

ALGORITHMS FOR ESTIMATION OF POSSIBLY NONSTATIONARY VECTOR TIME SERIES

GUOFU ZHOU

Washington University

First version received January 1991

Abstract. This paper presents efficient algorithms for evaluating the likelihood function and its gradient of possibly nonstationary vector autoregressive moving-average (VARMA) processes.

Keywords. Forecasting; likelihood function; innovation algorithm; Whittle algorithm; Bayesian analysis.

1. INTRODUCTION

Time series models are powerful tools in the study and understanding of economic dynamics, and yet there are still many gaps between theory and practice, one of which is the lack of efficient algorithms for implementing the necessary computations in the estimation of the parameters in a vector time series model. This is particularly true for possibly nonstationary models which have been of enormous interest in recent economic studies for at least two reasons. First, it is believed that the nonstationarity (unit-root) question is closely linked to important questions in economic theory. Second, econometric inferences could be misleading if this question is ignored (e.g. Nelson and Kang, 1981, 1984). New approaches and tests for nonstationary models are developed by Dickey and Fuller (1979), Chan (1988), Phillips (1988) and Priestley (1988), among others.

Since no analytical formulae are available for the maximum likelihood estimators in a possibly nonstationary vector autoregressive moving-average (VARMA) model, an iterative procedure has to be used in practice. This inevitably requires many computations of the likelihood function and its gradient. Therefore it is of importance to develop efficient algorithms to solve this problem. In addition, such algorithms are useful in Bayesian analysis. For example, as shown by Zhou (1990), the Monte Carlo integration approach is a solution to the computational problem of posterior densities which are analytically intractable, but to use this approach it is necessary to evaluate the likelihood function thousands of times.

The paper is organized as follows. In Section 2 we synthesize the prediction problem of a stochastic process because forecasts of time series have wide applications. Both the innovation algorithm of Brockwell and Davis (1987, 1988) and Whittle's algorithm (Whittle, 1963; Morf *et al.*, 1978) are

suggested for evaluating the best linear predictors. In Section 3, these algorithms are developed to evaluate the likelihood function and its gradient efficiently. In Section 4 a few remarks conclude the paper.

2. PREDICTIONS

Since forecasts of time series are of general interest, we discuss in some detail the 'forward' and 'backward' best linear predictors. Both the algorithm of Brockwell and Davis (1987, pp. 412–13) and the Whittle algorithm are provided to compute them efficiently. Since these two algorithms will be the foundations of the results in the next section, a comparison of their merits is made. However, it should be pointed out at the outset that there are important alternative algorithms for the computation of predictions, e.g. the lattice algorithms (Friedlander, 1982) and the methods of Kailath based on approximation by combinations of Toeplitz forms. These algorithms are not reviewed here because it seems that further studies are required to apply them to our present case of evaluation of the likelihood function and its gradient.

Let $\{Y_t : t = 0, \pm 1, \pm 2, \dots\}$ be an m -dimensional random process defined on some probability space (Ω, \mathcal{F}, P) . Without loss of generality, we assume that it has zero mean. We further assume that it has finite variance and covariances. Notice that, unless explicitly stated otherwise, the process Y_t is neither necessarily Gaussian nor stationary.

Consider the problem of predicting Y_{n+1} conditional on Y_1, \dots, Y_n . There are, presumably, many ways to construct various type of predictors based on different criteria. However, the best linear predictor is the most popular one. This may be due, at least theoretically, to its tractability. It is defined as a linear function of Y_1, \dots, Y_n ,

$$\hat{Y}_{n+1} = \Phi_{n,1} Y_n + \dots + \Phi_{n,n} Y_1, \quad (1)$$

where the coefficient matrices $\Phi_{i,j}$ are chosen to minimize the mean-squared error:

$$\min_{\Phi_{n,1}, \dots, \Phi_{n,n}} E(Y_{n+1} - \hat{Y}_{n+1})' Q (Y_{n+1} - \hat{Y}_{n+1}) \quad (2)$$

where E is the expectation operator and Q is a given $m \times m$ positive definite matrix. Thus, \hat{Y}_{n+1} is best in the sense that it has the least mean-squared error. Notice that this error is weighted by the matrix Q which reflects the importance of forecasting error. For example, if Q is chosen to be diagonal with the first diagonal element 10 and the rest 1, then the forecast of the first element of Y_{n+1} is measured as 10 times more important than the rest in minimizing the mean-squared error. It is clear that the prediction with a different Q will be different. Following the usual practice, we assume throughout that $Q = I$, the identity matrix. This will not lose any generality

of the theory, because the coefficient matrices for forecasting a series $\{Z_t\}$ with a general positive definite matrix Q can be obtained by a transformation of the $\Phi_{i,j}$ with $Q = I$. To see this, let $Y_t = RZ_t$ where R is a nonsingular matrix such that $Q = R'R$. The minimization problem

$$\min_{\Phi_{n,1}, \dots, \Phi_{n,n}} E(Z_{n+1} - \hat{Z}_{n+1})' Q (Z_{n+1} - \hat{Z}_{n+1})$$

is then

$$\min_{\Phi_{n,1}, \dots, \Phi_{n,n}} E(Y_{n+1} - \hat{Y}_{n+1})' (Y_{n+1} - \hat{Y}_{n+1}),$$

which is the same as (2) with $Q = I$. It follows from (1) that $R^{-1}\Phi_{n,1}R, \dots, R^{-1}\Phi_{n,n}R$ are the coefficient matrices for obtaining the best linear forecast of Z_{n+1} .

An alternative representation of the best linear predictor is the innovation form

$$\hat{Y}_{n+1} = \Theta_{n,1}(Y_n - \hat{Y}_n) + \dots + \Theta_{n,n}(Y_1 - \hat{Y}_1), \tag{3}$$

which expresses \hat{Y}_{n+1} as a linear function of the innovations $Y_k - \hat{Y}_k$ ($k = 1, \dots, n$). There are at least two advantages of considering predictors of this form. One is that the innovation algorithm (discussed later) is suitable for computing \hat{Y}_{n+1} in this way. The second advantage, as we shall see in Section 3, is that it has a close relationship to the fact that a density function can be written as a product of conditional densities. It is this role that makes the innovation algorithm very useful.

