

Modelling multiple time series via common factors

BY JIAZHU PAN

*Department of Statistics and Modelling Science, University of Strathclyde, Livingstone Tower,
Richmond Street, Glasgow G1 1XH, U.K.*

jiazhu@stams.strath.ac.uk

AND QIWEI YAO

*Department of Statistics, The London School of Economics and Political Science, Houghton
Street, London, WC2A 2AE, U.K.*

q.yao@lse.ac.uk

SUMMARY

We propose a new method for estimating common factors of multiple time series. One distinctive feature of the new approach is that it is applicable to some nonstationary time series. The unobservable, nonstationary factors are identified by expanding the white noise space step by step, thereby solving a high-dimensional optimization problem by several low-dimensional subproblems. Asymptotic properties of the estimation are investigated. The proposed methodology is illustrated with both simulated and real datasets.

Some key words: Cross-correlation function; Dimension reduction; Factor model; Multivariate time series; Non-stationarity; Portmanteau test; White noise.

1. INTRODUCTION

An important problem in modelling multivariate time series is to reduce the number of parameters involved. For example, a vector autoregressive and moving average model (VARMA) with moderately large order (p, q) is viable in practice only if a parsimonious representation is identified, by imposing constraints on the coefficient matrices; see Tiao & Tsay (1989), Reinsel (1997) and references therein. An alternative strategy is to reduce the dimensionality. Attempts along this line include, among others, approaches based on principal component analysis (Priestley et al., 1974; Brillinger, 1981; Stock & Watson, 2002), canonical correlation analysis (Box & Tiao, 1977; Geweke, 1977; Geweke & Singleton, 1981; Tiao & Tsay, 1989; Anderson, 2002), reduced rank regression (Ahn, 1997; Reinsel & Velu, 1998) and factor models (Engle & Watson, 1981; Peña & Box, 1987; Forni et al., 2000; Bai & Ng, 2002).

In this paper, we revisit factor models. Although the form of the model concerned is the same as that in, for example, Peña & Box (1987), our approach is novel in three respects. First, we allow factors to be nonstationary and the nonstationarity is not necessarily driven by unit roots. The latter was investigated in the context of factor models by, for example, Ahn (1997) and Peña & Poncela (2006). Secondly, our estimation method is new and it identifies the unobserved factors by expanding the white noise space step by step, thereby solving a high-dimensional optimization problem by several low-dimensional subproblems. Finally, we allow dependence between the factors and the white noise in the model.

We do not impose distributional assumptions in the model. Instead we use the portmanteau test to identify the white noise space. The key assumption in the theoretical exploration is that the sample cross-covariance functions converge in probability to constant limits; see Condition 1 in § 3. This may be implied by the ergodicity of stationary processes, and may also be fulfilled for some nonstationary mixing processes, purely deterministic trends and random walks; see Remark 2 in § 3.

2. MODELS AND METHODOLOGY

2.1. Factor models

Let $\{Y_t\}$ be a $d \times 1$ time series satisfying

$$Y_t = AX_t + \varepsilon_t, \quad (1)$$

where X_t is an $r \times 1$ time series with finite second moments, $r \leq d$ is unknown, A is a $d \times r$ unknown constant matrix, and $\{\varepsilon_t\}$ is a sequence of vector white noise processes with mean μ_ε and covariance matrix Σ_ε ; that is, ε_t and ε_s are uncorrelated for any $t \neq s$. Furthermore, we assume that there exists no linear combination of X_t , which is a white noise process; otherwise such a linear combination should be part of ε_t . We only observe Y_1, \dots, Y_n from model (1). To simplify the presentation, we assume that

$$S_0 \equiv n^{-1} \sum_{t=1}^n (Y_t - \bar{Y})(Y_t - \bar{Y})^T = I_d, \quad (2)$$

where $\bar{Y} = n^{-1} \sum_{1 \leq t \leq n} Y_t$, and I_d denotes the $d \times d$ identity matrix. In practice this amounts to replacing Y_t by $S_0^{-1/2} Y_t$ before the analysis.

The component variables of the unobserved X_t are called the factors, and A is called the factor loading matrix. We may assume that the rank of A is r ; otherwise (1) may be expressed equivalently in terms of a smaller number of factors. Model (1) is unchanged if we replace (A, X_t) by $(AH, H^{-1}X_t)$ for any invertible $r \times r$ matrix H , so we may assume that the column vectors of $A = (a_1, \dots, a_r)$ are orthonormal:

$$A^T A = I_r. \quad (3)$$

Even with the constraint (3), A and X_t are not uniquely determined in (1), as the aforementioned replacement is still applicable for any orthogonal H . However, the linear space spanned by the columns of A , denoted by $\mathcal{M}(A)$ and called the factor loading space, is a uniquely defined r -dimensional subspace in \mathcal{R}^d .

