

Practical methods for modelling weak VARMA processes: identification, estimation and specification with a macroeconomic application*

Jean-Marie Dufour[†]
McGill University

Denis Pelletier[‡]
North Carolina State University

May 2014
Compiled: May 19, 2014

*The authors thank Marine Carasco, John Galbraith, Nour Meddahi and Rui Castro for several useful comments. The second author gratefully acknowledges financial assistance from the Social Sciences and Humanities Research Council of Canada, the Government of Québec (fonds FCAR), the CRDE and CIRANO. Earlier versions of the paper circulated under the title *Linear Estimation of Weak VARMA Models With a Macroeconomic Application*. This work was supported by the Social Sciences and Humanities Research Council of Canada, the Natural Sciences and Engineering Research Council of Canada, the Canadian Network of Centres of Excellence [program on *Mathematics of Information Technology and Complex Systems* (MITACS)], the Canada Council for the Arts (Killam Fellowship), the CIREQ, the CIRANO, and the Fonds FCAR (Government of Québec).

[†]William Dow Professor of Economics, McGill University, Centre interuniversitaire de recherche en analyse des organisations (CIRANO), and Centre interuniversitaire de recherche en économie quantitative (CIREQ). Mailing address: Department of Economics, McGill University, Leacock Building, Room 519, 855 Sherbrooke Street West, Montréal, Québec H3A 2T7, Canada. TEL: (1) 514 398 8879; FAX: (1) 514 398 4938; e-mail: jean-marie.dufour@mcgill.ca . Web page: <http://www.jeanmariedufour.com>

[‡]Department of Economics, Box 8110, North Carolina State University, Raleigh, NC 27695-8110, USA. Email: denis_pelletier@ncsu.edu. Web page: <http://www4.ncsu.edu/dpellet>

ABSTRACT

We consider the problem of developing practical methods for modelling weak VARMA processes. In a first part, we propose new identified VARMA representations, the *diagonal MA equation form* and the *final MA equation form*, where the MA operator is diagonal and scalar respectively. Both of these representations have the important feature that they constitute relatively simple modifications of a VAR model (in contrast with the echelon representation). In a second part, we study the problem of estimating VARMA models by relatively simple methods which only require linear regressions. We consider a generalization of the regression-based estimation method proposed by Hannan and Rissanen (1982). The asymptotic properties of the estimator are derived under weak hypotheses on the innovations (uncorrelated and strong mixing) so as to broaden the class of models to which it can be applied. In a third part, we present a modified information criterion which gives consistent estimates of the orders under the proposed representations. To demonstrate the importance of using VARMA models to study multivariate time series we compare the impulse-response functions and the out-of-sample forecasts generated by VARMA and VAR models.

Key words: linear regression; VARMA; final equation form; information criterion; weak representation; strong mixing condition; impulse-response function.

Journal of Economic Literature Classification: C13, C32, C51, E0.

1. Introduction

In time series analysis and econometrics, VARMA models are scarcely used to represent multivariate time series. VAR models are much more widely employed because they are easier to implement. The latter models can be estimated by least squares methods, while VARMA models typically require nonlinear methods (such as maximum likelihood). Specification is also easier for VAR models since only one lag order must be chosen.

VAR models, however, have important drawbacks. First, they are typically less parsimonious than VARMA models [*e.g.*, see Lütkepohl and Poskitt (1996b)]. Second, the family of VAR models is not closed under marginalization and temporal aggregation [see Lütkepohl (1991)]. The truth cannot always be a VAR. If a vector satisfies a VAR model, subvectors do not typically satisfy VAR models (but VARMA models). Similarly, if the variables of a VAR process are observed at a different frequency, the resulting process is not a VAR process. In contrast, the class of VARMA models is closed under such operations. Furthermore, Athanasopoulos and Vahid (2008) argue that there is no compelling reason for restricting macroeconomic forecasting to VAR models and show that VARMA models can forecast macroeconomic variables more accurately than VARs. Chen, Choi, and Escanciano (2012) refers to many examples in macroeconomics where the models contain an MA component.

The importance of nonlinear models has been growing in the time series literature. These models are interesting and useful but may be hard to use. Because of this and the fact that many important classes of nonlinear processes admit an ARMA representation [*e.g.*, see Francq and Zakoïan (1998), Francq, Roy, and Zakoïan (2005)] many researchers and practitioners still have an interest in linear ARMA models. However, the innovations in these ARMA representations do not have the usual i.i.d. or m.d.s. property, although they are uncorrelated. One must then be careful before applying usual results to the estimation of ARMA models because they usually rely on the above strong assumptions [*e.g.*, see Brockwell and Davis (1991) and Lütkepohl (1991)]. We refer to these as strong and semi-strong ARMA models respectively, by opposition to weak ARMA models where the innovations are only uncorrelated. The i.i.d. and m.d.s. properties are also not robust to aggregation (the i.i.d. Gaussian case being an exception); see Francq and Zakoïan (1998), Francq, Roy, and Zakoïan (2005), Palm and Nijman (1984), Nijman and Palm (1990), Drost (1993). In fact, the Wold decomposition only guarantees that the innovations are uncorrelated.

It follows that (weak) VARMA models appear to be preferable from a theoretical viewpoint, but their adoption is complicated by identification and estimation difficulties. The direct multivariate generalization of ARMA models does not give an identified representation [see Lütkepohl (1991, Section 7.1.1)]. It follows that one has to decide on a set of constraints to impose so as to achieve identification. Standard estimation methods for VARMA models (maximum likelihood, nonlinear least squares) require nonlinear optimization which may not be feasible as soon as the model involves a few time series, because the number of parameters can increase quickly.

In this paper, we consider the problem of modeling weak VARMA processes. Our goal is to develop a procedure which will ease the use of these models. It will cover three basic modelling operations: identification, estimation and specification.

First, in order to avoid identification problems and to further ease the use of VARMA models,

we introduce three new identified VARMA representations, the *diagonal MA equation form*, the *final MA equation form* and the *diagonal AR equation form*. Under the diagonal MA equation form (diagonal AR equation form) representation, the MA (AR) operator is diagonal and each lag operator may have a different order. Under the final MA equation form representation the MA operator is scalar, *i.e.* the operators are equal across equations. The diagonal and final MA equation form representations can be interpreted as simple extensions of the VAR model, which should be appealing to practitioners who prefer to employ VAR models due to their ease of use. The identified VARMA representation which is the most widely employed in the empirical literature is the *echelon form*. Specification of VARMA models in echelon form does not amount to specifying the order p and q as with ARMA models. Under this representation, VARMA models are specified by as many parameters, called Kronecker indices, as the number of time series studied. These indices determine the order of the elements of the AR and MA operators in a non trivial way. The complicated nature of the echelon form representation is a major reason why practitioners are not using VARMA models, so the introduction of a simpler identified representation is interesting.

Second, we consider the problem of estimating VARMA models by relatively simple methods which only require linear regressions. For that purpose, we consider a multivariate generalization of the regression-based estimation method proposed by Hannan and Rissanen (1982) for univariate ARMA models. The method is performed in three steps. In a first step, a long autoregression is fitted to the data. In the second step, the lagged innovations in the ARMA model are replaced by the corresponding residuals from the long autoregression and a regression is performed. In a third step, the data from the second step are filtered so as to give estimates that have the same asymptotic covariance matrix than one would get with the maximum likelihood [claimed in Hannan and Rissanen (1982), proven in Zhao-Guo (1985)]. Extension of this innovation-substitution method to VARMA models was also proposed by Hannan and Kavalieris (1984a) and Koreisha and Pukkila (1989), under the assumption that the innovations are a m.d.s.

Here, we extend these results by showing that the linear regression-based estimators are consistent under weak hypotheses on the innovations and how filtering in the third step gives estimators that have the same asymptotic distribution as their nonlinear counterparts (maximum likelihood if the innovations are i.i.d., or nonlinear least squares if they are merely uncorrelated). In the non i.i.d. case, we consider strong mixing conditions [Doukhan (1995), Bosq (1998)], rather than the usual m.d.s. assumption. By using weaker assumptions for the process of the innovations, we broaden the class of processes to which our method can be applied.

Third, we suggest a modified information criterion to choose the orders of VARMA models under these representations. This criterion is to be minimized in the second step of the estimation method over the orders of the AR and MA operators and gives consistent estimates of these orders. Our criterion is a generalization of the information criterion proposed by Hannan and Rissanen (1982), which was later corrected by Hannan and Rissanen (1983) and Hannan and Kavalieris (1984b), for choosing the orders p and q in ARMA models. The idea of generalizing this information criterion is mentioned in Koreisha and Pukkila (1989) but a specific generalization and theoretical properties are not presented.

Fourth, the method is applied to U.S. macroeconomic data previously studied by Bernanke and Mihov (1998) and McMillin (2001). To illustrate the impact of using VARMA models instead of

VAR models to study multivariate time series we compare the impulse-response functions generated by each model. We show that we can obtain much more precise estimates of the impulse-response function by using VARMA models instead of VAR models.

The rest of the paper is organized as follows. Our framework and notation are described in section 2. The new identified representations are presented in section 3. In section 4, we present the estimation method. In section 5, we describe the information criterion used for choosing the orders of VARMA models under the representation proposed in our work. Section 6 contains results of Monte Carlo simulations which illustrate the properties of our method. Section 7 presents the macroeconomic application where we compare the impulse-response functions from a VAR model and VARMA models. Section 8 contains a few concluding remarks. Finally, proofs are in the appendix.

2. Framework

Consider the following K -variate zero mean VARMA(p, q) model in standard representation for a real-valued series Y_t :

$$Y_t = \sum_{i=1}^p \Phi_i Y_{t-i} + U_t - \sum_{j=1}^q \Theta_j U_{t-j} \quad (2.1)$$

where U_t is a sequence of uncorrelated random variables defined on some probability space $(\Omega, \mathcal{A}, \mathcal{P})$. The vectors Y_t and U_t contain the K univariate time series: $Y_t = [y_{1t}, \dots, y_{Kt}]'$ and $U_t = [u_{1t}, \dots, u_{Kt}]'$. We can also write the previous equation with lag operators:

$$\Phi(L)Y_t = \Theta(L)U_t \quad (2.2)$$

where

$$\Phi(L) = I_K - \Phi_1 L - \dots - \Phi_p L^p, \quad \Theta(L) = I_K - \Theta_1 L - \dots - \Theta_q L^q. \quad (2.3)$$

Let H_t be the Hilbert space generated by $(Y_s, s < t)$. The process U_t can be interpreted as the linear innovation of Y_t :

$$U_t = Y_t - E_L[Y_t | H_t]. \quad (2.4)$$

We assume that Y_t is a strictly stationary and ergodic sequence and that the process U_t has common variance ($\text{Var}[U_t] = \Sigma_U$) and finite fourth moment ($E[|u_{it}|^{4+2\varepsilon}] < \infty$, for all i and t , where $\varepsilon > 0$). The case of I(1) and cointegrated variables is left for future work. We make the zero mean-mean hypothesis only to simplify notation.

Assuming that the process Y_t is stable,

$$\det[\Phi(z)] \neq 0 \text{ for all } |z| \leq 1, \quad (2.5)$$

and invertible,

$$\det[\Theta(z)] \neq 0 \text{ for all } |z| \leq 1, \quad (2.6)$$

it can be represented as an infinite VAR,

$$\Pi(L)Y_t = U_t, \quad (2.7)$$

where

$$\Pi(L) = \Theta(L)^{-1}\Phi(L) = I_K - \sum_{i=1}^{\infty} \Pi_i L^i, \quad (2.8)$$

or an infinite VMA

$$Y_t = \Psi(L)U_t, \quad (2.9)$$

where

$$\Psi(L) = \Phi(L)^{-1}\Theta(L) = I_K - \sum_{j=1}^{\infty} \Psi_j L^j. \quad (2.10)$$

We will denote by $\varphi_{ik}(L)$ the polynomial in row i and column k of $\Phi(L)$, and the row i or column k of $\Phi(L)$ by

$$\Phi_{i\bullet}(L) = [\varphi_{i1}(L), \dots, \varphi_{iK}(L)], \quad (2.11)$$

$$\Phi_{\bullet k}(L) = [\varphi_{1k}(L), \dots, \varphi_{Kk}(L)]'. \quad (2.12)$$

The diag operator creates a diagonal matrix,

$$\text{diag}[\varphi_{ii}(L)] = \text{diag}[\varphi_{11}(L), \dots, \varphi_{KK}(L)] = \begin{bmatrix} \varphi_{11}(L) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \varphi_{KK}(L) \end{bmatrix} \quad (2.13)$$

where

$$\varphi_{ii}(L) = 1 - \varphi_{ii,1}L - \dots - \varphi_{ii,p}L^p. \quad (2.14)$$

The function $\deg[\varphi(L)]$ returns the degree of the polynomial $\varphi(L)$ and the function $\dim(\gamma)$ gives the length of the vector γ .

We need to impose some structure on the process U_t . The typical hypothesis which is imposed in the time series literature is that the U_t 's are either independent and identically distributed (i.i.d.) or a martingale difference sequence (m.d.s.). In this work, we do not impose such strong assumptions because we want to broaden the class of models to which it can be applied. We only assume that it satisfies a strong mixing condition [Doukhan (1995), Bosq (1998)]. Let U_t be a strictly stationary process, and

$$\alpha(h) = \sup_{\substack{B \in \sigma(U_s, s \leq t) \\ C \in \sigma(U_s, s \geq t+h)}} |\Pr(B \cap C) - \Pr(B)\Pr(C)| \quad (2.15)$$

the α -mixing coefficient of order $h \geq 1$, where $\sigma(U_s, s \leq t)$ and $\sigma(U_s, s \geq t+h)$ the σ -algebras associated with $\{U_s : s \leq t\}$ and $\sigma(U_s : s \geq t+h)$ respectively. We suppose that U_t is strong mixing, *i.e.*

$$\sum_{h=1}^{\infty} \alpha(h)^{\varepsilon/(2+\varepsilon)} < \infty \quad \text{for some } \varepsilon > 0. \quad (2.16)$$

This is a fairly minimal condition that will be satisfied by many processes of interest.

3. Identification and diagonal VARMA representations

It is important to note that we cannot work with the standard representation (2.1) because it is not identified. To help us gain intuition on the identification of VARMA models we can consider a more general representation where Φ_0 and Θ_0 are not identity matrices:

$$\Phi_0 Y_t = \Phi_1 Y_{t-1} + \cdots + \Phi_p Y_{t-p} + \Theta_0 U_t - \Theta_1 U_{t-1} + \cdots + \Theta_q U_{t-q}. \quad (3.1)$$

By this specification, we mean the well-defined process

$$Y_t = (\Phi_0 - \Phi_1 L - \cdots - \Phi_p L^p)^{-1} (\Theta_0 + \Theta_1 L + \cdots + \Theta_q L^q) U_t. \quad (3.2)$$

But we can see this such process has a standard representation if Φ_0 and Θ_0 are nonsingular. To see this, we premultiply (3.1) by Φ_0^{-1} and define $\bar{U}_t = \Phi_0^{-1} \Theta_0 U_t$:

$$\begin{aligned} Y_t = & \Phi_0^{-1} \Phi_1 Y_{t-1} + \cdots + \Phi_0^{-1} \Phi_p Y_{t-p} \\ & + \bar{U}_t - \Phi_0^{-1} \Theta_1 \Theta_0^{-1} \Phi_0 \bar{U}_{t-1} - \cdots - \Phi_0^{-1} \Theta_q \Theta_0^{-1} \Phi_0 \bar{U}_{t-q}. \end{aligned} \quad (3.3)$$

Redefining the matrices, we get a representation of type (2.1). As long as Φ_0 and Θ_0 are nonsingular, we can transform a non-standard VARMA into a standard one.

We say that two VARMA representations are equivalent if $\Phi(L)^{-1} \Theta(L)$ results in the same operator $\Psi(L)$. Thus, to ensure uniqueness of a VARMA representation, we must impose restrictions on the AR and MA operators such that for a given $\Psi(L)$ there is one and only one set of operators $\Phi(L)$ and $\Theta(L)$ that can generate this infinite MA representation.

A first restriction that we impose is a multivariate equivalent of the coprime property in the univariate case. We do not want factors of $\Phi(L)$ and $\Theta(L)$ to “cancel out” when $\Phi(L)^{-1} \Theta(L)$ is computed. This feature is called the *left-coprime* property [see Hannan (1969), Hannan and Deistler (1988) and Lütkepohl (1993)]. There exist more than one representation which guarantee the uniqueness of the left-coprime operators. The predominant representation in the economics literature is the *echelon form* [see Deistler and Hannan (1981), Hannan and Kavalieris (1984b), Lütkepohl (1993), Lütkepohl and Poskitt (1996a)]. It requires the selection of Kronecker indices, which conceptually is not as easy as selecting the orders p and q of an ARMA model.¹ This might be a reason why practitioners are reluctant to employ VARMA models.

