



MONASH University

Australia

Department of Econometrics
and Business Statistics

<http://www.buseco.monash.edu.au/depts/ebs/pubs/wpapers/>

**Autoregressive Approximation in Nonstandard Situations:
The Non-Invertible and Fractionally Integrated Cases**

D.S. Poskitt

June 2005

Working Paper 16/05

Autoregressive Approximation in Nonstandard Situations: The Non-Invertible and Fractionally Integrated Cases

D. S. Poskitt*

Department of Econometrics and Business Statistics, Monash University

Abstract

Autoregressive models are commonly employed to analyze empirical time series. In practice, however, any autoregressive model will only be an approximation to reality and in order to achieve a reasonable approximation and allow for full generality the order of the autoregression, h say, must be allowed to go to infinity with T , the sample size. Although results are available on the estimation of autoregressive models when h increases indefinitely with T such results are usually predicated on assumptions that exclude (i) non-invertible processes and (ii) fractionally integrated processes. In this paper we will investigate the consequences of fitting long autoregressions under regularity conditions that allow for these two situations and where an infinite autoregressive representation of the process need not exist. Uniform convergence rates for the sample autocovariances are derived and corresponding convergence rates for the estimates of $AR(h)$ approximations are established. A central limit theorem for the coefficient estimates is also obtained. An extension of a result on the predictive optimality of AIC to fractional and non-invertible processes is obtained.

JEL subject classifications : Primary C14, C32; Secondary C53

Key words and phrases : autoregression, autoregressive approximation, fractional process, non-invertibility, order selection, asymptotic efficiency.

Version of June 3, 2005

*Corresponding address: Don Poskitt, Department of Econometrics and Business Statistics, Monash University, Victoria 3800, Australia Tel.:+61-3-9905-9378; fax:+61-3-9905-5474.

E-mail address: Don.Poskitt@Buseco.monash.edu.au

1 Introduction

The use of autoregressive models has a long history that can be traced back to the early papers of Akaike (1969) and Parzen (1974) and beyond to the prescient work of Yule (1921). It is not surprising given this long history that there is a substantial literature dealing with such models :- using the *Google* web browser with the search word autoregression produced 17,600 sites, the word autoregressive produced 89,700! Nevertheless, there are still some gaps in the theory of autoregressive approximation that need to be filled if autoregressive modelling is to be routinely extended to the type of long memory processes that have come to play such an important role in time series analysis. Many empirical time series exhibit long-term persistence and Beran (1992, 1994) and Baillie (1996) provide a brief history of the application of long memory processes and a review of various statistical procedures for analyzing such processes. This paper extends the theory of autoregressive approximation to both long memory and non-invertible processes.

In order to set the scene let $y(t)$ for $t \in \mathbb{Z}$ denote a linearly regular, covariance-stationary process,

$$y(t) = \sum_{j=0}^{\infty} k(j)\varepsilon(t-j) \quad (1.1)$$

where $\varepsilon(t)$, $t \in \mathbb{Z}$, is a zero mean white noise process with variance σ^2 and the impulse response coefficients satisfy the conditions $k(0) = 1$ and $\sum_{j \geq 0} k(j)^2 < \infty$.

Assumption 1 Let \mathcal{E}_t denote the σ -algebra of events determined by $\varepsilon(s)$, $s \leq t$. It will be supposed throughout the paper that $\varepsilon(t)$ is ergodic and that

$$E[\varepsilon(t) | \mathcal{E}_{t-1}] = 0 \quad \text{and} \quad E[\varepsilon(t)^2 | \mathcal{E}_{t-1}] = \sigma^2. \quad (1.2)$$

Furthermore, $E[\varepsilon(t)^4] < \infty$.

Assumption 1 imposes a classical martingale difference structure on the innovations $\varepsilon(t)$. The significance of this assumption here is that it implies that the minimum mean squared error predictor of $y(t)$ given \mathcal{E}_{t-1} , $\bar{y}_{(t|t-1, \dots, \infty)}$ say, is the linear predictor, Hannan and Deistler (1988, Theorem 1.4.2).

Consider now the best linear predictor of $y(t)$ based on $y(t-j)$, $j = 1, \dots, h$. Let $\gamma(\tau) = \gamma(-\tau) = E[y(t)y(t+\tau)] = \sigma^2 \sum_{r \geq 0} k(r)k(\tau+r)$, $\tau = 0, 1, \dots$, denote the autocovariance function of the process $y(t)$. The coefficients of the minimum mean squared predictor of $y(t)$ based only on the finite past $y(t-1), \dots, y(t-h)$, denoted $\phi_h(j)$, $j = 0, \dots, h$, are obtained by solving the Yule-Walker equations

$$\sum_{j=0}^h \phi_h(j)\gamma(j-k) = \delta_{0k}\sigma_h^2, \quad k = 0, 1, \dots, h, \quad (1.3)$$

where δ_{0k} is Kronecker's delta, $\phi_h(0) = 1$ and

$$\sigma_h^2 = E[\epsilon_h(t)^2] \quad (1.4)$$

is the minimising value of the prediction error variance associated with the prediction error

$$\epsilon_h(t) = \sum_{j=0}^h \phi_h(j)y(t-j). \quad (1.5)$$

If h is chosen appropriately then it seems reasonable to suppose that the optimal predictor $\bar{y}_{\langle t|t-1, \dots, t-h \rangle} = \phi_h(1)y(t-1) + \dots + \phi_h(h)y(t-h)$ determined from the autoregressive model of order h implicit in its calculation will form a good approximation to the best predictor $\bar{y}_{\langle t|t-1, \dots, \infty \rangle}$. It is this notion, in part, that constitutes the background to the use of autoregressive models as a means of analysing observed time series.

Heuristically speaking it is clear that the order of the autoregression must be allowed to go to infinity in order to capture the influence of effects in the remote past and from a theoretical perspective h must be allowed to increase indefinitely in order to achieve full generality. Although results are available on the estimation of autoregressive models when $h \rightarrow \infty$ with the sample size T , such results are usually predicated on the assumption that the process admits an infinite autoregressive representation with coefficients that tend to zero at an appropriate rate. Such assumptions are often expressed in terms of particular summability conditions on the autoregressive coefficients themselves, or equivalently the Wold representation. Thus it is commonly assumed that (i) the transfer function associated with Wold's representation,

$$k(z) = \sum_{j=0}^{\infty} k(j)z^j,$$

is invertible, which following common practice is defined to mean $k(z) \neq 0$, $|z| \leq 1$, and, (ii) a summability condition such as $\sum_{j \geq 0} |k(j)| < \infty$, or $\sum_{j \geq 0} j|k(j)|^2 < \infty$, or $\sum_{j \geq 0} j^{\frac{1}{2}}|k(j)| < \infty$ holds. See Hannan and Deistler (1988, Section 7.4) for example. There are two critical cases that do not meet such conditions (i) non-invertible processes, of course, and (ii) fractionally integrated processes. One contribution of this paper is to show that such assumptions can be relaxed and that results on the statistical properties of AR approximations can be extended to allow for fractional and non-invertible processes.

Examination of non-invertible processes is motivated by the observation that, although it might be argued that processes observed in the real world are unlikely to exhibit spectral zeroes, lack of invertibility might be induced by the actions of the practitioner, by over-differencing for example. The consequences of such over-differencing for the subsequent analysis of an observed time series are then of interest.

Fractionally integrated processes are of interest in their own right. Fractional processes,

for which the notation $I(d)$ is commonly used, were introduced by Granger and Joyeux (1980), and were independently described in Hosking (1980), and it has been shown that they exhibit dynamic behaviour very similar that observed with many empirical time series. The class of fractionally integrated $I(d)$ processes can be characterized by the specification

$$y(t) = \sum_{j \geq 0} k(j)\varepsilon(t-j) = k(z)\varepsilon(t) = \frac{\kappa(z)}{(1-z)^d}\varepsilon(t)$$

wherein, as will be done henceforth in expressions of this type, the indeterminate z is interpreted as the lag operator, that is $z\varepsilon(t) = \varepsilon(t-1)$. For any $b > -1$ the operator $(1-z)^b$ is defined via the binomial expansion

$$(1-z)^b = 1 - bz + \frac{b(b-1)z^2}{2!} - \frac{b(b-1)(b-2)z^3}{3!} + \dots,$$

which yields the result that

$$\frac{1}{(1-z)^d} = \sum_{j=0}^{\infty} \frac{\Gamma(j+d)z^j}{\Gamma(j+1)\Gamma(d)},$$

where the gamma function $\Gamma(x) = \int_0^{\infty} t^{x-1}e^{-t}dt$ for $x \geq 0$ and the relation $\Gamma(x+1) = x\Gamma(x)$ defines $\Gamma(x)$ for $x < 0$. Hence

$$k(j) = \sum_{r=0}^j \frac{\kappa(j-r)\Gamma(r+d)}{\Gamma(r+1)\Gamma(d)} \quad j = 1, 2, \dots$$

where $\kappa(z) = \sum_{j \geq 0} \kappa(j)z^j$. If $\kappa(z)$ is such that $\sum_{j \geq 0} |\kappa(j)| < \infty$, $\kappa(z)$ might be the transfer function of a stable and invertible autoregressive moving-average (*ARMA*) process for example, then using Sterling's approximation it can be shown that

$$k(j) \sim \frac{\kappa(1)}{\Gamma(d)}j^{d-1} \quad \text{as } j \rightarrow \infty. \tag{1.6}$$

From (1.6) it follows that $\sum_{j \geq 0} |k(j)|^2 < \infty$ if $|d| < 0.5$ and $y(t)$ is well-defined as the limit in mean square of a covariance-stationary process with spectral density

$$f(\omega) = \frac{\sigma^2|k(e^{i\omega})|^2}{2\pi} = \frac{\sigma^2|\kappa(e^{i\omega})|^2}{2\pi|1 - e^{i\omega}|^{2d}}.$$

Using the result that $|1 - e^{i\omega}|^{2d} = |2 \sin(\omega/2)|^{2d}$ and $\sin(\omega/2) \sim \omega/2$ as $\omega \rightarrow 0$ it can be shown that the spectral density obeys the inverse power law $f(\omega) \sim \sigma^2|\kappa(1)|^2/2\pi\omega^{2d}$ as ω approaches zero. Similarly, the autocovariance function declines at a hyperbolic rate, $\gamma(\tau) \sim C\tau^{2d-1}$, $C \neq 0$, as $\tau \rightarrow \infty$, and not at an exponential rate as it would for a stable and invertible *ARMA* process. Throughout the paper C will stand for a universal, though not the same, constant. Note that the impulse response coefficients of $k(z)$ are not absolutely

summable if $d > 0$ and $k(1) = 0$ if $d < 0$. For a more detailed examination of the properties outlined above see Beran (1994).

The use of fractional models depends on the practitioner being able to conduct appropriate inference. Likelihood based methods have been studied in Fox and Taqqu (1986), Sowell (1992) and Beran (1995), for example, and non-parametric and semi-parametric techniques have also been advanced, as in Robinson (1995). Tieslau, Schmidt and Baillie (1996) proposed a minimum-distance estimator that is structured in terms of the autocorrelations and Martin and Wilkinson (1999) consider an efficient method of moments, or indirect, estimator constructed using an autoregression as an auxiliary model. One potential use of autoregressions, beyond their use as approximations in their own right, is to develop similar indirect methods of estimation. Beran, Bhansali and Ocker (1998) discuss the modelling of finite order autoregressive processes driven by fractional Gaussian noise; here we consider long autoregressive approximations to general processes.

We establish uniform convergence rates for the sample autocovariances and derive corresponding convergence rates for the estimates of $AR(h)$ approximations, where $h \rightarrow \infty$ with T , under regularity conditions that allow for both non-invertible and fractionally integrated processes. A central limit theorem for the coefficient estimates is also obtained. All these results are, to the authors knowledge, new to the literature. A major contribution of this paper is to provide a verification of a conjecture of Beran (1992, p. 410) concerning the extension of a result on the predictive optimality of AIC due to Shibata (1980) to fractional and non-invertible processes.

The paper proceeds as follows. In Section 2 results from the theory of stochastic processes that provide a rationale for a consideration of AR approximations in more general settings than are currently considered are reviewed. Section 3 outlines the estimation techniques to be discussed. These two sections provide the background, establish notation and present the basic assumptions. Section 4 lists some of the fundamental results that underly the statistical properties of the estimators considered in Section 3. The properties of autoregressive approximations are discussed in detail in Section 5. Section 6 of the paper presents a central limit result for the autoregressive estimator. Section 7 closes the paper with a small simulation study illustrating the (finite sample) practical impact of the (asymptotic) results obtained. Proofs are assembled together in the appendix.

