p \phi_k) & \text{if } j > p \end{cases} \tag{3.1}$$

with  $\phi_k = 0$  if  $k < 1$ . If we denote as  $r_i^{-1}$ ,  $i = 1, \dots, p$ , the roots of the characteristic equation  $\phi(B) = 0$ , then

$$\left(1 - \sum_{k=1}^p \phi_k B^k\right) = \prod_{i=1}^p (1 - r_i B). \tag{3.2}$$

Therefore, negative values of  $r_i$  increase the value of  $\pi_j$ ,  $j > p$ , and decrease the similarity of  $w_t$  and an  $AR(p)$ .

Thus, the definition of a nearly nonstationary process needs (i) a parameterization that converges to the unit root with the sample size and (ii)

a constant term that can reflect the influence of the remaining roots in finite samples. Phillips (1987) and Chan and Wei (1987) define a nearly nonstationary process for the AR(1) case by reparameterizing  $\rho = \exp(-c/T) = 1 - c/T + o(T^{-1})$ , where  $c$  is a fixed constant. In this definition, the convergence rate to unity is fixed to be  $O(T^{-1})$ . These authors use this definition to provide asymptotic theory for the estimation of  $\rho$ . The formulation is justified by Phillips (1987) because this is the order of consistency of the least squares estimator, and by Chan and Wei (1987) because this is the order of the observed Fisher information of  $\rho$  under normality. In order to analyze the consequences of overdifferencing with different convergence rates we will define  $\rho$  as

$$\rho = \exp\left(-\frac{c}{T^\beta}\right) \tag{3.3}$$

with  $c$  and  $\beta$  being fixed constants. We deal only with stationary processes, and hence  $c, \beta > 0$ . Time series generated by (2.1) and (3.3) formally constitute a triangular array of the type  $\{y_{iT}: t = 1, \dots, T; T = 1, 2, \dots\}$ . Since this formulation is not essential in this paper, we will still use the notation  $\{y_t\}$  to refer to this process. It has to be noted that, since  $\alpha = E(y_t)(1 - \rho)\phi(1)$ , the process has no constant term if  $\rho = 1$ .

Given a sample from a process generated by (2.1) and (3.3), the analyst has to decide whether to estimate  $\rho$  or to impose the value  $\rho = 1$ . By the properties of least squares estimators it can be proved that the least squares estimator of  $\rho$  satisfies  $\hat{\rho} = \rho + O_p(T^{-(\beta+1)/2})$ , whereas imposing unity has the property  $1 = \rho + O(T^{-\beta})$ . Then, for  $\beta > 1$ , the convergence rate when imposing unity is faster than estimating by least squares. This result helps to explain why processes with  $\beta > 1$  are, for some purposes, better modeled in differences.

#### 4. PROPERTIES OF ESTIMATORS IN THE OVERDIFFERENCED PROCESS

##### 4.1. Root- $T$ consistency

Let  $\{w_{t|p}\}$  be the true AR( $p$ ) process  $\phi(B)w_{t|p} = a_t$ . This process follows the Markovian representation  $W_{t|p} = A_p W_{t-1|p} + U_{t,p}$ . The  $p \times p$  matrix  $A_p$  has the same structure as  $A_0$  with the coefficients  $(\phi_1, \dots, \phi_p)$  in the first row and  $W_{t|p} = (w_{t|p}, \dots, w_{t-p+1|p})'$ . Then, from (2.3),

$$w_t = \phi^{-1}(B) \left\{ 1 - \frac{(1 - \rho)B}{1 - \rho B} \right\} a_t = w_{t|p} - \sum_{j=0}^{\infty} \psi_j (1 - \rho) z_{t-1-j} \tag{4.1}$$

where  $\psi_j$  are the coefficients of  $\phi^{-1}(B)$ , and  $(1 - \rho B)z_t = a_t$ . Let us denote  $\Gamma_{w|p} = E(W_{t|p} W_{t|p}')$  and  $\gamma_{w|p} = E(W_{t|p} w_{t+1|p})$ . We define the sampling autocovariances as  $\hat{\Gamma}_{w|p} = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} W_{j|p} W_{j|p}'$ ,  $\hat{\gamma}_{w|p} = (T - p - 1)^{-1} \sum_{j=p+1}^{T-1} W_{j|p} w_{j+1|p}$ , and also make the following assumption.

A4.  $E(\|\hat{\Gamma}_{w|p}^{-1}\|^{2k})$  ( $k = 1, 2, \dots, k_0$ ) is bounded for all finite and sufficiently large  $T$  and some  $k_0$ .

The distance between the sampling second-order moments of  $w_t$  and  $w_{t|p}$  is determined in the following theorem.

**THEOREM 1.** *Let  $\{w_t\}$  be the process (2.3) and let  $w_1, \dots, w_T$  be a sample from this process. Let  $\rho$  be defined as in (3.3) with  $\beta \geq 1$ . Then*

$$\begin{aligned} (a) \quad & \hat{\Gamma}_w = \hat{\Gamma}_{w|p} + O_p(T^{-1/2}) \\ (b) \quad & \hat{\gamma}_w = \hat{\gamma}_{w|p} + O_p(T^{-1/2}). \end{aligned}$$

The proof is given in the Appendix. Since  $w_{t|p}$  is a stationary process, then  $\hat{\gamma}_{w|p} = \gamma_{w|p} + O_p(T^{-1/2})$ . Applying this result and Theorem 1, the following corollary holds.

**COROLLARY 1.** *Assume the conditions of Theorem 1 hold. Then*

$$\begin{aligned} (a) \quad & \hat{\Gamma}_w = \Gamma_{w|p} + O_p(T^{-1/2}) \\ (b) \quad & \hat{\gamma}_w = \gamma_{w|p} + O_p(T^{-1/2}). \end{aligned}$$

We can now prove root- $T$  consistency of  $\hat{\phi}$ . See the proof in the Appendix.

**THEOREM 2.** *Assume the conditions of Theorem 1 hold. Then*

$$\hat{\phi} = \phi + O_p(T^{-1/2}).$$

#### 4.2. Bias and mean squared error

Let  $\hat{\phi}_{|p}$  be the least squares estimator of  $\phi$  from a sample from a true  $AR(p)$  process. The bias and mean squared error (MSE) of this estimator, of a properly specified autoregression, have been widely investigated (see, for instance, Bhansali, 1981; Kunitomo and Yamamoto, 1985; Shaman and Stine, 1988; and references therein). Since the similarity between the estimator  $\hat{\phi}$ , of the  $ARIMA(p+1, 1, 1)$  misspecified as an  $AR(p)$ , and  $\hat{\phi}_{|p}$  depends on the near nonstationarity hypothesis, we will express their differences in terms of  $\rho$ . The following theorems formulate the first- and second-order moments of the least squares estimator  $\hat{\phi}$  around the true parameter  $\phi$  as the respective moments of  $\hat{\phi}_{|p}$  plus an error term depending on  $\rho$ .

**THEOREM 3.** *Assume A1 (with  $s_0 = 8$ ), A2, A3 and A4. Then*

$$E(\hat{\phi} - \phi) = E(\hat{\phi}_{|p} - \phi) + O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{1/2}\right\}. \quad (4.2)$$

The proof is given in the Appendix. Since  $(1-\rho)/(1+\rho) = O(T^{-\beta})$  and given that  $E(\hat{\phi}_{|p} - \phi) = O(T^{-1})$  (see, for instance, Bhansali, 1981) we need a value

$\beta > 2$  for the biases to be equal up to terms of order  $O(T^{-1})$ , whereas for root- $T$  consistency we only need  $\beta \geq 1$ .

**THEOREM 4.** *Assume A1 (with  $s_0 = 8$ ), A2, A3 and A4. Then*

$$E\{(\hat{\phi} - \phi)(\hat{\phi} - \phi)'\} = E\{(\hat{\phi}_{|p} - \phi)(\hat{\phi}_{|p} - \phi)'\} + O\left[\max\left\{\left(\frac{1-\rho}{1+\rho}\right)^{1/2} T^{-1/2}, \frac{1-\rho}{1+\rho}\right\}\right].$$

See the proof in the Appendix. We can see from this theorem that the MSEs are closer to each other than the biases. If  $\rho$  is such that  $\beta > 1$  then the two expressions for the MSE are equal up to terms  $O(T^{-1})$ .