The predictor \hat{Y}_{n+1} is a forecast of the future obtained by using the information of the past. Thus it may be called a 'forward' predictor. In contrast, the backward predictor is predicting 'backwards' in the sense that it predicts the value of Y_1 conditional on the observations of Y_{n+1}, \dots, Y_2 . Notice that no special properties of $\{Y_1, \dots, Y_n\}$ have been used in the definition of \hat{Y}_{n+1} other than the fact Y_1, \dots, Y_n are random variables. Therefore, by replacing $\{Y_1, \dots, Y_n\}$ with $\{Y_{n+1}, \dots, Y_2\}$ the backward predictor can be similarly defined and written as

$$\tilde{Y}_1 \equiv \tilde{\Phi}_{n,1}Y_2 + \dots + \tilde{\Phi}_{n,n}Y_{n+1}. \tag{4}$$

The usefulness of this type of predictor will be clarified later.

From the perspective of projections in the Hilbert space $\mathcal{L}^2(\Omega, \mathcal{F}, P)$, the i th element of \hat{Y}_{n+1} is exactly the usual projection of Y_{n+1}^i on the closed subspace generated by the elements of Y_1, \dots, Y_n . Hence, the existence of the forward predictors or the existence of the coefficient matrices is guaranteed by the standard prediction theorems in a Hilbert space, but the uniqueness of the coefficient matrices can only be assured by further assuming the nonsingularity of the $nm \times nm$ covariance matrix $\Gamma \equiv (K(i, j))$, where $K(i, j) \equiv EY_iY_j'$, $m \times m$. This follows from Brockwell and Davis (1987), especially Chapters 2, 5 and 11. Similar assertions are also valid for the backward predictors. The nonsingularity assumption is equivalent to the

requirement that none of the elements of Y_1, \dots, Y_{n+1} is an exact linear combination of the rest. It thus excludes redundancy of the information in the forecast. Since it is trivially satisfied in applications, we henceforth make this assumption throughout so that there will be no confusion about obtaining possibly different coefficient matrices. Therefore the remaining key issue is to present computationally efficient ways to obtain these unique matrices.

Before taking up the computation problem, we provide an example of the predictors which are familiar to econometricians. Let the process be Gaussian and put $n = t$. The forward prediction given by (1) is then the familiar conditional expectation, $\hat{Y}_{t+1} = E(Y_{t+1} | Y_1, \dots, Y_t)$. In a univariate case, it is the well-known best mean-square linear estimation of Y_{t+1} conditional on the available information up to time t . In a rational expectation economy, the coefficient matrices in (1) reflect the weights of past information on agents' expectations, while those in (3) catch the effects of 'surprises' on the expectations. This example also illustrates that both the forward and backward predictions can be of interest to econometricians.

Of much practical value is the innovation algorithm that evaluates the Θ_{ij} s in (3). Let V_{n-1} be the covariance matrix of the prediction errors, i.e.

$$V_{n-1} \equiv E(Y_n - \hat{Y}_n)(Y_n - \hat{Y}_n)' \quad (5)$$

Then we have the following algorithm.

INNOVATION ALGORITHM. If $\{Y_t\}$ has zero mean and finite covariances $K(i, j) = EY_i Y_j'$, then the forward predictors are given by

$$\hat{Y}_{n+1} \equiv \begin{cases} 0 & \text{if } n = 0, \\ \sum_{k=1}^n \Theta_{n,k} (Y_{n+1-k} - \hat{Y}_{n+1-k}) & \text{otherwise} \end{cases} \quad (6)$$

where the coefficients and error covariance matrices are recursively determined from

$$V_0 = K(1, 1),$$

$$\Theta_{n,n-k} \left[K(n+1, k+1) - \sum_{j=0}^{k-1} \Theta_{n,n-j} V_j \Theta_{k,k-j}' \right] V_k^{-1}$$

$$k = 0, 1, \dots, n-1, \quad (7)$$

$$V = K(n+1, n+1) - \sum_{j=0}^{n-1} \Theta_{n,n-j} V_j \Theta_{n,n-j}'$$

PROOF. See Brockwell and Davis (1987, pp. 412–13).

In a univariate case, the innovation algorithm is a refinement of the square-root-free Cholesky decomposition procedure. Similar ideas can be traced back to Kailath (1968, 1970), Rissanen and Barbosa (1969) and Ansley

(1979), to name just a few. The univariate innovation algorithm is also analyzed by Wincek and Reinsel (1986). However, it appears that it is Brockwell and Davis (1987) (see also Brockwell and Davis, 1988) who first give a thorough treatment and generalize it to evaluate predictors and likelihood functions of multivariate Gaussian stationary processes.

It is interesting to know how the innovation algorithm works in terms of matrices. Recall that every positive definite matrix can be uniquely decomposed as LDL' , a product of three matrices, where D is a diagonal matrix and L is lower triangular with identity diagonal elements. It is easy to show that, in a similar fashion, the innovation algorithm decomposes the covariance matrix Γ into

$$\Gamma = \Theta V \Theta'$$

where

$$\Theta = \begin{bmatrix} I & \mathbf{0} & \mathbf{0} & & \mathbf{0} \\ \Theta_{1,1} & I & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \Theta_{n,n} & \Theta_{n,n-1} & \Theta_{n,n-2} & & I \end{bmatrix}$$

$$V = \begin{bmatrix} V_0 & \mathbf{0} & \mathbf{0} & & \mathbf{0} \\ \mathbf{0} & V_1 & \mathbf{0} & & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & & V_n \end{bmatrix}.$$

Notice that V is *not* a diagonal matrix, but a matrix having $m \times m$ matrices as its diagonal elements and zeros elsewhere.