2 Linear Prediction and Autoregressive Approximation

Since by assumption $y(t)$ is a regular process then we know from a famous result due to Szegö (1939) and Kolmogorov (1941) that

$$\int_{-\pi}^{\pi} \log\{f(\omega)\} d\omega > -\infty,$$

and it is not possible to determine $y(t + 1)$ precisely from its own history up to time t , i.e

$$\sigma^2 = 2\pi \exp\left\{\frac{1}{2\pi} \int_{-\pi}^{\pi} \log\{f(\omega)\} d\omega\right\} > 0. \quad (2.1)$$

The transfer function $k(z)$ has no zeroes inside the unit circle and $|k(e^{i\omega})|^2 > 0$ almost everywhere (a.e.) where $|k(e^{i\omega})|^2 = \lim_{\rho \uparrow 1} |k(\rho e^{i\omega})|^2$, the radial limit of $k(z)$ on the boundary of the unit circle $|z| = 1$, see Grenander and Rosenblatt (1957). In the context of autoregressive modelling it is now standard practice to strengthen the mini-phase condition, $k(z) \neq 0$, $|z| < 1$, by adding the condition that $k(z)$ has no zeroes on the unit circle. If $k(z)$ has neither zeroes nor singularities for $|z| \leq 1$ then the same is true of

$$\phi(z) = \sum_{j=0}^{\infty} \phi(j)z^j = \frac{1}{k(z)}$$

where the $k(j)$ and $\phi(j)$ are related by the recursions

$$\phi(0) = k(0) = 1, \quad \sum_{i=0}^j k(i)\phi(j-i) = 0, \quad j = 1, 2, \dots \quad (2.2)$$

Inversion of the operator $k(z)$ results in an infinite autoregressive, $AR(\infty)$, representation of $y(t)$ equivalent to (1.1),

$$\sum_{j=0}^{\infty} \phi(j)y(t-j) = \varepsilon(t), \quad (2.3)$$

to which the $AR(h)$ model in (1.5) yields an obvious approximation, the restriction that $k(z) \neq 0$, $|z| \leq 1$, thereby being motivated by the ease of interpretation that it provides.

If $y(t)$ is regular, however, it is not necessary for $k(z)$ to be invertible in order for there to be an autoregression that yields an appropriate approximation to the process. Rewriting the Yule-Walker equations in matrix-vector notation yields $\mathbf{\Gamma}_h \boldsymbol{\phi}_h = -\boldsymbol{\gamma}_h$ where $\mathbf{\Gamma}_h = [\gamma(i-j)]_{i,j=1,\dots,h}$, $\boldsymbol{\phi}_h = (\phi_h(1), \dots, \phi_h(h))'$ and $\boldsymbol{\gamma}_h = (\gamma(1), \dots, \gamma(h))'$. Note that regularity of $y(t)$ implies that $\mathbf{\Gamma}_h$ is nonsingular for all h . Otherwise there would exist a non-null h element unit vector \mathbf{x} such that

$$\mathbf{x}'\mathbf{\Gamma}_h\mathbf{x} = \int_{-\pi}^{\pi} \left| \sum_{s=1}^h x_s \exp(-i\omega s) \right|^2 f(\omega) d\omega = 0,$$

which implies, Munroe (1953, Theorem 25.7), that $|\sum_{j=1}^h x_j \exp(-i\omega s)|^2 = 0$ a.e. since $f(\omega) > 0$ a.e.. Bessel's inequality now implies that $x_s = 0$ for $s = 1, \dots, h$, leading to the conclusion that $\mathbf{\Gamma}_h > 0$ *reductio ad absurdum*. It follows that $\boldsymbol{\phi}_h$ is unique and $\phi_h(z) = \sum_{j=0}^h \phi_h(j)z^j \neq 0$, $|z| \leq 1$. Solving (1.3) using the Levinson (1947)-Durbin (1960) algorithm

$$\phi_h(j) = \phi_{h-1}(j) + \phi_h(h)\phi_{h-1}(h-j), \quad \phi_h(0) = 1, \quad j = 1, \dots, h-1$$

$$\begin{aligned}\phi_h(h) &= \sum_{j=0}^{h-1} \phi_{h-1}(j)\gamma(h-j)/\sigma_{h-1}^2 \\ \sigma_h^2 &= \sigma_{h-1}^2(1 - \phi_h(h)^2)\end{aligned}\tag{2.4}$$

initiated at $\phi_0(0) = 1$ and $\sigma_0^2 = \gamma(0)$, and using the relationship $\sigma_h^2 = \det(\mathbf{\Gamma}_{h+1})/\det(\mathbf{\Gamma}_h)$, which leads to the conclusion that $|\phi_h(h)| < 1$ for all h , we can see that σ_h^2 is monotonically decreasing in h . Basic Hilbert space arguments can also be used to show that $\lim_{h \rightarrow \infty} \sigma_h^2 = \sigma^2$. The later follows from the following result.

Lemma 2.1 *If $y(t)$ is a linearly regular, covariance-stationary process then*

$$\lim_{h \rightarrow \infty} E[(\epsilon_h(t) - \varepsilon(t))^2] = 0.$$

Thus, heuristically at least, we may think of (2.3) as providing an expression for the innovation in terms of $y(t), y(t-1), \dots$ where the coefficients $\phi(j), j = 1, \dots$, are determined not by inverting $k(z)$ and solving (2.2) but as the limit of the stable autoregressive operators $\phi_h(z)$ as $h \rightarrow \infty$. Indeed, Wold (1938) first derived (1.1) by fitting autoregressions of ever increasing order.

Example(i): Suppose that $y(t) = \varepsilon(t) - \varepsilon(t-1)$. Then $y(t)$ is regular but not invertible. Substituting $k(z) = 1 - z$ into (2.2) gives coefficients $\phi(j) = 1$ for all j but $\sum_{j>0} y(t-j)$ is not convergent since $E[(\sum_{j \geq n} y(t-j))^2] = E[\varepsilon(t-n)^2] = \sigma^2$ for any n . The best predictor $\bar{y}_{\langle t|t-1, \dots, t-h \rangle}$ is nevertheless well defined. Solving the Yule-Walker equations with

$$\mathbf{\Gamma}_h = \sigma^2 \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & \dots & 0 \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & -1 \\ 0 & \dots & 0 & -1 & 2 \end{bmatrix} \quad \text{and} \quad \boldsymbol{\gamma}_h = \sigma^2 \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

it is easily verified that the solution is

$$\phi_h = \left(\frac{h}{h+1}, \frac{h-1}{h+1}, \dots, \frac{1}{h+1} \right)' \quad \text{and} \quad \sigma_h^2 = \sigma^2 \left\{ 1 + \frac{1}{h+1} \right\}.$$

Expanding the residual process

$$\epsilon_h(t) = \sum_{j=0}^h \left\{ 1 - \frac{j}{h+1} \right\} y(t-j)$$

as a function of the innovations and rearranging terms gives

$$\epsilon_h(t) = \varepsilon(t) - \frac{1}{h+1} \sum_{j=1}^{h+1} \varepsilon(t-j).$$

Since $\varepsilon(t)$ are martingale differences the second term in this last expression obeys the law of the iterated logarithm in h and $\varepsilon_h(t)$ converges to $\varepsilon(t)$ almost surely as $h \rightarrow \infty$, not just in mean square. \square

Example(ii): Now suppose that $y(t)$ is a fractional noise process, $y(t) = (1 - z)^{-d}\varepsilon(t)$, $|d| < 0.5$. Then it can be shown that $y(t)$ is the solution to the stochastic difference equation $\sum_{j \geq 0} \psi(j)y(t - j) = \varepsilon(t)$ where

$$\psi(j) = \frac{\Gamma(j - d)}{\Gamma(j + 1)\Gamma(-d)},$$

the coefficients in the binomial expansion of $(1 - z)^d$. Thus $y(t)$ admits an infinite autoregressive representation for all $d \in (-0.5, 0.5)$ even though $k(z) = (1 - z)^{-d}$ is not invertible if $-0.5 < d < 0$ since $k(z) = 0$ when $z = 1$. Inserting the recursion $\gamma(h) = \gamma(h - 1)(h + d - 1)/(h + d)$, $h = 1, 2, \dots$, $\gamma(0) = \sigma^2\Gamma(1 - 2d)/\Gamma^2(1 - d)$, into the Levinson-Durbin algorithm we find that

$$\phi_h(j) = \psi(j) \left\{ \frac{\Gamma(h + 1)\Gamma(h + 1 - d - j)}{\Gamma(h + 1 - j)\Gamma(h + 1 - d)} \right\}$$

for $j = 1, \dots, h$ and

$$\sigma_h^2 = \sigma^2 \frac{\Gamma(h + 1)\Gamma(h + 1 - 2d)}{\Gamma^2(h + 1 - d)}.$$

Now, from Sterling's approximation it follows that $\Gamma(x + 1 + a)/\Gamma(x + 1) = x^a\{1 + o(1)\}$ for $|a| < 1$ as $x \rightarrow \infty$ and from this it is straightforward to show that

$$\frac{\Gamma(h + 1)\Gamma(h + 1 - 2d)}{\Gamma^2(h + 1 - d)} = \left\{ 1 + \frac{(1 - d)}{h} \right\}^d \{1 + o(1)\}$$

and $\sigma_h^2 \rightarrow \sigma^2$ as $h \rightarrow \infty$, illustrating directly the consequence of Lemma 2.1 in this case. It also follows that $|\phi_h(j) - \psi(j)| \rightarrow 0$ for all $j = 1, \dots, h$ as $h \rightarrow \infty$ for, as might have been anticipated, the sequence of autoregressions characterized by $\phi_h(z)$, $h = 1, 2, \dots$, converge in mean square to the infinite autoregressive representation $(1 - z)^d y(t) = \varepsilon(t)$. \square

From the preceding discussion it is apparent that it is the regularity of $y(t)$ that is important in the context of autoregressive modelling rather than invertibility in the conventional sense that $k(z) \neq 0$ for $|z| \leq 1$. This observation gives rise to the following:

Assumption 2 *The series $y(t)$ is a linearly regular, covariance-stationary process with Wold representation $y(t) = \sum_{j \geq 0} k(j)\varepsilon(t - j)$ where $k(z) = \kappa(z)/(1 - z)^d$ for $|d| < 0.5$ and $\kappa(z)$ is a causal transfer function with impulse response coefficients satisfying $\sum_{j \geq 0} |\kappa(j)| < \infty$.*

3 Model Fitting

Let $y(t)$, $t = 1, \dots, T$ denote a realisation of T observations on an observed process and set

$$c_T(r) = c_T(-r) = T^{-1} \sum_{t=r+1}^T y(t-r)y(t), \quad r = 0, 1, \dots, T-1, \quad (3.1)$$

the sample autocovariance function. Substituting $c_T(r)$ for $\gamma(r)$ in the Yule-Walker equations and solving for $\phi_h(j)$, $j = 1, \dots, h$ and σ_h yields estimates of the parameters in the $AR(h)$ model. Noting the correspondence with the method of moments we will denote the Yule-Walker estimator and its associated estimates by the use of an over-bar. This estimator has the advantage that it can be readily calculated via the Levinson-Durbin recursions, and being based on Toeplitz calculations $\bar{\phi}_h(z)$ will be stable. The variance estimate $\bar{\sigma}_h^2 = c_T(0) + \sum_{j=1}^h \bar{\phi}_h(j)c_T(j)$ need not minimize the empirical mean squared error however.

Estimating the parameters by directly minimizing the observed mean squared error $T^{-1} \sum_{t=1}^T (y(t) - \phi_h(1)y(t-1) + \dots + \phi_h(h)y(t-h))^2$ leads to the least squares estimates of course, which we shall denote by use of a caré. The least squares estimator is obtained by solving the normal equations

$$\mathbf{M}_h \hat{\phi}_h = -\mathbf{m}_h$$

where

$$\mathbf{M}_h = T^{-1} \sum_{t=1}^T \begin{bmatrix} y(t-1) \\ \vdots \\ y(t-h) \end{bmatrix} (y(t-1), \dots, y(t-h))$$

and

$$\mathbf{m}_h = T^{-1} \sum_{t=1}^T y(t) \begin{bmatrix} y(t-1) \\ \vdots \\ y(t-h) \end{bmatrix}.$$

In the above expressions the pre-sample values $y(1-h), \dots, y(0)$ are assumed to be equal to zero. For ease of exposition and notational simplicity summations will continue to be expressed in this manner in what follows. In practice the range of summation for the least squares estimator is often taken as $t = h+1, \dots, T$. The effects of the elimination of the initial terms will, for given h , be asymptotically negligible. Efficient numerical methods for the solution of least squares problems of this type are readily available of course. By way of contrast with the Yule-Walker estimator,

$$\hat{\sigma}^2 = T^{-1} \sum_{t=1}^T (y(t) - \hat{\phi}_h(1)y(t-1) + \dots + \hat{\phi}_h(h)y(t-h))^2$$

minimizes the observed mean squared error but there is no guarantee that $\hat{\phi}_h(z)$ will be

stable.