Model (1) has been studied by Peña & Box (1987) who assume that ε_t and X_{t+k} are uncorrelated for any integers t and k , and Y_t is stationary. Under those conditions, the number of factors r is the maximum rank of the autocovariance matrices of Y_t over all nonzero lags. Furthermore, both A and r may be estimated via standard eigenanalysis; see Peña & Box (1987).

2.2. Estimation of A and r

Our goal is to estimate $\mathcal{M}(A)$, or its orthogonal complement $\mathcal{M}(B)$, where $B = (b_1, \dots, b_{d-r})$ is a $d \times (d-r)$ matrix for which (A, B) forms a $d \times d$ orthogonal matrix, i.e. $B^T A = 0$ and $B^T B = I_{d-r}$; see also (3). It follows from (1) that

$$B^T Y_t = B^T \varepsilon_t, \quad (4)$$

and hence $\{B^T Y_t, t = 0, \pm 1, \dots\}$ is a $(d - r) \times 1$ white noise process. Therefore,

$$\text{corr}(b_i^T Y_t, b_j^T Y_{t-k}) = 0, \quad (5)$$

for any $i, j = 1, \dots, d - r$ and $k = 1, \dots, p$, where $p \geq 1$ is an arbitrary integer. Under assumption (2), $b_i^T S_k b_j$ is the sample correlation coefficient between $b_i^T Y_t$ and $b_j^T Y_{t-k}$, where

$$S_k = n^{-1} \sum_{t=k+1}^n (Y_t - \bar{Y})(Y_{t-k} - \bar{Y})^T. \quad (6)$$

This suggests that we may estimate B by minimizing

$$\Psi_n(B) \equiv \sum_{k=1}^p \|B^T S_k B\|^2 = \sum_{k=1}^p \sum_{1 \leq i, j \leq d-r} \rho_k(b_i, b_j)^2, \quad (7)$$

where the matrix norm $\|H\|$ is defined as $\{\text{tr}(H^T H)\}^{1/2}$, and $\rho_k(b, a) = b^T S_k a$.

Minimizing (7) leads to a constrained optimization problem with $d \times (d - r)$ variables. Furthermore, r is unknown. Below we present a stepwise expansion algorithm for estimating the columns of B as well as the number of columns r . Put

$$\psi(b) = \sum_{k=1}^p \rho_k(b, b)^2, \quad \psi_m(b) = \sum_{k=1}^p \sum_{i=1}^{m-1} \{\rho_k(b, \hat{b}_i)^2 + \rho_k(\hat{b}_i, b)^2\},$$

and let $\alpha \in (0, 1)$ be the level of significance tests.

Step 1. Let \hat{b}_1 be a unit vector that minimizes $\psi(b)$. Compute the Ljung–Box–Pierce portman-teau test statistic:

$$L_{p,1} = n(n+2) \sum_{k=1}^p \frac{\rho_k(\hat{b}_1, \hat{b}_1)^2}{n-k}. \quad (8)$$

Terminate the algorithm with $\hat{r} = d$ and $\hat{B} = 0$ if $L_{p,1}$ is greater than the upper α -point of the χ_p^2 distribution. Otherwise proceed to Step 2.

Step 2. For $m = 2, \dots, d$, let \hat{b}_m minimize $\psi(b) + \psi_m(b)$ subject to the constraints

$$\|b\| = 1, \quad b^T \hat{b}_i = 0 \quad \text{for } i = 1, \dots, m-1. \quad (9)$$

Terminate the algorithm with $\hat{r} = d - m + 1$ and $\hat{B} = (\hat{b}_1, \dots, \hat{b}_{m-1})$ if

$$L_{p,m} \equiv n^2 \sum_{k=1}^p \frac{1}{n-k} \left[\rho_k(\hat{b}_m, \hat{b}_m)^2 + \sum_{j=1}^{m-1} \{\rho_k(\hat{b}_m, \hat{b}_j)^2 + \rho_k(\hat{b}_j, \hat{b}_m)^2\} \right] \quad (10)$$

is greater than the upper α -point of the $\chi_{p(2m-1)}^2$ distribution (Reinsel, 1997, pp. 149–50).

Step 3. If $L_{p,m}$ never exceeds the critical value for all $1 \leq m \leq d$, let $\hat{r} = 0$ and $\hat{B} = I_d$.