In this work, to ease the use of VARMA models we present new VARMA representations which can be seen as a simple extensions of the VAR model. To introduce them, we first review another identified representation, the *final equation form*, which will refer to as the *final AR equation form*,

¹Specification of VARMA models in echelon form is discussed for example in Hannan and Kavalieris (1984b), Lütkepohl and Claessen (1997), Poskitt (1992), Nsiri and Roy (1992, 1996), Lütkepohl and Poskitt (1996b), Bartel and Lütkepohl (1998). A more general and in-depth discussion of identification of VARMA models can be found in Hannan and Deistler (1988, Chapter 2).

under which the AR operator is scalar [see Zellner and Palm (1974), Hannan (1976), Wallis (1977), Lütkepohl (1993)].

Definition 3.1 (Final AR equation form) *The VARMA representation (2.1) is said to be in final AR equation form if $\Phi(L) = \varphi(L)I_K$, where $\varphi(L) = 1 - \varphi_1L - \dots - \varphi_pL^p$ is a scalar polynomial with $\varphi_p \neq 0$.*

To see how we can obtain a VARMA model with a final AR equation form representation, we can proceed as follows. By standard linear algebra, we have

$$\Phi^*(L)\Phi(L) = \Phi(L)\Phi^*(L) = \det[\Phi(L)]I_K \quad (3.4)$$

where $\Phi^*(L)$ is the adjoint matrix of $\Phi(L)$. On multiplying both sides of (2.2) by $\Phi^*(L)$, we get:

$$\det[\Phi(L)]Y_t = \Phi(L)^*\Theta(L)U_t. \quad (3.5)$$

This representation may not be attractive for several reasons. First, it is quite far from usual VAR models by excluding lagged values of other variables in each equation (*e.g.*, the AR part of the first equation includes lagged values of y_{1t} but no lagged values of y_{2t}, \dots, y_{Kt}). Further, the AR coefficients are the same in all the equations, which will require a polynomial of higher order pK . Second, the interaction between the different variables is modeled through the MA part of the model, which may have to be quite complex.

However, more convenient alternative representations can be derived through analogous manipulations. Upon multiplying both sides of (2.2) by $\Theta^*(L)$, we get:

$$\Theta(L)^*\Phi(L)Y_t = \det[\Theta(L)]U_t \quad (3.6)$$

where $\Theta(L)^*$ is the adjoint matrix of $\Theta(L)$. We refer to VARMA models in (3.6) as being in *final MA equation form*.

Definition 3.2 (Final MA equation form) *The VARMA representation (2.1) is said to be in final MA equation form if $\Theta(L) = \theta(L)I_K$, where $\theta(L) = 1 - \theta_1L - \dots - \theta_qL^q$ is a scalar operator with $\theta_q \neq 0$.*

The same criticism regarding the parsimony of the final equation form would apply but it is possible to get a more parsimonious representation by looking at common structures across equations. Suppose there are common roots across rows for some columns of $\Theta(L)$, so that starting from (2.1) we can write:

$$\Phi(L)Y_t = \bar{\Theta}(L)D(L)U_t, \quad (3.7)$$

$$\bar{\Theta}^*(L)\Phi(L)Y_t = \det[\bar{\Theta}(L)]D(L)U_t, \quad (3.8)$$

where $D(L) = \text{diag}[d_1(L), \dots, d_K(L)]$ and $d_j(L)$ is a polynomial common to $\theta_{ij}(L)$, $\forall i = 1, \dots, K$. We see that allowing non-equal diagonal polynomials in the moving average as in equation (3.8)

may yield a more parsimonious representation than in (3.6). We will call the representation (3.8) *diagonal MA equation form* representation.

Definition 3.3 (Diagonal MA equation form) *The VARMA representation (2.1) is said to be in diagonal MA equation form if $\Theta(L) = \text{diag}[\theta_{ii}(L)] = I_K - \Theta_1 L - \dots - \Theta_q L^q$ where $\theta_{ii}(L) = 1 - \theta_{ii,1} L - \dots - \theta_{ii,q_i} L^{q_i}$, $\theta_{ii,q_i} \neq 0$, and $q = \max_{1 \leq i \leq K} (q_i)$.*

This representation is interesting because contrary to the echelon form it is easy to specify. We don't have to deal with rules for the orders of the off-diagonal elements in the AR and MA operators. The fact that it can be seen as a simple extension of the VAR model is also appealing. Practitioners are comfortable using VAR models, so simply adding lags of u_{it} to equation i is a natural extension of the VAR model which could give a more parsimonious representation. It also has the advantage of putting the simple structure on the MA polynomials, the part which complicates the estimation, rather than the AR part as in the final AR equation form. Notice that in VARMA models, it is not necessary to include lags of all the innovations u_{1t}, \dots, u_{Kt} in every equation. This could entice practitioners to consider VARMA models if it is combined with a simple regression-based estimation method. For this representation to be useful, it needs to be identified. This is demonstrated in Theorem 3.8 below under the following assumptions and using Lemma 3.7 below.

Assumption 3.4 *The matrices $\Phi(z)$ and $\Theta(z)$ have the following form:*

$$\Phi(z) = I_K - \Phi_1 z - \dots - \Phi_p z^p, \quad \Theta(z) = I_K - \Theta_1 z - \dots - \Theta_q z^q.$$

Assumption 3.5 *$\Theta(z)$ is diagonal:*

$$\Theta(z) = \text{diag}[\theta_{ii}(z)]$$

where $\theta_{ii}(z) = 1 - \theta_{ii,1} z - \dots - \theta_{ii,q_i} z^{q_i}$ and $\theta_{ii,q_i} \neq 0$.

Assumption 3.6 *For each $i = 1, \dots, K$, there are no roots common to $\Phi_{i\bullet}(z)$ and $\theta_{ii}(z)$, i.e. there is no value z^* such that $\Phi_{i\bullet}(z^*) = 0$ and $\theta_{ii}(z^*) = 0$.*

Lemma 3.7 *Let $[\Phi(z), \Theta(z)]$ and $[\bar{\Phi}(z), \bar{\Theta}(z)]$ be two pairs of polynomial matrices which satisfy the Assumptions 3.4 to 3.6. If R_0 is a positive constant such that*

$$\Phi(z)^{-1} \Theta(z) = \bar{\Phi}(z)^{-1} \bar{\Theta}(z)$$

for $0 \leq |z| < R_0$, then

$$\Phi(z) = \bar{\Phi}(z) \text{ and } \Theta(z) = \bar{\Theta}(z), \forall z.$$

The proof of this lemma as well as other propositions appear in the Appendix. In Lemma 3.7, the condition

$$\Phi(z)^{-1} \Theta(z) = \bar{\Phi}(z)^{-1} \bar{\Theta}(z) \tag{3.9}$$

could be replaced by

$$\Theta(z)^{-1} \Phi(z) = \bar{\Theta}(z)^{-1} \bar{\Phi}(z) \tag{3.10}$$

since by assumption the inverse of $\Theta(z)$ and $\bar{\Theta}(z)$ exist. The assumptions **3.4** to **3.6** and conditions in Lemma **3.7** allow $\det[\Phi(z)]$ and $\det[\Theta(z)]$ to have roots on or inside the unit circle $|z| = 1$. It should be noted that Assumption **3.6** is weaker than the hypothesis that $\det[\Phi(L)]$ and $\det[\Theta(L)]$ have no common roots, which would be a generalization of the usual identification condition for ARMA models.

Theorem 3.8 (Identification of diagonal MA equation form representation) *Let the VARMA model be defined by equations (2.1) - (2.6) and let Assumptions **3.4-3.6** hold. If the VARMA model is in diagonal MA equation form, then it is identified.*

Similarly, we can demonstrate that the final MA equation form representation is identified under the following assumption.

Assumption 3.9 *There are no roots common to $\Phi(z)$ and $\theta(z)$, i.e. there is no value z^* such that $\Phi(z^*) = 0$ and $\theta(z^*) = 0$.*

Theorem 3.10 (Identification of final MA equation form representation) *Let the VARMA model be defined by equations (2.1)-(2.6) and let Assumptions **3.4** and **3.9** hold. If the VARMA model is in final MA equation form, then it is identified.*

From equation (3.6), we see that it is always possible to obtain a diagonal MA equation form representation starting from any VARMA representation. One case where we would obtain a diagonal and not final MA representation is when there are common factors across rows of columns of $\Theta(L)$ as in (3.8).

A strong appeal of the diagonal and final MA equation form representations is that it is easy to get the equivalent (in term of autocovariances) invertible MA representation of a non-invertible representation. With ARMA models, we simply have to invert the roots of the MA polynomial which are inside the unit circle and adjust the standard deviation of the innovations (divide it by the square of these roots): see Hamilton (1994, Section 3.7). The same procedure could be applied to VARMA models in diagonal or final MA equation form.

For VARMA representations where no particular simple structure is imposed on the MA part, at the moment we are not aware of an algorithm to go from the non-invertible to the invertible representation though theoretically this invertible representation exist and is unique as long as $\det[\Theta(z)] \neq 0$ for $|z| = 1$; see Hannan and Deistler (1988, chapter 1, section 3). So it might be troublesome to use a nonlinear optimization with these VARMA representations since we don't know how to go from the non-invertible to the invertible representation.

We can also consider the following natural generalization of the final AR equation form, where we simply replace the scalar AR operator by a diagonal operator.

Definition 3.11 (Diagonal AR equation form) *The VARMA representation (2.1) is said to be in diagonal AR equation form if $\Phi(L) = \text{diag}[\varphi_{ii}(L)] = I_K - \Phi_1 L - \dots - \Phi_p L^p$ where $\varphi_{ii}(L) = 1 - \varphi_{ii,1} L - \dots - \varphi_{ii,p_i} L^{p_i}$ and $p = \max_{1 \leq i \leq K} (p_i)$.*

Assumption 3.12 For each $i = 1, \dots, K$, there are no roots common to $\varphi_{ii}(z)$ and $\Theta_{i\bullet}(z)$, i.e. there is no value z^* such that $\varphi_{ii}(z^*) = 0$ and $\Theta_{i\bullet}(z^*) = 0$.

Theorem 3.13 (Identification of diagonal AR equation form representation) Let the VARMA model be defined by equations (2.1)-(2.6) and let Assumptions 3.4 and 3.12 hold. If the VARMA model is in diagonal AR equation form, then it is identified.

From Theorem 3.8, we can see that one way to ensure identification is to impose constraints on the MA operator. This is an alternative approach to the ones developed for example in Hannan (1971, 1976) where the identification is obtained by restricting the autoregressive part to be lower triangular with $\deg[\varphi_{ik}(L)] \leq \deg[\varphi_{ii}(L)]$ for $k > i$, or in the final AR equation form where $\Phi(L)$ is scalar. It may be more interesting to impose constraints on the moving average part instead because it is this part which causes problems in the estimation of VARMA models. Other identified representations which do not have a simple MA operator include the reversed echelon canonical form [see Poskitt (1992)] where the rows of the VARMA model in echelon form are permuted so that the Kronecker indices are ordered from smallest to largest, and the scalar component model [see Tiao and Tsay (1989)] where contemporaneous linear transformations of the vector process are considered. A general treatment of algebraic and topological structure underlying VARMA models is given in Hannan and Kavalieris (1984b). For the maximum likelihood estimation of linear state space models, data driven local coordinates are often used. See e.g. Ribarits, Deistler, and McKelvey (2004) and McKelvey, Helmersson, and Ribarits (2004). Theorem 2.7.1 in Hannan and Deistler (1988) provides general conditions for a class of ARMAX models to be identifiable. These conditions are satisfied by the proposed representations.

4. Estimation

We next introduce elements of notation for the parameters of our model. First, irrespective of the VARMA representation employed, we split the whole vector of parameters γ in two parts γ_1 (the parameters for the AR part) and γ_2 (MA part):

$$\gamma = \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}. \quad (4.1)$$

For a VARMA model in diagonal MA equation form, γ_1 and γ_2 are

$$\gamma_1 = [\varphi_{1\bullet,1}, \dots, \varphi_{1\bullet,p}, \dots, \varphi_{K\bullet,1}, \dots, \varphi_{K\bullet,p}]', \quad (4.2)$$

$$\gamma_2 = [\theta_{11,1}, \dots, \theta_{11,q_1}, \dots, \theta_{KK,1}, \dots, \theta_{KK,q_K}]', \quad (4.3)$$

while for a VARMA model in final MA equation form, γ_2 is

$$\gamma_2 = [\theta_1, \dots, \theta_q]'$$

For VARMA models in diagonal AR equation form, we simply invert γ_1 and γ_2 :

$$\gamma_1 = [\varphi_{11,1}, \dots, \varphi_{11,p_1}, \dots, \varphi_{KK,1}, \dots, \varphi_{KK,p_K}]', \quad (4.4)$$

$$\gamma_2 = [\theta_{1\bullet,1}, \dots, \theta_{1\bullet,q}, \dots, \theta_{K\bullet,1}, \dots, \theta_{K\bullet,q}]', \quad (4.5)$$

while for a VARMA model in final AR equation form,

$$\gamma_1 = [\varphi_1, \dots, \varphi_p]'. \quad (4.6)$$

The estimation method involves three steps. The observations go from $t = 1, \dots, T$.

Step 1. Estimate a VAR(n_T) to approximate the VARMA(p, q) and keep the residuals that we will call \hat{U}_t :

$$\hat{U}_t = Y_t - \sum_{j=1}^{n_T} \hat{\Pi}_j^{(n_T)} Y_{t-j} \quad \text{for } t = n_T + 1, \dots, T, \quad (4.7)$$

with $T > (K + 1)n_T$.

Step 2. With the residuals from step 1, compute an estimate of the covariance matrix of U_t , $\hat{\Sigma}_U = \frac{1}{T} \sum_{t=n_T+1}^T \hat{U}_t \hat{U}_t'$ and estimate by GLS the multivariate regression

$$\Phi(L)Y_t = [\Theta(L) - I_K] \hat{U}_t + e_t, \quad (4.8)$$

to get estimates $\tilde{\Phi}(L)$ and $\tilde{\Theta}(L)$ of $\Phi(L)$ and $\Theta(L)$. The estimator is

$$\tilde{\gamma} = \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1} \right]^{-1} \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} Y_t \right] \quad (4.9)$$

where $l = n_T + \max(p, q) + 1$. Setting

$$\mathbf{Y}_{t-1}^{(p)} = [y_{1,t-1}, \dots, y_{K,t-1}, \dots, y_{1,t-p}, \dots, y_{K,t-p}], \quad (4.10)$$

$$\hat{\mathbf{U}}_{t-1}^{(q)} = [\hat{u}_{1,t-1}, \dots, \hat{u}_{K,t-1}, \dots, \hat{u}_{1,t-q}, \dots, \hat{u}_{K,t-q}], \quad (4.11)$$

$$\mathbf{y}_{t-1}^{(k)} = [y_{k,t-1}, \dots, y_{k,t-p_k}], \quad (4.12)$$

$$\hat{\mathbf{u}}_{t-1}^{(k)} = [\hat{u}_{k,t-1}, \dots, \hat{u}_{k,t-q_k}], \quad (4.13)$$

the matrix \hat{Z}_{t-1} for the various representations is:

$$\hat{Z}_{t-1}^{DMA} = \begin{bmatrix} \mathbf{Y}_{t-1}^{(p)} & \cdots & 0 & -\hat{\mathbf{u}}_{t-1}^{(1)} & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \mathbf{Y}_{t-1}^{(p)} & 0 & \cdots & -\hat{\mathbf{u}}_{t-1}^{(K)} \end{bmatrix}, \quad (4.14)$$

$$\hat{Z}_{t-1}^{FMA} = \begin{bmatrix} \mathbf{Y}_{t-1}^{(p)} & \cdots & 0 & -\hat{\mathbf{u}}_{t-1}^{(1)} \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & \mathbf{Y}_{t-1}^{(p)} & -\hat{\mathbf{u}}_{t-1}^{(K)} \end{bmatrix}, \quad (4.15)$$

$$\hat{Z}_{t-1}^{DAR} = \begin{bmatrix} \mathbf{y}_{t-1}^{(1)} & \cdots & 0 & -\hat{\mathbf{U}}_{t-1}^{(q)} & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \mathbf{y}_{t-1}^{(K)} & 0 & 0 & -\hat{\mathbf{U}}_{t-1}^{(q)} \end{bmatrix}, \quad (4.16)$$

$$\hat{Z}_{t-1}^{FAR} = \begin{bmatrix} \mathbf{y}_{t-1}^{(1)} & -\hat{\mathbf{U}}_{t-1}^{(q)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{y}_{t-1}^{(K)} & 0 & 0 & -\hat{\mathbf{U}}_{t-1}^{(q)} \end{bmatrix}, \quad (4.17)$$

where *DMA*, *FMA*, *DAR* and *FAR* respectively stands for Diagonal MA, Final MA, Diagonal AR and Final AR equation form.