It will be shown below that $(\hat{\phi}'_h, \hat{\sigma}^2)$ and $(\bar{\phi}'_h, \bar{\sigma}^2)$ are asymptotically equivalent under the regularity conditions employed here, but they may exhibit quite different finite sample behaviour. Indeed, the stability of $\bar{\phi}_h(z)$ is known to give rise to significant biases in finite samples, biases that are not present with the least squares estimator. Tjøstheim and Paulsen (1983) present theoretical and empirical evidence of this phenomenon and show that when $y(t)$ is a finite autoregression then the first term in an asymptotic expansion of the bias of $\bar{\phi}_h$ has order of magnitude $O(T^{-1})$ but the size of the constant varies inversely with the distance of the zeroes of the true autoregressive operator from the unit circle. Hence, when the data generating mechanism shows strong autocorrelation it is possible for the bias in the Yule-Walker coefficient estimates to be substantial. This bias is known to feed through to other quantities of interest such as the prediction error variance Paulsen and Tjøstheim (1985) and estimates of power spectra Lysne and Tjøstheim (1987). Given that fractional processes can display long-range dependence with autocovariances that decay much slower than exponentially, similar effects are likely to be manifest when employing the Yule-Walker estimates under the current scenario *a-fortiori*. This suggests that the least squares estimator is to be preferred. Some empirical evidence illustrating these ideas is given below.

4 Some Asymptotic Theory

We begin with some asymptotic properties of the basic statistics that form the building blocks of the Yule-Walker and least squares estimators.

Theorem 4.1 *Suppose that $y(t)$ is a covariance-stationary process that satisfies Assumption 1 and Assumption 2 and that $H_T = o\{(T/\log T)^{\frac{1}{2}-d'}\}$ where $d' = \max\{0, d\}$. Then*

$$\max_{0 \leq \tau \leq H_T} |c_T(\tau) - \gamma(\tau)| = O\{(\log T/T)^{\frac{1}{2}-d'}\}.$$

This result is of interest in its own right for it indicates that the convergence rate of the autocovariance estimates of a fractional process equals that that obtains in the standard stationary case if $d < 0$ and $y(t)$ is anti-persistent but if $d > 0$ and $y(t)$ exhibits long memory then the convergence can be much slower.

Let

$$\begin{aligned} c_T(j, k) &= T^{-1} \sum_{t=1}^T y(t-j)y(t-k) \\ &= T^{-1} \sum_{t=\max\{j,k\}+1}^T y(t-j)y(t-k) \quad j, k = 0, 1, \dots, H_T. \end{aligned}$$

Whereas the autocovariance estimates $c_T(\tau)$ are used to calculate $\bar{\phi}_h$, it is the lag covariances $c_T(j, k)$ that determine the normal equations that define $\hat{\phi}_h$.

Theorem 4.2 *Under the same conditions as for Theorem 4.1*

$$\max_{0 \leq \tau \leq H_T} \max_{|j-k|=\tau} |c_T(j, k) - \gamma(\tau)| = O\{(\log T/T)^{\frac{1}{2}-d'}\}$$

uniformly in $j, k = 0, 1, \dots, H_T$.

Combining Theorem 4.1 with Theorem 4.2 gives rise to the following corollary.

Corollary 4.1 *If $y(t)$ satisfies assumptions 1 and 2 then the Yule-Walker and least squares autoregressive estimators $\bar{\phi}_h$ and $\hat{\phi}_h$ are asymptotically equivalent and*

$$\|\hat{\phi}_h - \bar{\phi}_h\|^2 = O\left\{\left(\frac{h^{1+4d}}{\lambda_{\min}(\mathbf{\Gamma}_h^4)}\right)\left(\frac{\log T}{T}\right)^{1-2d'}\right\} + O\left\{\left(\frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h^2)}\right)\left(\frac{\log T}{T}\right)^{1-2d'}\right\}.$$

In the light of Corollary 4.1 the results that follow will be expressed and proven in terms of the least squares or the Yule-Walker estimates, whichever is most convenient, it being understood that equivalent asymptotic properties will hold for both estimators.

Specific reference has not been made to the non-invertible case where $k(z) = \kappa(z)/(1 - z)^d = 0$, $|z| = 1$. This is because the existence of unit roots does not invalidate the conditions of Assumptions 1 and 2 and the results presented in this section will hold regardless, although they may not give the best rates of convergence possible. It is apparent from Corollary 4.1, however, that the presence of spectral zeroes could have an important impact via it's influence on the proximity of $\lambda_{\min}(\mathbf{\Gamma}_h)$ to zero. That this is so will be seen in Section 5.

In what follows consideration will be given to the properties of the estimates obtained by fitting an $AR(h)$ model where the order h is allowed to increase with T . Following the arguments in Section 2 we interpret this as a two step process. First, for any given h the parameters of the associated $AR(h)$ approximation must be estimated, and second, a "reasonable" value of h must be chosen. In the conventional case where an $AR(\infty)$ representation exists it is common practice to use Berk's inequality, Berk ('s 1974), to analyse the effects that the truncation due to using an $AR(h)$ approximation has on these two steps. Since under present assumptions an infinite autoregressive representation as in (2.3) is not guaranteed to exist this technique is not available to us. We can, nevertheless, handle the consequences of using an $AR(h)$ approximation if we know something of the relationship between the statistical properties of realizations of the innovations $\varepsilon(t)$ and realizations of the prediction errors, or residuals, $\epsilon_h(t)$.

Theorem 4.3 *Let $\varepsilon(t)$ and $\epsilon_h(t)$ denote the innovations and prediction errors associated with the minimum mean squared error predictors $\bar{y}_{\langle t|t-1, \dots, \infty \rangle}$ and $\bar{y}_{\langle t|t-1, \dots, t-h \rangle}$ of $y(t)$ where*

$y(t)$ satisfies Assumptions 1 and 2. Then

$$T^{-1} \sum_{t=1}^T \varepsilon(t) \{\epsilon_h(t) - \varepsilon(t)\} = O\{(\log \log T/T)^{\frac{1}{2}}\}$$

uniformly in h .

Theorem 4.3 implies that

$$T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 - T^{-1} \sum_{t=1}^T \varepsilon(t)^2 = T^{-1} \sum_{t=1}^T \{\epsilon_h(t) - \varepsilon(t)\}^2 + O\{(\log \log T/T)^{\frac{1}{2}}\}, \quad (4.1)$$

which provides an empirical counterpart to the result that $\sigma_h^2 \geq \sigma^2$ in that the first term on the right hand side of (4.1) will converge to $E[(\epsilon_h(t) - \varepsilon(t))^2] \geq 0$, by ergodicity, and thus for T sufficiently large $T^{-1} \sum_{t=1}^T \epsilon_h(t)^2$ will be bounded below by $T^{-1} \sum_{t=1}^T \varepsilon(t)^2$, with the difference converging to zero as h increases, see Lemma 2.1. It will be seen that (4.1) plays an important role in determining the behaviour of model selection devices for large T , as does the following result.

Theorem 4.4 *Let $y(t)$ and $\epsilon_h(t)$ be as in Theorem 4.3. Then uniformly in $h \leq H_T$*

$$\max_{1 \leq j \leq h} T^{-1} \sum_{t=1}^T \epsilon_h(t) y(t-j) = O\{(\log T/T)^{\frac{1}{2}-d'}\}.$$

Theorem 4.4 is the empirical counterpart of the result that the prediction error $\epsilon_h(t)$ is, by construction, orthogonal to $y(t-1), \dots, y(t-h)$, that is $E[\epsilon_h(t)y(t-j)] = 0$, $j = 1, \dots, h$.

5 Autoregressive Modelling

In practice, of course, neither $\varepsilon(t)$ nor $\epsilon_h(t)$ can be observed and their properties will have to be deduced by fitting autoregressive models to the data. We begin, therefore, by first establishing the consistency of the coefficient estimates of the $AR(h)$ model to those of the $AR(h)$ approximation to the process.

Theorem 5.1 *If $y(t)$ is a stationary process that satisfies Assumption 1 and Assumption 2 then uniformly in $h \leq H_T$*

$$\sum_{j=1}^h |\hat{\phi}_h(j) - \phi_h(j)|^2 = O\left\{ \left(\frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h^2)} \right) \left(\frac{\log T}{T} \right)^{1-2d'} \right\}.$$

The following theorem relates to the residuals

$$\hat{\epsilon}_h(t) = \sum_{j=0}^h \hat{\phi}_h(j) y(t-j)$$

as estimates of the prediction errors $\epsilon_h(t)$.

Theorem 5.2 *Under the same assumptions as for Theorem 5.1*

$$T^{-1} \sum_{t=1}^T \epsilon_h(t) \{\hat{\epsilon}_h(t) - \epsilon_h(t)\} = O \left\{ \left(\frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h)} \right) \left(\frac{\log T}{T} \right)^{1-2d'} \right\}.$$

uniformly in $h \leq H_T$.

Comparison of Theorem 5.2 with Theorem 4.3 indicates that whereas the deviation of $\epsilon_h(t)$ from $\varepsilon(t)$ relative to the magnitude of $\varepsilon(t)$, as measured by their covariation, converges to zero at a rate that is independent of d the same is not true of the corresponding relationship between $\hat{\epsilon}_h(t)$ and $\epsilon_h(t)$. The relevance of this observation stems from the fact that it is common practice to determine the order of the model to be employed by minimizing a model selection criterion of the form

$$SC_T(h) = \log(\hat{\sigma}_h^2) + hC_T/T$$

over the range $h = 0, 1, \dots, M_T$ where $\hat{\sigma}_h^2 = T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t)^2$ and $C_T > 0$ is chosen by the practitioner such that $C_T/T \rightarrow 0$ as $T \rightarrow \infty$, as is $M_T < H_T$. If $C_T = 2$ we have *AIC*, if $C_T = \log T$ we have *BIC*, Schwarz (1978), and setting $C_T = \log \log T$ we obtain the criterion advanced in Hannan and Quin (1979).

Consider the function

$$L_T(h) = (\sigma_h^2 - \sigma^2) + h\sigma^2/T.$$

Shibata (1980) introduced $L_T(h)$ as a figure of merit in the context of fitting autoregressive models to a truly infinite-order process. Shibata shows that if an $AR(h)$ model is fitted to a stationary Gaussian process that has an $AR(\infty)$ representation and it is used to predict an independent realization of the same process then the difference between the mean squared prediction error of the fitted model and the innovation variance converges in probability to $L_T(h)$. Thus, if $y(t)'$ denotes an independent realization of the process $y(t)$ then

$$E[(y(t)' - \sum_{i=1}^h \hat{\phi}_h(i)y(t-i)')^2] = \sigma_h^2 + E[\sum_{j=1}^h \sum_{i=1}^h (\hat{\phi}_h(i) - \phi_h(i))(\hat{\phi}_h(j) - \phi_h(j))\gamma(j-i)]$$

and given that $T^{\frac{1}{2}}(\hat{\phi}_h - \phi_h) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \sigma^2 \mathbf{\Gamma}_h^{-1})$ then the asymptotic expectation and probability limit of the second term is $h\sigma^2/T$. Noting that $E[(y(t)' - \bar{y}'_{(t|t-1, \dots, t-h)})^2] = \sigma_h^2 \geq \lim_{h \rightarrow \infty} E[\epsilon_h(t)^2] = \sigma^2$ we can see that the first term of $L_T(h)$ measures the fit of the model and the second reflects the inaccuracy or uncertainty in the determination of the parameters of $\bar{y}'_{(t|t-1, \dots, t-h)} = \phi_h(1)y(t-1)' + \dots + \phi_h(h)y(t-h)'$. Now, $L_T(h)$ is bounded below by $L_T(h_T^*)$ in the range $h = 0, 1, \dots, M_T$ where $L_T(h_T^*) = \min_{h=1, \dots, M_T} L_T(h)$ and Shibata defines a sequence of selected orders h'_T as being efficient if $\lim_{T \rightarrow \infty} L_T(h'_T)/L_T(h_T^*) = 1$.

Although the regularity conditions imposed by Shibata (1980) are too restrictive to be applicable here a similar rationale for consideration of $L_T(h)$ can be given. Observe also that by Theorem 5.2 the empirical difference $T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t)^2 - T^{-1} \sum_{t=1}^T \epsilon(t)^2$ equals

$$T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 - T^{-1} \sum_{t=1}^T \epsilon(t)^2 + T^{-1} \sum_{t=1}^T (\hat{\epsilon}_h(t) - \epsilon_h(t))^2 + O\{h(\log T/T)^{1-2d'}\}.$$

The limit of $T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 - T^{-1} \sum_{t=1}^T \epsilon(t)^2$ is $\sigma_h^2 - \sigma^2$ and the third term

$$T^{-1} \sum_{t=1}^T (\hat{\epsilon}_h(t) - \epsilon_h(t))^2 = T^{-1} \sum_{t=1}^T \sum_{j=1}^h \sum_{i=1}^h (\hat{\phi}_h(i) - \phi_h(i))(\hat{\phi}_h(j) - \phi_h(j))y(t-i)y(t-j).$$

is a consistent estimate of $\sum_{j=1}^h \sum_{i=1}^h (\hat{\phi}_h(i) - \phi_h(i))(\hat{\phi}_h(j) - \phi_h(j))\gamma(j-i)$ by Theorem 4.2. Thus $L_T(h)$ can be viewed as providing a limiting bound to the empirical difference in the mean squared prediction error and the innovation variance.