5. MSE OF  $H$ -STEPS-AHEAD PREDICTION

In this section, we obtain the MSE of predicting  $y_{T+H}$  for  $t = T$ . The PMSE of a properly specified autoregression is (see, for instance, Kunitomo and Yamamoto, 1985)

$$\begin{aligned} \text{PMSE}(\check{y}_{T+H}) &= \sigma^2 \sum_{h=0}^{H-1} (e'_{p+2} A_\alpha^h e_{p+2})^2 + \frac{\sigma^2}{T} \sum_{h=0}^{H-1} \sum_{k=0}^{H-1} (e'_{p+2} A_\alpha^h e_{p+2})(e'_{p+2} A_\alpha^k e_{p+2}) \\ &\quad \times \text{tr}(A_\alpha^{H-1-h} \Gamma_y A_\alpha'^{H-1-k} \Gamma_y^{-1}) + O(T^{-3/2}). \end{aligned} \tag{5.1}$$

This expression is inconvenient however, to compare the PMSE of the  $AR(p + 1)$  model ( $\text{PMSE}(\check{y}_{T+H})$ ) with the PMSE of the misspecified  $ARIMA(p, 1, 0)$  model ( $\text{PMSE}(\check{y}_{T+H})$ ). We will rewrite the estimated  $H$ -steps-ahead predictions in terms of their estimated increments ( $\check{w}_t$  and  $\hat{w}_t$ , respectively). Hence,  $\text{PMSE}(\check{y}_{T+H}) = \sum_{h=1}^H \text{PMSE}(\check{w}_{T+h}) + 2 \sum_{h=1}^H \sum_{k=h+1}^H E\{(w_{T+h} - \check{w}_{T+h})(w_{T+k} - \check{w}_{T+k})\}$ , where  $\check{w}_t = \check{y}_t - \check{y}_{t-1}$ . A similar expression applies for  $\text{PMSE}(\hat{y}_{T+H})$ .

5.1. PMSE of the properly specified  $AR(p + 1)$  predictor

Let  $\check{A}_\alpha$  be the least squares estimator of  $A_\alpha$  using the properly specified model (2.2). The estimated increment  $\check{w}_{T+h}$  defined as a function of the estimated coefficients  $\check{A}_\alpha$  is

$$\check{w}_{T+h} = e'_{p+2} \check{A}_\alpha^{h-1} (\check{A}_\alpha - I_{p+2}) Y_T \tag{5.2}$$

where  $I_{p+2}$  is the identity matrix. The observed value  $w_{T+h}$  is

$$w_{T+h} = e'_{p+2} A_\alpha^{h-1} (A_\alpha - I_{p+2}) Y_T + L_h$$

where  $L_h = L_1 - L_2$ , with  $L_1 = \sum_{k=0}^{h-1} e'_{p+2} A_\alpha^k U_{T+h-k, p+2}$  and  $L_2 = \sum_{k=1}^{h-1} e'_{p+2} A_\alpha^{k-1} U_{T+h-k, p+2}$ .

The  $\text{PMSE}(\tilde{w}_{T+h})$  and  $E\{(\tilde{w}_{T+h} - w_{T+h})(\tilde{w}_{T+k} - w_{T+k})\}$  are shown in the following theorem (see the proof in the Appendix). The assumptions about  $s_0$  in Theorems 5 and 6 are needed in order to apply the results of Kunitomo and Yamamoto (1985) in the proof of the theorems.

**THEOREM 5.** *Let  $w_t$  follow (2.3), where  $\rho = \exp(-c/T^\beta)$  and  $\beta > 1$ . Assume  $A_2, A_3, A_4$  and  $A_1$  with  $s_0 = 32$  when  $h = 1, 2$ , and  $s_0 = 16h$  when  $h \geq 3$ . Then*

$$\begin{aligned} \text{PMSE}(\tilde{w}_{T+h}) &= \sigma^2 \sum_{j=0}^{h-1} (e'_{p+2} A_1^j c_{p+2})^2 + \frac{\sigma^2}{T} \sum_{j=0}^{h-1} \sum_{k=0}^{h-1} (e'_p A_p^j e_p)(e'_p A_p^k e_p) \\ &\quad \times \text{tr}(A_\alpha^{h-1-j} \Gamma_y A_\alpha'^{h-1-k} \Gamma_y^{-1}) + O(T^{-3/2}) \end{aligned} \quad (5.3)$$

and, for  $k \geq h$ ,

$$\begin{aligned} E\{(\tilde{w}_{T+h} - w_{T+h})(\tilde{w}_{T+k} - w_{T+k})\} &= \sigma^2 \sum_{i=0}^{h-1} (e'_{p+2} A_1^i c_{p+2})(e'_{p+2} A_1^{i+(k-h)} c_{p+2}) \\ &\quad + \frac{\sigma^2}{T} \sum_{n=0}^{k-1} \sum_{i=0}^{h-1} (e'_p A_p^n e_p)(e'_p A_p^i e_p) \\ &\quad \times \text{tr}(A_\alpha^{h-1-i} \Gamma_y A_\alpha'^{k-1-n} \Gamma_y^{-1}) + O(T^{-3/2}) \end{aligned} \quad (5.4)$$

where  $c_{p+2} = (1, 0, \dots, 0, 1)'$ .

The terms on the right-hand side of (5.3) and (5.4) have two components. The first component includes the variance of the prediction errors and the covariance between prediction errors at different horizons, respectively, of the noninvertible  $\text{ARMA}(p+1, 1)$  process. The second component is the sampling error, due to the estimation of the  $p+2$  parameters of the vector  $\phi$ .

## 5.2. PMSE of the overdifferenced $\text{ARIMA}(p, 1, 0)$ predictor

Assume that we predict  $w_{T+h}$  with the predictor derived from the estimated  $\text{AR}(p)$ , i.e.  $\hat{w}_{T+h} = e'_p \hat{A}_p^h W_T$ , where  $\hat{A}_p$  is the least squares estimator of  $A_p$ . Then

$$\hat{w}_{T+h} = e'_p A_p^h W_T + e'_p (\hat{A}_p^h - A_p^h) W_T = E(w_{T+h}|T) + e'_p (\hat{A}_p^h - A_p^h) W_T.$$

The true value  $w_{T+h}$  is, from (2.4),  $w_{T+h} = e'_{p+2} A_1^h Z_T + L_h = E(w_{T+h}|T) + L_h$ .



Then the  $h$ -steps-ahead prediction error is  $w_{T+h} - \hat{w}_{T+h} = L_h - e'_p(\hat{A}_p^h - A_p^h)W_T - v_t$ , where, by (4.1),

$$v_t = E(w_{T+h} - \hat{w}_{T+h}|p|T) = \sum_{j=h-1}^{\infty} \psi_j(1 - \rho)z_{T+h-1-j} + \sum_{j=0}^{h-2} \psi_j(1 - \rho)\rho^{h-1-j}z_T. \tag{5.5}$$

The following theorem gives an approximation of order  $o(T^{-1})$  of the expectation of the lead- $h$  mean squared prediction error (see the proof in the Appendix).

**THEOREM 6.** *Let  $w_t$  follow (2.3), where  $\rho = \exp(-c/T^\beta)$  and  $\beta > 1$ . Assume A2, A3, A4 and A1 with  $s_0 = 32$  when  $h = 1, 2$  and  $s_0 = 16h$  when  $h \geq 3$ . Then*

$$\begin{aligned} \text{PMSE}(\hat{w}_{T+h}) &= \sigma^2 \sum_{k=0}^{h-1} (e'_{p+2} A_1^k c_{p+2})^2 + \frac{\sigma^2}{T} \sum_{j=0}^{h-1} \sum_{k=0}^{h-1} (e'_p A_p^j e_p)(e'_p A_p^k e_p) \\ &\quad \times \text{tr}(A_p^{h-1-j} \Gamma_{w|p} A_p'^{h-1-k} \Gamma_{w|p}^{-1}) + o(T^{-1}) \end{aligned} \tag{5.6}$$

and, for  $k \geq h$ ,

$$\begin{aligned} E\{(\hat{w}_{T+h} - w_{T+h})(\hat{w}_{T+k} - w_{T+k})\} &= \sigma^2 \sum_{i=0}^{h-1} (e'_{p+2} A_1^i c_{p+2})(e'_{p+2} A_1^{i+(k-h)} c_{p+2}) \\ &\quad + \frac{\sigma^2}{T} \sum_{n=0}^{k-1} \sum_{i=0}^{h-1} (e'_p A_p^n e_p)(e'_p A_p^i e_p) \\ &\quad \times \text{tr}(A_p^{h-1-i} \Gamma_{w|p} A_p^{k-1-n} \Gamma_{w|p}^{-1}) + o(T^{-1}) \end{aligned} \tag{5.7}$$

where  $c_{p+2} = (1, 0, \dots, 0, 1)'$ .