The efficiency of an algorithm is often judged by how many multiplications it takes. Given the coefficient matrices $\Theta_{i,j}$, the predictors \hat{Y}_1 to \hat{Y}_{n+1} can be computed recursively from (6), where the number of multiplications is easily seen to be

$$\mu_1 = m^2 + m^2 2 + \dots + m^2 n = \frac{n(n+1)m^2}{2}$$

This is the straightforward step. The key difficulty is to evaluate the $\Theta_{i,j}$, and this takes

$$\rho_1 = \frac{n(n+1)(2n+1)}{6} m^* + n(n+2)m^*$$

where m^* is the number of multiplications for inverting an $m \times m$ matrix. Notice that Γ is an $nm \times nm$ matrix and therefore, if one uses the standard method of linear algebra to find the decomposition (8), it is known to take $O\{(mn)^3\}$ multiplications. The contribution of the innovation algorithm is that it reduces the number of multiplications to $\rho_1 = O\{(mn)^3/3\}$, only about a third of what the traditional algorithm needs for large n . This fact is important in practice because $n+1$ is often the sample size, which can be very large, and thousands of computations are required. Moreover, the

innovation algorithm becomes far better than any available approach (to the author's knowledge) when the positive matrix $\Gamma = K(i, j)$ is in a band form, i.e. $K(i, j) = \mathbf{0}$ whenever $|i - j| \geq q$ for some positive integer q , $q < n$. In this case, it is easy to check that $\Theta_{n,k} = \mathbf{0}$ for $k > q$, and as a result there is a drastic reduction in both the storage and the amount of computation time. In fact, the number of multiplications is now only

$$\rho_2 = (n - q)q(q + 2)m^* + nm^* + \frac{q(q + 1)(2q + 2)}{6}m^* + q(q + 1)m^*. \quad (10)$$

For a given q , Equation (10) suggests that the innovation algorithm in the current case takes a multiplication of only first order in n (if n is large), whereas ρ_1 and the multiplications in the standard matrix inverting method are of third order.

Perhaps the best way to see how the recursions work is to examine one example. Consider the VARMA(0, 1) or VMA(1) process

$$Y_t = \mathcal{E}_t + \Theta_1 \mathcal{E}_{t-1}$$

where $\mathcal{E}_t \sim \text{IIN}(0, \Sigma)$. It is easy to check that

$$K(i, i) = EY_i Y_i = \Sigma + \Theta_1 \Sigma \Theta_1' \quad K(i, i + 1) = EY_i Y_{i+1} = \Sigma \Theta_1'$$

and

$$K(i, j) = EY_i Y_j' = 0 \quad \text{for } |i - j| \geq 2.$$

An application of the innovation algorithm thus gives, after the first recursion,

$$V_0 = \Sigma + \Theta_1 \Sigma \Theta_1' \quad \Theta_{1,1} = \Sigma \Theta_1' V_0^{-1} \quad V_1 = V_0 - \Theta_{1,1} V_0 \Theta_{1,1}'.$$

Generally,

$$\Theta_{n,1} = \Sigma \Theta_1' V_{n-1}^{-1} \quad \Theta_{n,j} = 0 \quad \text{for } 2 \leq j \leq n$$

and

$$V_n = V_0 - \Theta_{n,1} V_{n-1} \Theta_{n,1}'.$$

With the computed $\{\Theta_{i,1}, i = 1, \dots, n\}$, the prediction is straightforward to compute:

$$\hat{Y}_0 = 0 \quad \hat{Y}_{i+1} = \Theta_{i,1}(Y_i - \hat{Y}_i) \quad \text{for } i = 1, \dots, n.$$

Now consider how to evaluate the coefficient matrices Φ_{ij} in (1). It seems that they can be retrieved by the innovation algorithm. However, direct calculation is much simpler by well-known Whittle algorithm (Whittle, 1963; 1983), which is a generalization of the well-known Durbin-Levinson algorithm. However, in his important paper, Whittle failed to observe that the evaluation of the prediction error matrices (V_n and \tilde{V}_n) can be computed more efficiently. This is the version presented below.

A process $\{Y_t\}$ is said to be stationary (weakly stationary in time series literature) if Y_1, \dots, Y_k have the same covariance matrices as Y_{t+1}, \dots, Y_{t+k} for any t and k , and thus we can denote $K(t+k, t)$ by $\Gamma(k)$. Let \tilde{V}_{n-1} be the error covariance matrix of the backward predictor. It has exactly the same form as (5) with Y_n and \hat{Y}_n replaced by Y_1 and \tilde{Y}_1 . Since the following algorithm computes both the forward and the backward predictor, we name it the Whittle algorithm.

WHITTLE ALGORITHM. If $\{Y_t\}$ is zero mean stationary, then the forward and backward predictors are given by (1) and (3) respectively with the coefficient matrices recursively determined by

$$V_0 = \tilde{V}_0 = \Gamma(0),$$

$$\Phi_{n,n} = \left[\Gamma(n) - \sum_{j=1}^{n-1} \Gamma(n-j) \tilde{\Phi}'_{n-1,j} \right] \tilde{V}_{n-1}^{-1},$$

$$\tilde{\Phi}_{n,n} = \left[\Gamma(n)' - \sum_{j=1}^{n-1} \Gamma(n-j)' \Phi'_{n-1,j} \right] V_{n-1}^{-1},$$

$$\Phi_{n,k} = \Phi_{n-1,k} - \Phi_{n,n} \tilde{\Phi}_{n-1,n-k} \quad k = 1, \dots, n-1, \quad (11)$$

$$\tilde{\Phi}_{n,k} = \tilde{\Phi}_{n-1,k} - \tilde{\Phi}_{n,n} \Phi_{n-1,n-k} \quad k = 1, \dots, n-1,$$

$$V_n = (I_m - \Phi_{n,n} \tilde{\Phi}_{n,n}) V_{n-1}, \quad (12)$$

$$\tilde{V}_n = (I_m - \tilde{\Phi}_{n,n} \Phi_{n,n}) \tilde{V}_{n-1}.$$

PROOF. See Morf *et al.* (1978) or Zhou (1990).

It is easy to show that the number of multiplications in the Whittle algorithm is

$$\rho_3 = (2n^2 + 6n - 1)m^*. \quad (14)$$

The striking difference between the innovation algorithm and the Whittle algorithm occurs when both of them are considered from a matrix perspective. As shown in (8), the innovation algorithm essentially produces the block LDL' decomposition of the covariance matrix Γ . In contrast, what the Whittle algorithm does is to decompose the inverse of this covariance matrix into a similar form:

$$\Gamma^{-1} = \Phi' V^{-1} \Phi,$$

where

$$\Phi \equiv \begin{bmatrix} I & 0 & 0 & \dots & 0 \\ -\Phi_{1,1} & I & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -\Phi_{n,1} & -\Phi_{n,n-1} & -\Phi_{n,n-2} & \dots & I \end{bmatrix}$$

$$V^{-1} = \begin{bmatrix} V_0^{-1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & V_1^{-1} & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & V_n^{-1} \end{bmatrix}$$

which is well known and easy to derive. The representation (15) is a block form of the usual UDU' decomposition of matrix theory. Since (15) implies $\Phi\Gamma\Phi' = V$, recalling the elementary transformations in matrix theory, we find that Φ captures all the elementary transformations that are necessary to make the covariance Γ block diagonal.