Set

$$\bar{L}_T(h) = \log \left(1 + \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \epsilon(t)^2}{\sum_{t=1}^T \epsilon(t)^2} \right) + \frac{h}{T}$$

and let \bar{h}_T^* denote a sequence of non-negative integers at each of which the minimum of $\bar{L}_T(h)$ with respect to h is attained, that is

$$\bar{L}_T(\bar{h}_T^*) = \min_{0 \leq h \leq M_T} \bar{L}_T(h)$$

or equivalently $\bar{h}_T^* = \operatorname{argmin}_{0,1,\dots,M_T} \bar{L}_T(h)$.

Theorem 5.3 *If $y(t)$ is a covariance-stationary process that satisfies Assumptions 1 and 2 then*

$$\lim_{T \rightarrow \infty} \left| \frac{\sigma^2 \bar{L}_T(\bar{h}_T^*)}{L_T(h_T^*)} - 1 \right| = 0$$

almost surely where $h_T^ = \operatorname{argmin}_{0,1,\dots,M_T} L_T(h)$.*

The criterion $\bar{L}_T(h)$ is unfeasible, but letting $AIC_T(h)$ denote the criterion $SC_T(h)$ when $C_T = 2$ we can deduce from Theorem 5.4 presented immediately below that

$$\begin{aligned} \max_{1 \leq h \leq M_T} |AIC_T(h) - \bar{L}_T(h)| &\leq \left| \log(T^{-1} \sum_{t=1}^T \epsilon(t)^2) \right| + \frac{M_T}{T} \\ &\quad + O \left\{ \left(\frac{M_T}{\lambda_{\min}(\mathbf{\Gamma}_{M_T})} \right) \left(\frac{\log T}{T} \right)^{1-2d'} \right\}. \end{aligned} \quad (5.1)$$

Theorem 5.4 *Under the same assumptions as for Theorem 5.3*

$$SC_T(h) = \log(T^{-1} \sum_{t=1}^T \epsilon(t)^2) + \log \left(1 + \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \epsilon(t)^2}{\sum_{t=1}^T \epsilon(t)^2} \right)$$

$$+\frac{hC_T}{T} + O\left\{\left(\frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h)}\right)\left(\frac{\log T}{T}\right)^{1-2d'}\right\}$$

uniformly in $h = 0, 1, \dots, H_T$.

Note that if $(M_T/\lambda_{\min}(\mathbf{\Gamma}_{M_T}))(\log T/T)^{1-2d'}$ converges to zero as $T \rightarrow \infty$ then the only non-vanishing term on the right hand side of (5.1) (the first term) is independent of both d and h . We can therefore conclude that $h_T^{AIC}/\bar{h}_T^* \rightarrow 1$ as $T \rightarrow \infty$ where h_T^{AIC} is the autoregressive order determined by $AIC_T(h)$ provided that $M_T/\lambda_{\min}(\mathbf{\Gamma}_{M_T}) = o\left\{(T/\log T)^{1-2d'}\right\}$.

Theorem 5.5 *If $y(t)$ is a covariance-stationary process that satisfies Assumptions 1 and 2 and $h_T^{AIC} = \operatorname{argmin}_{0,1,\dots,M_T} AIC_T(h)$ where $\lim_{T \rightarrow \infty} (M_T/\lambda_{\min}(\mathbf{\Gamma}_{M_T}))(\log T/T)^{1-2d'} = 0$ then the $AR(h_T^{AIC})$ model is asymptotically efficient in the sense that*

$$L_T(h_T^{AIC}) = L_T(h_T^*)\{1 + o(1)\}$$

almost surely as $T \rightarrow \infty$.

Alternative methods of autoregressive order determination that do not share the same structure as $SC_T(h)$ above have been proposed in the literature. The criterion autoregressive transfer function suggested by Parzen (1974) and the mean squared prediction error criterion of Mallows (1973), for example. Parzen's criterion can be expressed as

$$CAT_T(h) = 1 - \frac{(T-h)\tilde{\sigma}^2}{T\hat{\sigma}_h^{-2}} + \frac{h}{T}$$

and Mallows's statistic

$$MC_T(h) = T\left(\frac{\hat{\sigma}_h^2}{\tilde{\sigma}^2} - 1\right) + 2h$$

where

$$\tilde{\sigma}^2 = 2\pi \exp\left\{(2\pi N)^{-1} \sum_{j=1}^N \sum_{\tau=1-T}^{T-1} c_T(\tau) \cos(2\pi j\tau/T) + \gamma'\right\},$$

$\gamma' = 0.57721$ (Eulers constant) and $N = \lfloor (T-1)/2 \rfloor$, a nonparametric estimate of the innovation variance constructed from the periodogram by analogy with (2.1). Simple algebra shows that

$$CAT_T(h) - CAT_T(h-1) = \left\{\frac{T^{-1}(T-h-1)\hat{\sigma}_h^{-2} - T^{-1}(T-h)\hat{\sigma}_{h-1}^{-2}}{\hat{\sigma}_h^{-2}\hat{\sigma}_{h-1}^{-2}}\right\}\tilde{\sigma}^2 + \frac{1}{T}$$

and

$$MC_T(h) - MC_T(h-1) = T\left(\frac{\hat{\sigma}_h^2 - \hat{\sigma}_{h-1}^2}{\tilde{\sigma}^2}\right) + 2,$$

while from Theorem 5.2 and expression (4.1) it follows that

$$AIC_T(h) - AIC_T(h-1) = \frac{\hat{\sigma}_h^2 - \hat{\sigma}_{h-1}^2}{T^{-1} \sum_{t=1}^T \varepsilon(t)^2} + \frac{2}{T} + o(\hat{\sigma}_h^2 - \hat{\sigma}_{h-1}^2).$$

Similarly, it is straightforward to show that the final prediction error criterion

$$FPE_T(h) = \left(\frac{T+h}{T-h} \right) \hat{\sigma}_h^2$$

introduced by Akaike (1970) satisfies $\log FPE_T(h) = AIC_T(h) + O(T^{-2})$. Thus, bare remainder terms, we can anticipate that these criteria will move together and will be minimized at the same value of h . This suggests, and it can be shown, that $CAT_T(h)$, $MC_T(h)$ and $FPE_T(h)$ will also be asymptotically efficient selection criteria.

Finally, we must verify that a value of M_T that satisfies the requirements of Theorem 5.5 can be found. To investigate this in further detail it is necessary to give explicit structure to the spectral zeroes of the process and this is done by modifying Assumption 2.

Assumption 3 *The series $y(t)$ is a linearly regular, covariance-stationary process with Wold representation $y(t) = \sum_{j \geq 0} k(j)\varepsilon(t-j)$ where $k(z) = \kappa(z)/(1-z)^d$ for $|d| < 0.5$ and $\kappa(z) = u(z)\mu(z)$ where*

$$u(z) = (1-z)^{\nu_1}(1+z)^{\nu_2} \prod_{j=3}^n (1 - 2\cos(\theta_j)z + z^2)^{\nu_j},$$

with $0 < \theta_j < \pi$, $j = 3, \dots, n$, $\nu_j \geq 0$, $j = 1, \dots, n$, and $\mu(z)$ is a causal transfer function such that $|\mu(z)| \neq 0$, $|z| \leq 1$, $\sum_{j \geq 0} |\mu(j)| < \infty$.

The factor $u(z)$ is a possibly fractional operator whose roots lie on the unit circle. By appropriate choice of n and the theta it can be thought of as modelling spectral zeroes or troughs, as might occur after the incorrect application of seasonal differencing for example. The following result indicates how the zeroes of $u(z)$ impact on the behaviour of $\lambda_{\min}(\mathbf{\Gamma}_h)$ and it is this behaviour that determines appropriate bounds on the autoregressive order h and M_T .

Lemma 5.7 *Under Assumption 3*

$$\lambda_{\min}(\mathbf{\Gamma}_h) \begin{cases} \geq \inf_{\omega} \sigma^2 |\mu(e^{i\omega})|^2 / 2\pi |1 - e^{i\omega}|^{2d} > 0 & \text{if } u(z) \equiv 1 \text{ and } d \geq 0, \\ = O(h^{2d}) & \text{if } u(z) \equiv 1 \text{ and } d < 0, \\ = O(h^{-2 \max\{\nu_1-d, \nu_2, \nu_3, \dots, \nu_n\}}) & \text{otherwise,} \end{cases}$$

wherein $u(z) \equiv 1$ means that $\nu_1 = \dots = \nu_n = 0$.

Lemma 5.7 indicates that the presence of spectral zeroes of the type characterized by Assumption 3 leads to a consideration of terms of order $O\{h^{1+4q}(\log T/T)^{1-2d'}\}$, or smaller,

where $q \geq 0$. In order to operationalize the above estimation procedures a value for M_T must be chosen that ensures that $M_T^{1+4q}(\log T/T)^{1-2d'} \rightarrow 0$ as $T \rightarrow \infty$ for all possible values of q and d , both of which are unknown to the practitioner of course. One such choice is $M_T = [c(\log T)^a]$, the integer part of $c(\log T)^a$ for some $a \geq 1$ and $c > 0$.

6 A Central Limit Theorem

We now wish to establish the asymptotic distribution of the autoregressive estimator $\hat{\phi}_h$, or equivalently $\bar{\phi}_h$, under the regularity conditions considered in this paper. The difficulty is that the convergence rate of the autocovariance estimates upon which the coefficient estimators are based depends on the value of d . If $-0.5 < d < 0.25$ then the asymptotic distribution of $T^{\frac{1}{2}}(c_T(\tau) - \gamma(\tau))$ is normal, but when $d \geq 0.25$ the autocovariances are no longer \sqrt{T} consistent. See Hosking (1996) for details. Given that in practice the value of d will not be known, we seek a transformation that will lead to a conventional \sqrt{T} consistent, asymptotic normal approximation in which the parameters of the approximating distribution can be determined without explicit knowledge of d .

An interesting feature of the autocovariances noted by Hosking (1996, p. 268) is that when $d \in [0.25, 0.5)$ they contain a common slowly varying component that can be removed by differencing. Indeed, in this case $T^{\frac{1}{2}}(c_T(\tau) - \gamma(\tau)) = T^{2d-\frac{1}{2}}\varrho_T - \zeta_T(\tau)$ where ϱ_T and $\zeta_T(\tau)$ have non-degenerate limiting distributions, a Rosenblatt process and a Normal distribution respectively. Thus, if $v(\tau) = \gamma(\tau) - \gamma(0)$ and $u_T(\tau) = c_T(\tau) - c_T(0)$ then from Hosking (1996, Theorem 5) it follows that $T^{\frac{1}{2}}\{u_T(\tau) - v(\tau)\}$, for $\tau = 1, \dots, h$, have a non-degenerate multivariate Normal limiting distribution with mean zero and covariance matrix

$$\Delta_h = \left[\frac{1}{2} \sum_{s=-\infty}^{\infty} (\gamma(s) - \gamma(s-k) - \gamma(s+l) + \gamma(s-k+l))^2 + K_4 v(k)v(l) \right]_{k,l=1,\dots,h}$$

where K_4 is the fourth cumulant of $\varepsilon(t)$, whatever is the value of d . This suggests that some form of differencing, or centering, may be necessary to achieve our desired outcome and ultimately gives rise to the following result.

Theorem 6.1 *Let $\mathbf{C}_h = \mathbf{I}_h - h^{-1}\mathbf{1}\mathbf{1}'$ denote the h th order centering matrix where $\mathbf{1} = (1, 1, \dots, 1)'$ is the h element sum vector. Set $\Phi_h = \mathbf{I}_h + \mathbf{P}_h$ where \mathbf{P}_h equals*

$$\begin{bmatrix} \phi_h(2) & \phi_h(3) & \cdots & \cdots & \phi_h(h) & 0 \\ \phi_h(3) & \phi_h(4) & \cdots & \phi_h(h) & 0 & 0 \\ \vdots & \vdots & & 0 & 0 & 0 \\ \phi_h(h-1) & \phi_h(h) & & & & \\ \phi_h(h) & 0 & \cdots & \cdots & 0 & 0 \\ 0 & 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & \cdots & \cdots & 0 & 0 \\ \phi_h(1) & 0 & \cdots & \cdots & 0 & 0 \\ \phi_h(2) & \phi_h(1) & 0 & \cdots & \cdots & 0 \\ \vdots & & & & & \vdots \\ \phi_h(h-2) & \phi_h(h-3) & \cdots & \phi_h(1) & 0 & 0 \\ \phi_h(h-1) & \phi_h(h-2) & \cdots & \cdots & \phi_h(1) & 0 \end{bmatrix}.$$

Then for any h component vectors λ_h , where $1 \leq h \leq M_T$, $M_T = [c(\log T)^a]$, $a \geq 1$, $c > 0$, such that $0 < \|\lambda_h\| < \infty$ the scalars $T^{\frac{1}{2}}\lambda_h' \mathbf{C}_h \Gamma_h (\bar{\phi}_h - \phi_h)$ form a triangular array equal to $\beta_{h,T} + \rho_{h,T}$ where $\rho_{h,T} = o_p(1)$ and $\beta_{h,T}/\eta_h \xrightarrow{\mathcal{L}} N(0, 1)$ where $\eta_h^2 = \lambda_h' (\mathbf{C}_h \Phi_h \Delta_h \Phi_h' \mathbf{C}_h) \lambda_h$.