The terms on the right-hand side of (5.6) and (5.7) have two components. The first one, the variance of the prediction errors and their covariance between different horizons of the true ARIMA( $p + 1, 1, 1$ ) process, is the same as in Theorem 5. The second one is the sampling error due to the estimation of the  $p$  parameters  $\phi$ , in contrast with the estimation of the  $p + 2$  parameters of the AR( $p + 1$ ) model. It should be observed that this second component differs from that in the previous subsection only in the elements inside the trace operators.

6. COMPARING PREDICTION ACCURACY

In this section, we compare the PMSEs found in the last section for the two models. We prove that, under the assumption of near nonstationarity exposed in (3.3), with  $\beta > 1$ , overdifferencing may produce lower PMSE (to terms of small order). The expressions in Theorems 5 and 6 reveal that the only difference between  $\text{PMSE}(\hat{y}_{T+H})$  and  $\text{PMSE}(\hat{y}_{T+H})$  is in the elements inside the trace operators. These traces can be compared using the following two lemmas: Lemma 1 compares such a trace in processes with and without a constant term; Lemma 2 compares the trace in nearly nonstationary processes with no constant term and the trace in the overdifferenced model. The proofs of these lemmas can be found in the Appendix.

LEMMA 1. *Let  $y_t$  follow process (2.1). Then*

$$\text{tr}(A_\alpha^i \Gamma_y A_\alpha'^j \Gamma_y^{-1}) = 1 + \text{tr}(A_0^i \Gamma_{\bar{y}} A_0'^j \Gamma_{\bar{y}}^{-1}).$$

LEMMA 2. *Let  $y_t$  follow process (2.1) with  $\rho = \exp(-c/T^\beta)$  and  $\beta > 1$ . Then*

$$\text{tr}(A_0^i \Gamma_{\bar{y}} A_0'^j \Gamma_{\bar{y}}^{-1}) = \rho^{i+j} + \text{tr}(A_p^i \Gamma_{w|p} A_p^j \Gamma_{w|p}^{-1}) + o(T^{-1}).$$

Now we can prove the advantage of overdifferencing when the process is nearly nonstationary.

THEOREM 7. *Let  $y_t$  follow process (2.1) with  $\rho = \exp(-c/T^\beta)$  and  $\beta > 1$ , and let the conditions of Theorems 5 and 6 hold. Then, for  $H \geq 1$ ,*

$$\text{PMSE}(\check{y}_{T+H}) - \text{PMSE}(\hat{y}_{T+H}) = \nu_H + o(H^2 T^{-1})$$

where

$$\nu_H = \frac{\sigma^2}{T} \left( \sum_{h=1}^H \sum_{j=0}^{h-1} \psi_j \right)^2 + \frac{\sigma^2}{T} \left( \sum_{h=1}^H \sum_{j=0}^{h-1} \psi_j \rho^{h-1-j} \right)^2 > 0 \tag{6.1}$$

with  $\psi_j = (e_p' A_p^j e_p)$ ,  $j = 1, \dots, H$ .

The proof is a direct application of Lemmas 1 and 2 to the differences between (5.3) and (5.6) and between expressions (5.7) and (5.4).

Expression (6.1) shows that the advantage of the overdifferenced model can be decomposed into two parts. The first term at the right-hand side of (6.1) is the result of applying Lemma 1 and therefore is due to the MSE of estimating the constant term  $\alpha$  in the  $\text{AR}(p+1)$  model. The second term is the result of applying Lemma 2 and then is due to the MSE of estimating an extra parameter in the  $\text{AR}(p+1)$ . Thus, the superior forecasting performance of the model  $\text{ARIMA}(p, 1, 0)$  is due to its more parsimonious representation. For  $H = 1$  the difference is  $2\sigma^2/T$  if a constant is needed, and  $\sigma^2/T$  if  $\alpha = 0$  and

no constant is estimated. This result is similar to that of Ledolter and Abraham (1981) for overspecified models, where they state that each unnecessary estimated parameter increases the one-step-ahead PMSE by  $\sigma^2/T$ .

Although these results are applicable to a general stationary autoregression, it is interesting to analyze the AR(1) case. First, its simplicity avoids the use of some asymptotic approximations. Second, the results will not be affected by any other root, as shown in (3.2), and they can be considered as a neutral benchmark. The PMSE of the proper predictor in this case can be evaluated with (5.1), whereas the PMSE in the overdifferenced model is easily evaluated using a random walk as predictor. The following remarks summarize the results for both the AR(1) case with no intercept (AR(1)) and the AR(1) case with intercept (AR(1,  $\mu$ )).

REMARK 1. Let  $y_t$  follow the process  $y_t = \rho y_{t-1} + a_t$ ,  $|\rho| < 1$ . Then  $PMSE(\hat{y}_{T+H}) - PMSE(\hat{y}_{T+H}) = \nu_{H|AR(1)} + o(H^2 T^{-3/2})$ , where

$$\nu_{H|AR(1)} = \sigma^2 \left\{ \frac{H^2 \rho^{2(H-1)}}{T} - \frac{(1 - \rho^H)^2}{1 - \rho^2} \right\}. \tag{6.2}$$

Table I shows the values of  $\rho$  that make  $\nu_{H|AR(1)} = 0$ . Larger values will produce  $\nu_{H|AR(1)} > 0$ . These values of  $\rho$  increase with  $H$ . Therefore, as the horizon grows, the process needs to be closer to the unit root in order to get some gain when differencing. The advantage of overdifferencing tends, then, to decrease when the horizon is large. It can also be seen that as  $H \rightarrow \infty$  the limit of (6.2) is negative. Then, the advantage of the overdifferenced predictor eventually disappears. If  $\rho$  is close enough to unity, this will happen at a horizon of no practical interest. This result has an interpretation in terms of the mean reversion of the true process. Since the process is stationary, its long-term prediction is the unconditional mean, which in this case is known. Therefore, the AR(1) predictor will forecast the long term with no error, whereas the random walk will not. Manipulating (6.2), we can conclude that, up to terms of small order, overdifferencing can produce better forecasts if

TABLE I  
VALUES OF  $\rho$  TO OBTAIN  $\nu_{H|AR(1)} = 0$  AND  $\nu_{H|AR(1, \mu)} = 0$

T	AR(1), Horizon					AR(1, $\mu$ ), Horizon				
	1	2	5	10	20	1	2	5	10	20
25	0.923	0.937	0.940	0.951	0.963	0.852	0.862	0.881	0.898	0.913
50	0.961	0.965	0.966	0.970	0.976	0.923	0.926	0.932	0.940	0.948
75	0.974	0.976	0.976	0.978	0.982	0.948	0.949	0.953	0.957	0.962
100	0.980	0.981	0.982	0.983	0.985	0.961	0.962	0.964	0.966	0.970
150	0.987	0.987	0.987	0.988	0.989	0.974	0.974	0.975	0.976	0.978
300	0.993	0.994	0.994	0.994	0.994	0.987	0.987	0.987	0.988	0.988

$$\rho > \exp\left(-\frac{2}{T+4H}\right). \quad (6.3)$$

This expression can be approximated, omitting the influence of  $H$ , as  $\rho > \exp(-2/T)$ . This value of  $c = 2$  agrees with the empirical work of Stock (1996).

REMARK 2. Let  $y_t$  follow the process  $y_t = \alpha + \rho y_{t-1} + a_t$ ,  $|\rho| < 1$ . Then  $\text{PMSE}(\check{y}_{T+H}) - \text{PMSE}(\hat{y}_{T+H}) = v_{H|\text{AR}(1,\mu)} + o(H^2 T^{-3/2})$ , where

$$v_{H|\text{AR}(1,\mu)} = \sigma^2 \left\{ \frac{H^2 \rho^{2(H-1)}}{T} + \frac{(1-\rho^H)^2}{T(1-\rho)^2} - \frac{(1-\rho^H)^2}{1-\rho^2} \right\}. \quad (6.4)$$

Table I shows the values of  $\rho$  that make  $v_{H|\text{AR}(1,\mu)} = 0$ . From (6.4) it can be verified that the overdifferenced predictor produces better forecasts, up to terms of small order, if

$$\rho > \exp\left(-\frac{4}{T+4H}\right) \quad (6.5)$$

which can be simplified as  $\rho > \exp(-4/T)$ . In this case, the limit of (6.4) as  $H \rightarrow \infty$  is still positive if  $\rho > \exp(-2/T)$ .

## 7. A SIMULATION STUDY

In this section, we illustrate the preceding results with a simulation exercise. We consider three different AR(2) models: M1,  $(1 - 0.5B)(1 - \rho B)y_t = 10 + a_t$ ; M2,  $(1 - 0.5B)(1 - \rho B)y_t = a_t$ ; and M3,  $(1 + 0.8B)(1 - \rho B)y_t = 10 + a_t$ , with  $\rho = 0.9, 0.92, 0.94, 0.96, 0.98, 0.99$ . Sample sizes are  $T = 50, 100$ . Real series usually have non-zero mean, and models M1 and M3 can illustrate the consequences of overdifferencing in such series. Also, model M2 can arise when in doubt about a second difference.