In addition, although the innovation algorithm and the Whittle algorithm can both be used to evaluate the forward predictions, there are important differences. First, the Whittle algorithm requires the stationarity of the process, whereas the innovation algorithm does not. Second, the Whittle algorithm is uniquely suited for computing the backward predictions. Third, the innovation algorithm is far better than direct inversion of the Γ matrix, which takes $O(m^3T^3)$ multiplications, but it still consumes $\rho_1 = O(m^3T^3/3)$ multiplications and hence is also a third-order algorithm, where $T = n + 1$ is the sample size. In contrast, the Whittle algorithm is a $\rho_3 = O(3m^3T^2)$ procedure. Fourth, when applied to pure VMA(q) processes, the innovation algorithm becomes very attractive and takes only $\rho_2 = O\{q(q+2)m^3T\}$ multiplications. However, the Whittle algorithm is better for fitting pure VAR models.

Since the determination of the covariance structure of a pure stationary VAR process (without specifying the initial conditions) is rather tedious, we do not give an example to show that the Whittle algorithm performs better in a pure VAR model (this is because $\rho_1 < \rho_3$ and $\rho_1/\rho_3 \rightarrow \infty$ as $n \rightarrow \infty$). Instead, we give an example that shows the comparative advantage of the innovation algorithm. Consider the following numerical experiment on a bivariate ARMA(2, 2) model with $T = 50$:

$$X_t = \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \epsilon_t + \Theta_1 \epsilon_{t-1} + \Phi_2 \epsilon_{t-2}.$$

Then predictions of X_t are easily obtained from the predictions of Y_t , where $Y_t = X_t - \Phi_1 X_{t-1} - \Phi_2 X_{t-2}$ is a pure VMA(2) process. On a Sun 3/50, the innovation algorithm takes 0.03 min for all the predictions. Although all the results agree up to 8 digits in double precision computations,¹ the Whittle algorithm and the direct inversion take relatively more time: 0.10 min and 1.46 min respectively. This is what one would expect from the foregoing theoretical analysis. To illustrate this, notice that $m = q = 2$ and $n = 49$ in this example. Thus, by (10) and (14), the innovation algorithm and the Whittle algorithm would require $\rho_2 = 3496$ and $\rho_3 = 58800$ multiplications respectively, whereas the standard direct inversion approach takes $\rho = (2 \times 49)^3 = 941192$. The actual computation time spent agrees with the relationship $\rho > \rho_3 > \rho_2$. But the proportion of $0.03/0.10 = 30\%$, which measures how fast the innovation algorithm is relative to the Whittle

algorithm, is different from $\rho_2/\rho_3 = 6\%$ (this is the corresponding theoretical measure). This may be explained by the fact that some fixed amount of time has to be used in any of the algorithms. So, if an algorithm takes longer, the fixed amount could be relatively small and less important. As a result, the proportions are likely to be more equal. Indeed, $\rho_2/\rho = 0.037\%$ is closer to $0.03/1.46 = 0.02\%$ than is any other proportion to the theoretical estimation.

3. LIKELIHOOD FUNCTION

The prediction algorithms discussed in the previous section are useful not only for prediction and simulation but also for likelihood function evaluations. Armed with them, we are able in this section to derive algorithms that compute the exact likelihood function and its first-order derivatives for a Gaussian VARMA(p, q) process:

$$X_t = \Phi_1 X_{t-1} + \dots + \Phi_p X_{t-p} + \mathcal{E}_t + \Theta_1 \mathcal{E}_{t-1} + \dots + \Theta_q \mathcal{E}_{t-q}, \quad (16)$$

where the initial values X_0, \dots, X_{1-p} are given constants, \mathcal{E}_t are normal and independent identically distributed disturbances ($\mathcal{E}_t \sim \text{IIN}(0, \Sigma)$, $t = 1, \dots, T$) and Σ is the $m \times m$ variance-covariance matrix of the disturbances that is positively definite. Unlike some other studies (Anderson, 1980; Reinsel, 1979a,b), we do not require the predetermination of the initial \mathcal{E}_s , i.e. $\mathcal{E}_0, \dots, \mathcal{E}_{1-q}$. However, since we allow the possibility of nonstationarity, it is necessary to specify the initial motion of the X . This is done here, for simplicity, by letting the p initial values be given.

We provide two methods for obtaining the maximum likelihood estimator (the set of parameter values which maximize the likelihood function). The first is a one-step minimization problem which has all the parameters of the model as choice variables. The second, however, is a two-step procedure which is simpler in terms of the computation task and storage. This is because in the first step the minimization problem to be solved involves a much smaller number of parameters, i.e. only those that appear in the moving average and disturbances, and in the second step the computation is totally analytical—the maximum likelihood estimator of the autoregressive parameters is obtained as an explicit function of the known estimator obtained from the first step. All the algorithms are computationally efficient and easily implemented.

To apply the prediction algorithms to the VARMA model (16), we let Y_t be a linear filter of the X_t :

$$Y_t = X_t - \Phi_1 X_{t-1} - \dots - \Phi_p X_{t-p}. \quad (17)$$

The model implies that the filter should admit a multivariate moving-average representation, and hence $\{Y_t\}$ is stationary, although $\{X_t\}$ may not be.