A corollary of Theorem 6.1, that follows from Bernstein's Lemma, is that if $\lambda_h = \mathbf{C}_h \lambda_h$ then a zero mean normal distribution with variance $\lambda_h' (\Phi_h \Delta_h \Phi_h') \lambda_h$ can be used as an asymptotic approximation to the large sample distribution of $T^{\frac{1}{2}}\lambda_h' \Gamma_h (\bar{\phi}_h - \phi_h)$. The condition that $\lambda_h = \mathbf{C}_h \lambda_h$ implies, of course, that the elements of λ_h must sum to zero.

7 Empirical Illustrations

This section of the paper reports the outcome of some simulation experiments designed to illustrate the theoretical results and properties discussed above. The experiments are based on three data generating mechanisms, the non-invertible moving average process $y(t) = \varepsilon(t) - \varepsilon(t-1)$ of Example (i) and two cases of the fractional noise process $y(t) = \varepsilon(t)/(1-z)^d$ of Example (ii) with $d = 0.125$ and 0.375 . For all three processes $\varepsilon(t)$ is standardized, Gaussian white noise. For each process the sample sizes $T = 100, 200, 500, 1000$ were considered and the values and figures presented here are all based on $R = 1000$ replications. In light of the discussion in Section 3, the behaviour of both the Yule-Walker and least squares estimates is examined.¹

Figure 1 presents the relative frequency of occurrence of the different orders given by h_T^{AIC} when $T = 100$ and the value $M_T = 2\sqrt{T} = [(\log T)^{1.962}] = 20$ is employed. At this sample size the dispersion of h_T^{AIC} about h_T^* is quite large, for all three processes, and there are no obvious differences in the observed performance of the two estimators. As T increases, however, first, the orders chosen by h_T^{AIC} become more concentrated around h_T^* , in accord with the predictions of Section 5, and secondly, the values of h_T^{AIC} produced by the Yule-Walker estimator are generally smaller than those given by least squares. The latter feature is shown in Table 1, which gives the average value of h_T^{AIC} compared to h_T^* for each model and sample size. Some indication of why the Yule-Walker procedure produces smaller values of h_T^{AIC} than least squares can be found in Table 2, which presents the empirical variance and the empirical bias of the estimates of the partial autocorrelation $\phi_h(h)$ for $h = h_T^*$. For all three processes and at all sample sizes the Yule-Walker estimate exhibits a larger bias than does least squares. The bias pushes the estimate of $\phi_h(h)$ towards the origin, leading to smaller values of h being selected. This behaviour is most noticeable in the case of the

¹The properties of the estimation procedure proposed by Burg (1968) were also investigated. This produces an estimate of $\phi_h(z)$ that is, like the Yule-Walker estimate, stable, but the finite sample properties of Burg's estimator mimic those of the least squares estimate. Detailed particulars of Burg's algorithm and other features of the simulations reported here can be found in Grose and Poskitt (2005), *Empirical Evidence on Nonstandard Autoregressive Approximations*, URL: <http://www-personal.buseco.monash.edu.au/~sgrose/papers/AREmpirical.pdf>, where a more extensive range of simulation experiments are documented.

Figure 1: Relative frequency of occurrence of h_T^{AIC} , $T = 100$, for (a) $y(t) = \varepsilon(t) - \varepsilon(t - 1)$, (b) $y(t) = \varepsilon(t)/(1 - z)^{0.125}$ and (c) $y(t) = \varepsilon(t)/(1 - z)^{0.375}$.

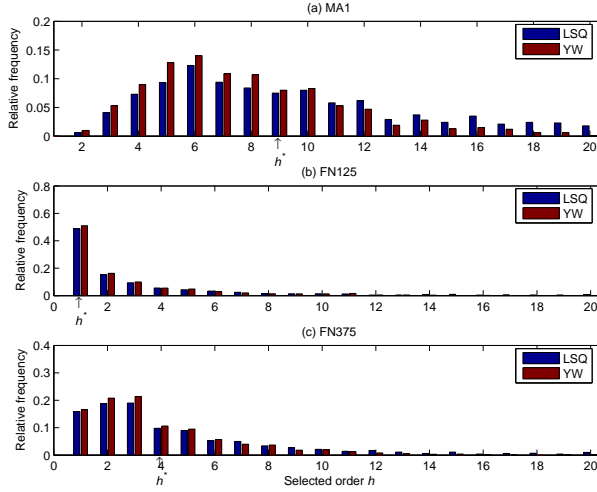


Table 1: Average value of h_T^{AIC} compared to h_T^* .

Process	T	h_T^*	h_T^{AIC} (LS)	h_T^{AIC} (YW)
$y(t) = \varepsilon(t) - \varepsilon(t - 1)$	100	9	9.218	7.855
	200	13	13.209	11.936
	500	21	21.208	19.904
	1000	31	30.993	29.29
$(1 - z)^{0.125}y(t) = \varepsilon(t)$	100	1	3.27	2.716
	200	1	3.181	2.982
	500	2	4.224	4.091
	1000	4	5.422	5.331
$(1 - z)^{0.375}y(t) = \varepsilon(t)$	100	4	4.658	3.916
	200	5	6.078	5.484
	500	8	8.694	8.326
	1000	12	11.992	11.72

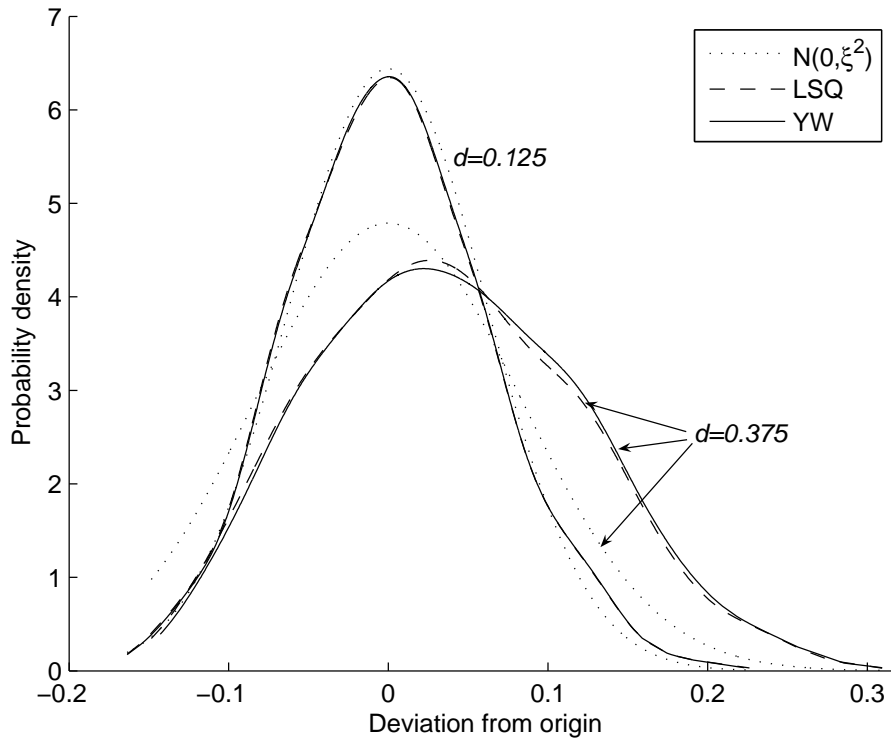
non-invertible moving average process, where the bias of the Yule-Walker estimates exceeds that of least squares by an order of magnitude even when $T = 1000$.

Figures 2 and 3 illustrate the impact of the distributional properties discussed in Section 6 for the two fractional noise processes. Figure 2 plots the empirical distribution of $h^{-1} \sum_{j=1}^h (\bar{\phi}_h(j) - \phi(j))$, the average deviation or coefficient error of the Yule-Walker estimator, together with the empirical distribution of $h^{-1} \sum_{j=1}^h (\hat{\phi}_h(j) - \phi(j))$, the average deviation of the least squares estimator, when $h = h_T^*$ and $T = 500$. The density estimates are constructed from the simulated values using a Gaussian kernel with bandwidth equal to 75% of the over-smoothed bandwidth i.e., $0.75 \xi \sqrt{(243/35R)}$ where ξ is the standard

Table 2: Partial Autocorrelation Estimates for $h = h_T^*$.

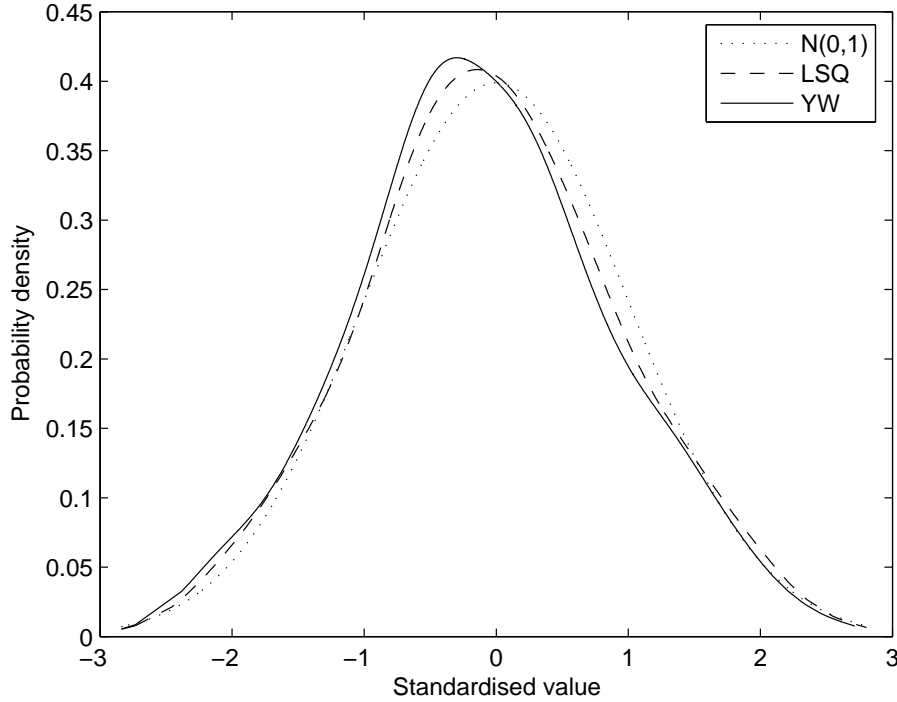
Process	T	$\phi_h(h)$	LS		YW	
			Variance	Bias	Variance	Bias
$y(t) = \varepsilon(t) - \varepsilon(t - 1)$	100	0.1	0.010335	-0.00278	0.008645	-0.016423
	200	0.071429	0.004958	-0.00151	0.004601	-0.011105
	500	0.045455	0.001991	-0.00028	0.001871	-0.005715
	1000	0.031250	0.001059	-0.00087	0.001041	-0.004680
$(1 - z)^{0.125}y(t) = \varepsilon(t)$	100	-0.142857	0.011214	0.001727	0.010982	0.003122
	200	-0.142857	0.005617	0.000256	0.005558	0.000974
	500	-0.066667	0.002027	0.002196	0.002014	0.002440
	1000	-0.032258	0.001051	0.002046	0.001043	0.002185
$(1 - z)^{0.375}y(t) = \varepsilon(t)$	100	-0.103448	0.012071	0.028433	0.011374	0.035746
	200	-0.081081	0.005793	0.009894	0.005540	0.013757
	500	-0.049180	0.002156	0.007040	0.002104	0.008661
	1000	-0.032258	0.001024	0.003841	0.001007	0.004590

Figure 2: Empirical distribution of $h^{-1} \sum_{j=1}^h (\bar{\phi}_h(j) - \phi(j))$ for fractional noise processes $y(t) = \varepsilon(t)/(1 - z)^d$ with $d = 0.125$ and $d = 0.375$, $h = h_T^*$ and $T = 500$.



deviation observed over the replications, see Wand and Jones (1995). Comparison of the estimated distributions to a normal curve of error with zero mean and variance ξ^2 indicates that when $d = 0.125$ the distribution of the average coefficient error is reasonably close

Figure 3: Observed distribution of $\varphi_{\lambda,T}$ for $y(t) = \varepsilon(t)/(1-z)^{0.375}$, when $\lambda'_h = (1, 0, \dots, 0, -1)$, $T = 1000$.



to normal for both estimators. When $d = 0.375$, however, the presence of the Rosenblatt process in the limiting behaviour of the underlying statistics is manifest in a marked distortion in the distribution relative to the shape anticipated of a normal random variable, particularly in the right hand tail of the distribution. This distortion is still present when $T = 1000$ and does not disappear asymptotically. By way of contrast, Figure 3 plots the observed distributions of $\bar{\varphi}_{\lambda,T} = T^{\frac{1}{2}} \lambda'_h \Gamma_h (\bar{\phi}_h - \phi_h) / (\lambda'_h (\Phi_h \Delta_h \Phi'_h) \lambda_h)^{\frac{1}{2}}$, $h = h_T^*$, and, to use an obvious notation, $\hat{\varphi}_{\lambda,T}$, obtained from realizations of the process $y(t) = \varepsilon(t)/(1-z)^{0.375}$ when $\lambda'_h = (1, 0, \dots, 0, -1)$ and $T = 1000$. The empirical distributions are overlaid with a standard normal density. Although some bias is still apparent even at this sample size, more so for $\bar{\phi}_h$ than $\hat{\phi}_h$, kurtosis and skewness of the type observed previously with this process has now gone and the operation of Theorem 6.1 is apparent.