Step 3. Using the second step estimates, we first form new residuals

$$\tilde{U}_t = Y_t - \sum_{i=1}^p \tilde{\Phi}_i Y_{t-i} + \sum_{j=1}^q \tilde{\Theta}_j \tilde{U}_{t-j} \quad (4.18)$$

initiating with $\tilde{U}_t = 0$, $t \leq \max(p, q)$, and we define

$$X_t = \sum_{j=1}^q \tilde{\Theta}_j X_{t-j} + Y_t, \quad (4.19)$$

$$W_t = \sum_{j=1}^q \tilde{\Theta}_j W_{t-j} + \tilde{U}_t, \quad (4.20)$$

initiating with $X_t = W_t = 0$ for $t \leq \max(p, q)$. We also compute a new estimate of Σ_U , $\tilde{\Sigma}_U = \frac{1}{T} \sum_{t=l'}^T \tilde{U}_t \tilde{U}_t'$, with $l' = \max(p, q) + 1$. Then we regress by GLS $\tilde{U}_t + X_t - W_t$ on \tilde{V}_{t-1} with

$$\tilde{V}_t = \sum_{j=1}^q \tilde{\Theta}_j \tilde{V}_{t-j} + \tilde{Z}_t \quad (4.21)$$

where \tilde{Z}_t is just like \hat{Z}_t from step 2 except that it is computed with \tilde{U}_t instead of \hat{U}_t to obtain regression coefficients that we call $\hat{\Phi}_i$ and $\hat{\Theta}_j$:

$$\hat{\gamma} = \left[\sum_{t=l'}^T \tilde{V}_{t-1}' \tilde{\Sigma}_U^{-1} \tilde{V}_{t-1} \right]^{-1} \left[\sum_{t=l'}^T \tilde{V}_{t-1}' \tilde{\Sigma}_U^{-1} [\tilde{U}_t + X_t - W_t] \right]. \quad (4.22)$$

The properties of the above estimators are summarized in the following three theorems. Theorem 4.1 is a generalization of results from Lewis and Reinsel (1985) where convergence is demonstrated for mixing rather than i.i.d. innovations. We denote the Euclidean norm by $\|B\|^2 = \text{tr}(B'B)$.

Theorem 4.1 (VARMA first step estimates) *Let (i) the VARMA model be defined by equations (2.1)-(2.6); (ii) for some $\varepsilon > 0$ the strong mixing condition (2.16) holds and $E[|u_{it}|^{4+2\varepsilon}] < \infty$, $\forall i$. If*

$n_T/\log(T) \rightarrow \infty$ and $n_T^2/T \rightarrow 0$ as $T \rightarrow \infty$, then for the first stage estimates

$$\sum_{j=1}^{n_T} \|\hat{\Pi}_j^{(n_T)} - \Pi_j\| = O_p(n_T T^{-1/2}). \quad (4.23)$$

Theorem 4.2 (VARMA second step estimates) *Under the assumptions of Theorem 4.1 and the assumption that the VARMA model is identified, then the second stage estimator converge in probability to the true value and*

$$\sqrt{T}(\tilde{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, \tilde{J}^{-1} \tilde{I} \tilde{J}^{-1})$$

where

$$\tilde{I} = \sum_{j=-\infty}^{\infty} E \left[\{Z'_{t-1} \Sigma_U^{-1} U_t\} \{Z'_{t-1-j} \Sigma_U^{-1} U_{t-j}\}' \right], \quad \tilde{J} = E [Z'_{t-1} \Sigma_U^{-1} Z_{t-1}], \quad (4.24)$$

and Z_{t-1} is equal to the matrix \hat{Z}_{t-1} where \hat{U}_t is replaced by U_t . Further, if $m_T^4/T \rightarrow 0$ with $m_T \rightarrow \infty$ then the matrix \tilde{I} and \tilde{J} can be consistently estimated in probability respectively by

$$\tilde{I}_T = \frac{1}{T} \sum_{j=-m_T}^{m_T} \omega(j, m_T) \sum_{t=l+|j|}^T \{\hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \tilde{U}_t\} \{\hat{Z}'_{t-1-j} \hat{\Sigma}_U^{-1} \tilde{U}_{t-j}\}', \quad (4.25)$$

$$\tilde{J}_T = \frac{1}{T} \sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1}, \quad (4.26)$$

with $\omega(j, m_T) = 1 - |j|/(m_T + 1)$.

Theorem 4.3 (VARMA third step estimates) *Under the assumptions of Theorem 4.2, the third stage estimator converge in probability to the true value, and*

$$\sqrt{T}(\hat{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, \hat{J}^{-1} \hat{I} \hat{J}^{-1}) \quad (4.27)$$

with

$$\hat{I} = \sum_{j=-\infty}^{\infty} E \left[\{V'_{t-1} \Sigma_U^{-1} U_t\} \{V'_{t-1-j} \Sigma_U^{-1} U_{t-j}\}' \right], \quad \hat{J} = E [V'_{t-1} \Sigma_U^{-1} V_{t-1}] \quad (4.28)$$

and V_{t-1} is equal to the matrix \tilde{V}_{t-1} where \tilde{U}_t is replaced by U_t . Further, if $m_T^4/T \rightarrow 0$ with $m_T \rightarrow \infty$ then the matrix \hat{I} and \hat{J} can be consistently estimated in probability respectively by

$$\hat{I}_T = \frac{1}{T} \sum_{j=-m_T}^{m_T} \omega(j, m_T) \sum_{t=l'+|j|}^T \{\tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} \tilde{U}_t\} \{\tilde{V}'_{t-1-j} \tilde{\Sigma}_U^{-1} \tilde{U}_{t-j}\}', \quad (4.29)$$

$$\hat{J}_T = \frac{1}{T} \sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} \tilde{V}_{t-1}, \quad (4.30)$$

where \tilde{U}_t are the filtered residuals computed with $\hat{\gamma}$.

Notice the simplicity of this estimation method. Only three regressions are needed so we can

avoid all the caveats associated with nonlinear optimizations. This is an important problem with VARMA models where one typically deals with a high number of parameters and numerical convergence may be hard to obtain. This is especially important when we consider the fact that the asymptotic distribution of our estimators, on which we would base our inference, may be a bad approximation to the finite-sample distribution in high-dimensional dynamic models. Because of this, an estimation procedure which only requires linear methods is interesting since it suggests that simulation-based procedures – bootstrap techniques for example – should be used, something that would be impractical if the estimation is based on non-linear optimizations.

It is also important to mention that this procedure is not specific to the representations considered in this work. The expressions can be easily adapted to other identified representation, *e.g.* the echelon form. Since our estimation method is only based on regressions we can afford to use a less parsimonious representation whereas for nonlinear method it is highly important to keep the number of parameters to a minimum. An advantage of the proposed diagonal MA and final MA representations is that if the second step estimates do not correspond to an invertible MA representation (roots inside the unit circle), it is easy to get the corresponding invertible representation² to be able to perform Step 3.

For the estimation of VARMA models the emphasis has been on maximizing the likelihood (minimizing by nonlinear least squares) quickly. There are two ways of doing this. The first is having quick and efficient algorithm to evaluate the likelihood [*e.g.* Luceño (1994) and the reference therein, Mauricio (2002), Shea (1989)]. The second is to find preliminary consistent estimates that can be computed quickly to initialize the optimization algorithm. We are not the first to present a generalization to VARMA models of the Hannan and Rissanen (1982) estimation procedure for ARMA models [whose asymptotic properties are further studied in Zhao-Guo (1985) and Saikkonen (1986)]; see also Durbin (1960), Hannan and Kavalieris (1984a), Hannan, Kavalieris, and Mackisack (1986), Poskitt (1987), Koreisha and Pukkila (1990a, 1990b, 1995), Pukkila, Koreisha, and Kallinen (1990), Galbraith and Zinde-Walsh (1994, 1997). A similar method in three steps is also presented in Hannan and Kavalieris (1984a) where the third step is presented as a correction to the second step estimates.

There are many variations around the innovation-substitution approach for the estimation of VARMA models but with the exception of Hannan and Kavalieris (1984b)³ and us, none use a third step to get efficient estimators, surely because these procedures are often seen as a way to get initial values to start up a nonlinear optimization [*e.g.*, see Poskitt (1992), Koreisha and Pukkila (1989), Lütkepohl and Claessen (1997)]. In one of them, Koreisha and Pukkila (1989), lagged and current innovations are replaced by the corresponding residuals and a regression is performed. This is asymptotically the same as the first two steps of our method. Other variations are described in Hannan and Kavalieris (1986), Hannan and Deistler (1988), Huang and Guo (1990), Spliid (1983), Reinsel, Basu, and Yap (1992), Poskitt and Lütkepohl (1995), Lütkepohl and Poskitt (1996b) and Flores de Frutos and Serrano (2002). Another approach is to use the link that exist between the VARMA parameters and the infinite VAR and VMA representations. See Galbraith, Ullah, and

²See the comments after Theorem 3.10.

³They use a similar third step that is presented as a correction to the second step estimator but suggest that the third step should be iterated. They assume that the $\{U_i\}$ is a m.d.s.

Zinde-Walsh (2000) for the estimation of VMA models using a VAR. VARMA models can also be estimated with subspace methods, which is based on multiple regressions and a weighted singular value decomposition [see Bauer and Wagner (2002, 2008), Bauer (2005a, 2005b)].

Here, however, we supply a distributional theory which holds under much weaker assumptions. In the articles cited above, the data generating processes considered have innovations that are either i.i.d. or at a minimum form a martingale difference sequence. This allow us to study a broader class of models, *e.g.* temporally aggregated processes, marginalized processes, weak representation of nonlinear models.

We can ask ourselves what is the cost of not doing the nonlinear estimation. For a given sample size we will certainly lose some efficiency because of the first step estimation. We can nonetheless compare the asymptotic variance matrix of our estimator with the corresponding nonlinear estimator. We first can see that if the innovations are a m.d.s., then the asymptotic variance of our linear estimator is the same as the variance of maximum likelihood estimates under Gaussianity. The variance of MLE for i.i.d. Gaussian innovations is given in Lütkepohl (1993):

$$I = plim \left[\frac{1}{T} \sum_{t=1}^T \frac{\partial U_t'}{\partial \gamma} \Sigma^{-1} \frac{\partial U_t}{\partial \gamma'} \right]^{-1}. \quad (4.31)$$

We can transform this expression so as to obtain an equation more closely related to our previous results. First, we split γ in the same two vectors γ_1 (the AR parameters) and γ_2 (the MA parameters), then we compute $\partial U_t / \partial \gamma_1'$ and $\partial U_t / \partial \gamma_2'$. We know that

$$U_t = Y_t - \Phi_1 Y_{t-1} - \dots - \Phi_p Y_{t-p} + \Theta_1 U_{t-1} + \dots + \Theta_q U_{t-q}. \quad (4.32)$$

So taking the derivative with respect to γ_1' :

$$\frac{\partial U_t}{\partial \gamma_1'} = -Z_{\bullet 1: \dim(\gamma_1), t-1} + \Theta_1 \frac{\partial U_{t-1}}{\partial \gamma_1'} + \dots + \Theta_q \frac{\partial U_{t-q}}{\partial \gamma_1'}, \quad (4.33)$$

$$\Theta(L) \frac{\partial U_t}{\partial \gamma_1'} = -Z_{\bullet 1: \dim(\gamma_1), t-1}, \quad (4.34)$$

$$\frac{\partial U_t}{\partial \gamma_1'} = -\Theta(L)^{-1} Z_{\bullet 1: \dim(\gamma_1), t-1}, \quad (4.35)$$

where $Z_{\bullet 1: \dim(\gamma_1), t-1}$ is the first $\dim(\gamma_1)$ columns of Z_{t-1} . Similarly, the derivative with respect to γ_2' is

$$\begin{aligned} \frac{\partial U_t}{\partial \gamma_2'} &= -Z_{\bullet \dim(\gamma_1)+1: \dim(\gamma), t-1} + \Theta_1 \frac{\partial U_{t-1}}{\partial \gamma_2'} + \dots + \Theta_q \frac{\partial U_{t-q}}{\partial \gamma_2'} \\ &= -\Theta(L)^{-1} Z_{\bullet \dim(\gamma_1)+1: \dim(\gamma), t-1} \end{aligned} \quad (4.36)$$

Combining the two expressions we see that

$$\frac{\partial U_t}{\partial \gamma'} = -V_{t-1} \quad (4.37)$$

so the variance matrix for maximum likelihood estimates I is equal to the matrix J^{-1} from the third step estimation. Moreover if U_t is a m.d.s. we see that we have the equality $J = I$ so that the asymptotic variance matrix that we get in the third step of our method is the same as one would get by doing the maximum likelihood. For the weak VARMA case, from the results in Boubacar Maïnassara and Francq (2011) we know that the asymptotic covariance matrix of the QMLE estimator of γ is equal to $J^{-1}IJ^{-1}$ with

$$I = 4 \sum_{k=-\infty}^{\infty} \text{Cov} \left[U_t \Sigma^{-1} \frac{\partial U_t}{\partial \gamma'} ; U_{t-k} \Sigma^{-1} \frac{\partial U_{t-k}}{\partial \gamma'} \right], \quad J = 2E \left[\frac{\partial U_t'}{\partial \gamma} \Sigma^{-1} \frac{\partial U_t}{\partial \gamma'} \right] \quad (4.38)$$

In our previous results we saw that $\partial U_t / \partial \gamma' = V_{t-1}$. From this we see that $J = 2\hat{J}$, $I = 4\hat{I}$ and our third-step estimator have the same asymptotic variance-covariance matrix as maximum likelihood or non-linear least squares estimators depending on the properties of the innovations. To get a feel for the loss of efficiency in finite samples due to replacing the true innovations by residuals from a long VAR we performed Monte Carlo simulations and report the results in section 6. As pointed out by a referee, from the result in (4.38), we can interpret the third step as a Newton-Raphson step in the minimization of the sum of squared residuals $\sum_{t=1}^T U_t' U_t$. It can also be seen as a GLS correction induced by the MA structure in the error we make when we replace the true error term U_t by the first step residual \hat{U}_t [see Reinsel, Basu, and Yap (1992)].

5. Order selection

We still have unknowns in our model, the orders of the AR and MA operators. If no theory specifies these parameters, we have to use a statistical procedure to choose them. We propose the following information criterion method to choose the orders for VARMA models in the different identified representations proposed in Section 3. In the second step of the estimation, we compute for all $p \leq P$ and $q \leq Q$ the following information criterion:

$$\log(\det \tilde{\Sigma}_U) + \dim(\gamma) \frac{(\log T)^{1+\delta}}{T}. \quad (5.1)$$

We then choose \hat{p} and \hat{q} as the set which minimizes the information criterion. We assume that the upper bound P and Q on the order of the AR and MA part are bigger than the true values of p and q (or that they slowly grow with the sample size). The properties of \hat{p} and \hat{q} are summarized in the following theorem.

Theorem 5.1 (Estimation of the order p and q in VARMA models) *Under the assumptions of Theorem 4.1, if $n_T = O((\log T)^{1+\delta_1})$ from some $\delta_1 > 0$ with $\delta_1 < \delta$ and the orders are chosen by minimizing (5.1), then \hat{p} and \hat{q} converge in probability to their true value.*

In practice, this procedure can lead to a search over too many models for the diagonal representations. A valid alternative is to search for the true orders by proceeding equation by equation. In the second step of the estimation, instead of doing a simultaneous estimation, just perform univariate regressions. For a VARMA model in diagonal MA equation form, regress

$$y_{it} = \sum_{j=1}^{p_i} \sum_{k=1}^K \varphi_{ik,j} y_{k,t-j} - \sum_{j=1}^{q_i} \theta_{ii,j} \hat{u}_{i,t-j} + e_{it}, \quad (5.2)$$

for $i = 1, \dots, K$, while for a VARMA models in diagonal AR equation form we regress

$$y_{it} = \sum_{j=1}^{p_i} \varphi_{ii,j} y_{i,t-j} - \sum_{j=1}^{q_i} \sum_{k=1}^K \theta_{ik,j} \hat{u}_{k,t-j} + e_{it}. \quad (5.3)$$

We then chose \hat{p}_i and \hat{q}_i as the orders which minimize the following information criterion:

$$\log(\hat{\sigma}_i^2) + g(p_i, q_i) \frac{(\log T)^{1+\delta}}{T} \quad (5.4)$$

where $\delta > 0$ and $g(p_i, q_i) = p_i K + q_i$ or $g(p_i, q_i) = p_i + q_i K$ for the diagonal MA or AR equation form representation respectively. The global order for the autoregressive operator is then $\hat{p} = \max(\hat{p}_1, \dots, \hat{p}_K)$ for the diagonal MA representation and, similarly for the diagonal AR representation, $\hat{q} = \max(\hat{q}_1, \dots, \hat{q}_K)$. We see that this equation by equation selection procedure is not only easier to apply, it can lead to more parsimonious representations by identifying rows of zeros coefficients in Φ_j or Θ_j .