An important aspect in the simulation exercise is the possibility of obtaining an explosive estimated predictor. There are two main reasons to avoid these explosive situations. First, they are of limited practical interest. A typical situation where a practitioner has doubts about differencing, for forecasting purposes, deals mainly with estimated  $\rho$  close to, but lower than, unity. Second, the explosive nature of the predictions generated with a predictor with  $\hat{\rho} > 1$  produces an excessive influence on the averages resulting from the simulations, because the explosive estimated predictor is easily worse than its overdifferenced counterpart, especially in the long term. Unreported simulations show that very few explosive estimated predictors can have an extremely large influence in the computations, given a too optimistic representation of the effect of overdifferencing. Therefore, in order to obtain a clearer picture of what can be expected from overdifferencing in a real situation, we have considered only

those replications in which estimated roots were outside the unit circle. The percentage of rejected replications is low. For instance, if  $\rho = 0.98$  and  $T = 100$  it is 1%, and with  $T = 50$  it is 2.7%.

In each replication, we generate a random sample of the process of size  $500 + T + 30$  with random noise  $a_t \sim N(0, 1)$ . The first 500 observations were ignored to avoid the effect of initial values, and the last 30 were used to evaluate the prediction error. By averaging the prediction squared errors of 20 000 valid replications we obtain  $V_y(H)$  and  $V_w(H)$  as the sampling estimation of the PMSE of forecasting  $y_{T+H}$  using the forecasts generated by the correct AR(2) model or the overdifferenced ARIMA(1, 1, 0) model respectively. Figures 1–3 show the ratio  $\{V_y(H) - V_w(H)\}/V_y(H)$  for M1 and M3 as a function of  $T$  and  $\rho$ . This ratio represents the empirical expected gain (or loss if negative) of overdifferencing at each horizon. The figures reveal

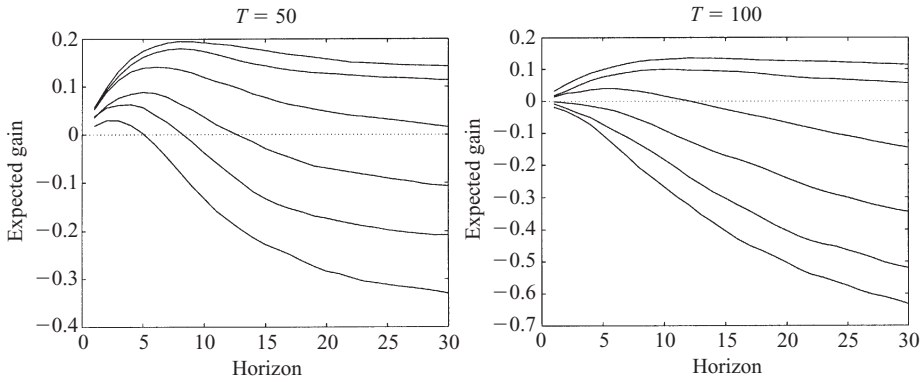


FIGURE 1.  $\{V_y(H) - V_w(H)\}/V_y(H)$  of Model M1 for horizon  $H = 1, \dots, 30$  and sample size  $T$ . The values of  $\rho$  are (from bottom to top) 0.90, 0.92, 0.94, 0.96, 0.98, 0.99.

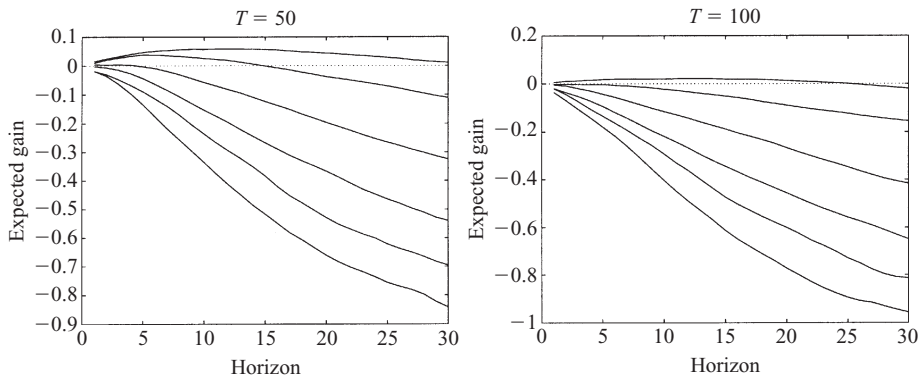


FIGURE 2.  $\{V_y(H) - V_w(H)\}/V_y(H)$  of Model M2 for horizon  $H = 1, \dots, 30$  and sample size  $T$ . The values of  $\rho$  are (from bottom to top) 0.90, 0.92, 0.94, 0.96, 0.98, 0.99.

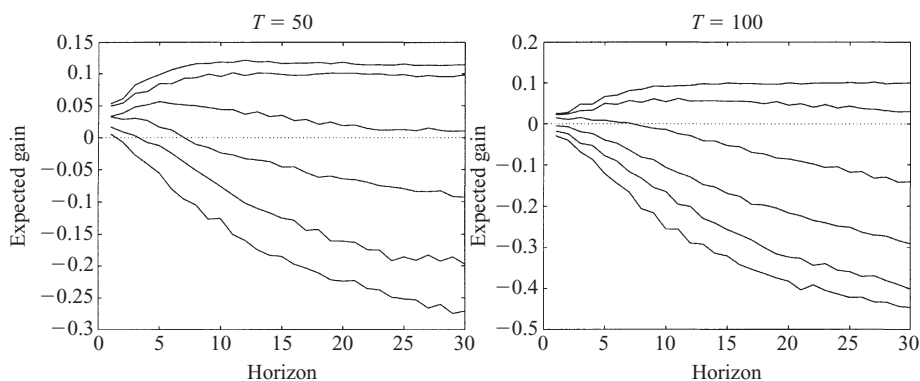


FIGURE 3.  $\{V_y(H) - V_w(H)\}/V_y(H)$  of Model M3 for horizon  $H = 1, \dots, 30$  and sample size  $T$ . The values of  $\rho$  are (from bottom to top) 0.90, 0.92, 0.94, 0.96, 0.98, 0.99.

that, as expected from the theoretical results, there are situations where overdifferencing outperformed the true model. The expected gain increases with the size of  $\rho$  and decreases with  $T$ . Also, in agreement with Equation (3.2), the gain is larger in the model with positive second root (M1) than in the model with negative root (M3). The gain substantially decreases if  $\alpha = 0$  (M2).

The main feature of these figures is the divergence of the curves as the horizon increases. In the very short term, the difference between the two predictors is very small, even negligible. Nevertheless, in the medium or long term the gain or loss can be important. The risk of falling into an important loss if  $\rho$  is not large enough can be diminished, however, if some efficient rule to decide about differencing is used. A second important aspect of these figures is that in the long run ( $H \gg T^{1/2}$ ) the gain decreases and can be negative. Also, as proved in the last section, the gain in the model with no constant always disappears at sufficiently large  $H$ .

Figures 4 and 5 show the absolute values of  $V_y(H)$  and  $V_w(H)$  for selected values of  $\rho$ . These figures also contain the population PMSE of the process (dotted lines). These population values can be obtained from the first term on the right-hand side of expression (5.1). The distance from these population curves to each solid line is the PMSE due to the estimation of the unknown parameters. It can be seen that the sampling variability of the nondifferenced predictor (line with symbol  $+$ ) increases noticeably when the number of parameters increases (models M1 and M3 with respect to M2). This increment of the PMSE due to the estimation of the parameters causes that the overdifferenced predictor (line with symbol  $\circ$ ) can outperform its competitor when the process approaches nonstationarity.

It can be seen that the theoretical results accurately explain this finite sample performance. Since these results depend mainly on the size of the roots rather than on their number, it is reasonable to foresee similar conclusions in larger autoregressions.