Acknowledgement: This research was partially supported by the Australian Research Council under Grant DP0452717. The computations leading to the results presented in Section 7 were carried out by Simone Grose using MATLAB. I am grateful to Simone Grose, Andy Tremayne and Jon Wellner for helpful and constructive comments on a previous version of this paper.

Appendix:Proofs

Proof of Lemma 2.1: See Anderson (1971, Theorem 7.6.6). ■

Proof of Theorem 4.1: Assume that $\frac{1}{4} < d < \frac{1}{2}$. By Theorem 3 of Hosking (1996) $E[(c_T(\tau) - \gamma(\tau))^2] = O(T^{-2(1-2d)})$ and from Chebychev's inequality

$$Pr \left(|c_T(\tau) - \gamma(\tau)| > \delta(\log T/T)^{\frac{1}{2}-d} \right) \leq \frac{C}{\delta^2} \frac{1}{(T \log T)^{1-2d}}.$$

Now set $\Delta_\tau(T) = (c_T(\tau) - \gamma(\tau))(\log T/T)^{d-\frac{1}{2}}$. Then for $T' = T^{4/(1-2d)}$

$$\sum_{T=1}^{\infty} Pr \left(\max_{|\tau| \leq H_{T'}} |\Delta_\tau(T')| > \delta \right) \leq \frac{C}{\delta^2} \sum_{T=1}^{\infty} \left(\frac{1-2d}{4 \log T} \right)^{3(1-2d)/2} \frac{1}{T^2} < \infty$$

and by the Borel-Cantelli lemma $\Delta_\tau(T') \rightarrow 0$ a.s. uniformly in τ , $|\tau| \leq H_{T'}$.

Let $N^2 = T'$. Then for all T such that $N^2 < T < (N+1)^2$

$$\begin{aligned} \Delta_\tau(T) &= \left(\frac{T}{\log T} \frac{\log(N+1)^2}{(N+1)^2} \right)^{\frac{1}{2}(1-2d)} \frac{(N+1)^2}{T} \Delta_\tau((N+1)^2) \\ &\quad - \left(\frac{T}{\log T} \right)^{\frac{1}{2}(1-2d)} \frac{1}{T} \sum_{t=T}^{(N+1)^2} (y(t)y(t-\tau) - \gamma(\tau)) \end{aligned}$$

and

$$\begin{aligned} \max_{|\tau| \leq H_T} |\Delta_\tau(T)| &\leq \max_{|\tau| \leq H_T} \left(\frac{T}{\log T} \frac{\log(N+1)^2}{(N+1)^2} \right)^{\frac{1}{2}(1-2d)} \left(1 + \frac{1}{N} \right)^2 |\Delta_\tau((N+1)^2)| \\ &\quad + \max_{|\tau| \leq H_T} \left(\frac{T}{\log T} \right)^{\frac{1}{2}(1-2d)} \frac{1}{T} \left| \sum_{t=T}^{(N+1)^2} (y(t)y(t-\tau) - \gamma(\tau)) \right|. \end{aligned} \quad (\text{A.1})$$

But

$$\frac{T}{\log T} \frac{\log(N+1)^2}{(N+1)^2} \leq \frac{\log(N+1)}{\log N} \rightarrow 1 \text{ as } N \rightarrow \infty$$

and by what has already been shown it follows that the first term on the right hand side of (A.1) converges to zero a.s..

Moreover, using Chebychev's inequality once more we can bound

$$Pr \left(\max_{|\tau| \leq H_T} \left| \sum_{t=T}^{(N+1)^2} (y(t)y(t-\tau) - \gamma(\tau)) \right| \geq \delta(\log T)^{\frac{1}{2}(1-2d)} T^{\frac{1}{2}(1+2d)} \right) \quad (\text{A.2})$$

by

$$\left(\frac{(N+1)^2}{\log(N+1)^2} \right)^{\frac{1}{2}(1-2d)} \cdot \frac{C}{\delta^2(\log T)^{(1-2d)} T^{(1+2d)}} \cdot (2N+1)^{4d}$$

where the first factor accounts for the maximum being taken over $H_T < H_{(N+1)^2}$ terms and

the last factor arises because the sum contains $(N + 1)^2 - T + 1 < (2N + 1)$ summands. Thus we can deduce that the probability in (A.2) is less than

$$\frac{(2N + 1)^{4d}(N + 1)^{(1-2d)}}{(2 \log N)^{3(1-2d)/2} N^{2+4d}} \leq \frac{18}{N^{1+2d}}$$

and hence, via the Borel-Cantelli lemma, that the second term of (A.1) converges to zero since the series $\{N^{-(1+2d)}\}$ is convergent.

A similar proof using the method of subsequences can be employed to establish the result for the remaining cases, $d = \frac{1}{4}$, when $E[(c_T(\tau) - \gamma(\tau))^2] = O(\log T/T)$, and $d \in (-\frac{1}{2}, \frac{1}{4})$, when $E[(c_T(\tau) - \gamma(\tau))^2] = O(T^{-1})$. \blacksquare

Proof of Theorem 4.2: The following relationship exists between the elements of the sequence $c_T(\tau)$ and the $c_T(j, k)$ for $j - k = \tau = 0, \pm 1, \dots, \pm H_T$:

$$T\{c_T(\tau) - c_T(j, k)\} = \sum_{s=B_T(j,k)}^T y(s - |\tau|)y(s) = D_T(j, k, \tau) \quad (\text{A.3})$$

where $B_T(j, k) = T + 1 - \min\{j, k\}$. Note that $D_T(j, k, \tau)$ contains $\min\{j, k\}$, or at most H_T , summands. Now, since $\varepsilon(t)$ has finite fourth moment

$$E[y(t)^4] = \sigma^4 \left(\sum_{j=0}^{\infty} k(j)^2 \right)^2 + K_4 \sum_{j=0}^{\infty} k(j)^4 < \infty,$$

where K_4 is the fourth cumulant of $\varepsilon(t)$, and the variance of $D_T(j, k, \tau)$ is dominated by CH_T^2 uniformly in $j, k = 0, 1, \dots, H_T$. Thus

$$\begin{aligned} Pr \left(\max_{|\tau| \leq H_T} |D_T(j, k, \tau)| \geq \delta T \right) &< H_T \frac{CH_T^2}{\delta^2 T^2} \\ &\leq \frac{C}{(\log T)^{3/2} T^{\frac{1}{2}}} \end{aligned} \quad (\text{A.4})$$

where the final inequality follows since for $0 < d < \frac{1}{2}$ $H_T = o\{(T/\log T)^{\frac{1}{2}}(\log T/T)^d\}$ and $(\log T/T)^d < 1$ and $H_T = o\{(T/\log T)^{\frac{1}{2}}\}$ for $-\frac{1}{2} < d \leq 0$.

Along the subsequence $T' = T^4$ it follows that $\lim_{T \rightarrow \infty} \max_{|\tau| \leq H_{T'}} T'^{-1} |D_{T'}(j, k, \tau)| = 0$ a.s. because $\{T^{-2}\}$ is a convergent series. Furthermore, letting $N^2 = T'$, then for all T in between N^2 and $(N + 1)^2$ we can bound $|N^{-2}D_{N^2}(j, k, \tau) - T^{-1}D_T(j, k, \tau)|$ by

$$\left| \frac{(T - N^2)D_{N^2}(j, k, \tau)}{TN^2} \right| + \left| \frac{\sum_{s=B_{N^2}(j,k)}^{B_T(j,k)} y(s - |\tau|)y(s) - \sum_{s=N^2+1}^T y(s - |\tau|)y(s)}{T} \right|. \quad (\text{A.5})$$

The first term in (A.5) converges to zero uniformly in j, k and τ by what has already been proved since $(T - N^2)/TN^2 \leq (2N + 1)/N^4$ and the second term converges similarly via an application of Chebychev's inequality and the Borel-Cantelli lemma.

To show the latter, consider the case $d \in (-\frac{1}{2}, \frac{1}{4})$ for example. By Theorems 4.1 and 4.2 and Theorem 3 of Hosking (1996) the variance of the numerator can be bounded by $C(T - N^2) \leq C(2N + 1)$ uniformly in $j, k = 0, 1, \dots, H_T$ so

$$Pr \left(\max_{|\tau| \leq H_T} \left| \sum_{s=B_{N^2}(j,k)}^{B_T(j,k)} y(s - |\tau|)y(s) - \sum_{s=N^2+1}^T y(s - |\tau|)y(s) \right| \geq \delta T \right) < H_T \cdot \frac{C(2N + 1)}{\delta^2 T^2}$$

and $H_T(2N + 1)/T^2 \leq 2(N + 1)(2N + 1)/N^4 \log(N + 1) < 6/N^2$. \blacksquare

For convenience and completeness we now state a result taken from Poskitt (2000).

Lemma A.1 : Let \mathbf{A}_T and \mathbf{B}_T denote two $h \times h$ (stochastic) matrices such that $\|\mathbf{A}_T - \mathbf{B}_T\|$ equals $O(C_T)$ where $C_T \rightarrow 0$ as $T \rightarrow \infty$ and suppose that $\liminf_{T \rightarrow \infty} \lambda_{\min}[\mathbf{B}_T] \geq \delta_h > 0$. Then \mathbf{A}_T is nonsingular for all T sufficiently large and $\|\mathbf{A}_T^{-1} - \mathbf{B}_T^{-1}\| = (\delta_h)^{-2}O(C_T)$.

Proof of Corollary 4.1: Let $\bar{\Gamma}_h = [c_T(i - j)]_{i,j=1,\dots,h}$ and $\bar{\gamma}_h = (c_T(1), \dots, c_T(h))'$. Then

$$\begin{aligned} \bar{\phi}_h - \hat{\phi}_h &= \bar{\Gamma}_h^{-1} \bar{\gamma}_h - \mathbf{M}_h^{-1} \mathbf{m}_h \\ &= (\bar{\Gamma}_h^{-1} - \mathbf{M}_h^{-1}) \mathbf{m}_h + \bar{\Gamma}_h^{-1} (\bar{\gamma}_h - \mathbf{m}_h). \end{aligned} \quad (\text{A.6})$$

From Theorem 4.1 it follows that

$$\limsup_{T \rightarrow \infty} \|\Gamma_h - \bar{\Gamma}_h\|^2 = O\{h^2(\log T/T)^{1-2d'}\} = o(1)$$

and hence that

$$\begin{aligned} \liminf_{T \rightarrow \infty} \lambda_{\min}(\bar{\Gamma}_h) &\geq \lambda_{\min}(\Gamma_h) - \limsup_{T \rightarrow \infty} \|\Gamma_h - \bar{\Gamma}_h\| \\ &= \lambda_{\min}(\Gamma_h) > 0. \end{aligned}$$

Using the same argument in conjunction with Theorem 4.2 it can also be shown that $\liminf_{T \rightarrow \infty} \lambda_{\min}(\mathbf{M}_h) \geq \lambda_{\min}(\Gamma_h)$. From Lemma A.1 it follows that $\|\bar{\Gamma}_h^{-1} - \mathbf{M}_h^{-1}\| = O\{(h/\lambda_{\min}(\Gamma_h^2))(\log T/T)^{\frac{1}{2}(1-2d')}\}$. Now the first term on the right hand side of (A.6) can be bounded in norm by $\|(\bar{\Gamma}_h^{-1} - \mathbf{M}_h^{-1})\| \cdot (\|\gamma_h\| + \|\mathbf{m}_h - \gamma_h\|)$, which equals

$$O\{(h/\lambda_{\min}(\Gamma_h^2))(\log T/T)^{\frac{1}{2}(1-2d')}(h^{\frac{1}{2}(4d-1)} + h^{\frac{1}{2}}(\log T/T)^{\frac{1}{2}(1-2d')})\},$$

and the norm of the second term of (A.6) is $O\{(h^{\frac{1}{2}}/\lambda_{\min}(\Gamma_h))(\log T/T)^{\frac{1}{2}-d'}\}$. \blacksquare