Theorem 5.2 (Estimation of the order p and q in diagonal VARMA models) *Under the assumption of Theorem 5.1, if the VARMA model is in either the diagonal MA or AR equation form and the orders are chosen by minimizing (5.4), then \hat{p}_i and \hat{q}_i , $i = 1, \dots, K$, converge in probability to their true value.*

The criterion in equation (5.1) is a generalization of the information criterion proposed by Hannan and Rissanen (1982) which the authors acknowledged that it must in fact be modified to provide consistent estimates of the order, p and q . The original criterion was

$$\log \tilde{\sigma}^2 + (p + q) \frac{(\log T)^\delta}{T} \quad (5.5)$$

with $\delta > 0$. But in Hannan and Rissanen (1983) they acknowledged that $\tilde{\sigma}^2 - \frac{1}{T} \sum_{t=1}^T u_t^2$ is $O_p(n_T T^{-1})$ and not $O_p(T^{-1})$ so the penalty $(\log T)^\delta / T$ is not strong enough. The authors proposed two possible modifications to their procedure. The simpler is to take $(\log T)^{1+\delta}$ instead of $(\log T)^\delta$ in the information criterion so that the penalty on $p + q$ will dominate $\log \hat{\sigma}^2$ in the criterion. The second, which they favored and was used in later work [see Hannan and Kavalieris (1984b)], is to modify the first step of the procedure. Instead of taking $n_T = O(\log T)$ they used another information criterion to choose the order of the long autoregression and they iterated the

whole procedure picking a potentially different p and q at every iteration. A similar approach is also proposed in Poskitt (1987). In this work we prefer the first solution so as to keep the procedure as simple as possible.

For the identification of the order of VARMA models, it all depends on the representation used. Although it was one of the first representation studied, not much work has been done with the final AR equation form. People felt that this representation gives VARMA models with too many parameters. A complete procedure to fit VARMA models under this representation is given in Lütkepohl (1993): One would first fit an ARMA(p_i, q_i) model to every univariate time series, using maybe the procedure of Hannan and Rissanen (1982). To build the VARMA representation in final AR equation form, knowing that the VAR operator is the same for every equation we would take it to be the product of all the univariate AR polynomials. This would give a VAR operator of order $p = \sum_{i=1}^K p_i$. Accordingly, for the VMA part we would take $q = \max_k [q_k + \sum_{i=1, i \neq k}^K p_i]$. It is no wonder that people feel that the final equation form uses too many parameters.

For VARMA models in echelon form, there has been a lot more work done on the identification of Kronecker indices. The problem has been studied by, among others, Hannan and Kavalieris (1984b), Poskitt (1992) and Lütkepohl and Poskitt (1996b). Non-stationary or cointegrated systems are considered by Huang and Guo (1990), Bartel and Lütkepohl (1998), and Lütkepohl and Claessen (1997). Additional references are given in Lütkepohl (1993, Chapter 8).

As for weak VARMA models estimated by QMLE, Boubacar Maïnassara (2012) propose a modified Akaike's information criteria for selecting the orders p and q .

A complementing approach to specify VARMA models, which is based on Cooper and Wood (1982), aims at finding simplifying structures via some combinations of the different series to obtain more parsimonious models. It includes Tiao and Tsay (1989), Tsay (1989a, 1989b, 1991) and Nsiri and Roy (1992, 1996).

The final stage of ARMA model specification usually involve analyzing the residuals, *i.e.* checking if they are uncorrelated. Popular tools include portmanteau tests such as Box-Pierce [Box and Pierce (1970)] and Ljung-Box [Ljung and Box (1978)] tests, and their multivariate generalization [Lütkepohl (1993, Section 5.2.9)]. Those tests are not directly applicable in our case because they are derived under strong assumptions for the innovations (independence or martingale difference). But recent developments for weak ARMA and VARMA models are applicable. They include Francq, Roy, and Zakoïan (2005) and Shao (2011) (weak ARMA), Francq and Raïssi (2007) (weak VAR), Boubacar Maïnassara (2011) and Katayama (2012) (weak VARMA).

6. Monte Carlo simulations

To illustrate the performance of our estimation method we ran two types of simulations. For the first type, weak VARMA models were simulated where the innovations are not independent nor a m.d.s. but merely uncorrelated. The second type of simulations involves strong VARMA models (VARMA models with i.i.d. Gaussian innovations). All the simulated models are bivariate so the results are easier to analyze. The results are generated using Ox version 3.30 on Linux [see Doornik (1999)]. We performed 1000 simulations for each model. The results with strong VARMA models being comparable to those for weak VARMA models, we only report results for the latter.

We simulate weak VARMA processes by directly simulating weak innovations, from which we build the simulated series. From the results in Drost and Nijman (1993), we know that the temporal aggregation of a strong GARCH process (where the standardized innovations are i.i.d.) will give a weak process⁴. Suppose \tilde{U}_t is given by the following bivariate ARCH model:

$$\tilde{U}_t = H_t^{1/2} \varepsilon_t, \quad H_t = \Omega + \alpha \tilde{U}_{t-1} \tilde{U}_{t-1}' \quad (6.1)$$

where ε_t is i.i.d. $N(0, I_2)$, $H_t^{1/2}$ is the Cholesky decomposition of H_t and α is a scalar. If we consider \tilde{U}_t as a stock variable, then temporal aggregation of \tilde{U}_t over two periods, *i.e.*

$$U_t = \tilde{U}_{2t} \quad (6.2)$$

will give a weak process. The series U_t will be uncorrelated but not a m.d.s., its mean will be zero and the variance will be $\Omega(1 - \alpha^2)/(1 - \alpha)$.

In these examples, because the innovations are not a m.d.s., we cannot do maximum likelihood. We instead employ nonlinear generalized least-squares (GLS), *i.e.* we minimize the nonlinear least squares, compute an estimate of the variance matrix of the innovations and then do nonlinear GLS. We did not apply this procedure, partly to reduce the estimation time in our Monte Carlo study, partly because there is no asymptotic gain in iterating.

In these simulations the sample size is 250 observations, which represent about 20 years of monthly data, a reasonable sample size for macroeconomic data. Tables 1 gives results for a VARMA model in final MA equation form [VARMA(1,1)], while results for VARMA models in diagonal MA equation form are given in Tables 2 and 3 [VARMA(1,1) with $q = (1, 1)$ and VARMA(2,1) with $q = (1, 1)$ respectively]. We present the results (mean, standard deviations, root mean square error, 5% quantile, 95% quantile and median) for the second (when the number of parameters does not exceed five) and third step estimates, and the nonlinear GLS estimates (using the true value of the parameters as initial values). Samples for which the optimization algorithm did not converge were dropped (this happened for less than 1% of the simulations). In our simulations, we took

$$\Omega = \begin{bmatrix} 1.0 & 0.7 \\ 0.7 & 1.0 \end{bmatrix}, \quad \alpha = 0.3. \quad (6.3)$$

From looking at the RMSE, a first thing to notice is that there can be sizable improvement in doing the third step. Some of the third step RMSEs in Tables 1 and 2 are more than 50% smaller than for the second step. This is an interesting observation considering that the third step basically involve only one extra regression. Comparing the third step RMSEs and the RMSEs for the nonlinear GLS estimates, we see that the former are usually no more than 15% bigger. This is also an interesting observation. The cost of avoiding a numerical optimization, which can become quite challenging as the number of time series studied or order of the operators increases, appears to be small.

In the top part of these tables we also present the results for the selection of the order of the

⁴Another way of simulating a weak VARMA process is to time-aggregate a strong VARMA process with innovations that have skewed marginal distributions (*e.g.*, a mixture of two Gaussian distributions with different means but mean zero unconditionally). We can appeal to the results of Francq and Zakoian (1998, Section 2.2.1) to claim that the resulting VARMA is only weak.

operators using our proposed information criterion. For models in final MA equation form, we have to select the orders p and q , and for models in diagonal MA equation, the selection is over p , q_1 and q_2 . In Table 1, we see that for VARMA models in final MA equation form the most frequently chosen orders are the true ones, and the criterion will tend to pick a higher value for q than for p . This result might partially be skewed by the fact that the simulated models have a highly persistent moving average ($\theta_1 = 0.9$). For VARMA models in diagonal equation form (Tables 2 and 3), we get similar results, the orders selected with the highest frequency are the true ones.

7. Application to a macroeconomics model of the U.S. monetary policy

To illustrate our estimation method and the gains that can be obtained from using a more parsimonious representation, we fit VARMA and VAR models to six macroeconomic time series and compute the impulse-response functions generated by each model. What people typically do to get the impulse-response functions is first fit a VAR to their multiple time series and then get the implied infinite VMA representation. The order of the VAR required for macro series is usually high. For example, Bernanke and Mihov (1998) use a VAR(13) to model six monthly macroeconomic time series when about 30 years of data are available. The resulting standard errors for the impulse-response functions are very large, like in most macroeconomic study. We can ask ourselves how much of this is due to the fact that so many parameters are estimated. To try to answer this we will study the impulse-response functions generated by VARMA models estimated on the same data. We will concentrate on VARMA models in final MA equation form.

Our example is based on McMillin (2001) who compare numerous identification restrictions for the structural effects of monetary policy shocks using the same dataset as Bernanke and Mihov (1998).⁵ The series are plotted in Figure 1. One of the model studied is a VAR applied to the first difference of the series, in order, gdp_m , $(pssc_m - pgdp_m)$, $fyff$, $nbrec1$, $tr1$, $pssc_m$. With an argument based on Keating (2002), the author state that using this ordering of the variables the Cholesky decomposition, based on long-run macroeconomic restrictions, which are described in an appendix, of the variance matrix of the innovations will identify the structural effects of the policy variable $nbrec1$ without imposing any contemporaneous restrictions among the variables. Since the model is in first difference, the impulse-response at a given order is the cumulative shocks up to that order.

By fitting a VAR(12) to these series we get basically the same impulse-response functions and confidence bands as in McMillin (2001) They are plotted in Figure 2. The impulse-response function for the output and federal funds rate tends to zero as the order increases which is consistent with the notion that a monetary variable does not have a long term impact on real variables. The

⁵The dataset consist of the log of the real GDP (gdp_m), total bank reserves ($tr1$), nonborrowed reserves ($nbrec1$), federal funds rate ($fyff$), log of the GDP deflator ($pgdp_m$), log of the Dow-Jones index of spot commodity prices ($pssc_m$). These are monthly data and cover the period January 1962 to December 1996. The monthly data for real GDP and the GDP deflator were constructed by state space methods, using a list of monthly interpolator variables and assuming that the interpolation error is describable as an AR(1) process. Both total reserves and nonborrowed reserves are normalized by a 36-month moving average of total reserves.

Table 1. Estimation of a weak final MA equation form VARMA(1,1).

The simulated model is a weak VARMA(1,1) in final MA equation form with $\varphi_{11,1} = 0.5$, $\varphi_{21,1} = 0.7$, $\varphi_{12,1} = -0.6$, $\varphi_{22,1} = 0.3$ and $\theta_1 = 0.9$. The variance of the innovations is 1.3 and the covariance is 0.91. Sample size is 250, the length of the long AR is $n_T = 15$, the number of repetition is 1000. The parameter in the criterion is $\delta = 0.5$.

Frequencies of selection of (\hat{p}, \hat{q}) using the information criteria.

$p \setminus q$	0	1	2	3	4	5
0	0.000	0.000	0.000	0.000	0.000	0.000
1	0.000	0.953	0.025	0.002	0.001	0.001
2	0.000	0.001	0.014	0.003	0.000	0.000
3	0.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	0.000	0.000	0.000

	Value	Average	Std. dev.	RMSE	5%	95%	Median
Second step							
$\varphi_{11,1}$	0.5	0.441	0.065	0.088	0.328	0.544	0.444
$\varphi_{21,1}$	0.7	0.675	0.060	0.065	0.576	0.770	0.677
$\varphi_{12,1}$	-0.6	-0.631	0.054	0.062	-0.717	-0.540	-0.632
$\varphi_{22,1}$	0.3	0.229	0.057	0.091	0.134	0.321	0.230
θ_1	0.9	0.825	0.057	0.095	0.725	0.917	0.826
Third step							
$\varphi_{11,1}$	0.5	0.491	0.055	0.055	0.399	0.580	0.494
$\varphi_{21,1}$	0.7	0.695	0.053	0.054	0.603	0.779	0.698
$\varphi_{12,1}$	-0.6	-0.601	0.049	0.049	-0.680	-0.519	-0.600
$\varphi_{22,1}$	0.3	0.294	0.050	0.051	0.204	0.375	0.294
θ_1	0.9	0.887	0.034	0.037	0.830	0.940	0.886
NLLS							
$\varphi_{11,1}$	0.5	0.495	0.050	0.051	0.412	0.579	0.496
$\varphi_{21,1}$	0.7	0.702	0.049	0.049	0.621	0.781	0.702
$\varphi_{12,1}$	-0.6	-0.609	0.043	0.044	-0.681	-0.540	-0.609
$\varphi_{22,1}$	0.3	0.288	0.046	0.048	0.214	0.366	0.287
θ_1	0.9	0.887	0.028	0.031	0.838	0.929	0.888

impulse response of the price level increases as we let the order grow and does not revert to zero.

We next estimate VARMA models for the four representations proposed in this work. The information criterion picked a VARMA(3,10) for the final MA representation. The impulse-response functions for this model are plotted in Figure 3. The behavior of the impulse-response function for GDP, the federal funds rate and the price level from the VARMA models are similar to what we obtained with a VAR. The most notable differences are that the initial decrease in the federal funds rate is smaller (0.20 versus 0.32 percentage point) and the GDP is peaking earlier with the VARMA.

It is not surprising that VAR and VARMA models are giving similar impulse-response functions since they both are a way of getting an infinite MA representation. What is more interesting is the comparison of the width of the confidence bands for the VAR and VARMA's impulse-response functions.⁶ For GDP and the federal funds rate, we see that the bands are much smaller for the VARMA model and they shrink more quickly as the horizon increases. The confidence bands for these two variables should be collapsing around their IRF since there should be no long-term effect of the policy variable so the uncertainty should decrease as the horizon increases. The situation is different for the price level. For this variable the confidence band grows with the order. Again this is not so surprising because we expect that a change in the non-borrowed reserves should have a long-term impact on the price level. With a non-dying impact it is natural that the uncertainty about this impact can grow as time passes.

The result that the confidence bands around IRFs can be shorter with a VARMA than with a VAR could be expected since these models are simple extensions of the VAR approach. The introduction of a simple MA operator allows the reduction of the required AR order so we can get more precise estimates, which translate into more precise impulse-response functions.

Another way of comparing the performance of VAR and VARMA models is to compare their out-of-sample forecasts using a metric (*e.g.*, RMSE as in our example). Employing the same dataset as above, we recursively estimated the models and computed the out-of-sample forecasts, starting at observation 300 until the end of the sample. The orders of the different models are chosen by minimizing the RMSE over the possible values⁷. The results for the VAR, VARMA diagonal MA and VARMA final MA representations are presented in Table 4. We see that reduction of up to 12% of the RMSE can be obtained by using a VARMA model instead of a VAR, the greatest gain being for one-step ahead VARMA in final MA representation.

8. Conclusion

In this paper, we proposed a modeling and estimation method which ease the use of VARMA models. We first propose new identified VARMA representations, the final MA equation form and the diagonal MA equation form. These two representations are simple extensions of the class of VAR models where we add a simple MA operator, either a scalar or a diagonal operator. The addition of

⁶The confidence bands are computed by performing a parametric bootstrap using Gaussian innovations.

⁷For the VARMA diagonal MA representation we don't search over all the possible orders because it would involve the estimation of too many models. We instead proceed in two steps. We first impose that all the q_i orders are equal which gives us an upper bound for the value of MA orders. In a second step, one q_i after the other we check to see if a lower order for the MA order of that equation would lower the RMSE.

Figure 1. Macroeconomic series.