Consistent with Section 2, denote by $\Gamma(k)$ the covariance matrix of Y_{t+k} with Y_t . It follows that

$$\Gamma(k) = EY_{t+k}Y_t' = \begin{cases} \sum_{j=0}^{q-k} \Theta_{j+k} \Sigma \Theta_j' & \text{if } k \leq q \\ 0 & \text{otherwise} \end{cases}$$

with $\Theta_0 \equiv I$, the identity matrix. By the normality assumption on the model errors, Y_{n+1} given Y_1, \dots, Y_n must be normal. Furthermore, this normal distribution has the forward predictor as its mean and the error covariance matrix as its covariance matrix. Notice that a density can be written as a recursive product of a series of conditional densities:

$$P(Y_1, \dots, Y_T) = P(Y_1)P(Y_2|Y_1) \dots P(Y_T|Y_1, \dots, Y_{T-1}).$$

We thus derive the log likelihood (up to a constant):

$$\log \mathcal{L}(\boldsymbol{\theta}) = -\frac{1}{2} \sum_{j=1}^T \log |V_{j-1}| - \frac{1}{2} \sum_{j=1}^T (Y_j - \hat{Y}_j)' V_{j-1}^{-1} (Y_j - \hat{Y}_j),$$

where $\boldsymbol{\theta}$ stands for all the parameters. This formula expresses the likelihood explicitly as a function of the filters $\{Y_j\}$ and their forward predictions. It is straightforward to obtain Y_j as a function of the observations on X_t from (17). \hat{Y}_j is easily computed using either the innovation algorithm or the Whittle algorithm of the previous section. In our present case, the innovation algorithm is faster than the Whittle algorithm because Y_t has the special property that $\Gamma(k) = \mathbf{0}$ whenever $k > q$. Nevertheless, the Whittle algorithm is useful for checking the code and for other purposes, some of which we have already mentioned in the previous section and more will be discussed later. In short, there are two efficient methods for the likelihood function evaluation. The first method uses the innovation algorithm, and the second is based on the Whittle algorithm.

With either of the algorithms, the maximum likelihood estimator can be obtained by applying a derivative-free method to maximize the likelihood function. Equivalently, one solves a minimization problem with the objective function

$$\mathbf{\Pi}(\boldsymbol{\theta}) \equiv \sum_{j=1}^T \log |V_{j-1}| + \sum_{j=1}^T (Y_j - \hat{Y}_j)' V_{j-1}^{-1} (Y_j - \hat{Y}_j)$$

over all the parameters. It should be noticed that the minimization problem (20) is not unconstrained because the choice of the elements of Σ must be made in a subspace of all possible values such that Σ is a positive definite matrix. Since unconstrained problems are much easier to solve, we eliminate the constraint by making a parameter transformation: $\Sigma = LL'$, where L is lower triangular with positive diagonal elements. This is the Cholesky

decomposition of the positive definite Σ matrix and L exists uniquely. Then an unconstrained minimization² of (20) over the parameters

$$\Phi_i, \Theta_j, L \quad (i = 1, \dots, p; j = 1, \dots, q) \quad (21)$$

will produce the maximum likelihood estimator with Σ being replaced by the decomposition wherever it appears in the evaluation of $\Pi(\theta)$.

In order to facilitate solving the minimization problem and to obtain a better estimation of the Hessian matrix, it is of practical importance to be able to evaluate not only the log likelihood function but also its gradient. As in the case of log likelihood evaluation, two algorithms, which are based upon either the innovation algorithm or the Whittle algorithm, can be used to this end. Because the derivations are similar, and it may be sufficient in practice to use just the first, we derive only those formulae that depend solely on the innovation algorithm.

The basic idea is to differentiate the recursive relationships and write the results in as simple a form as possible. It is along these lines that Wincek and Reinsel (1986) derive the gradient formula for a univariate time series. Fortunately, the procedure can be generalized to the multivariate model. Recalling the standard results in matrix theory, we have the differential of the objective function (20):

$$\begin{aligned} d\Pi(\theta) = & \sum_{j=1}^T \text{tr} V_{j-1}^{-1} dV_{j-1} - \sum_{j=1}^T Y_j^{*'} V_{j-1}^{-1} dV_{j-1} V_{j-1}^{-1} Y_j^* \\ & + 2 \sum_{j=1}^T Y_j^{*'} V_{j-1}^{-1} dY_j^* \quad Y_j^* \equiv Y_j - \hat{Y}_j, \end{aligned} \quad (22)$$

where tr is the trace operator. Starting from here, we obtain the derivatives in the following three steps with respect to (w.r.t.) the three different sets of parameters in (21).

3.1. Derivatives with respect to Φ_l ($l = 1, \dots, p$)

Let $\theta_{i,j}^l$ be the (i, j) th element of the matrix Φ_l . Holding all other parameters constant, then the derivative w.r.t. $\theta_{i,j}^l$ is just the differential (divided by $d\theta_{i,j}^l$), which is seen from (22) to be

$$d\Pi(\theta) = 2 \sum_{j=1}^T Y_j^{*'} V_{j-1}^{-1} dY_j^*$$

where the V have already been obtained by using the innovation algorithm. The dY^* can be calculated recursively from

$$dY_{n+1}^* = dY_{n+1} - \sum_{k=1}^{\min(n,q)} \Theta_{n,k} dY_{n+1-k}^*$$

where dY_{n+1} is given by the $m \times 1$ vector

$$dY_{n+1} = - \begin{bmatrix} 0 \\ \vdots \\ 0 \\ x_j^{n+1-p} \\ 0 \\ \vdots \\ 0 \end{bmatrix} d\theta_{i,j}^l \quad (\text{ith row})$$

Both of the last two formulae are easily derived from their definitions ((22), (3) and (17)).

3.2. Derivatives with respect to L

In order to obtain the partials, we need to compute all variables in the differential (22). This can be accomplished by evaluating the differentials of the Σ first, the Γ s second, the prediction coefficient matrices third and the Y s and Y^* s last. The partial derivative of the objective function w.r.t. l_{ij} ($j \leq i$) is then readily given by (22).