Proof of Theorem 4.3: Let $\rho(z) = \sum_{j \geq 1} \rho(j)z^j = \phi_h(z)k(z) - 1$. Then $\epsilon_h(t) - \varepsilon(t) = \sum_{j \geq 1} \rho(j)\varepsilon(t - j)$. From Parseval's relation

$$\sum_{j \geq 1} \rho(j)^2 = \int_{-\pi}^{\pi} |\phi_h(e^{i\omega})k(e^{i\omega}) - 1|^2 d\omega = 2\pi\sigma^{-2}E[(\epsilon_h(t) - \varepsilon(t))^2] < \infty$$

and therefore we can conclude that $T^{-1} \sum_{t=1}^T \varepsilon(t)\{\epsilon_h(t) - \varepsilon(t)\} = O\{(\log \log T/T)^{\frac{1}{2}}\}$ by

Theorem 5.3.5. of Hannan and Deistler (1988). ■

Proof of Theorem 4.4: By definition $\epsilon_h(t) = \sum_{j=0}^h \phi_h(j)y(t-j)$. Simple substitution now gives us

$$T^{-1} \sum_{t=1}^T \epsilon_h(t)y(t-r) = \sum_{j=0}^h \phi_h(j)T^{-1} \sum_{t=1}^T y(t-j)y(t-r) = \sum_{j=0}^h \phi_h(j)c_T(j,r),$$

which by Theorem 4.2 equals

$$\sum_{j=0}^h \phi_h(j)[\gamma(j-r) + O\{(\log T/T)^{\frac{1}{2}-d'}\}].$$

Since $\phi_h(j)$, $j = 1, \dots, h$, solve the Yule-Walker equations $\sum_{j=0}^h \phi_h(j)\gamma(j-r) = 0$ for $r = 1, \dots, h$. Moreover, $\phi_h(z) \neq 0$, $|z| \leq 1$, and there exists constants $C < \infty$ and $\zeta < 1$ such that $|\phi_h(j)| < C\zeta$ and $\sum_{j=0}^h |\phi_h(j)| < C(1 - \zeta^{h+1})/(1 - \zeta) < C(1 - \zeta)^{-1}$ so that $\sum_{j=0}^h \phi_h(j)O\{(\log T/T)^{\frac{1}{2}-d'}\} = O\{(\log T/T)^{\frac{1}{2}-d'}\}$. Hence the desired result. ■

Proof of Theorem 5.1: Substituting $\epsilon_h(t) = \sum_{j=0}^h \phi_h(j)y(t-j)$ into the normal equations yields the expression

$$\mathbf{M}_h(\hat{\phi}_h - \phi_h) = T^{-1} \sum_{t=1}^T \epsilon_h(t) \begin{bmatrix} y(t-1) \\ \vdots \\ y(t-h) \end{bmatrix}.$$

It follows that

$$\|\hat{\phi}_h - \phi_h\|^2 \leq \frac{1}{\lambda_{\min}(\mathbf{\Gamma}_h^2)} \sum_{j=1}^h \left(T^{-1} \sum_{t=1}^T \epsilon_h(t)y(t-j) \right)^2$$

and hence that $\|\hat{\phi}_h - \phi_h\|^2 = (1/\lambda_{\min}(\mathbf{\Gamma}_h^2))O\{h(\log T/T)^{1-2d'}\}$ by Theorem 4.2. ■

Proof of Theorem 5.2: From the definition of $\hat{\epsilon}_h(t)$ and $\epsilon_h(t)$ we get

$$\hat{\epsilon}_h(t) - \epsilon_h(t) = \sum_{j=1}^h \{\hat{\phi}_h(j) - \phi_h(j)\}y(t-j)$$

and from the Cauchy-Schwartz inequality, Theorem 4.4 and Theorem 5.1 we have

$$\begin{aligned} |T^{-1} \sum_{t=1}^T \epsilon_h(t)\{\hat{\epsilon}_h(t) - \epsilon_h(t)\}| &= |T^{-1} \sum_{t=1}^T \sum_{j=1}^h \{\hat{\phi}_h(j) - \phi_h(j)\}\epsilon_h(t)y(t-j)| \\ &\leq \left[\|\hat{\phi}_h - \phi_h\|^2 \cdot \sum_{j=1}^h \left(T^{-1} \sum_{t=1}^T \epsilon_h(t)y(t-j) \right)^2 \right]^{\frac{1}{2}} \\ &= O \left\{ \frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h)} \left(\frac{\log T}{T} \right)^{1-2d'} \right\}, \end{aligned} \quad (\text{A.7})$$

giving the result of the theorem. \blacksquare

Proof of Theorem 5.3: Since $h/T \rightarrow 0$ as $T \rightarrow \infty$ and $\sigma_h^2 - \sigma^2$ is monotonically decreasing in h it follows that $h_T^* \rightarrow \infty$ as $T \rightarrow \infty$. Similarly, for T sufficiently large the behaviour of

$$\log \left(1 + \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sum_{t=1}^T \varepsilon(t)^2} \right)$$

will be determined by that of $\log(1 + (\sigma_h^2 - \sigma^2)/\sigma^2)$. The latter is decreasing in h and it follows that $\bar{h}_T^* \rightarrow \infty$ as $T \rightarrow \infty$. Indeed, expanding $\bar{L}_T(h)$ using $\log(1+x) = \sum_{r \geq 1} (-)^{r-1} x^r / r$ and recognizing from Lemma 2.1 and equation (4.3) that $T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 - T^{-1} \sum_{t=1}^T \varepsilon(t)^2 = E[(\epsilon_h(t) - \varepsilon(t))^2] + o(1)$ will converge to zero as h increases we find that

$$\bar{L}_T(h) = \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sum_{t=1}^T \varepsilon(t)^2} + \frac{h}{T} + o \left\{ \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sum_{t=1}^T \varepsilon(t)^2} \right\}$$

and

$$\begin{aligned} \left| \bar{L}_T(h) - \frac{L_T(h)}{\sigma^2} \right| &\leq \left| \left(\frac{\sigma^2}{\sum_{t=1}^T \varepsilon(t)^2} \right) \left(\frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sigma_h^2 - \sigma^2} \right) - 1 \right| \left(\frac{\sigma_h^2 - \sigma^2}{\sigma^2} \right) \\ &\quad + o \left\{ \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sum_{t=1}^T \varepsilon(t)^2} \right\} \\ &= o \{ (\sigma_h^2 - \sigma^2) / \sigma^2 \}. \end{aligned} \tag{A.8}$$

From (A.8) we conclude that

$$\left| \frac{\sigma^2 \bar{L}_T(h)}{L_T(h)} - 1 \right| = \left(\frac{\sigma^2}{(\sigma_h^2 - \sigma^2) + h/T} \right) o \{ (\sigma_h^2 - \sigma^2) / \sigma^2 \} = o(1).$$

By definition of \bar{h}_T^* and h_T^* as the minimizing values of, respectively, $\bar{L}_T(h)$ and $L_T(h)$ over the common range $h = 0, 1, \dots, M_T$ it now follows that

$$\frac{\sigma^2 \bar{L}_T(\bar{h}_T^*)}{L_T(h_T^*)} = \frac{\bar{L}_T(\bar{h}_T^*)}{\bar{L}_T(h_T^*) \{1 + o(1)\}} \leq 1 + o(1)$$

and

$$\frac{\sigma^2 \bar{L}_T(\bar{h}_T^*)}{L_T(h_T^*)} = \frac{L_T(\bar{h}_T^*) \{1 + o(1)\}}{L_T(h_T^*)} \geq 1 + o(1),$$

which implies that

$$\left| \frac{\sigma^2 \bar{L}_T(\bar{h}_T^*)}{L_T(h_T^*)} - 1 \right| = o(1),$$

as required. \blacksquare

Proof of Theorem 5.4: The least squares residual $\hat{\epsilon}_h(t)$ is by construction orthogonal to

$y(t-1), \dots, y(t-h)$ for $t = 1, \dots, T$ and thus

$$T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t) \{\hat{\epsilon}_h(t) - \epsilon_h(t)\} = \sum_{j=1}^h \{\hat{\phi}_h(j) - \phi_h(j)\} T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t) y(t-j) = 0.$$

The residual mean square can therefore be re-expressed as

$$T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t)^2 = T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 + T^{-1} \sum_{t=1}^T \epsilon_h(t) \{\hat{\epsilon}_h(t) - \epsilon_h(t)\}$$

and the right hand side equals

$$T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 + O \left\{ \frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h)} \left(\frac{\log T}{T} \right)^{1-2d'} \right\}$$

by Theorem 5.2. A trivial re-expression of $T^{-1} \sum_{t=1}^T \epsilon_h(t)^2$ as the sum of $T^{-1} \sum_{t=1}^T \varepsilon(t)^2$ and $T^{-1} \sum_{t=1}^T \epsilon_h(t)^2 - T^{-1} \sum_{t=1}^T \varepsilon(t)^2$, used in conjunction with the usual McLaurin expansion of $\log(1+x)$ as given above, now yields the result that

$$\begin{aligned} \log T^{-1} \sum_{t=1}^T \hat{\epsilon}_h(t)^2 &= \log T^{-1} \sum_{t=1}^T \varepsilon(t)^2 + \log \left(1 + \frac{\sum_{t=1}^T \epsilon_h(t)^2 - \sum_{t=1}^T \varepsilon(t)^2}{\sum_{t=1}^T \varepsilon(t)^2} \right) \\ &\quad + O \left\{ \left(\frac{\sum_{t=1}^T \varepsilon(t)^2}{T} \right)^{-1} \frac{h}{\lambda_{\min}(\mathbf{\Gamma}_h)} \left(\frac{\log T}{T} \right)^{1-2d'} \right\}. \end{aligned}$$

But $T^{-1} \sum_{t=1}^T \varepsilon(t)^2$ converges to σ^2 a.s., giving the result as stated in the theorem. \blacksquare

Proof of Lemma 5.7: Let \mathbf{x} denote a unit eigenvector associated with the eigenvalue $\lambda_{\min}(\mathbf{\Gamma}_h)$. Then $\lambda_{\min}(\mathbf{\Gamma}_h) = \mathbf{x}' \mathbf{\Gamma}_h \mathbf{x}$. Now consider:

Case (i), $u(z) \equiv 1$ and $d \geq 0$. Then $f(\omega) = \sigma^2 |\mu(e^{i\omega})|^2 / 2\pi |1 - e^{i\omega}|^{2d}$ and

$$\begin{aligned} \lambda_{\min}(\mathbf{\Gamma}_h) &= \int_{-\pi}^{\pi} \left| \sum_{s=1}^h x_s \exp(-i\omega s) \right|^2 f(\omega) d\omega \\ &\geq \inf_{\omega} f(\omega) \int_{-\pi}^{\pi} \left| \sum_{s=1}^h x_s \exp(-i\omega s) \right|^2 d\omega = \inf_{\omega} f(\omega) > 0. \end{aligned}$$

Case (ii), $u(z) \equiv 1$ and $d < 0$. An adaptation of the circulant imbedding argument underlying the simulation technique of Davies and Harte (1987) yields the result that $\mathbf{\Gamma}_h = \mathbf{U}^* \mathbf{\Lambda} \mathbf{U}$ where the $(2h+1) \times h$ matrix

$$\mathbf{U} = [(2h+1)^{-\frac{1}{2}} \exp(-i2\pi(j-1)(k-1)/(2h+1))]_{j=1, \dots, 2h+1, k=1, \dots, h},$$

and $\mathbf{\Lambda} = 2\pi \text{diag}\{f_h(\omega_0), \dots, f_h(\omega_{2h})\}$,

$$f_h(\omega) = \frac{1}{2\pi} \sum_{\tau=-h}^h \gamma(\tau) \exp(-i\omega\tau), \quad \omega_j = 2\pi j/(2h+1), \quad j = 0, \dots, 2h,$$

as can be readily verified via straightforward, if somewhat tedious, algebra. Hence $\lambda_{\min}(\mathbf{\Gamma}_h) = \mathbf{x}'\mathbf{\Gamma}_h\mathbf{x} = \mathbf{w}^*\mathbf{\Lambda}\mathbf{w}$ where $\mathbf{w} = \mathbf{U}\mathbf{x}$ and $\|\mathbf{w}\| = 1$ since $\|\mathbf{x}\| = 1$ and $\mathbf{U}^*\mathbf{U} = \mathbf{I}$.

From the Rayleigh-Ritz theorem it follows that

$$\lambda_{\min}(\mathbf{\Gamma}_h) \geq \min\{2\pi f_h(\omega_0), \dots, 2\pi f_h(\omega_{2h})\}. \quad (\text{A.9})$$

But

$$f_h(\omega) = f(\omega) - \int_{-\pi}^{\pi} \{f(\omega) - f(\theta)\} D_h(\omega - \theta) d\theta$$

where $D_h(\theta) = \sin((h + \frac{1}{2})\theta)/2\pi \sin(\theta/2)$, Dirichlet's kernel, and since $f(\cdot)$ is absolutely integrable and continuous a.e. it follows from the Riemann-Lebesgue lemma that

$$\lim_{h \rightarrow \infty} \sup_{0 \leq \omega \leq \pi} \int_{-\pi}^{\pi} \{f(\omega) - f(\theta)\} D_h(\omega - \theta) d\theta = 0.$$

We can therefore conclude that for every $\delta > 0$

$$|\text{argmin}_{\omega \in \{\omega_0, \dots, \omega_{2h}\}} f_h(\omega) - \text{argmin}_{\omega \in \{\omega_0, \dots, \omega_{2h}\}} f(\omega)| < \delta \quad (\text{A.10})$$

for h sufficiently large. Now, $f(\omega_0) = 0$ and therefore

$$\min\{2\pi f_h(\omega_0), \dots, 2\pi f_h(\omega_{2h})\} = 2\pi f_h(\omega_0) = \sum_{\tau=-h}^h \gamma(\tau) = O(h^{2d})$$

for all h sufficiently large.