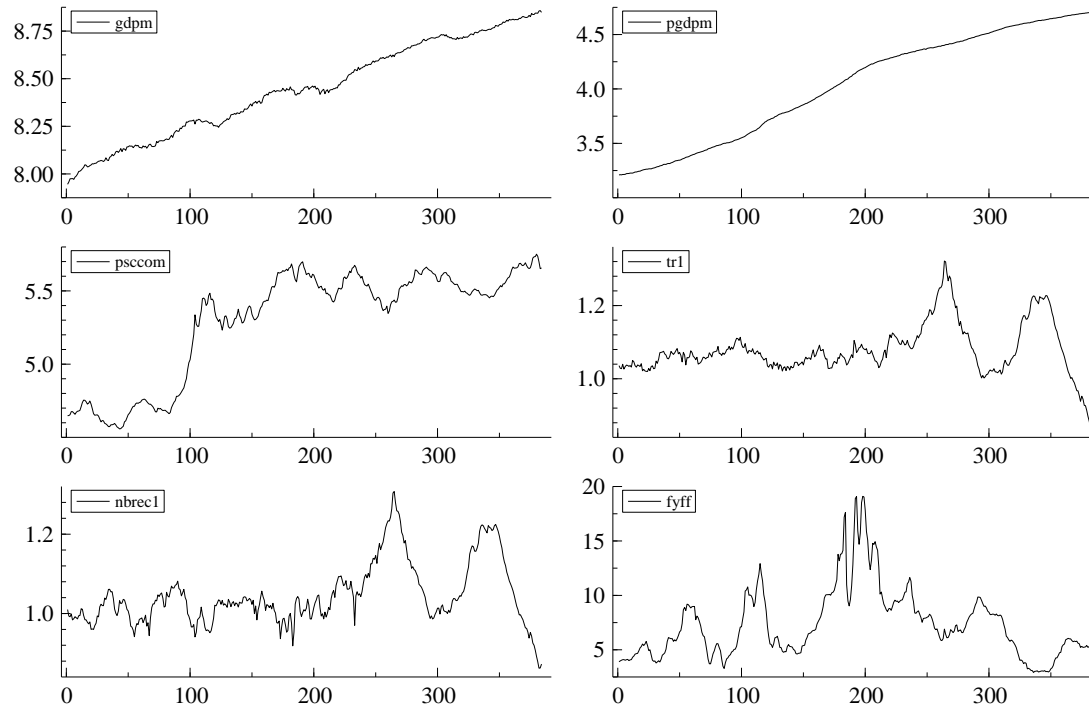


Figure 2. Impulse-response functions for VAR model.

A VAR(12) is fitted to the first difference of the six time series. The confidence band represent a one standard deviation. The standard deviations are derived from a parametric bootstrap using Gaussian innovations.

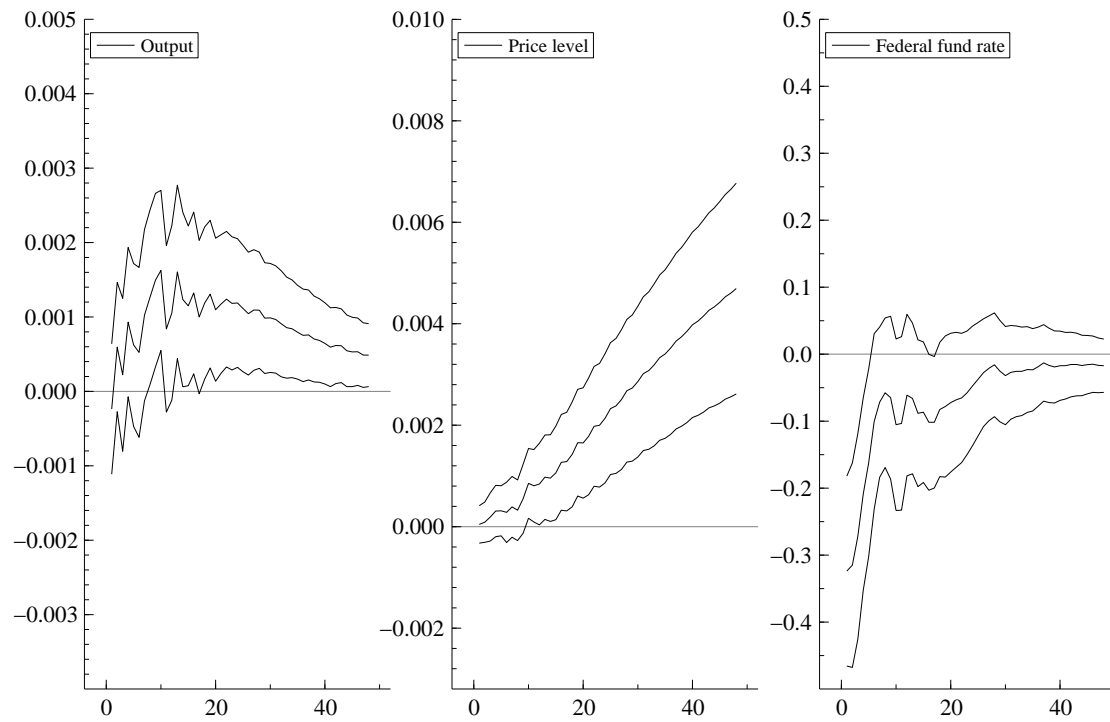


Figure 3. Impulse-response functions for VARMA model in final MA equation form.

A VARMA(3,10) is fitted to the first difference of the six time series. The confidence band represent a one standard deviation. The standard deviations are derived from a parametric bootstrap using Gaussian innovations.

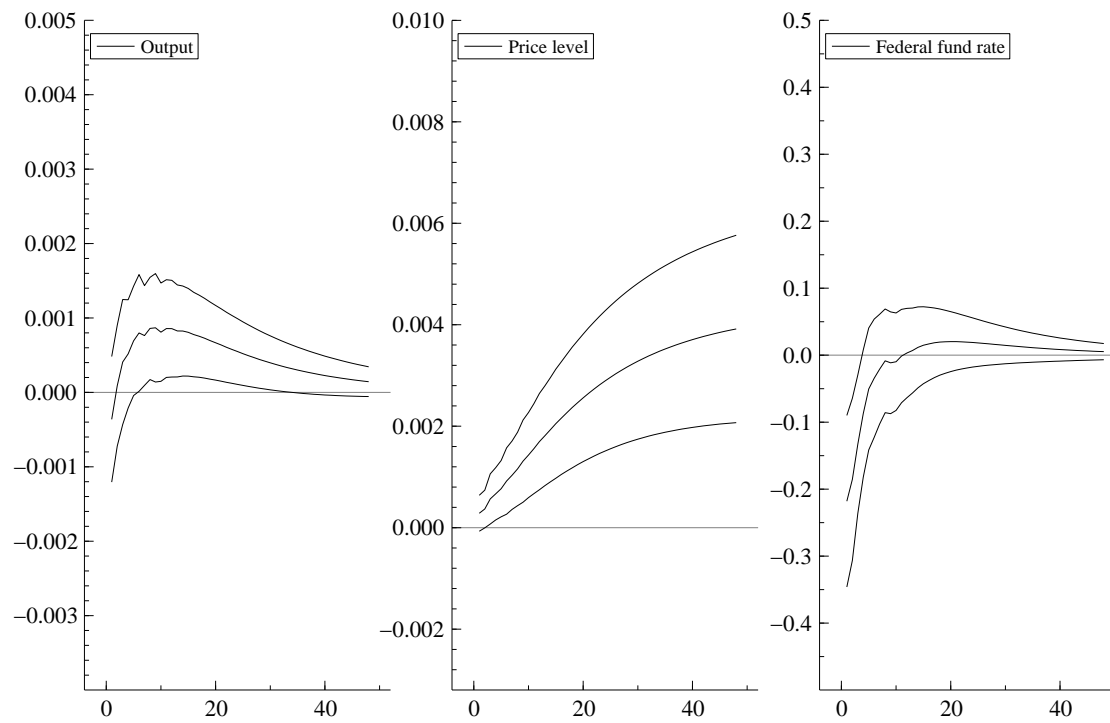


Table 4. RMSE for VAR and VARMA models]RMSE for VAR and VARMA models

Step ahead	VAR	VARMA diag. MA	VARMA final MA
1	0.0829 $p = 1$	0.0778 $p = 0$ $q = (2, 1, 2, 1, 0, 1)$	0.0725 $p = 0$ $q = 18$
3	0.0794 $p = 1$	0.0775 $p = 1$ $q = (1, 1, 1, 1, 0, 1)$	0.0729 $p = 1$ $q = 15$
6	0.0822 $p = 7$	0.0826 $p = 1$ $q = (0, 1, 2, 2, 2, 1)$	0.0773 $p = 0$ $q = 18$
9	0.0867 $p = 2$	0.0803 $p = 3$ $q = (8, 11, 11, 1, 11, 11)$	0.0798 $p = 0$ $q = 18$
12	0.0836 $p = 2$	0.0805 $p = 3$ $q = (6, 6, 6, 3, 1, 6)$	0.0807 $p = 0$ $q = 18$

a MA part can give more parsimonious representations, yet the simple form of the MA operators does not introduce undue complications.

To ease the estimation we studied the problem of estimating VARMA models by relatively simple methods which only require linear regressions. For that purpose, we considered a generalization of the regression-based estimation method proposed by Hannan and Rissanen (1982) for univariate ARMA models. Our method is in three steps. In a first step a long VAR is fitted to the data. In the second step, the lagged innovations in the VARMA model are replaced by the corresponding lagged residuals from the first step and a regression is performed. In a third step, the data from the second step are filtered and another regression is performed. We showed that the third-step estimators have the same asymptotic variance as their nonlinear counterpart (Gaussian maximum likelihood if the innovations are i.i.d., or generalized nonlinear least squares if they are merely uncorrelated). In the non i.i.d. case, we consider strong mixing conditions, rather than the usual martingale difference sequence assumption. We make these minimal assumptions on the innovations to broaden the class of models to which this method can be applied.

We also proposed a modified information criterion that gives consistent estimates of the orders of the AR and MA operators of the proposed VARMA representations. This criterion is to be minimized in the second step of the estimation method over a set of possible values for the different orders.

Monte Carlo simulation results indicates that the estimation method works well for small sample sizes and the information criterion picks the true value of the order p and q most of the time. These results holds for sample sizes commonly used in macroeconomics, *i.e.* 20 years of monthly data or 250 sample points. To demonstrate the importance of using VARMA models to study multivariate

time series we compare the impulse-response functions and the out-of-sample forecasts generated by VARMA and VAR models when these models are applied to the dataset of macroeconomic time series used by Bernanke and Mihov (1998).

A. Proofs

Lemma A.1 Let U and V be random variables measurable with respect to $\mathcal{F}_{-\infty}^0$ and \mathcal{F}_n^∞ , respectively where \mathcal{F}_a^b is the σ -algebra of events generated by sets of the form $\{(X_{i_1}, X_{i_2}, \dots, X_{i_n}) \in E_n\}$ with $a \leq i_1 < i_2 < \dots < i_n \leq b$, and E_n is some n -dimensional Borel set. Let r_1, r_2, r_3 be positive numbers. Assume that $\|U\|_{r_1} < \infty$ and $\|V\|_{r_2} < \infty$ where $\|U\|_r = (E[|U|^r])^{1/r}$. If $r_1^{-1} + r_2^{-1} + r_3^{-1} = 1$, then there exists a positive constant c_0 independent of U, V and n , such that

$$|E[UV] - E[U]E[V]| \leq c_0 \|U\|_{r_1} \|V\|_{r_2} \alpha(n)^{1/r_3}.$$

where $\alpha(n)$ is defined in equation (2.15).

Proof. See Davydov (1968).

Lemma A.2 If the random process $\{y_t\}$ is stationary and satisfies the strong mixing condition (2.15), with $E|y_t|^{2+\varepsilon} < \infty$ for some $\varepsilon > 0$, and if $\sum_{j=1}^{\infty} \alpha(j)^{\varepsilon/(2+\varepsilon)} < \infty$, then

$$\begin{aligned} \sigma^2 &\equiv \lim_{T \rightarrow \infty} \text{Var}[y_1 + \dots + y_T] \\ &= E[(y_t - E[y_t])^2] + 2 \sum_{j=1}^{\infty} E[(y_t - E[y_t])(y_{t+j} - E[y_{t+j}])]. \end{aligned}$$

Moreover, if $\sigma \neq 0$ and $E[y_t] = 0$, then

$$\Pr \left[\frac{y_1 + \dots + y_T}{\sigma \sqrt{T}} < z \right] \xrightarrow{T \rightarrow \infty} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du.$$

Proof. See Ibragimov (1962).

Proof of Lemma 3.7

Clearly, $\Phi(0) = \Theta(0) = I_K$ and $\det[\Phi(0)] = \det[\Theta(0)] = 1 \neq 0$. The polynomials $\det[\Phi(z)]$ and $\det[\Theta(z)]$ are different from zero in a neighborhood of zero. So we can choose $R_0 > 0$ such that $\det[\Phi(z)] \neq 0$ and $\det[\Theta(z)] \neq 0$ for $0 \leq |z| < R_0$. It follows that the matrices $\Phi(z)$ and $\Theta(z)$ are invertible for $0 \leq |z| < R_0$.

Let

$$C_0 = \{ |z| < R_0 \}$$

and

$$\Psi(z) = \Phi(z)^{-1} \Theta(z)$$

for $z \in C_0$. Since

$$\Phi(z)^{-1} = \frac{1}{\det[\Phi(z)]} \Phi^*(z), \quad \Theta(z)^{-1} = \frac{1}{\det[\Theta(z)]} \Theta^*(z),$$

where $\Phi^*(z)$ and $\Theta^*(z)$ are matrices of polynomials, it follows that, for $z \in C_0$, each element of $\Phi(z)^{-1}$ and $\Theta(z)^{-1}$ is a rational function whose denominator is different from zero. Thus, for

$z \in C_0$, $\Phi(z)^{-1}$ and $\Theta(z)^{-1}$ are matrices of analytic functions, and the function

$$\Psi(z) = \Phi(z)^{-1}\Theta(z)$$

is analytic in the circle $0 \leq |z| < R_0$. Hence, it has a unique representation of the form

$$\Psi(z) = \sum_{k=0}^{\infty} \Psi_k z^k, \quad z \in C_0.$$

By assumption,

$$\Psi(z) = \Phi(z)^{-1}\Theta(z) = \bar{\Phi}(z)^{-1}\bar{\Theta}(z)$$

for $z \in C_0$. Hence, for $z \in C_0$,

$$\begin{aligned} \bar{\Phi}(z)\Phi(z)^{-1}\Theta(z) &= \bar{\Theta}(z), \\ \bar{\Phi}(z)\Phi(z)^{-1} &= \bar{\Theta}(z)\Theta(z)^{-1} \equiv \Delta(z), \end{aligned} \tag{A.1}$$

where $\Delta(z)$ is a diagonal matrix because $\Theta(z)$ and $\bar{\Theta}(z)$ are both diagonal,

$$\Delta(z) = \text{diag}[\delta_{ii}(z)],$$

where

$$\delta_{ii}(z) = \frac{\bar{\theta}_{ii}(z)}{\theta_{ii}(z)}, \quad \theta_{ii}(0) = 1, \quad \delta_{ii}(0) = \bar{\theta}_{ii}(0), \quad i = 1, \dots, K. \tag{A.2}$$

From (A.2), it follows that each $\delta_{ii}(z)$ is rational with no pole in C_0 such that $\delta_{ii}(0) = 1$, so it can be written in the form

$$\delta_{ii}(z) = \frac{e_i(z)}{f_i(z)}$$

where $e_i(z)$ and $f_i(z)$ have no common roots, $f_i(z) \neq 0$ for $z \in C_0$ and $\delta_{ii}(0) = e_i(0) = 1$. From (A.1), it follows that for $z \in C_0$

$$\bar{\theta}_{ii}(z) = \delta_{ii}(z)\theta_{ii}(z), \quad \bar{\varphi}_{ij}(z) = \delta_{ii}(z)\varphi_{ij}(z), \quad i, j = 1, \dots, K.$$

We first show that $\delta_{ii}(z)$ must be a polynomial. If $f_i(z) \neq 1$, then its order cannot be greater than the order $q_i \equiv \deg[\theta_{ii}(z)]$ for otherwise $\bar{\theta}_{ii}(z)$ would not be a polynomial. Similarly, if $f_i(z) \neq 1$ and is a polynomial of order less or equal to q_i , then all its roots must be roots of $\theta_{ii}(z)$ and $\varphi_{ij}(z)$, for otherwise $\bar{\theta}_{ii}(z)$ or $\bar{\varphi}_{ij}(z)$ would be a rational function. If $q_i \geq 1$, these roots are then common to $\theta_{ii}(z)$ and $\varphi_{ij}(z)$, $j = 1, \dots, K$, which is in contradiction with Assumption 3.6. Thus the degree of $f_i(z)$ must be zero, and $\delta_{ii}(z)$ is a polynomial.

If $\delta_{ii}(z)$ is a polynomial of degree greater than zero, this would entail that $\bar{\theta}_{ii}(z)$ and $\bar{\varphi}_{ij}(z)$ have roots in common, in contradiction with Assumption 3.6. Thus $\delta_{ii}(z)$ must be a constant. Further, $\delta_{ii}(0) = 1$ so that for $i = 1, \dots, K$,

$$\bar{\theta}_{ii}(z) = \theta_{ii}(z), \quad \bar{\varphi}_{ij}(z) = \varphi_{ij}(z), \quad j = 1, \dots, K,$$

hence

$$\bar{\Theta}(z) = \Theta(z), \quad \bar{\Phi}(z) = \Phi(z).$$

Proof of Theorem 3.8. Under the assumption that the VARMA process is invertible, we can write

$$\Theta(L)^{-1}\Phi(L)Y_t = U_t.$$

Now suppose by contradiction that there exist operators $\bar{\Phi}(L)$ and $\bar{\Theta}(L)$, with $\bar{\Theta}(L)$ diagonal and invertible, and $\bar{\Phi}(L) \neq \Phi(L)$ or $\bar{\Theta}(L) \neq \Theta(L)$, such that

$$\bar{\Theta}(L)^{-1}\bar{\Phi}(L) = \Theta(L)^{-1}\Phi(L),$$

If the above equality hold, then it must also be the case that

$$\bar{\Theta}(z)^{-1}\bar{\Phi}(z) = \Theta(z)^{-1}\Phi(z), \quad \forall z \in C_0,$$

where $C_0 = \{z \in \mathbb{C} \mid 0 \leq |z| < R_0\}$ and $R_0 > 0$. By Lemma 3.7, it follows that

$$\bar{\Phi}(z) = \Phi(z), \quad \bar{\Theta}(z) = \Theta(z) \quad \forall z.$$

Hence, the representation is unique.