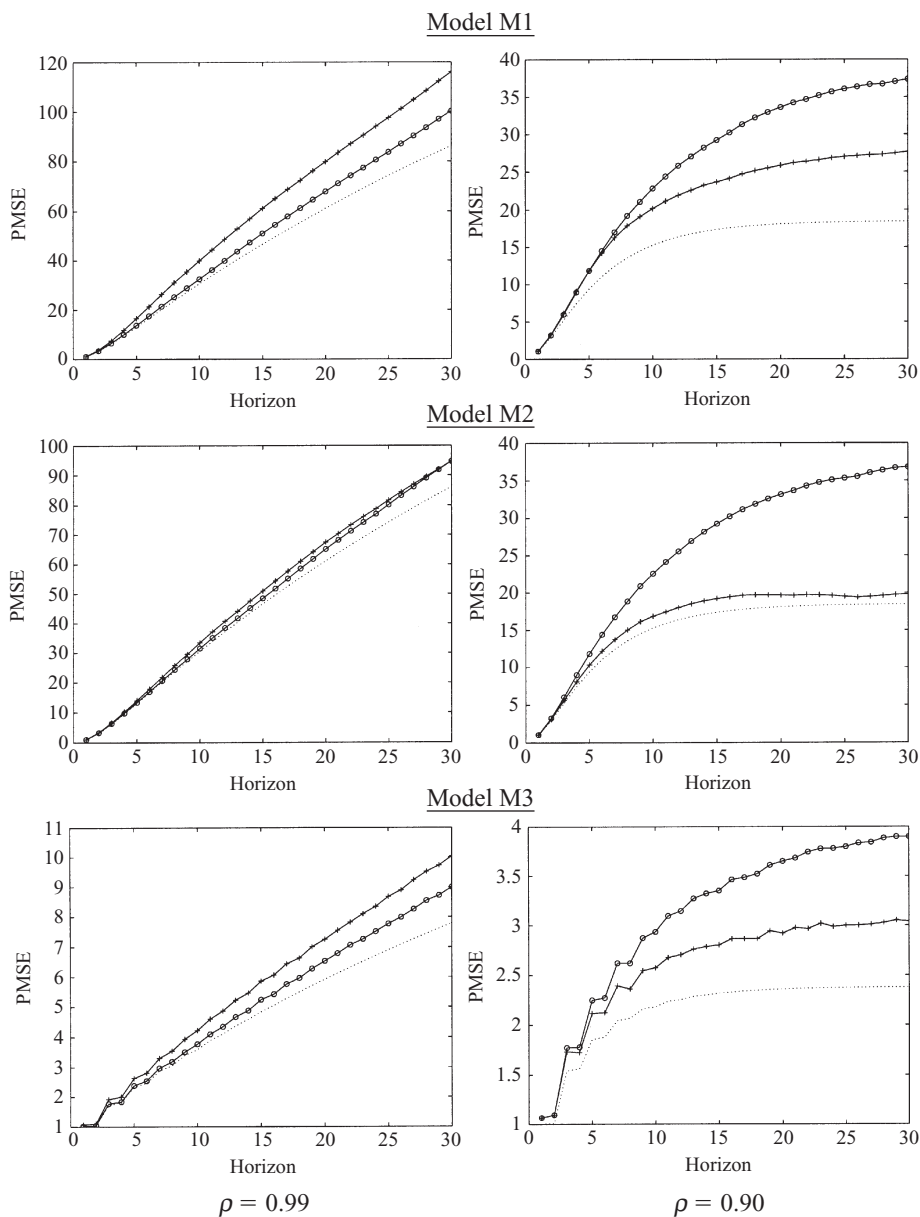


FIGURE 4. Values of  $V_y$  (+),  $V_w$  (o) and population PMSE (.....). Sample size  $T = 50$ .

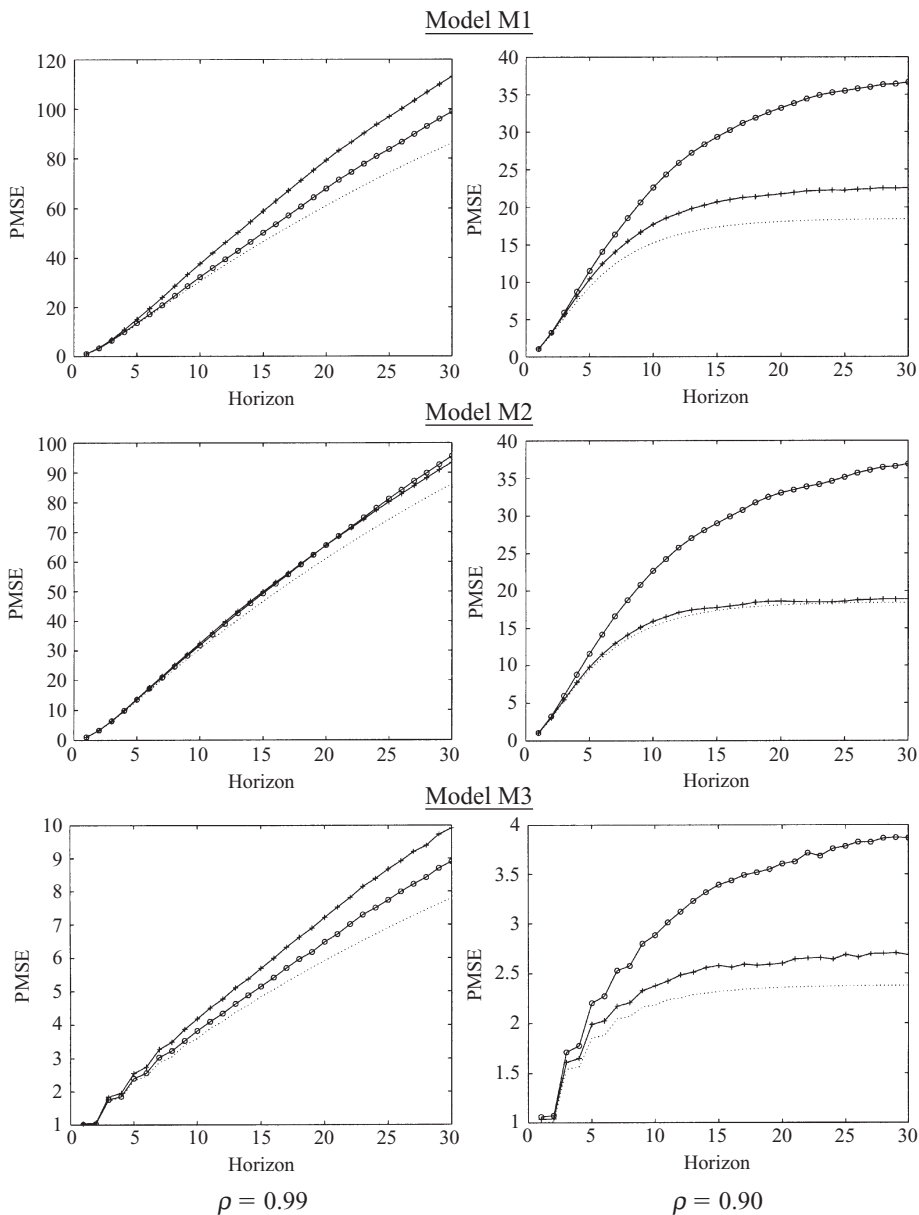


FIGURE 5. Values of  $V_y$  (+),  $V_w$  (o) and population PMSE (.....). Sample size  $T = 100$ .



APPENDIX

*Lemmas*

We present some lemmas used for the proof of theorems in subsequent sections. For an arbitrary  $p \times 1$  vector  $x$  and a  $p \times p$  matrix  $M$ , let  $\|x\| = (x'x)^{1/2}$  be the Euclidean norm of  $x$  and  $\|M\| = \sup_{\|x\| \leq 1} (x'M'x)^{1/2}$  be the matrix norm of  $M$ .

LEMMA A1. Assume A1 and A2, with  $s_0 = 2k$  and  $k \geq 1$ . Then, as  $T \rightarrow \infty$ ,

$$E(\|\hat{\Gamma}_w - \hat{\Gamma}_{w|p}\|^k) = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}\right\}$$

and

$$E(\|\hat{\gamma}_w - \hat{\gamma}_{w|p}\|^k) = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}\right\}. \tag{A1}$$