Notice that all other parameters are held constant. It is straightforward to see that the differential of the Σ is given by the $m \times m$ matrix

$$d\Sigma = \begin{bmatrix} & & & l_{ij} \\ & & & \vdots \\ & & & \vdots \\ l_{jj} & \dots & 2l_{ij} & l_{mj} \\ & & \vdots & \\ & & l_{mj} & \end{bmatrix} dl_{ij} \quad (\text{ith row}).$$

(ith column)

The differential of the covariance matrix is then computed from

$$d\Gamma(k) = \begin{cases} \sum_{j=1}^{q-k} \Theta_{j+k} d\Sigma \Theta_j' & \text{if } k \leq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Now we can compute the differentials of the prediction coefficient matrices. Similar to the likelihood evaluations, the differentials are recursively determined from the differential form of the innovation algorithm:

$$\begin{aligned} dV_0 &= d\Gamma(0), \\ d\Theta_{n, \min(n,q)} &= [d\Gamma\{\min(n, q)\} - \Theta_{n, \min(n,q)} V_{\max(0, n-q)}^{-1}], \\ d\Theta_{n, n-k} &= \{ [d\Gamma(n-k) - \Theta_{n, n-k} dV_k] \\ &\quad + \sum_{j=\max(0, n-q)}^{k-1} (d\Theta_{n, n-j} V_j \Theta_{k, k-j}' \\ &\quad + \Theta_{n, n-j} dV_j \Theta_{k, k-j}' + \Theta_{n, n-j} V_j d\Theta_{k, k-j}') \} V_k^{-1} \end{aligned}$$

$$k = \max(1, n - q), \max(1, n - q) + 1, \dots, n -$$

$$dV_n = d\Gamma(0) -$$

$$\sum_{j=\max(0, n-q)}^{n-1} (d\Theta_{n, n-j} V_j \Theta'_{n, n-j} + \Theta_{n, n-j} dV_j \Theta'_{n, n-j} + \Theta_{n, n-j} V_j d\Theta'_{n, n-j}),$$

$$n = 1, \dots, T - 1.$$

The remaining calculations for obtaining the partial w.r.t. l_{ij} ($j \leq i$) follow easily.

3.3. Derivatives with respect to Θ_l ($l = 1, \dots, q$)

Again, it is sufficient to find the differential $d\Pi(\theta)$ when all parameters are held constant except θ_{ij}^l , the (i, j) th element of the matrix Θ_l . Now the differentials of the Γ should be computed from

$$d\Gamma(k) = \begin{cases} \Theta_{l+k} \Sigma d\Theta'_l & \text{if } l < k, l \leq q - k, \\ \mathbf{0} & \text{if } l < k, l > q - k, \\ \Theta_{l+k} \Sigma d\Theta'_l + d\Theta_l \Sigma & \text{if } l = k, l \leq q - k, \\ d\Theta_l \Sigma & \text{if } l = k, l > q - k, \\ \Theta_{l+k} \Sigma d\Theta'_l + d\Theta_l \Sigma \Theta'_{l-k} & \text{if } l > k, l \leq q - k, \\ d\Theta_l \Sigma \Theta'_{l-k} & \text{if } l > k, l > q - k. \end{cases}$$

The remaining procedures are exactly the same as in step 2.

The above procedures are the first method mentioned for obtaining the maximum likelihood estimators. The second method is motivated from the idea of generalized least squares (GLS). It turns out that the maximum likelihood estimator of the autoregressive parameters, conditional on Σ (or L) and the moving-average parameters, can be obtained analytically as a function of other parameters of the model. Therefore the maximum likelihood estimator can be obtained by a two-step procedure. In the first step, the conditional likelihood function is maximized to get the maximum likelihood estimator of the moving-average component and covariance parameters. Then, in the second step, the autoregressive parameters are obtained in closed form by using GLS.

Taking a transpose on both sides of the model (16) and stacking together, we obtain its matrix form:

$$\begin{bmatrix} X'_1 \\ \vdots \\ X'_T \end{bmatrix} = \begin{bmatrix} X'_0 & & & X'_{1-p} \\ \vdots & & & \vdots \\ X'_{T-1} & \dots & & X'_{T-p} \end{bmatrix} \begin{bmatrix} \Phi'_1 \\ \vdots \\ \Phi'_p \end{bmatrix} + \begin{bmatrix} \mathcal{E}'_1 \\ \vdots \\ \mathcal{E}'_T \end{bmatrix}$$

$$+ \begin{bmatrix} \mathcal{E}'_0 & \dots & \mathcal{E}'_{1-p} \\ \vdots & & \vdots \\ \mathcal{E}'_{T-1} & \dots & \mathcal{E}'_{T-p} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Theta}'_1 \\ \vdots \\ \boldsymbol{\Theta}'_p \end{bmatrix}$$

To apply the standard GLS results, we need to use the vector form of the model. In fact, it is enough to do this for the dependent vector. By definition (17) we have

$$\begin{bmatrix} \mathbf{Y}'_1 \\ \vdots \\ \mathbf{Y}'_T \end{bmatrix} \equiv \begin{bmatrix} \mathbf{X}'_1 \\ \vdots \\ \mathbf{X}'_T \end{bmatrix} - \begin{bmatrix} \mathbf{X}'_0 & \dots & \mathbf{X}'_{1-p} \\ \vdots & & \vdots \\ \mathbf{X}'_{T-1} & \dots & \mathbf{X}'_{T-p} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Phi}'_1 \\ \vdots \\ \boldsymbol{\Phi}'_p \end{bmatrix}.$$

Vectorizing this expression, we write the result

$$\mathbf{y} = \mathbf{x} - \bar{\mathcal{X}}\boldsymbol{\phi},$$

where $\boldsymbol{\phi}$ is the vector form of the stacked parameters in the autoregressive component and $\bar{\mathcal{X}}$ is the $mT \times m^2p$ matrix given by

$$\bar{\mathcal{X}} = \mathbf{I} \otimes \begin{bmatrix} \mathbf{X}'_0 & \mathbf{X}'_{1-p} \\ \vdots & \vdots \\ \mathbf{X}'_{T-1} & \mathbf{X}'_{T-p} \end{bmatrix},$$

where \mathbf{I} is the $m \times m$ identity matrix and \otimes is the Kronecker operator.