Proof of Theorem 3.10. The proof can be easily adapted from the proof of Theorem 3.8 once we replace Assumption 3.6 by Assumption 3.9.

Lemma A.3 (Infinite VAR convergence) *If the VARMA model is invertible and if $n_T/\log(T) \rightarrow \infty$ as $T \rightarrow \infty$, then $\sum_{k=1}^K \sum_{j=n_T+1}^{\infty} |\pi_{ik,j}| = o(T^{-1})$ for $i = 1, \dots, K$, where $\pi_{ik,j}$ represent the parameters in $\Pi(L) = \Theta(L)^{-1}\Phi(L)$.*

Proof of Lemma A.3. The matrix $\Theta(L)^{-1}$ can be seen has its adjoint matrix divided by its determinant. Since Y_t is invertible, the roots of $\det\Theta(L)$ are outside the unit circle and so the elements of $\Pi(L) = \Theta(L)^{-1}\Phi(L)$ decrease exponentially:

$$|\pi_{ik,j}| \leq c\rho^j, \quad \forall i, m,$$

with $c > 0$ and $0 < \rho < 1$. From this

$$\begin{aligned} T \sum_{k=1}^K \sum_{j=n_T+1}^T |\pi_{ik,j}| &\leq T \sum_{k=1}^K \sum_{j=n_T+1}^T c\rho^j \\ &\leq cKT \frac{\rho^{n_T+1}}{1-\rho} \\ &\rightarrow 0 \end{aligned}$$

as $T \rightarrow \infty$ if $n_T/\log(T) \rightarrow \infty$ because $|\rho| < 1$.

From the proof of Lemma A.3, we see that the condition $n_T/\log T \rightarrow \infty$ could be replaced by a weaker condition like $n_T = \kappa \log(T)$ with $\kappa > 1/\log(\rho)$ where ρ is the value given the upper bound at which the parameters $\pi_{ik,j}$ are declining to zero. A drawback if this assumption is that it would depend on the persistence of the process.

Lemma A.4 (Covariance estimation) *If the process $\{Y_t\}$ is a strictly stationary VARMA process with $\{U_t\}$ uncorrelated, $E[|u_{it}|^{4+2\varepsilon}] < \infty$ for some $\varepsilon > 0$, α -mixing with $\sum_{h=1}^{\infty} \alpha(h)^{\varepsilon/(2+\varepsilon)} < \infty$ then*

$$\frac{1}{T} \sum_{t=1}^T y_{i,t-r} y_{i',t-s} - E[y_{i,t-r} y_{i',t-s}] = O_{ms}(T^{-1/2}) \quad \forall i, i', r, s.$$

Proof of Lemma A.4. In a preliminary step, let us prove that the following result holds (assuming that $s > r$ without loss of generality):

$$\frac{1}{T^2} \sum_{t=1}^T \sum_{t'=1}^T \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] = O(1/T). \quad (\text{A.3})$$

We start by breaking this sum in the following parts:

$$\begin{aligned} & \frac{1}{T^2} \sum_{t=1}^T \sum_{t'=1}^T \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] \\ = & \frac{1}{T^2} \sum_{t=1}^{T-(s-r)-1} \sum_{t'=t+1+(s-r)}^T \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] \\ & + \frac{1}{T^2} \sum_{t'=1}^{T-(s-r)-1} \sum_{t=t'+1+(s-r)}^T \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] \\ & + \frac{1}{T^2} \sum_{t=1+(s-r)}^{T-(s-r)} \sum_{t'=t-(s-r)}^{t+(s-r)} \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] \\ & + \frac{1}{T^2} \sum_{t=1}^{1+(s-r)} \sum_{t'=1}^{t+(s-r)} \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] \\ & + \frac{1}{T^2} \sum_{t'=T-(s-r)}^T \sum_{t=T-(s-r)-(T-t')}^T \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}]. \end{aligned} \quad (\text{A.4})$$

The last three double sums are $O(1/T)$ since the covariances are finite and the number of terms is of order T . For the first two double sums, using Davydov's inequality (lemma A.1), the strong mixing hypothesis and the finite fourth moment we know that

$$\lim_{T \rightarrow \infty} \sum_{t'=t+1+(s-r)}^T | \text{Cov} [u_{i,t-r} u_{i',t-s}; u_{i,t'-r} u_{i',t'-s}] |$$

$$\begin{aligned}
&\leq \lim_{T \rightarrow \infty} \sum_{t'=t+1+(s-r)}^T c_0 \|u_{i,t-r} u_{i',t-s}\|_{2+\varepsilon} \|u_{i,t'-r} u_{i',t'-s}\|_{2+\varepsilon} \alpha (t' - t - (s-r))^{\varepsilon/(2+\varepsilon)} \\
&< \infty,
\end{aligned}$$

from which we conclude that the first two terms converge to zero at rate $1/T$.

Now that have the result in Equation (A.3), we first notice that by stationarity of the process,

$$E \left[\frac{1}{T} \sum_{t=1}^T y_{i,t-r} y_{i',t-s} \right] - E[y_{i,t-r} y_{i',t-s}] = 0.$$

Now taking the variance and writing the covariances in terms of the innovations U_t :

$$\begin{aligned}
\text{Var} \left[\frac{1}{T} \sum_{t=1}^T y_{i,t-r} y_{i',t-s} \right] &= \frac{1}{T^2} \sum_{t=1}^T \sum_{t'=1}^T \text{Cov} [y_{i,t-r} y_{i',t-s}; y_{i,t'-r} y_{i',t'-s}] \\
&\leq \sum_{j_1=0}^{\infty} \sum_{j'_1=0}^{\infty} \sum_{j_2=0}^{\infty} \sum_{j'_2=0}^{\infty} \sum_{k_1=1}^K \sum_{k'_1=1}^K \sum_{k_2=1}^K \sum_{k'_2=1}^K |\psi_{ik_1, j_1}| |\psi_{i'k'_1, j'_1}| |\psi_{ik_2, j_2}| |\psi_{i'k'_2, j'_2}| \\
&\quad \frac{1}{T^2} \sum_{t=1}^T \sum_{t'=1}^T |\text{Cov} [u_{k_1, t-r-j_1} u_{k'_1, t-s-j'_1}; u_{k_2, t'-r-j_2} u_{k'_2, t'-s-j'_2}]|. \quad (\text{A.5})
\end{aligned}$$

From the assumption of stationarity we know that the ψ 's are decreasing exponentially, and from Equation (A.3) we get that the right-hand side of Equation (A.5) is $O(1/T)$. Hence,

$$\frac{1}{T} \sum_{t=1}^T y_{i,t-r} y_{i',t-s} - E[y_{i,t-r} y_{i',t-s}] = O_{ms}(T^{-1/2}) \quad \forall i, i', r, s.$$

Corollary A.5 (Moment estimation) *Under the assumption of Lemma A.4,*

$$\frac{1}{T} \sum_{t=1}^T y_{i,t-r} u_{i',t-s} - E[y_{i,t-r} u_{i',t-s}] = O_{ms}(T^{-1/2}) \quad \forall i, i', r, s.$$

Proof of Lemma A.5. The proof is very similar to the proof of Lemma A.4 where in Equation (A.5) some of the ψ 's would be zero.

Proof of Theorem 4.1. We first introduce some additional matrix norms:

$$\|B\|_2^2 = \sup_{l \neq 0} \frac{l' B' B l}{l' l}, \quad (\text{A.6})$$

$$\|B\|_1 = \max_{i \leq j \leq n} \sum_{i=1}^n |b_{ij}|, \quad (\text{A.7})$$

$$\|B\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |b_{ij}|, \quad (\text{A.8})$$

where (A.6) is the largest eigenvalue of $B'B$. Useful inequalities relating these norms are given in Horn and Johnson (1985, p. 313):

$$\|AB\|^2 \leq \|A\|_2^2 \|B\|^2, \quad \|AB\|^2 \leq \|A\|^2 \|B\|_2^2, \quad \|B\|_2^2 \leq \|B\|_1 \|B\|_\infty. \quad (\text{A.9})$$

In the first step estimation, we regress

$$y_{it} = \sum_{j=1}^{n_T} \sum_{k=1}^K \pi_{ik,j} y_{k,t-j} + e_{it}, \quad (\text{A.10})$$

when in fact

$$y_{it} = \sum_{j=1}^{\infty} \sum_{k=1}^K \pi_{ik,j} y_{k,t-j} + u_{it}.$$

If we let

$$\hat{B}(n_T) = \sum_{t=n_T+1}^T \frac{\mathbf{Y}'_{t-1}(n_T) \mathbf{Y}_{t-1}(n_T)}{T - n_T},$$

then OLS applied to (A.10) yields:

$$\begin{aligned} \hat{\Pi}_{i\bullet}^{(n_T)} &= [\hat{\pi}_{i\bullet,1}, \dots, \hat{\pi}_{i\bullet,n_T}]' \\ &= \hat{B}(n_T)^{-1} \sum_{t=n_T+1}^T \frac{\mathbf{Y}_{t-1}^{(n_T)'} y_{it}}{T - n_T} \\ &= \hat{B}(n_T)^{-1} \sum_{t=n_T+1}^T \frac{\mathbf{Y}_{t-1}^{(n_T)'}}{T - n_T} \left\{ \sum_{j=1}^{\infty} \pi_{i\bullet,j} Y_{t-j} + u_{it} \right\} \\ &= \Pi_{i\bullet}^{(n_T)} + \hat{B}(n_T)^{-1} \sum_{t=n_T+1}^T \frac{\mathbf{Y}_{t-1}^{(n_T)'}}{T - n_T} \left\{ \sum_{j=n_T+1}^{\infty} \pi_{i\bullet,j} Y_{t-j} + u_{it} \right\}. \end{aligned}$$

Rearranging the elements,

$$\hat{\Pi}_{i\bullet}^{(n_T)} - \Pi_{i\bullet}^{(n_T)} = \hat{B}(n_T)^{-1} \sum_{t=n_T+1}^T \frac{\mathbf{Y}_{t-1}^{(n_T)'}}{T - n_T} \left\{ \sum_{j=n_T+1}^{\infty} \pi_{i\bullet,j} Y_{t-j} \right\} + \hat{B}(n_T)^{-1} \sum_{t=n_T+1}^T \frac{\mathbf{Y}_{t-1}^{(n_T)'}}{T - n_T} u_{it}.$$

Using inequalities (A.9) and the fact that $\hat{B}(n_T)$ is symmetric,

$$\|\hat{\Pi}_{i\bullet}^{(n_T)} - \Pi_{i\bullet}^{(n_T)}\| \leq \|\hat{B}(n_T)^{-1}\|_2 \|V_{1T}\| + \|\hat{B}(n_T)^{-1}\|_2 \|V_{2T}\|, \quad (\text{A.11})$$

where

$$\begin{aligned} V_{1T} &= \frac{1}{T-n_T} \sum_{t=n_T+1}^T \mathbf{Y}_{t-1}^{(n_T)'} \sum_{j=n_T+1}^{\infty} \boldsymbol{\pi}_{i_{\bullet},j} Y_{t-j}, \\ V_{2T} &= \frac{1}{T-n_T} \sum_{t=n_T+1}^T \mathbf{Y}_{t-1}^{(n_T)'} \mathbf{u}_{it}. \end{aligned}$$

Firstly, $\|V_{2T}\|^2$ can be expanded into

$$\begin{aligned} \|V_{2T}\|^2 &= \text{tr}(V_{2T}' V_{2T}) \\ &= \sum_{k=1}^K \sum_{j=1}^{n_T} \left(\frac{\sum_{t=n_T+1}^T y_{k,t-j} \mathbf{u}_{it}}{T-n_T} \right)^2 \\ &= \sum_{k=1}^K \sum_{j=1}^{n_T} \left(\underbrace{E[y_{k,t-j} \mathbf{u}_{it}]}_{=0} + O_p(T^{-1/2}) \right)^2. \end{aligned}$$

It follows that $\|V_{2T}\|^2 = O_p(n_T T^{-1/2})$. Similarly, for $\|V_{1T}\|^2$

$$\begin{aligned} \|V_{1T}\|^2 &= \text{tr}(V_{1T}' V_{1T}) \\ &= \sum_{k=1}^K \sum_{j=1}^{n_T} \left(\frac{\sum_{t=n_T+1}^T y_{k,t-j} [\sum_{j'=n_T+1}^{\infty} \sum_{k'=1}^K \boldsymbol{\pi}_{ik',j'} y_{k',t-j'}]}{T-n_T} \right)^2 \\ &= \sum_{k=1}^K \sum_{j=1}^{n_T} \left(\sum_{k'=1}^K \sum_{j'=n_T+1}^{\infty} \boldsymbol{\pi}_{ik',j'} \frac{1}{T-n_T} \sum_{t=n_T+1}^T y_{k,t-j} y_{k',t-j'} \right)^2 \\ &= \sum_{k=1}^K \sum_{j=1}^{n_T} \left(\sum_{k'=1}^K \sum_{j'=n_T+1}^{\infty} \boldsymbol{\pi}_{ik',j'} \left[\text{Cov}[y_{k,t-j}; y_{k',t-j'}] + O_p(T^{-1/2}) \right] \right)^2. \end{aligned}$$

From Lemma **A.3**, we know that $\sum_{j=n_T+1}^{\infty} |\boldsymbol{\pi}_{ik,j}| = o(T^{-1})$ and it follows that $\sum_{j'=n_T+1}^{\infty} \boldsymbol{\pi}_{ik',j'} \text{Cov}[y_{k,t-j}; y_{k',t-j'}] = o(T^{-1})$. Hence, $\|V_{1T}\|^2 = o_p(n_T T^{-1})$.

For $\|\hat{\mathbf{B}}(n_T)^{-1}\|_1$, the existence of $\hat{\mathbf{B}}(n_T)^{-1}$ is guaranteed by a lemma that can be found in Tiao and Tsay (1983). The argument is the following. It is clear that $\hat{\mathbf{B}}(n_T)$ is a symmetric non-negative definite matrix. To show that it is positive definite take any arbitrary vector $c = [c_1, \dots, c_{Kn_T}]'$ and consider

$$c' \hat{\mathbf{B}}(n_T) c = \frac{1}{(T-n_T)^2} \sum_{t=n_T+1}^T \left(\sum_{j=1}^{n_T} \sum_{k=1}^K c_{(j-1)K+k} y_{k,t-j} \right)^2.$$

If $c' \hat{B}(n_T)c = 0$, then

$$\sum_{j=1}^{n_T} \sum_{k=1}^K c_{(j-1)K+k} y_{k,t-j} = 0 \quad \text{for } t = n_T + 1, \dots, T,$$

which, since $T > (K+1)n_T$, is a system of linear equations of Kn_T unknowns and more than Kn_T equations. Since Y_t is real-valued and non deterministic, this implies that $c = 0$ (except for a set with measure zero). This proves that $\hat{B}(n_T)$ is positive definite.

The final step is to show that $\|\hat{B}(n_T)^{-1}\|_2$ is bounded. We first see that

$$\|\hat{B}(n_T)^{-1}\|_2 \leq \|B(n_T)^{-1}\|_2 + \|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2$$

where $B(n_T)$ denotes the $(Kn_T \times Kn_T)$ matrix of the corresponding covariances instead of the empirical covariances. As in the univariate case Berk (1974, p. 491), $\|B(n_T)^{-1}\|_2$ is uniformly bounded above by a positive constant F for all n_T since Y_t is stationary and invertible. Next, using a similar argument as in the proof of Theorem 1 in Lewis and Reinsel (1985), we show that $\|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2 \xrightarrow{p} 0$. From previous results, $E[\|\hat{B}(n_T) - B(n_T)\|_2^2] \leq E[\|\hat{B}(n_T) - B(n_T)\|^2] \leq c_0 \frac{n_T^2}{T} \rightarrow 0$ as $T \rightarrow \infty$ for some positive constant c_0 . Then, from

$$\begin{aligned} \|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2 &= \|\hat{B}(n_T)^{-1}[\hat{B}(n_T) - B(n_T)]B(n_T)^{-1}\|_2 \\ &\leq F(\|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2 + F)\|\hat{B}(n_T) - B(n_T)\|_2, \end{aligned}$$

we have

$$0 \leq \Xi_{n_T} = \frac{\|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2}{F(\|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2 + F)} \leq \|\hat{B}(n_T) - B(n_T)\|_2$$

so that as $T \rightarrow \infty$, $\Xi_{n_T} \xrightarrow{p} 0$ and $\|\hat{B}(n_T)^{-1} - B(n_T)^{-1}\|_2 = F^2 \Xi_{n_T} / (1 - F \Xi_{n_T}) \xrightarrow{p} 0$. Hence, $\|\hat{\Pi}_{i_\bullet}^{(n_T)} - \Pi_{i_\bullet}(n_T)\| = O_p(n_T T^{-1/2})$.