Remark 1. (i) The algorithm increases the dimension of $\mathcal{M}(B)$ by 1 each time until a newly selected direction \hat{b}_m does not lead to a white noise process. Note that (9) ensures that all those \hat{b}_j are orthogonal to each other.

(ii) The minimization problem in Step 2 is d -dimensional subject to constraint (9). It may be reduced to an unconstrained optimization problem with $d - m$ free variables. Note that the vector

b satisfying (9) is of the form

$$b = D_m u, \quad (11)$$

where u is any $(d - m + 1) \times 1$ unit vector, D_m is a $d \times (d - m + 1)$ matrix with columns given by the $(d - m + 1)$ orthonormal eigenvectors of the matrix $I_d - B_{m-1} B_{m-1}^T$ corresponding to the $(d - m + 1)$ -fold eigenvalue 1, where $B_m = (\hat{b}_1, \dots, \hat{b}_m)$. Also note that any $k \times 1$ unit vector is of the form $u^T = (u_1, \dots, u_k)$, where

$$u_1 = \prod_{j=1}^{k-1} \cos \theta_j, \quad u_i = \sin \theta_{i-1} \prod_{j=i}^{k-1} \cos \theta_j \quad (i = 2, \dots, k-1), \quad u_k = \sin \theta_{k-1}.$$

In the above expressions, $\theta_1, \dots, \theta_{k-1}$ are $(k - 1)$ free parameters.

(iii) Note that $\hat{B}^T \hat{B} = I_{d-\hat{r}}$. We may let the columns of \hat{A} be the \hat{r} orthonormal eigenvectors of $I_d - \hat{B} \hat{B}^T$ corresponding to the common eigenvalue 1. It holds that $\hat{A}^T \hat{A} = I_{\hat{r}}$.

(iv) The multivariate portmanteau test statistic $L_{p,m}$ given in (10) has a normalized constant n^2 , which is different from $n(n + 2)$ used in the univariate case (8). For the univariate case, the modified constant $n(n + 2)$ was suggested to improve the finite-sample accuracy; see Ljung & Box (1978). For multivariate cases, a different suggestion, proposed by Li & McLeod (1981), uses

$$L_{p,m}^* = L_{p,m} + \frac{p(p + 1)(2m - 1)}{2n} \quad (12)$$

instead of $L_{p,m}$ as the test statistic. Our numerical experiment indicates that both $L_{p,m}$ and $L_{p,m}^*$ work reasonably well with moderately large sample sizes, unless $d \gg r$. For the latter cases, both $L_{p,m}$ and $L_{p,m}^*$ may lead to substantially overestimated r . In our context, an obvious alternative is to use a more stable univariate version,

$$L'_{p,m} = n(n + 2) \sum_{k=1}^p \frac{\rho_k(\hat{b}_m, \hat{b}_m)^2}{n - k}, \quad (13)$$

instead of $L_{p,m}$ in Step 2. Then the critical value of the test is the upper α -point of the χ_p^2 distribution.

(v) Although we do not require the processes $\{Y_t\}$ and $\{X_t\}$ to be stationary, our method rests on the fact that there is no autocorrelation in the white noise process $\{\varepsilon_t\}$. Furthermore, the asymptotic χ^2 distributions of the portmanteau tests used in determining r typically rely on the assumption that $\{\varepsilon_t\}$ be independent and identically distributed. An early investigation of these tests in more general settings is given by Francq et al. (2005).

(vi) When Y_t is nonstationary, the sample cross-covariance function S_k is no longer a meaningful covariance measure. However, since ε_t is white noise and is stationary, $c_1^T S_k c_2$ is the proper sample covariance between $c_1^T Y_t$ and $c_2^T Y_{t-k}$ for any vectors $c_1, c_2 \in \mathcal{M}(B)$. In fact our method relies on the fact that $c_1^T S_k c_2$ is close to 0 for any $1 \leq k \leq p$. This also indicates that in practice we should not use large p as, for example, $c_1^T S_p c_2$ is a poor estimator of $\text{cov}(c_1^T Y_t, c_2^T Y_{t-p})$ when p is too large.

(vii) When the number of factors r is given, we may omit all the test steps, and stop the algorithm after obtaining $\hat{b}_1, \dots, \hat{b}_r$ by solving the r optimization problems.

2.3. Modelling with estimated factors

Note that $\hat{A} \hat{A}^T + \hat{B} \hat{B}^T = I_d$. Once we have obtained \hat{A} , it follows from (1) that

$$Y_t = \hat{A} \xi_t + e_t, \quad (14)$$

where

$$\xi_t = \hat{A}^T Y_t = \hat{A}^T A X_t + \hat{A}^T \varepsilon_t, \quad e_t = \hat{B} \hat{B}^T Y_t. \quad (15)$$

We treat e_t as a white noise process, and estimate $\text{var}(e_t)$ by the sample variance of $\hat{B} \hat{B}^T Y_t$.