For brevity, we introduce the notation

$$\mathbf{Y} \equiv \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_T \end{bmatrix} \equiv \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_T \end{bmatrix}$$

and denote by $\boldsymbol{\Gamma}$ the covariance matrix of the random vector \mathbf{Y} . Decompose this matrix as $\boldsymbol{\Gamma}^{-1} = \mathbf{F}'\mathbf{F}$ (this matrix, as a tool in the proof, does not actually need to be computed), and let \mathbf{G} be an arrangement of the \mathbf{F} matrix,

$$\mathbf{G} = (\mathbf{F}_1, \mathbf{F}_{m+1}, \dots, \mathbf{F}_{(T-1)m+1}; \dots; \mathbf{F}_m, \dots, \mathbf{F}_{(T-1)m+m})$$

where \mathbf{F}_i is the i th column of \mathbf{F} . Then the log conditional likelihood function is (up to a conditioning constant)

$$\mathbf{Y}'\boldsymbol{\Gamma}^{-1}\mathbf{Y} = (\mathbf{Y}^* - \mathcal{X}^*\boldsymbol{\phi})'(\mathbf{Y}^* - \mathcal{X}^*\boldsymbol{\phi}),$$

where

$$\mathbf{Y}^* = \mathbf{G}\mathbf{x} \quad \mathcal{X}^* = \mathbf{G}\bar{\mathcal{X}}.$$

Therefore, conditioning on $\boldsymbol{\Theta}_1, \dots, \boldsymbol{\Theta}_q$ and $\boldsymbol{\Sigma}$, the maximum likelihood estimator of the vector $\boldsymbol{\phi}$ which contains all parameters of $\boldsymbol{\Phi}_1, \dots, \boldsymbol{\Phi}_p$ as its elements is

$$\boldsymbol{\phi}_{\text{ML}} = (\mathcal{X}^{*\prime}\mathcal{X}^*)^{-1}\mathcal{X}^{*\prime}\mathbf{Y}^* \quad (23)$$

and the conditional covariance matrix is $\text{var}(\boldsymbol{\phi}_{\text{ML}}) = (\mathcal{X}^{*\prime}\mathcal{X}^*)^{-1}$. This matrix is useful for generating the conditional samples from the posterior density in Bayesian analysis.

It is obvious that $Y^* = FX$ and $\mathcal{X}^* = F\mathcal{X}$, where \mathcal{X} is an appropriate arrangement of $\tilde{\mathcal{X}}$. Let D be the diagonal matrix $D = \text{diag}[V_0, V_1, \dots, V_{T-1}]$. Replacing ϕ by its conditional estimation (23) in (20), we obtain the reduced objective function

$$\Pi(\theta) = \sum_{j=1}^T \log |V_{j-1}| + \tilde{Y}' \tilde{M} \tilde{Y}, \quad (24)$$

where $\tilde{Y} = X - \hat{X}$ and

$$\tilde{M} = D^{-1} - D^{-1}(\mathcal{X} - \hat{\mathcal{X}})[(\mathcal{X} - \hat{\mathcal{X}})' D^{-1}(\mathcal{X} - \hat{\mathcal{X}})]^{-1}(\mathcal{X} - \hat{\mathcal{X}})' D^{-1}.$$

The vector \hat{X} denotes the projection of X , whose definition and computation are discussed in Section 2. Denoted analogously, the matrix $\hat{\mathcal{X}}$ is composed of the projections of the corresponding columns.

The reduced objective function depends only upon part of the parameters $\Theta_1, \dots, \Theta_q$ and Σ . Therefore the θ in (24) should be understood to represent these parameters. The maximum likelihood estimator can be obtained by minimizing this objective function. For practical purposes, it may be easier if it is computed from the following compact form:

$$\Pi(\theta) = \sum_{j=1}^T [\log |V_{j-1}| + \tilde{Y}'_j V_{j-1}^{-1} \tilde{Y}_j] - U' W^{-1} U \quad (25)$$

where the U and W are an $m^2 p$ vector and an $m^2 p \times m^2 p$ matrix, with elements defined respectively by

$$U_i = \sum_{k=1}^T Z_{i,k} V_k^{-1} \tilde{Y}_k \quad i = 1, \dots, m^2 p$$

and

$$W_{i,j} = \sum_{k=1}^T Z_{i,k} V_k^{-1} Z_{j,k} \quad i, j = 1, \dots, m^2 p.$$

The $Z_{i,k}$ ($1 \leq i \leq T$, $1 \leq k \leq m^2 p$) are vectors of m elements. All $Z_{i,k}$ ($1 \leq k \leq m^2 p$) stacked one on top of the other form the i th column of the \mathcal{X} matrix. Even though the Whittle algorithm can be used to evaluate the

$$-Y = CY \quad E(Y - Y)(Y - Y) =$$

$$Y' \Gamma^{-1} Y = (CY)' D^{-1} (CY) = (K\mathcal{Y})' D^{-1} (K\mathcal{Y}),$$

where K is an $mT \times mT$ matrix, obtained as a column arrangement of C in the same way as for G . Let $K = [K_1, \dots, K_m]$, with the K_i ($1 \leq i \leq m$) $mT \times T$ matrices. Thus, we get exactly the same formula for evaluating the reduced objective function as given by (25) except for the obvious differences that

$$U_i = (K_i Z)' D^{-1} \tilde{Y} \quad i = \underline{1}, \dots, m^2 p$$

$$W_{i,j} = (K_i Z)' D^{-1} (K_j Z) \quad i, j = \underline{1}, \dots, m^2 p,$$

where the $T \times mp$ matrix Z is defined by

$$Z = \begin{bmatrix} X'_0 & X'_{1-p} \\ \vdots & \vdots \\ X'_{T-1} & X'_{T-p} \end{bmatrix}.$$

Finally, we give the formulae for evaluation of the gradient of the reduced objective function. It is sufficient to provide only

$$\begin{aligned} d\Pi(\theta) &= \sum_{j=1}^T \text{tr} V_{j-1}^{-1} dV_{j-1} - (\tilde{Y} - \tilde{\mathcal{X}}\phi)' D^{-1} dDD^{-1} (\tilde{Y} - \tilde{\mathcal{X}}\phi) \\ &\quad + 2(\tilde{Y} - \tilde{\mathcal{X}}\phi)' D^{-1} (d\tilde{Y} - d\tilde{\mathcal{X}}\phi - \tilde{\mathcal{X}}d\phi), \end{aligned} \quad (26)$$

where $d\phi = W^{-1}[dU - dW\phi]$, with dW and dU given by

$$dW = d\tilde{\mathcal{X}}' D^{-1} \tilde{\mathcal{X}} + \tilde{\mathcal{X}}' D^{-1} [d\tilde{\mathcal{X}} - dDD^{-1} \tilde{\mathcal{X}}],$$

$$dU = d\tilde{\mathcal{X}}' D^{-1} \tilde{\mathcal{Y}} + \tilde{\mathcal{X}}' D^{-1} [d\tilde{\mathcal{Y}} - dDD^{-1} \tilde{\mathcal{Y}}]$$

respectively. The remaining differentials can be found in the same way as in the first method.