Proof of Theorem 4.2. If we denote by Z_{t-1} the equivalent of \hat{Z}_{t-1} which contains the true innovations u_{kt} instead of the residuals \hat{u}_{kt} ,

$$\begin{aligned} \hat{\gamma} &= \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1} \right]^{-1} \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} (Z_{t-1} \gamma + U_t) \right] \\ &= \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1} \right]^{-1} \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} Z_{t-1} \right] \gamma + \\ &\quad \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1} \right]^{-1} \left[\sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} U_t \right]. \end{aligned}$$

Firstly, we show that $\hat{\Sigma}_U \xrightarrow{p} \Sigma_U$. We can write the residual \hat{U}_t as

$$\begin{aligned}\hat{U}_t &= \hat{\Pi}^{nr}(L)Y_t \\ &= \hat{\Pi}^{nr}(L)\Psi(L)U_t \\ &= [I_K + (\hat{\Pi}^{nr}(L)\Psi(L) - I_K)]U_t \\ &= [I_K + (\hat{\Pi}^{nr}(L) - \Pi(L))\Psi(L)]U_t \\ &= U_t + (\hat{\Pi}^{nr}(L) - \Pi(L))Y_t.\end{aligned}$$

Using the results from Lemma A.4, Theorem 4.1 where we showed that $\sum_{l=1}^{n_T} \|\hat{\Pi}_l^{(nr)} - \Pi_l\| = O_p(n_T T^{-1/2})$, combined with $\sum_{l=n_T+1}^{\infty} \|\Pi_l\| = o(T^{-1})$ if $\log(T)/n_T \rightarrow 0$ as $T \rightarrow \infty$, we can conclude that

$$\hat{\Sigma}_U = \frac{1}{T - n_T} \sum_{t=n_T+1}^T \hat{U}_t \hat{U}_t' = \frac{1}{T - n_T} \sum_{t=n_T+1}^T U_t U_t' + o_p(T^{-1/2}) \xrightarrow{p} \Sigma_U.$$

To show that $\frac{1}{T} \sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} \hat{Z}_{t-1}$ converge to $\tilde{J} = E[Z'_{t-1} \Sigma_U^{-1} Z_{t-1}]$ in probability, since $\hat{\Sigma}_U \xrightarrow{p} \Sigma_U$ we only have to show:

- $\frac{1}{T} \sum_{t=l}^T y_{i,t-j} y_{k,t-j'} \xrightarrow{p} E[y_{i,t-j} y_{k,t-j'}]$,
- $\frac{1}{T} \sum_{t=l}^T y_{i,t-j} \hat{u}_{k,t-j'} \xrightarrow{p} E[y_{i,t-j} u_{k,t-j'}]$,
- $\frac{1}{T} \sum_{t=l}^T \hat{u}_{i,t-j} \hat{u}_{k,t-j'} \xrightarrow{p} E[u_{i,t-j} u_{k,t-j'}]$.

The first is proved in Lemma A.4. The second can be proved in a similar manner. Start by writing

$$\begin{aligned}\frac{1}{T} \sum_{t=l}^T y_{i,t-j} \hat{u}_{k,t-j'} &= \frac{1}{T} \sum_{t=l}^T y_{i,t-j} u_{k,t-j'} + \frac{1}{T} \sum_{t=l}^T y_{i,t-j} (\hat{u}_{k,t-j'} - u_{k,t-j'}) \\ &= \frac{1}{T} \sum_{t=l}^T y_{i,t-j} u_{k,t-j'} + \frac{1}{T} \sum_{t=l}^T \sum_{m=1}^{n_T} \sum_{k'=1}^K (\pi_{kk',m} - \hat{\pi}_{kk',m}) y_{i,t-j} y_{k',t-m} \\ &\quad + \frac{1}{T} \sum_{t=l}^T \sum_{m=n_T+1}^{\infty} \sum_{k'=1}^K \pi_{kk',m} y_{i,t-j} y_{k',t-m}\end{aligned}\tag{A.12}$$

Proving that the first term in (A.12), $\frac{1}{T} \sum_{t=l}^T y_{i,t-j} u_{k,t-j'}$, converges in quadratic mean to $E[y_{i,t-j} u_{k,t-j'}]$ is very similar to the proof in Lemma A.4 where we express $y_{i,t-j}$ as an infinite MA so we omit the details to shorten the exposition. Proving that the second and third term converge to zero in probability is also similar; combine the results of Lemma A.4 and Theorem 4.1 for the second, Lemmas A.3 and A.4 for the third. Combining all these results we can conclude that $\tilde{\gamma} \xrightarrow{ms} \gamma$.

For the asymptotic distribution, since $\hat{\Sigma}_U \xrightarrow{p} \Sigma_U$, the limit distribution of $\frac{1}{\sqrt{T}} \sum_{t=l}^T \hat{Z}'_{t-1} \hat{\Sigma}_U^{-1} U_t$ will be the same as that of $\frac{1}{\sqrt{T}} \sum_{t=l}^T \hat{Z}'_{t-1} \Sigma_U^{-1} U_t$. For the latter, we can prove the asymptotic normality

using an argument similar to the one used in Francq and Zakoian (1998, Lemma 4). The argument is the following. Neglecting the constants in Σ_U^{-1} , $\frac{1}{\sqrt{T}} \sum_{t=l}^T \hat{Z}'_{t-1} \Sigma_U^{-1} U_t$ contains terms such $\frac{1}{\sqrt{T}} \sum_{t=l}^T u_{i,t} y_{k,t-j}$ with $i, k = 1, \dots, K$ and $j = 1, \dots, \max(p, q)$. Using the MA(∞) representation of Y_t ,

$$\frac{1}{\sqrt{T}} \sum_{t=l}^T u_{i,t} y_{k,t-j} = \frac{1}{\sqrt{T}} \sum_{t=l}^T u_{i,t} \left(\sum_{k'=1}^K \sum_{j'=0}^{\infty} \psi_{kk',j'} u_{k',t-j-j'} \right) = \frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(1)} + \frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(2)}$$

where for any positive integer r ,

$$\begin{aligned} \mathbf{A}_{r,t}^{(1)} &= \sum_{j'=0}^r \sum_{k'=1}^K \psi_{kk',j'} u_{i,t} u_{k',t-j-j'}, \\ \mathbf{A}_{r,t}^{(2)} &= \sum_{j'=r+1}^{\infty} \sum_{k'=1}^K \psi_{kk',j'} u_{i,t} u_{k',t-j-j'}. \end{aligned}$$

First note that $\mathbf{A}_{r,t}^{(1)}$ is a function of a finite number of values from the process $\{U_t\}$. Therefore, the stationary process $\{\mathbf{A}_{r,t}^{(1)}\}$ satisfies a mixing property of the form (2.16). Lemma A.2 implies that $\frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(1)}$ has a limiting distribution $\mathcal{N}(0, \tilde{\tau}_r)$ and as $r \rightarrow \infty$, $\tilde{\tau}_r \rightarrow \tilde{\tau}$.

Now we will show that $E[\frac{1}{T} (\sum_{t=l}^T \mathbf{A}_{r,t}^{(2)})^2]$ converges to 0 uniformly in T as $r \rightarrow \infty$. It will follow that the limiting distribution of $\frac{1}{\sqrt{T}} \sum_{t=l}^T u_{i,t} y_{k,t-j}$ is the same as the limiting distribution of $\frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(1)}$ from a straightforward adaptation of a result given in Anderson (1971, Corollary 7.1.1, p. 426). We have

$$\begin{aligned} & \text{Var} \left[\frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(2)} \right] \\ &= \text{Var} \left[\frac{1}{\sqrt{T}} \sum_{t=l}^T \sum_{j'=r+1}^{\infty} \sum_{k'=1}^K \psi_{kk',j'} u_{i,t} u_{k',t-j-j'} \right] \\ &\leq \sum_{j_1=r+1}^{\infty} \sum_{j_2=r+1}^{\infty} \sum_{k_1=1}^K \sum_{k_2=1}^K |\psi_{kk_1,j_1}| |\psi_{kk_2,j_2}| \frac{1}{T} \sum_{t=l}^T \sum_{t'=l}^T |\text{cov}(u_{i,t} u_{k_1,t-j-j_1}; u_{i,t'} u_{k_2,t'-j-j_2})| \\ &\leq \sum_{j_1=r+1}^{\infty} \sum_{j_2=r+1}^{\infty} \sum_{k_1=1}^K \sum_{k_2=1}^K |\psi_{kk_1,j_1}| |\psi_{kk_2,j_2}| \frac{1}{T} \sum_{t=l}^T \sum_{t'=l}^{\infty} |\text{cov}(u_{i,t} u_{k_1,t-j-j_1}; u_{i,t'} u_{k_2,t'-j-j_2})| \\ &\leq C \sum_{j_1=r+1}^{\infty} \sum_{j_2=r+1}^{\infty} \sum_{k_1=1}^K \sum_{k_2=1}^K |\psi_{kk_1,j_1}| |\psi_{kk_2,j_2}| \end{aligned}$$

for some positive constant C following a similar argument as in the proof of Lemma A.4. Thus,

$$\sup_T \text{Var} \left[\frac{1}{\sqrt{T}} \sum_{t=l}^T \mathbf{A}_{r,t}^{(2)} \right] \rightarrow 0$$

as $r \rightarrow \infty$.

We can extend this asymptotic normality to all the elements of $\frac{1}{\sqrt{T}} \sum_{t=l}^T \hat{\mathbf{Z}}'_{t-1} \hat{\Sigma}_U^{-1} U_t$ to conclude that

$$\frac{1}{\sqrt{T}} \sum_{t=l}^T \hat{\mathbf{Z}}'_{t-1} \hat{\Sigma}_U^{-1} U_t \xrightarrow{d} \mathcal{N}(0, \hat{I})$$

with \tilde{I} defined in Equation (4.24). From this,

$$\sqrt{T}(\tilde{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, \tilde{J}^{-1} \tilde{I} \tilde{J}^{-1}).$$

From the preceding results, it is obvious that \tilde{J} can be consistently estimated by \tilde{J}_T as defined in Equation (4.26) and using Theorem 2 of Newey and West (1987) or more explicit results from Francq and Zakoian (2000) for weak ARMA models, we know that $\tilde{J}_T \xrightarrow{p} \tilde{J}$ if we take $m_T^4/T \rightarrow 0$ with $m_T \rightarrow \infty$ as $T \rightarrow \infty$.

Proof of Theorem 4.3. First we can rewrite X_t , W_t and \tilde{V}_t as

$$X_t = \hat{\Theta}(L)^{-1} Y_t, \quad W_t = \hat{\Theta}(L)^{-1} \tilde{U}_t, \quad \tilde{V}_t = \hat{\Theta}(L)^{-1} \tilde{Z}_t.$$

We can also rewrite $\tilde{U}_t + X_t - W_t$ as

$$\begin{aligned} \tilde{U}_t + X_t - W_t &= \hat{\Theta}(L)^{-1} Y_t + \tilde{U}_t - \hat{\Theta}(L)^{-1} \tilde{U}_t \\ &= \hat{\Theta}(L)^{-1} [Z_{t-1} \gamma + U_t] + \tilde{U}_t - \hat{\Theta}(L)^{-1} \tilde{U}_t \\ &= \hat{\Theta}(L)^{-1} Z_{t-1} \gamma + \hat{\Theta}(L)^{-1} U_t + \tilde{U}_t - \hat{\Theta}(L)^{-1} \tilde{U}_t \\ &= V_{t-1} \gamma + U_t + [(\tilde{U}_t - U_t) - \hat{\Theta}(L)^{-1} (\tilde{U}_t - U_t)] \\ &= V_{t-1} \gamma + U_t + o_p(T^{-1/2}). \end{aligned}$$

With this, the regression becomes

$$\begin{aligned} \hat{\gamma} &= \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} \tilde{V}_{t-1} \right]^{-1} \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} (\tilde{U}_t + X_t - W_t) \right] \\ &= \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} \tilde{V}_{t-1} \right]^{-1} \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} V_{t-1} \right] \gamma + \\ &\quad \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} \tilde{V}_{t-1} \right]^{-1} \left[\sum_{t=l'}^T \tilde{V}'_{t-1} \tilde{\Sigma}_U^{-1} U_t \right] + o_p(T^{-1/2}). \end{aligned}$$

Just like in the proof of theorem 4.2 we see that $\hat{\gamma} - \gamma = O_p(T^{-1/2})$. With a similar application of Ibragimov's central limit theorem as in the proof of Theorem 4.2, we conclude that

$$\sqrt{T}(\hat{\gamma} - \gamma) \xrightarrow{d} \mathcal{N}(0, \hat{J}^{-1} \hat{I} \hat{J})$$

where \hat{I} and \hat{J} are defined in Equation (4.28). As in the proof of theorem 4.2 the matrices \hat{I} and \hat{J} can be consistently estimated respectively by \hat{I}_T and \hat{J}_T as defined in Equations (4.29) and (4.30).

Proof of Theorem 5.1.

Let us denote by $\tilde{\Sigma}_U(p, q)$ the value taken by $\tilde{\Sigma}_u$ for given orders p and q . The true value of p and q is denoted by p_0 and q_0 . The difference between the information criterion for given values of the orders p and q , and the true values p_0, q_0 is equal to

$$\log(\det \tilde{\Sigma}_U(p, q)) - \log(\det \tilde{\Sigma}_U(p_0, q_0)) + [\dim \gamma(p, q) - \dim \gamma(p_0, q_0)] \frac{(\log T)^{1+\delta}}{T}. \quad (\text{A.13})$$

First, consider the case where $p < p_0$ or $q < q_0$. In this case, as T grows to infinity, eventually $\det \tilde{\Sigma}_U(p, q) > \det \tilde{\Sigma}_U(p_0, q_0)$ because of the left-coprime property while the penalty term is shrinking to zero. As a result, (A.13) would become positive as $T \rightarrow \infty$. So eventually we must have $p \geq p_0$ and $q \geq q_0$.

Next, to discuss the case where the $p \geq p_0$ and $q \geq q_0$, we can start by writing the residuals of the second step estimation as

$$\begin{aligned} \tilde{U}_t &= \tilde{\Phi}(L)Y_t - (\tilde{\Theta}(L) - I_K) \hat{U}_t \\ &= \tilde{\Phi}(L)Y_t - (\tilde{\Theta}(L) - I_K) \hat{\Pi}^{(nr)}(L)Y_t \\ &= \left[\tilde{\Phi}(L) - (\tilde{\Theta}(L) - I_K) \hat{\Pi}^{(nr)}(L) \right] Y_t \\ &= \left[\tilde{\Phi}(L) - (\tilde{\Theta}(L) - I_K) \hat{\Pi}^{(nr)}(L) \right] \Psi_0(L)U_t \\ &= \left[\tilde{\Phi}(L) - \tilde{\Theta}(L) \hat{\Pi}^{(nr)}(L) + \hat{\Pi}^{(nr)}(L) \right] \Psi_0(L)U_t \\ &= \left[(\tilde{\Phi}(L) - \Phi_0(L)) + \Phi_0(L) - \tilde{\Theta}(L) \hat{\Pi}^{(nr)}(L) + (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) + \Pi_0(L) \right] \Psi_0(L)U_t \\ &= \left[(\tilde{\Phi}(L) - \Phi_0(L)) + (\Theta_0(L) - \tilde{\Theta}(L)) \Pi_0(L) - \tilde{\Theta}(L) (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) + \right. \\ &\quad \left. (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) + \Pi_0(L) \right] \Psi_0(L)U_t \\ &= \left[(\tilde{\Phi}(L) - \Phi_0(L)) \Psi_0(L) + (\Theta_0(L) - \tilde{\Theta}(L)) - \tilde{\Theta}(L) (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) \Psi_0(L) + \right. \\ &\quad \left. (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) \Psi_0(L) + I_K \right] U_t \\ &= \left[(\tilde{\Phi}(L) - \Phi_0(L)) \Psi_0(L) + (\Theta_0(L) - \tilde{\Theta}(L)) - \chi(L) + C(L) + I_K \right] U_t, \end{aligned} \quad (\text{A.14})$$

where $\chi(L) = \tilde{\Theta}(L) (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) \Psi_0(L)$ and $C(L) = (\hat{\Pi}^{(nr)}(L) - \Pi_0(L)) \Psi_0(L)$.