PROOF. Let  $m_{ij}$  be a generic element of  $M$ . Since  $E(\|M\|^k) = O\{\max_{i,j} E(|m_{ij}|^k)\}$ ,  $i, j = 1, \dots, p$ , and by Minkowski's inequality,  $E(\|\hat{\Gamma}_w - \hat{\Gamma}_{w|p}\|^k) = O(\max_{t,s} E|w_t w_{t-s} - w_{t|p} w_{t-s|p}|^k)$ . A similar result applies to (A1). Using the decomposition (4.1), and by Minkowski's inequality,

$$E|w_t w_{t-s} - w_{t|p} w_{t-s|p}|^k \leq \{(E|w_{t|p} r_{t-s}|^k)^{1/k} + (E|w_{t-s|p} r_t|^k)^{1/k} + (E|r_t r_{t-s}|^k)^{1/k}\}^k \tag{A2}$$

where  $r_t = \sum_{j=0}^{\infty} \psi_j(1-\rho)z_{t-1-j}$ . By Hölder's inequality,  $E|w_{t|p} r_{t-s}|^k \leq (E|w_{t|p}|^{2k} E|r_{t-s}|^{2k})^{1/2}$ . Also, by Assumption A1,  $E|w_{t|p}|^{2k} = O(1)$ . Similarly,  $E|r_{t-s}|^{2k} \leq \{\sum_{j=0}^{\infty} |\psi_j(1-\rho)|(E|z_{t-s-1-j}|^{2k})^{1/2k}\}^{2k}$ , where it can be verified that

$$E|z_t|^{2k} \leq E\left(\sum_{j=0}^{\infty} |\rho^{2j} a_{t-j}^2|\right)^k \leq \left\{\sum_{j=0}^{\infty} (E|\rho^{2j} a_{t-j}^2|^k)^{1/k}\right\}^k = O\left\{\left(\frac{1}{1-\rho^2}\right)^k\right\}.$$

Therefore,

$$E|r_{t-s}|^{2k} = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^k\right\} \tag{A3}$$

and hence

$$E|w_{t|p} r_{t-s}|^k = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}\right\}.$$

A similar result applies to the second term in (A2). The third term in (A2) can also be solved following the previous arguments. It can be shown that

$$E|r_t r_{t-s}|^k = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^k\right\}.$$

Applying these results to (A2) proves the lemma.

LEMMA A2. Assume A1, A2, A3 and A4, with  $s_0 = 2k$ . Then, as  $T \rightarrow \infty$ ,

$$E(\|\hat{\Gamma}_w^{-1} - \hat{\Gamma}_{w|p}^{-1}\|^k) = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}\right\}.$$

PROOF. It can be verified that  $\|\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1}\|^k = \|\hat{T}_w^{-1}(\hat{T}_w - \hat{T}_{w|p})\hat{T}_{w|p}^{-1}\|^k$ . By Hölder's inequality and Lemma A1 the result follows.

LEMMA A3. Assume A1, A2, A3 and A4, with  $s_0 = 4k$ . Then, as  $T \rightarrow \infty$ ,

$$E(\|\hat{\phi} - \hat{\phi}_{|p}\|^k) = O\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}\right\} \quad (\text{A4})$$

$$E(\|\hat{\phi} - \phi\|^k) = O\left[\max\left\{\left(\frac{1-\rho}{1+\rho}\right)^{k/2}, T^{-k/2}\right\}\right]. \quad (\text{A5})$$

PROOF. The estimator  $\hat{\phi}$  can be expressed as

$$\hat{\phi} = (\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1})(\hat{\gamma}_w - \hat{\gamma}_{w|p}) + (\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1})\hat{\gamma}_{w|p} + \hat{T}_{w|p}^{-1}(\hat{\gamma}_w - \hat{\gamma}_{w|p}) + \hat{\phi}_{|p}$$

where  $\hat{\phi}_{|p} = \hat{T}_{w|p}^{-1}\hat{\gamma}_{w|p}$ . By Minkowski's inequality we obtain

$$\begin{aligned} E(\|\hat{\phi} - \hat{\phi}_{|p}\|^k) &\leq (E\{\|(\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1})(\hat{\gamma}_w - \hat{\gamma}_{w|p})\|^k\})^{1/k} \\ &\quad + [E\{\|(\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1})\hat{\gamma}_{w|p}\|^k\}]^{1/k} + [E\{\|\hat{T}_{w|p}^{-1}(\hat{\gamma}_w - \hat{\gamma}_{w|p})\|^k\}]^{1/k}. \end{aligned}$$

By Hölder's inequality and applying Lemmas A1 and A2, expression (A4) holds. In order to prove (A5) we use the decomposition  $\hat{\phi} - \phi = \Gamma_{w|p}^{-1}(\hat{\gamma}_w - \gamma_{w|p}) + (\hat{T}_w^{-1} - \Gamma_{w|p}^{-1})\hat{\gamma}_w$ , and also the decompositions  $\hat{\gamma}_w - \gamma_{w|p} = (\hat{\gamma}_w - \hat{\gamma}_{w|p}) + (\hat{\gamma}_{w|p} - \gamma_{w|p})$  and  $\hat{T}_w^{-1} - \Gamma_{w|p}^{-1} = (\hat{T}_w^{-1} - \hat{T}_{w|p}^{-1}) + (\hat{T}_{w|p}^{-1} - \Gamma_{w|p}^{-1})$ . Applying that  $E(\|\hat{\gamma}_{w|p} - \gamma_{w|p}\|^{2k}) = O(T^{-k})$  and  $E(\|\hat{T}_{w|p}^{-1} - \Gamma_{w|p}^{-1}\|^{2k}) = O(T^{-k})$  (see, for instance, Lemma 3.3 of Bhansali, 1981), and using the same arguments as before, completes the result.

#### Proofs of results in Section 4

PROOF OF THEOREM 1. Since  $E(z_t^2) = \sigma^2/(1-\rho^2)$ , and by Chebyshev's inequality, we obtain  $z_t = O_p\{(1-\rho^2)^{-1/2}\}$ . Hence,

$$r_t = O_p\left\{\left(\frac{1-\rho}{1+\rho}\right)^{1/2}\right\}. \quad (\text{A6})$$

Since  $(1-\rho)/(1+\rho) = O(T^{-\beta})$ , then  $w_t = w_{t|p} + o_p(T^{-1/2})$ .

The elements of  $\hat{T}_w$  and  $\hat{\gamma}_w$  can be decomposed as

$$\begin{aligned} \frac{\sum_{j=p+1}^{T-1} w_{j-t} w_{j-s}}{T-p-1} &= \frac{\sum_{j=p+1}^{T-1} w_{j-t|p} w_{j-s|p}}{T-p-1} - \frac{\sum_{j=p+1}^{T-1} w_{j-t|p} r_{j-s}}{T-p-1} \\ &\quad - \frac{\sum_{j=p+1}^{T-1} w_{j-s|p} r_{j-t}}{T-p-1} + \frac{\sum_{j=p+1}^{T-1} r_{j-t} r_{j-s}}{T-p-1}. \end{aligned}$$

Applying (A6) and the result that  $w_{t|p} = O_p(1)$ , it can be verified that

$$\frac{\sum_{j=p+1}^{T-1} w_{j-t} w_{j-s}}{T-p-1} = \frac{\sum_{j=p+1}^{T-1} w_{j-t|p} w_{j-s|p}}{T-p-1} + o_p(T^{-1/2})$$

and the theorem follows.

PROOF OF THEOREM 2. Using the decomposition  $\hat{\phi} - \phi = \Gamma_{w|p}^{-1}(\hat{\gamma}_w - \gamma_{w|p}) + (\hat{T}_w^{-1} - \Gamma_{w|p}^{-1})\hat{\gamma}_w$ , and by stationarity of  $\{w_{t|p}\}$ , we have  $\Gamma_{w|p}^{-1} = O(1)$ . Also, if

$\hat{\Gamma}_w^{-1}$  exists, we have  $\hat{\Gamma}_w^{-1} - \Gamma_{w|p}^{-1} = \hat{\Gamma}_w^{-1}(\Gamma_{w|p} - \hat{\Gamma}_w)\Gamma_{w|p}^{-1}$ . Therefore, applying Corollary 1,  $\hat{\phi} - \phi = O_p(T^{-1/2})$ .

PROOF OF THEOREM 3. It can be verified that  $E(\hat{\phi} - \hat{\phi}_{|p}) = E\{(\hat{\Gamma}_w^{-1} - \hat{\Gamma}_{w|p}^{-1})\hat{\gamma}_{w|p}\} + E\{\hat{\Gamma}_w^{-1}(\hat{\gamma}_w - \hat{\gamma}_{w|p})\}$ . Applying Hölder's inequality and Lemmas A2 and A1, the theorem follows.