We model the lower-dimensional process ξ_t by VARMA or state-space models. As we have pointed out, \hat{A} may be replaced by $\hat{A}H$ for any orthogonal H . We may choose \hat{A} appropriately such that ξ_t admits a simple model; see, for example, Tiao & Tsay (1989). Alternatively, we may apply principal components analysis to the factors; see Example 3 in § 4. Note that there is no need to update \hat{B} now since $\mathcal{M}(\hat{A}H) = \mathcal{M}(\hat{A})$, which is the orthogonal complement of \hat{B} .

3. THEORETICAL PROPERTIES

The factor loading matrix A is only identifiable up to $\mathcal{M}(A)$, a linear space spanned by its columns. We are effectively concerned with the estimation for the factor loading space $\mathcal{M}(A)$ rather than A itself. To make our statements clearer, we first introduce some notation.

For $r < d$, let \mathcal{H} be the set consisting of all $d \times (d - r)$ matrices H satisfying the condition $H^T H = I_{d-r}$. For $H_1, H_2 \in \mathcal{H}$, define

$$D(H_1, H_2) = \|(I_d - H_1 H_1^T) H_2\| = \{d - r - \text{tr}(H_1 H_1^T H_2 H_2^T)\}^{1/2}. \quad (16)$$

Note that $H_1 H_1^T$ is the projection matrix on to the linear space $\mathcal{M}(H_1)$, and $D(H_1, H_2) = 0$ if and only if $\mathcal{M}(H_1) = \mathcal{M}(H_2)$. Therefore, \mathcal{H} may be partitioned into the equivalent classes by D as follows: the D -distance between any two elements in each equivalent class is 0, and the D -distance between any two elements from two different classes is positive. Denote by $\mathcal{H}_D = \mathcal{H}/D$ the quotient space consisting of all those equivalent classes; that is, we treat H_1 and H_2 as the same element in \mathcal{H}_D if and only if $D(H_1, H_2) = 0$. Then (\mathcal{H}_D, D) forms a metric space in the sense that D is a well-defined distance measure on \mathcal{H}_D ; see Lemma A1(i) in the Appendix. Furthermore, the functions $\Psi_n(\cdot)$, defined in (7), and

$$\Psi(H) \equiv \sum_{k=1}^p \|H^T \Sigma_k H\|^2, \quad (17)$$

are well defined on \mathcal{H}_D ; see Lemma A1(ii) in the Appendix. In the above expression, the Σ_k are given in Condition 1 below.

We only consider the asymptotic properties for the estimation of the factor loading space with r known; it remains open how to establish the theoretical properties when r is unknown. Then the estimator for B may be defined as

$$\hat{B} = \arg \min_{H \in \mathcal{H}} \Psi_n(H). \quad (18)$$

We need the following regularity conditions.

Condition 1. As $n \rightarrow \infty$, $S_k \rightarrow \Sigma_k$ in probability for $k = 0, 1, \dots, p$, where Σ_k are nonnegative definite matrices, and $\Sigma_0 = I_d$.

Condition 2. The matrix B is the unique minimizer of $\Psi(\cdot)$ in the space \mathcal{H}_D ; that is, $\Psi(\cdot)$ reaches its minimum value at B' if and only if $D(B', B) = 0$, where B is specified at the beginning of § 2.2.

Condition 3. There exist constants $a > 0, c > 0$ for which $\Psi(H) - \Psi(B) \geq a\{D(H, B)\}^c$ for any $H \in \mathcal{H}$.

Remark 2. (i) Condition 1 does not require that the process Y_t be stationary. In fact it may hold when $ES_k \rightarrow \Sigma_k$ and Y_t is φ -mixing in the sense that $\varphi(m) \rightarrow 0$ as $m \rightarrow \infty$, where

$$\varphi(m) = \sup_{k \geq 1} \sup_{U \in \mathcal{F}_{-\infty}^k, V \in \mathcal{F}_{m+k}^\infty, \text{pr}(U) > 0} |\text{pr}(V|U) - \text{pr}(V)|, \quad (19)$$

and $\mathcal{F}_i^j = \sigma(Y_i, \dots, Y_j)$; see Lemma A2 in the Appendix. It also gives a sufficient condition, which ensures that the convergence in Condition 1 is almost sure. Examples of nonstationary φ -mixing processes include, among others, stationary φ -mixing processes plus nonconstant trends, and standardized random walks such as $Y_t = Y_{t-1} + n^{-1/2}\varepsilon_t$, $t = 1, \dots, n$, where $Y_0 \equiv 0$ and ε_t are independent and identically distributed with, for example, $E(\varepsilon_t^2) < \infty$. Condition 1 may also hold for some purely deterministic processes such as a linear trend $Y_t = t/n$, $t = 1, \dots, n$.