For the case where $p = p_0$ and $q = q_0$ and from the results of Theorems 4.1 and 4.2, it follows

that with an obvious abuse of notation⁸ $\|(\tilde{\Phi}(L) - \Phi_0(L))\Psi_0(L)\| = O_p(T^{-1/2})$, $\|(\Theta_0(L) - \tilde{\Theta}(L))\| = O_p(T^{-1/2})$, $\|\chi(L)\| = O_p(n_T T^{-1/2})$ and $\|C(L)\| = O_p(n_T T^{-1/2})$. Using the above representation of the residuals \tilde{U}_t , we get

$$\tilde{\Sigma}_U(p_0, q_0) = \frac{1}{T} \sum_{t=n_T+1}^T U_t U_t' + O_p(n_T T^{-1}),$$

Also, if $p \geq p_0$ and $q \geq 0$, if we rewrite (A.14) as

$$\tilde{U}_t = [\tilde{\Theta}(L) (\tilde{\Theta}(L)^{-1} \tilde{\Phi}(L) - \Pi_0(L)) \Psi_0(L)^{-1} - \chi(L) + C(L) + I_K] U_t, \quad (\text{A.15})$$

we can see that even if $p > p_0$ and $q > q_0$, in which case the VAR and MA operators are not identified, the estimated infinite VAR representation is converging to the true $\Pi_0(L)$ and as a result $\|\tilde{\Theta}(L)^{-1} \tilde{\Phi}(L) - \Pi_0(L)\| = O_p(T^{-1/2})$ and

$$\tilde{\Sigma}_U(p, q) = \frac{1}{T} \sum_{t=n_T+1}^T U_t U_t' + O_p(n_T T^{-1}) \quad (\text{A.16})$$

for $p \geq p_0$ and $q \geq q_0$.

It follows that if $n_T = O(\log(T)^{1+\delta_1})$ with $\delta_1 < \delta$, then the dominating term in (A.13) is the penalty term so as $T \rightarrow \infty$ with probability one $\hat{p} \rightarrow p_0$, $\hat{q} \rightarrow q_0$.

Proof of Theorem 5.2. The proof is very similar to the proof of Theorem 5.1.

⁸For example, by $\|(\Theta(L) - \tilde{\Theta}(L))\|^2$ we mean $\sum_{i=1}^K \sum_{k=1}^K \sum_{j=1}^q (\theta_{ik,j} - \tilde{\theta}_{ik,j})^2$.

References

- ANDERSON, T. (1971): *The Statistical Analysis of Time Series*. Wiley, New York.
- ATHANASOPOULOS, G., AND F. VAHID (2008): "VARMA versus VAR for macroeconomic forecasting," *Journal of Business and Economic Statistics*, 26(2), 237–252.
- BARTEL, H., AND LÜTKEPOHL (1998): "Estimating the Kronecker indices of cointegrated echelon-form VARMA models," *Econometrics Journal*, 1, C76–C99.
- BAUER, D. (2005a): "Comparing the CCA Subspace Method to Pseudo Maximum Likelihood Methods in the Case of No Exogenous Inputs," *Journal of Time Series Analysis*, 26(5), 631–668.
- (2005b): "Estimating Linear Dynamical Systems Using Subspace Methods," *Econometric Theory*, 21, 181–211.
- BAUER, D., AND M. WAGNER (2002): "Estimating cointegrated systems using subspace algorithms," *Journal of Econometrics*, 111, 47–84.
- (2008): "Using subspace algorithm cointegration analysis: Simulation performance and application to the term structure," *Computational Statistics and Data Analysis*, Forthcoming.
- BERK, K. N. (1974): "Consistent autoregressive spectral estimates," *The Annals of Statistics*, 2(3), 489–502.
- BERNANKE, B. S., AND I. MIHOV (1998): "Measuring Monetary Policy," *The Quarterly Journal of Economics*, 113(3), 869–902.
- BOSQ, D. (1998): *Nonparametric Statistics for Stochastic Processes - Estimation and Prediction*, no. 110 in Lecture Notes in statistics. Springer-Verlag, Berlin, second edn.
- BOUBACAR MAÏNASSARA, Y. (2011): "Multivariate portmanteau test for structural VARMA models with uncorrelated but non-independent error terms," *Journal of Statistical Planning and Inference*, 141, 2961–2975.
- (2012): "Selection of weak VARMA models by modified Akaike's information criteria," *Journal of Time Series Analysis*, 33, 121–130.
- BOUBACAR MAÏNASSARA, Y., AND C. FRANCQ (2011): "Estimating structural VARMA models with uncorrelated but non-independent error terms," *Journal of Multivariate Analysis*, 102, 496–505.
- BOX, G. E. P., AND D. A. PIERCE (1970): "Distribution of residual autocorrelations in autoregressive-integrated moving average time series models," *JASA*, 65, 1509–1526.
- BROCKWELL, P. J., AND R. A. DAVIS (1991): *Time series: theory and methods*, Springer Series in Statistics. Springer-Verlag, New York, second edn.
- CHEN, B., J. CHOI, AND J. C. ESCANCIANO (2012): "Testing for fundamental Vector Moving Average representations," University of Rochester.
- COOPER, D. M., AND E. F. WOOD (1982): "Identifying Multivariate Time Series Models," *Journal of Time Series Analysis*, 3, 153–164.
- DAVYDOV, Y. A. (1968): "Convergence of Distributions Generated by Stationary Stochastic Processes," *Theory of Probability and its Applications*, pp. 691–696.
- DEISTLER, M., AND E. J. HANNAN (1981): "Some Properties of the Parametrization of ARMA Systems with Unknown Order," *Journal of Multivariate Analysis*, 11, 474–484.

- DOORNIK, J. A. (1999): *Object-Oriented Matrix Programming Using Ox, 3rd ed.* Timberlake Consultants Press and Oxford, Oxford, U.K., www.nuff.ox.ac.uk/Users/Doornik.
- DOUKHAN, P. (1995): *Mixing - Properties and Examples*, no. 85 in *Lecture Notes in Statistics*. Springer-Verlag.
- DROST, F. C. (1993): "Temporal Aggregation of Time-Series," in *Econometric Analysis of Financial Markets*, ed. by J. Kaehler, and P. Kugler, pp. 11–21. Physica-Verlag, New York.
- DROST, F. C., AND T. E. NIJMAN (1993): "Temporal Aggregation of GARCH Processes," *Econometrica*, 61(4), 909–927.
- DURBIN, J. (1960): "The Fitting of Time Series Models," *Revue de l'Institut International de Statistique*, 28, 233.
- FLORES DE FRUTOS, R., AND G. R. SERRANO (2002): "A Generalized Least Squares Estimation Method for VARMA Models," *Statistics*, 36(4), 303–316.
- FRANCQ, C., AND H. RAÏSSI (2007): "Multivariate Portmanteau Test for Autoregressive Models with Uncorrelated but Nonindependent Errors," *Journal of Time Series Analysis*, 28(3), 454–470.
- FRANCQ, C., R. ROY, AND J.-M. ZAKOÏAN (2005): "Diagnostic checking in ARMA models with uncorrelated errors," *Journal of the American Statistical Association*, 100, 532–544.
- FRANCQ, C., AND J.-M. ZAKOÏAN (1998): "Estimating Linear Representations of Nonlinear Processes," *Journal of Statistical Planning and Inference*, 68, 145–165.
- (2000): "Covariance matrix estimation for estimators of mixing weak ARMA models," *Journal of Statistical Planning and Inference*, 83, 369–394.
- GALBRAITH, J. W., A. ULLAH, AND V. ZINDE-WALSH (2000): "Estimation of the Vector Moving Average Model by Vector Autoregression," *Econometric Reviews*, 21(2), 205–219.
- GALBRAITH, J. W., AND V. ZINDE-WALSH (1994): "A Simple Noniterative Estimator for Moving Average Models," *Biometrika*, 81(1), 143–155.
- (1997): "On Simple, Autoregression-Based Estimation and Identification Techniques for ARMA Models," *Biometrika*, 84(3), 685–696.
- HAMILTON, J. D. (1994): *Time Series Analysis*. Princeton University Press, Princeton, New Jersey.
- HANNAN, E. J. (1969): "The Identification of Vector Mixed Autoregressive- Moving Average Systems," *Biometrika*, 57, 223–225.
- (1971): "The Identification Problem for Multiple Equation System with Moving Average Errors," *Econometrica*, 39, 751–766.
- (1976): "The Identification and Parameterization of ARMAX and State Space Forms," *Econometrica*, 44(4), 713–723.
- HANNAN, E. J., AND M. DEISTLER (1988): *The Statistical Theory of Linear Systems*. John Wiley & Sons, New York.
- HANNAN, E. J., AND L. KAVALIERIS (1984a): "A Method for Autoregressive-Moving Average Estimation," *Biometrika*, 71(2), 273–280.
- (1984b): "Multivariate Linear Time Series Models," *Advances in Applied Probability*, 16, 492–561.

- (1986): “Regression, Autoregression Models,” *Journal of Time Series Analysis*, 7(1), 27–49.
- HANNAN, E. J., L. KAVALIERIS, AND M. MACKISACK (1986): “Recursive Estimation of Linear Systems,” *Biometrika*, 73(1), 119–133.
- HANNAN, E. J., AND J. RISSANEN (1982): “Recursive Estimation of Mixed Autoregressive-Moving-Average Order,” *Biometrika*, 69, 81–94, Errata 70 (1982), 303.
- (1983): “Errata: ”Recursive Estimation of Mixed Autoregressive-Moving Average Order”,” *Biometrika*, 70(1), 303.
- HORN, R. G., AND C. A. JOHNSON (1985): *Matrix Analysis*. Cambridge University Press, Cambridge, U.K.
- HUANG, D., AND L. GUO (1990): “Estimation of Nonstationary ARMAX Models Based on the Hannan-Rissanen Method,” *The Annals of Statistics*, 18(4), 1729–1756.
- IBRAGIMOV, I. A. (1962): “Some Limit Theorems for Stationary Processes,” *Theory of Probability and its Applications*, 7, 349–382.
- KATAYAMA, N. (2012): “Chi-squared portmanteau tests for structural VARMA models with uncorrelated errors,” *Journal of Time Series Analysis*, 33, 863–872.
- KEATING, J. W. (2002): “Structural Inference with Long-Run Recursive Empirical Models,” *Macroeconomic Dynamics*, 6(2), 266–283.
- KOREISHA, S. G., AND T. M. PUKKILA (1989): “Fast Linear Estimation Methods for Vector Autoregressive Moving-Average Models,” *Journal of Time Series Analysis*, 10(4), 325–339.
- (1990a): “A Generalized Least-Squares Approach for Estimation of Autoregressive-Moving-Average Models,” *Journal of Time Series Analysis*, 11(2), 139–151.
- (1990b): “Linear Methods for Estimating ARMA and Regression Models with Serial Correlation,” *Communications in Statistics, Part B -Simulation and Computation*, 19, 71–102.
- (1995): “A Comparison Between Different Order-Determination Criteria for Identification of ARIMA Models,” *Journal of Business and Economic Statistics*, 13, 127–131.
- LEWIS, R., AND G. C. REINSEL (1985): “Prediction of Multivariate Time Series by Autoregressive Model Fitting,” *Journal of Multivariate Analysis*, 16, 393–411.
- LJUNG, G. M., AND G. E. P. BOX (1978): “On a Measure of Lack of Fit in Time Series Models,” *Biometrika*, 65(2), 297–303.
- LUCEÑO, A. (1994): “A Fast Algorithm For The Exact Likelihood Of Stationary And Partially Nonstationary Vector Autoregressive-Moving Average Processes,” *Biometrika*, 81(3), 555–565.
- LÜTKEPOHL, H. (1991): *Introduction to Multiple Time Series Analysis*. Springer-Verlag, Berlin.
- (1993): *Introduction to Multiple Time Series Analysis*. Springer-Verlag, Berlin, second edn.
- LÜTKEPOHL, H., AND H. CLAESSEN (1997): “Analysis of Cointegrated VARMA Processes,” *Journal of Econometrics*, 80, 223–239.
- LÜTKEPOHL, H., AND D. S. POSKITT (1996a): “Consistent Estimation of the Number of Cointegration Relations in a Vector Autoregressive Model,” Discussion Paper 74, Institut für Statistik und Ökonometrie, Humboldt-Universität zu Berlin.
- LÜTKEPOHL, H., AND D. S. POSKITT (1996b): “Specification of Echelon-Form VARMA Models,” *Journal of Business and Economic Statistics*, 14, 69–80.

- MAURICIO, J. A. (2002): “An Algorithm for the Exact Likelihood of a Stationary Vector Autoregressive-Moving Average Model,” *Journal of Time Series Analysis*, 23(4), 473–486.
- MCKELVEY, T., A. HELMERSSON, AND T. RIBARITS (2004): “Data driven local coordinates for multivariable linear systems and their application to system identification,” *Automatica*, 40, 1629–1635.
- MCMILLIN, W. D. (2001): “The Effect of Monetary Policy Shocks: Comparing Contemporaneous versus Long-Run Identifying Restrictions,” *Southern Economic Journal*, 67(3), 618–636.
- NEWBY, W. K., AND K. D. WEST (1987): “A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix,” *Econometrica*, 55, 703–708.
- NIJMAN, T. H., AND F. PALM (1990): “Parameter Identification in ARMA Processes in the Presence of Regular But Incomplete Sampling,” *Journal of Time Series Analysis*, 11(3), 239–248.
- NSIRI, S., AND R. ROY (1992): “On the Identification of ARMA Echelon-Form Models,” *Canadian Journal of Statistics*, 20(4), 369–386.
- NSIRI, S., AND R. ROY (1996): “Identification of Refined ARMA Echelon Form Models for Multivariate Time Series,” *Journal of Multivariate Analysis*, 56, 207–231.
- PALM, F., AND T. H. NIJMAN (1984): “Missing Observations in the Dynamic Regression Model,” *Econometrica*, 52(6), 1415–1436.
- POSKITT, D. S. (1987): “A Modified Hannan-Rissanen Strategy for Mixed Autoregressive-Moving Average Oder Determination,” *Biometrika*, 74(4), 781–790.
- (1992): “Identification of Echelon Canonical Forms for Vector Linear Processes Using Least Squares,” *The Annals of Statistics*, 20, 195–215.
- POSKITT, D. S., AND H. LÜTKEPOHL (1995): “Consistent Specification of Cointegrated Autoregressive Moving-Average Systems,” Discussion Paper 54, Institut für Statistik und Ökonometrie, Humboldt-Universität zu Berlin.
- PUKKILA, T., S. KOREISHA, AND A. KALLINEN (1990): “The Identification of ARMA Models,” *Biometrika*, 77(3), 537–548.
- REINSEL, G. C., S. BASU, AND S. F. YAP (1992): “Maximum likelihood estimators in the multivariate autoregressive moving-average model from a generalized least squares viewpoint,” *Journal of Time Series Analysis*, 13(2), 133–145.
- RIBARITS, T., M. DEISTLER, AND T. MCKELVEY (2004): “An analysis of the parametrization by data driven local coordinates for multivariable linear systems,” *Automatica*, 40, 789–803.
- SAIKKONEN, P. (1986): “Asymptotic Properties of some Preliminary Estimators for Autoregressive Moving Average Time Series Models,” *Journal of Time Series Analysis*, 7, 133–155.
- SHAO, X. (2011): “Testing for white noise under unknown dependence and its applications to diagnostic checking for time series models,” *Econometric Theory*, 27, 312–343.
- SHEA, B. L. (1989): “The Exact Likelihood of a Vector Autoregressive Moving Average Model,” *Applied Statistics*, 38(1), 161–184.
- SPLIID, H. (1983): “A Fast Estimation Method for the Vector Autoregressive Moving Average Model with Exogenous Variables,” *Journal of the American Statistical Association*, 78(384), 843–849.
- TIAO, G. C., AND R. S. TSAY (1983): “Consistency Properties of Least Squares Estimates of Autoregressive Parameters in ARMA Models,” *The Annals of Statistics*, 11, 856–871.

- (1989): “Model Specification in Multivariate Time Series,” *Journal of the Royal Statistical Society, Series B*, 51(2), 157–213.
- TSAY, R. S. (1989a): “Identifying Multivariate Time Series Models,” *Journal of Time Series Analysis*, 10, 357–372.
- TSAY, R. S. (1989b): “Parsimonious Parameterization of Vector Autoregressive Moving Average Models,” *Journal of Business and Economic Statistics*, 7(3), 327–341.
- TSAY, R. S. (1991): “Two Canonical Forms for Vector VARMA Processes,” *Statistica Sinica*, 1, 247–269.
- WALLIS, K. F. (1977): “Multiple time series analysis and the final form of econometric models,” *Econometrica*, 45(6), 1481–1497.
- ZELLNER, A., AND F. PALM (1974): “Time series analysis and simultaneous equation econometric model,” *Journal of Econometrics*, 2(1), 17–54.
- ZHAO-GUO, C. (1985): “The Asymptotic Efficiency of a Linear Procedure of Estimation for ARMA Models,” *Journal of Time Series Analysis*, 6(1), 53–62.