POOF OF THEOREM 4. We can decompose

$$\begin{aligned} \text{MSE}(\hat{\phi}) &= \text{MSE}(\hat{\phi}_{|p}) + E\{(\hat{\phi} - \hat{\phi}_{|p})(\hat{\phi}_{|p} - \phi)'\} \\ &\quad + E\{(\hat{\phi}_{|p} - \phi)(\hat{\phi} - \hat{\phi}_{|p})'\} + E\{(\hat{\phi} - \hat{\phi}_{|p})(\hat{\phi}_{|p} - \phi)'\}. \end{aligned}$$

Since  $\|M\| \leq \{\text{tr}(M'M)\}^{1/2}$ , and applying Lemma A3,

$$E\{\|(\hat{\phi} - \hat{\phi}_{|p})(\hat{\phi} - \hat{\phi}_{|p})'\|\} \leq E(\|\hat{\phi} - \hat{\phi}_{|p}\|^2) = O\left(\frac{1 - \rho}{1 + \rho}\right).$$

Analogously, and applying the result that  $E(\|\hat{\phi}_{|p} - \phi\|^2) = O(T^{-1})$  (see, for instance, Bhansali, 1981), it can be verified that

$$\begin{aligned} E\{\|(\hat{\phi}_{|p} - \phi)(\hat{\phi} - \hat{\phi}_{|p})'\|\} &= O\left\{\left(\frac{1 - \rho}{1 + \rho}\right)^{1/2} T^{-1/2}\right\} \\ E\{\|(\hat{\phi} - \hat{\phi}_{|p})(\hat{\phi}_{|p} - \phi)'\|\} &= O\left\{\left(\frac{1 - \rho}{1 + \rho}\right)^{1/2} T^{-1/2}\right\} \end{aligned}$$

and the theorem follows.

*Proofs of results in Section 5*

PROOF OF THEOREM 5. The Taylor expansions of  $\check{A}_\alpha^h$  and  $\check{A}_\alpha^{h-1}$  around  $A_\alpha$  are

$$\check{A}_\alpha^k = A_\alpha^k \sum_{j=0}^{k-1} A_\alpha^j (\check{A}_\alpha - A_\alpha) A_\alpha^{k-1-j} + O_p(T^{-1}) \quad k = h, h - 1.$$

Then, using  $\sum_{j=0}^{h-2} A_\alpha^j (\check{A}_\alpha - A_\alpha) A_\alpha^{h-2-j} = \sum_{j=1}^{h-1} A_\alpha^{j-1} (\check{A}_\alpha - A_\alpha) A_\alpha^{h-1-j}$ , and given that  $A_\alpha - A_\alpha = e_{p+2}(\check{\phi} - \phi)'$ , we have

$$\begin{aligned} E(\check{w}_{T+h} - w_{T+h})^2 &= E(L_1 - L_2)^2 + E(C_{h,1}' Y_T Y_T' C_{h,1}) + E(C_{h,2}' Y_T Y_T' C_{h,2}) \\ &\quad + E(C_{h,1}' Y_T Y_T' C_{h,2}) + E(C_{h,2}' Y_T Y_T' C_{h,1}) + O(T^{-3/2}) \end{aligned} \tag{A7}$$

where  $C_{h,1}' = e_{p+2}' A_\alpha^0 e_{p+2} (\check{\phi} - \phi)' A_\alpha^{h-1}$  and  $C_{h,2}' = \sum_{j=1}^{h-1} e_{p+2}' A_\alpha^{j-1} (A_\alpha - I_{p+2}) e_{p+2} (\check{\phi} - \phi)' A_\alpha^{h-1-j}$ , and where we have used the result that  $E(\|A_\alpha - A_\alpha\|^k) = O(T^{-k/2})$  (see, for instance, Bhansali, 1981, or Kunitomo and Yamamoto, 1985).

If we denote the  $k$ th coefficient of  $\phi(B)^{-1}$  by  $\psi_{k[\text{AR}(p+1)]}$  and the  $k$ th coefficient of  $\phi(B)^{-1}(1 - B)$  by  $\psi_{k[\text{ARMA}(p+1,1)]}$ , then  $e_{p+2}' A_\alpha^{k-1} (A_\alpha - I_{p+2}) e_{p+2} = \psi_{h[\text{AR}(p+1)]} - \psi_{k-1[\text{AR}(p+1)]} = \psi_{k[\text{ARMA}(p+1,1)]} = e_{p+2}' A_1^k c_{p+2}$ , and hence

$$E(L_h^2) = E\{(L_1 - L_2)^2\} = \sigma^2 \sum_{k=0}^{h-1} (e_{p+2}' A_1^k c_{p+2})^2. \tag{A8}$$

Since the effect of the dependence between  $Y_T$  and  $\check{\phi}$  in the PMSE is  $O(T^{-3/2})$

(Kunitomo and Yamamoto, 1985) and applying that  $\text{MSE}(\hat{\phi}) = \sigma^2 \Gamma_y^{-1} / T + O(T^{-3/2})$ , we find

$$E(C'_{h,2} Y_T Y_T' C_{h,2}) = \frac{\sigma^2}{T} \sum_{j=1}^{h-1} \sum_{k=1}^{h-1} (e'_{p+2} A_1^{j-1} c_{p+2})(e'_{p+2} A_1^{k-1} c_{p+2}) \\ \times \text{tr}(A_\alpha^{h-1-j} \Gamma_y A_\alpha' A_\alpha^{h-1-k} \Gamma_y^{-1}) + O(T^{-3/2}).$$

Applying the same arguments to the remaining terms of (A8) we obtain

$$E\{(\hat{w}_{T+h} - w_{T+h})^2\} = \sigma^2 \sum_{k=0}^{h-1} (e'_{p+2} A_1^k c_{p+2})^2 + \frac{\sigma^2}{T} \sum_{j=0}^{h-1} \sum_{k=0}^{h-1} (e'_{p+2} A_1^j c_{p+2})(e'_{p+2} A_1^k c_{p+2}) \\ \times \text{tr}(A_\alpha^{h-1-j} \Gamma_y A_\alpha' A_\alpha^{h-1-k} \Gamma_y^{-1}) + O(T^{-3/2}).$$

If we denote the  $k$ th coefficient of  $\phi(B)^{-1}$  by  $\psi_{k[\text{AR}(p)]}$ , then  $\psi_{k[\text{ARMA}(p+1,1)]} = \psi_{k[\text{AR}(p)]} + O(1 - \rho)$ , and therefore  $e'_{p+2} A_1^k c_{p+2} = e'_p A_p^k e_p + O(1 - \rho)$ . Then, if  $\beta > 1$ , expression (5.3) holds. Similarly, using the previous arguments, the proof of (5.4) follows.

PROOF OF THEOREM 6. The expectation of the square of  $w_{T+h} - \hat{w}_{T+h}$  is

$$E\{(w_{T+h} - \hat{w}_{T+h})^2\} = E(L_h^2) + E\{e'_p(\hat{A}_p^h - A_p^h) W_T W_T' (\hat{A}_p^h - A_p^h)' e_p\} \\ + E(v_T^2) + 2E\{e'_p(\hat{A}_p^h - A_p^h) W_T v_T\} \quad (\text{A9})$$

where the term  $E(L_h^2)$  is the same as (A8). Applying (A3) with  $k = 1$  and Hölder's inequality, then  $E(v_T^2) = o(T^{-1})$ . In order to solve the remaining terms of (A9), we will use a Taylor expansion of  $\hat{A}_p$  around  $A_p$ . The magnitude of the remainder term is determined by the root- $T$  consistency of  $\hat{A}_p$ . Then

$$\hat{A}_p^h = A_p^h + \sum_{j=0}^{h-1} A_p^j (\hat{A}_p - A_p) A_p^{h-1-j} \\ + \sum_{j=1}^{h-1} \left\{ \sum_{k=0}^{j-1} A_p^k (\hat{A}_p - A_p) A_p^{j-1-k} \right\} \times (\hat{A}_p - A_p) A_p^{h-1-j} + O_p(T^{-3/2}).$$

Thus, by Lemma A3,  $E\{e'_p(\hat{A}_p^h - A_p^h) W_T v_T\} = O[E\{\|\hat{\phi} - \phi\|' W_T v_T\}] = o(T^{-1})$ . Let us denote  $B'_{h,1} = e'_p \sum_{j=0}^{h-1} A_p^j (\hat{A}_p - A_p) A_p^{h-1-j}$ . Then, by Hölder's inequality,  $E\{e'_p(\hat{A}_p^h - A_p^h) W_T W_T' (\hat{A}_p^h - A_p^h)' e_p\} = E(B'_{h,1} W_T W_T' B_{h,1}) + O(T^{-3/2})$ . Applying Theorem 4 and the result that the effect in the PMSE of the dependency between  $\hat{\phi}_{|p}$  and  $W_T$  is  $O(T^{-3/2})$  (Kunitomo and Yamamoto, 1985), it follows that

$$E(B'_{h,1} W_T W_T' B_{h,1}) = \frac{\sigma^2}{T} \sum_{j=0}^{h-1} \sum_{k=0}^{j-1} (e'_p A_p^j e_p)(e'_p A_p^k e_p) \\ \times \text{tr}(A_p^{h-1-j} \Gamma_w A_p' A_p^{h-1-k} \Gamma_w^{-1}) + o(T^{-1})$$

and the proof of (5.6) is completed. Similarly, by the same arguments, expression (5.7) can be obtained.