4. CONCLUDING REMARKS

We have provided computationally efficient methods for evaluating the likelihood function and its gradient and methods for obtaining the maximum likelihood estimator in a possibly nonstationary vector autoregressive moving-average (VARMA) model. These methods are clearly also applicable for VARMA models with regressors and regression models with VARMA errors. However, only the estimation problem of time series analysis has been dealt with here. For the distribution theory and hypothesis testing, the readers are referred to Dickey and Fuller (1979), Chan (1988), Phillips (1988) and Priestley (1988) and references therein. Further work on the algorithms includes generalizing them to other processes, for example, fractional VARMA and non-Gaussian processes.

NOTES

This paper is part of my Ph.D. thesis at Duke University. I am grateful to Jean-Francois Richard, George Tauchen and Mike West for serving on my thesis committee and offering helpful comments, and especially to John Geweke for his guidance and to an anonymous referee for his insightful comments on an earlier version of this paper. Naturally, the responsibility of any errors or difficulties is solely my own. Financial support from NSF Grant SES-8908365 and Fossett Foundation is also gratefully acknowledged.

¹ A FORTRAN code that implements all the algorithms is available from the author upon request.

² Strictly speaking, it still has the trivial constraint that L must have nonzero diagonal elements.

REFERENCES

- ANDERSON, T. W. (1980) Maximum likelihood estimation for vector autoregressive moving average models. In *Directions in Time Series* (eds D. R. Brillinger and G. C. Tiao). Institute of Mathematical Statistics, 49–59.
- ANSLEY, C. F. (1979) An algorithm for the exact likelihood of a mixed autoregressive–moving average process. *Biometrika* 66, 59–65.
- BOX, G. E. P. and JENKINS, G. M. (1976) *Time Series Analysis: Forecasting and Control*, revised edn. San Francisco, CA: Holden-Day.
- BROCKWELL, P. J. and DAVIS, R. A. (1987) *Time Series: Theory and Methods*. New York: Springer-Verlag.
- and — (1988) Applications of innovation representations in time series analysis. In *Probability and Statistics* (ed. J. N. Srivastara). Amsterdam: North-Holland.
- CHAN, N. H. (1988) Limiting distributions of least squares estimates of unstable autoregressive processes. *Ann. Statist.* 16, 367–401.
- DICKEY, D. A. and FULLER, W. A. (1979) Distribution of the estimators for autoregressive time series with a unit root. *J. Am. Statist. Assoc.* 74, 427–31.
- FRIEDLANDER, B. (1982) *Proc. IEEE* 70, 829–67.
- GEWEKE, J. (1984) Measures of conditional linear dependence and feedback. *J. Am. Statist. Assoc.* 79, 907–15.
- (1988) Antithetic acceleration of Monte Carlo integration in Bayesian inference. *J. Econometrics* 38, 73–90.
- (1989) The posterior distribution of roots in multivariate autoregressions. *1989 Proceedings of the Business and Economic Statistics Section—American Statistical Association*, forthcoming.
- HANNAN, E. J. (1970) *Multiple Time Series*. New York: Wiley.
- HILLMER, S. C. and TAO, G. C. (1979) Likelihood function of stationary multiple autoregressive moving average models. *J. Am. Statist. Assoc.* 74, 652–60.
- KAILATH, T. (1968) An innovation approach to least square estimation—Part I: Linear filtering in additive noise. *IEEE Trans. Autom. Control* 13, 646–54.
- (1970) The innovation approach to detection and estimation theory. *Proc. IEEE* 58, 680–95.
- KOHN, R. and ANSLEY, C. F. (1985) Computing the likelihood of and its derivatives for a Gaussian ARMA model. *J. Statist. Comput. Simulation* 22, 229–63.
- MELARD, G. (1984) A fast algorithm for the exact likelihood of autoregressive moving average models. *Appl. Statist.* 33, 104–14.
- MORF, M., VIEIRA, A. and KAILATH, T. (1978) Covariance characterization by partial autocorrelation matrices. *Ann. Statist.* 6, 643–8.
- NELSON, C. R. and KANG, H. (1981) Spurious periodicity in appropriately detrended time series. *Econometrica* 49, 741–51.

- and PLOSSER, C. I. (1982) Trends and random walks in macroeconomic time series. *J. Monet. Econ.* 8, 129–62.
- PHADKE, M. S. and KEDEM, G. (1978) Computation of the exact likelihood function of multivariate moving average models. *Biometrika* 65, 511–19.
- PHILLIP, P. C. B. (1987) Time series regression with unit roots. *Econometrica* 55, 277–302.
- PRIESTLEY, M. B. (1988) *Non-linear and Non-stationary Time Series Analysis*. New York: Academic Press.
- REINSEL, G. C. (1979) Maximum likelihood estimation of stochastic linear difference equations with autoregressive moving average errors. *Econometrica* 47, 129–51.
- RISSANEN, J. and BARBOSA, L. (1973) A fast algorithm for optimum linear predictors. *IEEE Trans. Autom. Control* 18, 555.
- TIAO, G. C. and BOX, G. E. P. (1981) Modelling multiple time series with applications. *J. Am. Statist. Assoc.* 76, 802–16.
- WHITTLE, P. (1963) On the fitting of multivariate autoregressions, and the approximate canonical factorization of a spectral density matrix. *Biometrika* 50, 129–34.
- (1983) *Prediction and Regulation by Linear Least-square Methods*, 2nd edn. Minneapolis, MN: University of Minnesota.
- WILSON, T. G. (1973) The estimation of parameters in multivariate time series models. *J. R. Statist. Soc. Ser. B* 35, 76–85.
- (1979) Some efficient computational procedures for high order ARMA models. *J. Statist. Comput. Simulation* 8, 301–9.
- WINCEK, M. A. and REINSEL, G. C. (1986) An exact maximum likelihood estimation procedure for regression-ARMA time series models with possibly nonconsecutive data. *J. R. Statist. Soc., B Ser.* 48, 303–13.
- ZHOU, G. (1990) A Bayesian analysis of time series with applications to stationarity and causality. Ph. D. Dissertation, Duke University, Durham, NC.