Case (iii). Assume for simplicity that $\nu_j > 0$ for all $j = 1, \dots, n$. Then $2\pi f(\omega)$ equals the product of $\sigma^2 |\mu(e^{i\omega})|^2 > 0$ and

$$|2 \sin(\omega/2)|^{2(\nu_1-d)} |2 \sin((\pi - \omega)/2)|^{2\nu_2} \prod_{j=3}^n |4 \sin((\theta_j + \omega)/2) \sin((\theta_j - \omega)/2)|^{2\nu_j}$$

and $f(\omega) = 0$ when $\omega = \theta_1 = 0$, when $\omega = \theta_2 = \pi$ and when $\omega = \theta_j$, $j = 3, \dots, n$.

From the argument employed in Case (ii) we know that (A.9) holds and that (A.10) obtains for h sufficiently large. Set $j_i(h) = [(2h+1)\theta_i/2\pi]$ for $i = 2, \dots, n$. Both $|\omega_{j_i(h)} - \theta_i|$ and $|\omega_{(j_i(h)+1)} - \theta_i|$ are less than $2\pi/(2h+1)$ and the Taylor expansion of $f(\cdot)$ about θ_i implies that $f(\omega_{j_i(h)}) = o(2\pi/(2h+1))$ and similarly $f(\omega_{(j_i(h)+1)}) = o(2\pi/(2h+1))$. Thus if

$\bar{\omega}_m = \operatorname{argmin}_{\omega \in \{\omega_0, \dots, \omega_{2h}\}} f_h(\omega)$ then $\bar{\omega}_m = \omega_0$ and

$$\min\{2\pi f_h(\omega_0), \dots, 2\pi f_h(\omega_{2h})\} = \sum_{\tau=-h}^h \gamma(\tau) = O(h^{2(d-\nu_1)})$$

or $\bar{\omega}_m$ equals either $\omega_{j_i(h)}$ or $\omega_{(j_i(h)+1)}$ for some $i \in \{2, \dots, n\}$ and

$$\min\{2\pi f_h(\omega_0), \dots, 2\pi f_h(\omega_{2h})\} = \sum_{\tau=-h}^h \int_{-\pi}^{\pi} g_i(\omega) e^{i\omega\tau} d\omega = O(h^{-2\nu_i})$$

where $g_i(\omega) = f(\bar{\omega}_m + \omega)$. It follows that $f_h(\bar{\omega}_m) = O(h^{-2m})$ where $m = \max\{\nu_1 - d, \nu_2, \nu_3, \dots, \nu_n\}$. \blacksquare

Proof of Theorem 6.1: By definition of $\bar{\phi}_h$

$$T^{\frac{1}{2}} \mathbf{C}_h \mathbf{\Gamma}_h (\bar{\phi}_h - \phi_h) = T^{\frac{1}{2}} \mathbf{C}_h (\bar{\gamma}_h - \gamma_h) - T^{\frac{1}{2}} \mathbf{C}_h (\bar{\mathbf{\Gamma}}_h - \mathbf{\Gamma}_h) \phi_h + T^{\frac{1}{2}} \mathbf{C}_h (\mathbf{\Gamma}_h - \bar{\mathbf{\Gamma}}_h) (\bar{\phi}_h - \phi_h).$$

Multiplying through by $\boldsymbol{\lambda}'_h$ and rearranging terms on the right hand side yields the result that $T^{\frac{1}{2}} \boldsymbol{\lambda}'_h \mathbf{C}_h \mathbf{\Gamma}_h (\bar{\phi}_h - \phi_h)$ equals $\boldsymbol{\beta}_{hT} + \boldsymbol{\rho}_{hT}$ where

$$\boldsymbol{\rho}_{hT} = T^{\frac{1}{2}} \boldsymbol{\lambda}'_h \mathbf{C}_h (\mathbf{\Gamma}_h - \bar{\mathbf{\Gamma}}_h) (\bar{\phi}_h - \phi_h)$$

and

$$\boldsymbol{\beta}_{hT} = T^{\frac{1}{2}} \boldsymbol{\lambda}'_h \mathbf{C}_h [(\bar{\gamma}_h - \gamma_h) - (\bar{\mathbf{\Gamma}}_h - \mathbf{\Gamma}_h) \phi_h].$$

Hence we are lead to consider the limiting behaviour of $\boldsymbol{\beta}_{hT}$ and $\boldsymbol{\rho}_{hT}$.

Let $\mathbf{D}_h = [u_T(i-j) - v(i-j)]_{i,j=1,\dots,h}$ and $\mathbf{d}_h = (u_T(1) - v(1), \dots, u_T(h) - v(h))'$. Then it is a simple exercise to show that $\mathbf{C}_h (\bar{\mathbf{\Gamma}}_h - \mathbf{\Gamma}_h) = \mathbf{C}_h \mathbf{D}_h$ and $\mathbf{C}_h (\bar{\gamma}_h - \gamma_h) = \mathbf{C}_h \mathbf{d}_h$ and it follows that $\boldsymbol{\beta}_{hT} = T^{\frac{1}{2}} \boldsymbol{\lambda}'_h \mathbf{C}_h [\mathbf{d}_h - \mathbf{D}_h \phi_h]$. Writing \mathbf{D}_h as $\mathbf{T}_h + \mathbf{T}'_h$ where \mathbf{T}_h is the lower triangular Toeplitz matrix with first column $(0, u_T(1) - v(1), \dots, u_T(h-1) - v(h-1))'$ we find, after some straightforward rearrangement, that $\mathbf{D}_h \phi_h = \mathbf{P}_h \mathbf{d}_h$ and hence that $\mathbf{d}_h - \mathbf{D}_h \phi_h = \boldsymbol{\Phi}_h \mathbf{d}_h$. Thus $\boldsymbol{\beta}_{h,T} = T^{\frac{1}{2}} \boldsymbol{\lambda}'_h \boldsymbol{\Phi}_h \mathbf{d}_h$. By Theorem 5 of Hosking (1996), however, $T^{\frac{1}{2}} \mathbf{d}_h$ converges in distribution to $N(\mathbf{0}, \boldsymbol{\Delta}_h)$. We can therefore conclude that $\boldsymbol{\beta}_{h,T} / \eta_h \xrightarrow{\mathcal{L}} N(\mathbf{0}, 1)$ where $\eta_h^2 = \boldsymbol{\lambda}'_h (\mathbf{C}_h \boldsymbol{\Phi}_h \boldsymbol{\Delta}_h \boldsymbol{\Phi}'_h \mathbf{C}_h) \boldsymbol{\lambda}_h$, as stated.

Similarly, $\boldsymbol{\rho}_{hT} = -T^{\frac{1}{2}} \boldsymbol{\lambda}'_h (\mathbf{C}_h \mathbf{D}_h) (\bar{\phi}_h - \phi_h)$. Corollary 4.1 and Theorem 5.1 imply that $\|(\bar{\phi}_h - \phi_h)\| = O(M_T^{1+4q} (\log T/T)^{1-2d'})$ where $q \geq 0$ and from Theorem 5 of Hosking (1996), once again, we have that $T^{\frac{1}{2}} \mathbf{D}_h = O_p(1)$. This leads to the conclusion that $\boldsymbol{\rho}_{hT} = o_p(1)$ and completes the proof. \blacksquare

References

- AKAIKE, H. (1969). Fitting autoregressive models for prediction. *Annals of Institute of Statistical Mathematics* **21** 243–247.
- AKAIKE, H. (1970). Statistical predictor identification. *Annals of Institute of Statistical Mathematics* **22** 203–217.
- ANDERSON, T. W. (1971). *The Statistical Analysis of Time Series*. J. Wiley, New York.
- BAILLIE, R. T. (1996). Long memory processes and fractional integration in econometrics. *Journal of Econometrics* **73** 5–59.
- BERAN, J. (1992). Long-range dependence. *Statistical Science* **7** 404–427.
- BERAN, J. (1994). *Statistics for Long Memory Processes*. Chapman and Hall, New York.
- BERAN, J. (1995). Maximum likelihood estimation of the differencing parameter for invertible short and long memory autoregressive integrated moving average models. *Journal of the Royal Statistical Society B* **57** 654–672.
- BERAN, J., BHANSALI, R. J. and OCKER, D. (1998). On unified model selection for stationary and nonstationary short- and long-memory autoregressive processes. *Biometrika* **85** 921–934.
- BERK, K. N. (1974). Consistent autoregressive spectral estimation. *Annals of Statistics* **2** 489–502.
- BURG, J. (1968). A new analysis technique for time series data. Tech. rep., Advanced Study Institute on Signal Processing, N.A.T.O., Enschede, Netherlands.
- DAVIES, R. B. and HARTE, D. S. (1987). Tests for hurst effect. *Biometrika* **74** 95–101.
- DURBIN, J. (1960). The fitting of time series models. *Review of International Statistical Institute* **28** 233–244.
- FOX, R. and TAQQU, M. S. (1986). Large sample properties of parameter estimates for strongly dependent stationary gaussian time series. *Annals of Statistics* **14** 517–532.
- GRANGER, C. W. J. and JOYEUX, R. (1980). An introduction to long-memory time series models and fractional differencing. *Journal of Time Series Analysis* **1** 15–29.
- GRENANDER, U. and ROSENBLATT, M. (1957). *Statistical Analysis of Stationary Times Series*. J. Wiley, New York.

- HANNAN, E. J. and DEISTLER, M. (1988). *The Statistical Theory of Linear Systems*. Wiley, New York.
- HANNAN, E. J. and QUIN, B. G. (1979). The determination of the order of an autoregression. *Journal of Royal Statistical Society B* **41** 190–195.
- HOSKING, J. R. M. (1980). Fractional differencing. *Biometrika* **68** 165–176.
- HOSKING, J. R. M. (1996). Asymptotic distributions of the sample mean, autocovariances, and autocorrelations of long memory time series. *Journal of Econometrics* **73** 261–284.
- KOLMORGOROV, A. N. (1941). Interpolation und extrapolation von stationären zufälligen folgen. *Bulletin Academy Science U. S. S. R., Mathematics Series* **5** 3–14.
- LEVINSON, N. (1947). The Wiener RMS (root mean square) error criterion in filter design and prediction. *Journal of Mathematical Physics* **25** 261–278.
- LYSNE, D. and TJØSTHEIM, D. (1987). Loss of spectral peaks in autoregressive spectral estimation. *Biometrika* **74** 200–206.
- MALLOWS, C. L. (1973). Some comments on C_p . *Technometrics* **15** 661–675.
- MARTIN, V. L. and WILKINSON, N. P. (1999). Indirect estimation of ARFIMA and VARFIMA models. *Journal of Econometrics* **93** 149–175.
- MUNROE, M. E. (1953). *Introduction to Measure and Integration*. Addison-Wesley, Reading.
- PARZEN, E. (1974). Some recent advances in time series modelling. *IEEE Transactions on Automatic Control* **AC-19** 723–730.
- PAULSEN, J. and TJØSTHEIM, D. (1985). On the estimation of residual variance and order in autoregressive time series. *Journal of the Royal Statistical Society B-47* 216–228.
- POSKITT, D. S. (2000). Strongly consistent determination of cointegrating rank via canonical correlations. *Journal of Business and Economic Statistics* **18** 71–90.
- ROBINSON, P. M. (1995). Log periodogram regression of time series with long memory. *Annals of Statistics* **23** 1048–1072.
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Annals of Statistics* **6** 461–464.
- SHIBATA, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *Annals of Statistics* **8** 147–164.
- SOWELL, F. (1992). Maximum likelihood estimation of stationary univariate fractionally integrated time series models. *Journal of Econometrics* **53** 165–188.

- SZEGÖ, G. (1939). *Orthogonal Polynomials*. American Mathematical Society Colloquium Publication.
- TIESLAU, M. A., SCHMIDT, P. and BAILLIE, R. T. (1996). A minimum-distance estimator for long memory processes. *Journal of Econometrics* **71** 249–264.
- TJØSTHEIM, D. and PAULSEN, J. (1983). Bias of some commonly-used time series estimators. *Biometrika* **70** 389–400.
- WAND, M. P. and JONES, M. C. (1995). *Kernel Smoothing*. Chapman and Hall.
- WOLD, H. (1938). *The Analysis of Stationary Time Series*. Almqvist and Wicksell, Uppsala, 2nd ed.
- YULE, G. U. (1921). On the time correlation problem. *Journal of the Royal Statistical Society* **84** 497–510.