*Proofs of Section 6*

PROOF OF LEMMA 1. Let us decompose  $Y_t$  as  $Y_t = (\tilde{Y}'_t, 0)' + \boldsymbol{\mu}$ , where  $\boldsymbol{\mu} = (\mu, \mu, \dots, \mu, 1)'$ . Since,  $\alpha = \mu(1 - \sum_{i=1}^{p+1} \phi_i)$ , it can be shown that  $A_\alpha^i \boldsymbol{\mu} A_\alpha^j = \bar{\boldsymbol{\mu}}$ , where  $\bar{\boldsymbol{\mu}} = \boldsymbol{\mu} \boldsymbol{\mu}'$ . Then  $A_\alpha^i \Gamma_y A_\alpha^j = A_\alpha^i \Gamma_y^* A_\alpha^j + \bar{\boldsymbol{\mu}}$ , where  $\Gamma_y^*$  is a  $(p+2) \times (p+2)$  matrix with  $\Gamma_y$  in the first  $(p+1) \times (p+1)$  submatrix and zero elsewhere. Also, the covariance matrix  $\Gamma_y$  has the following block structure:

$$\Gamma_y = \begin{pmatrix} \Gamma_0 & \boldsymbol{\mu}_0 \\ \boldsymbol{\mu}'_0 & 1 \end{pmatrix}$$

where  $\Gamma_0 = E(Y_{0t} Y'_{0t})$ , with  $Y_{0t} = (y_t, y_{t-1}, \dots, y_{t-p})'$  and  $\boldsymbol{\mu}_0 = E(Y_{0t})$ . Using the properties of the inverses of block matrices, we can partition  $\Gamma_y^{-1}$  as

$$\Gamma_y^{-1} = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$$

where  $B_{11} = (\Gamma_0 - \boldsymbol{\mu}_0 \boldsymbol{\mu}'_0)^{-1} = \Gamma_0^{-1}$ . Then it is verified that  $\text{tr}(A_\alpha^i \Gamma_y^* A_\alpha^j \Gamma_y^{-1}) = \text{tr}(A_\alpha^i \Gamma_y A_\alpha^j \Gamma_y^{-1})$ . Hence,  $\text{tr}(A_\alpha^i \Gamma_y A_\alpha^j \Gamma_y^{-1}) = \text{tr}(A_\alpha^i \Gamma_y A_\alpha^j \Gamma_y^{-1}) + \text{tr}(\boldsymbol{\mu} \Gamma_y^{-1})$ . Given that  $\text{tr}(\bar{\boldsymbol{\mu}} \Gamma_y^{-1}) = \boldsymbol{\mu}' \Gamma_y^{-1} \boldsymbol{\mu}$ , and applying a result of Searle (1984, p. 258), it can be seen that  $\boldsymbol{\mu}' \Gamma_y^{-1} \boldsymbol{\mu} = 1 - |\Gamma_y - \boldsymbol{\mu} \boldsymbol{\mu}'| / |\Gamma_y| = 1$ , since the last column and row of  $\Gamma_y - \boldsymbol{\mu} \boldsymbol{\mu}'$  are zero and  $\Gamma_y$  is invertible.

PROOF OF LEMMA 2. Let  $C$  be the following nonsingular matrix:

$$C = \begin{pmatrix} 1 & -\rho & 0 & \dots & 0 & 0 \\ 0 & 1 & -\rho & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & -\rho \\ 1 & -\phi_1 & -\phi_2 & \dots & -\phi_{p-1} & -\phi_p \end{pmatrix}$$

Then

$$D = CA_0 C^{-1} = \begin{pmatrix} A_p & 0 \\ 0 & \rho \end{pmatrix}$$

Let  $\lambda_k$  be an eigenvalue of the matrix  $Q = \Gamma_y^{-1} A_0^i \Gamma_y A_0^j$ . Then

$$|D^i \Gamma_C D^j - \lambda \Gamma_C| = 0 \tag{A10}$$

where  $\Gamma_C = C \Gamma_y C'$ . This matrix  $\Gamma_C$  can be considered as the covariance matrix of the transformed series  $Z_t = CY_t$ , where  $Z_t = (z_{1,t}, z_{1,t-1}, \dots, z_{1,t-p+1}, z_{2,t})'$  and

$$Z_t = DZ_{t-1} + a_t c_{p+1} \tag{A11}$$

Therefore, the first  $p \times p$  submatrix of  $\Gamma_C$  is the covariance matrix of a process  $z_{1,t}$  following the coefficient matrix  $A_p$  and noise  $a_t$ ; namely, the matrix  $\Gamma_{w|p}$ . Denoting by  $V_{12}, V_{21}$  and  $V_{22}$  the remaining submatrices of this partitioning, we can rewrite (A10) as

$$\begin{vmatrix} A_p^i \Gamma_{w|p} A_p^j - \lambda \Gamma_{w|p} & (A_p^i V_{12} \rho^j - \lambda V_{12}) V_{22}^{-1/2} \\ (\rho^i V_{21} A_p^i - \lambda V_{21}) V_{22}^{-1/2} & \rho^{i+j} - \lambda \end{vmatrix} = 0$$

From (A11), the term  $V_{22}$  is the variance of an AR(1) process with coefficient  $\rho$ . Therefore  $V_{22}^{-1} = O(1 - \rho)$ . Hence, using the rule to evaluate the determinant of a partitioned matrix (see, for instance, Searle, 1984)

$$|Q - \lambda I| = |A_p^i \Gamma_{w|p} A_p^j - \lambda \Gamma_{w|p}| \{ \rho^{i+j} + O(1 - \rho) - \lambda \} = 0$$

Since the trace of a matrix equals the sum of its eigenvalues, the lemma follows.

## ACKNOWLEDGMENT

This work has been partially supported by DGES grant PB96-0111, CICYT grant PB96-0339, IVIE, and Cátedra BBV de Métodos para la Mejora de la Calidad. The authors would like to thank Ngai Hang Chan, George Tiao, Mike Wiper and the participants of the NBER/NSF Time Series Seminar, Duke University, 1997, for helpful discussions and suggestions on this work. We also thank the referees and the editor for their valuable and constructive comments.

## REFERENCES

- BHANSALI, R. J. (1981) Effects of not knowing the order of an autoregressive process on the mean squared error of prediction—I. *J. Am. Stat. Assoc.* 76, 588–97.
- CAMPBELL, J. Y. and PERRON P. (1991) Pitfalls and opportunities: what macroeconomists should know about unit roots. *NBER Macroeconomics Annu.* 141–201.
- CHAN, N. H. and WEI, C. Z. (1987) Asymptotic inference for nearly nonstationary AR(1) processes. *Ann. Stat.* 15, 1050–63.
- DICKEY, D. A. and FULLER, W. A. (1979) Distribution of the estimators for autoregressive time series with a unit root. *J. Am. Stat. Assoc.* 74, 427–31.
- FULLER, W. A. (1976) *Introduction to Statistical Time Series*. New York: Wiley.
- and HASZA, D. P. (1981) Properties of predictors for autoregressive time series. *J. Am. Stat. Assoc.* 76, 155–61.
- GRANGER, C. W. J. and JOYEUX, R. (1980) An introduction to long-range time series models and fractional differencing. *J. Time Ser. Anal.* 1, 15–30.
- HARVEY, A. C. (1991) Finite sample prediction and overdifferencing. *J. Time Ser. Anal.* 2, 221–32.
- KUNITOMO, N. and YAMAMOTO, T. (1985) Properties of predictors in misspecified autoregressive time series models. *J. Am. Stat. Assoc.* 80, 941–50.
- LEDOLTER, J. and ABRAHAM, B. (1981) Parsimony and its importance in time series forecasting. *Technometrics* 23, 411–14.
- LÜTKEPOHL, H. (1991) *Introduction to Multiple Time Series Analysis*. Berlin: Springer.
- PHILLIPS, P. C. B. (1987) Towards a unified asymptotic theory for autoregression. *Biometrika* 74, 535–47.
- PLOSSER, C. I. and SCHWERT, G. W. (1977) Estimation of a non-invertible moving average process: the case of overdifferencing. *J. Economet.* 6, 199–224.
- and — (1978) Money, income and sunspots: measuring economic relationships and the effects of differencing. *J. Monet. Econ.* 4, 637–60.
- SEARLE, S. R. (1984) *Matrix Algebra useful for Statistics*. New York: Wiley.
- SHAMAN, P. and STINE, R. (1988) The bias of autoregressive coefficient estimators. *J. Am. Stat. Assoc.* 83, 842–48.
- STOCK, J. H. (1996) VAR, error correction and pretest forecast at long horizons. *Oxford Bull. Econ. Stat.* 58, 685–701.