(ii) Under model (1), $\Psi(B) = 0$. Condition 2 implies that $\Psi(C) \neq 0$ for any $C \in \mathcal{H}$ and $\mathcal{M}(C) \cap \mathcal{M}(A)$ is not an empty set.

THEOREM 1. *Under Conditions 1 and 2, $D(\hat{B}, B) \rightarrow 0$ in probability as $n \rightarrow \infty$. Furthermore, $D(\hat{B}, B) \rightarrow 0$ almost surely if the convergence in Condition 1 is also almost sure.*

THEOREM 2. *Let $n^{1/2}(ES_k - \Sigma_k) = O(1)$, and let Y_t be φ -mixing with $\varphi(m) = O(m^{-\lambda})$ for $\lambda > p/(p-2)$ and $\sup_{t \geq 1} E\|Y_t\|^p < \infty$ for some $p > 2$. Then*

$$\sup_{H \in \mathcal{H}} |\Psi_n(H) - \Psi(H)| = O_P(n^{-1/2}).$$

If, in addition, Condition 3 also holds, $D(\hat{B}, B) = O_P(n^{-1/(2c)})$.

Theorems 1 and 2 do not require Y_t to be a stationary process. Their proofs are given in the Appendix.

4. NUMERICAL PROPERTIES

We illustrate the methodology proposed in §2 with two simulated examples, one stationary and one nonstationary, and one real dataset. The numerical optimization was solved using the downhill simplex method; see §10.4 of Press et al. (1992). In the simulated examples, we set the significance level at 5% for the portmanteau tests used in our algorithm, and $p = 15$ in (8). The results with $p = 5, 10$ and 20 show similar patterns and, therefore, are not reported. We measure the errors in estimating the factor loading space $\mathcal{M}(A)$ by

$$D_1(A, \hat{A}) = ([\text{tr}\{\hat{A}^T(I_d - AA^T)\hat{A}\} + \text{tr}(\hat{B}^T AA^T \hat{B})]/d)^{1/2}.$$

It may be shown that $D_1(A, \hat{A}) \in [0, 1]$, and it equals 0 if and only if $\mathcal{M}(A) = \mathcal{M}(\hat{A})$, and 1 if and only if $\mathcal{M}(A) = \mathcal{M}(\hat{B})$.

Example 1. Let $Y_{ti} = X_{ti} + \varepsilon_{ti}$ for $i = 1, 2, 3$, and $Y_{ti} = \varepsilon_{ti}$ for $i = 4, \dots, d$, where

$$\begin{aligned} X_{t1} &= 0.8X_{t-1,1} + e_{t1}, & X_{t2} &= e_{t2} + 0.9e_{t-1,2} + 0.3e_{t-2,2}, \\ X_{t3} &= -0.5X_{t-1,3} - \varepsilon_{t3} + 0.8\varepsilon_{t-1,3}, \end{aligned}$$

and all ε_{ij} and e_{ij} are independent and standard normal. Because of the presence of ε_{t3} in the equation for X_{t3} , X_t and ε_t are mutually dependent. In this setting, the number of true factors is $r = 3$, and the factor loading matrix may be taken as $A = (I_3, 0)^T$, where 0 denotes the $3 \times (d-3)$ matrix with all elements equal to 0. We set the sample size at $n = 300, 600$ and 1000 , and the dimension of Y_t at $d = 5, 10$ and 20 . For each setting, we generated 1000 samples from this

Table 1. Relative frequencies for \hat{r} taking different values in Example 1.
The true value of r is 3

d	n	\hat{r}							
		0	1	2	3	4	5	≥ 6	
5	300	0.000	0.209	0.444	0.345	0.002	0.000		
	600	0.000	0.071	0.286	0.633	0.010	0.000		
	1000	0.000	0.004	0.051	0.933	0.120	0.000		
10	300	0.000	0.219	0.524	0.255	0.002	0.000	0.000	
	600	0.000	0.049	0.290	0.649	0.012	0.000	0.000	
	1000	0.000	0.007	0.062	0.898	0.033	0.000	0.000	
20	300	0.000	0.162	0.543	0.285	0.010	0.000	0.000	
	600	0.000	0.033	0.305	0.609	0.053	0.000	0.000	
	1000	0.000	0.004	0.066	0.822	0.103	0.005	0.000	

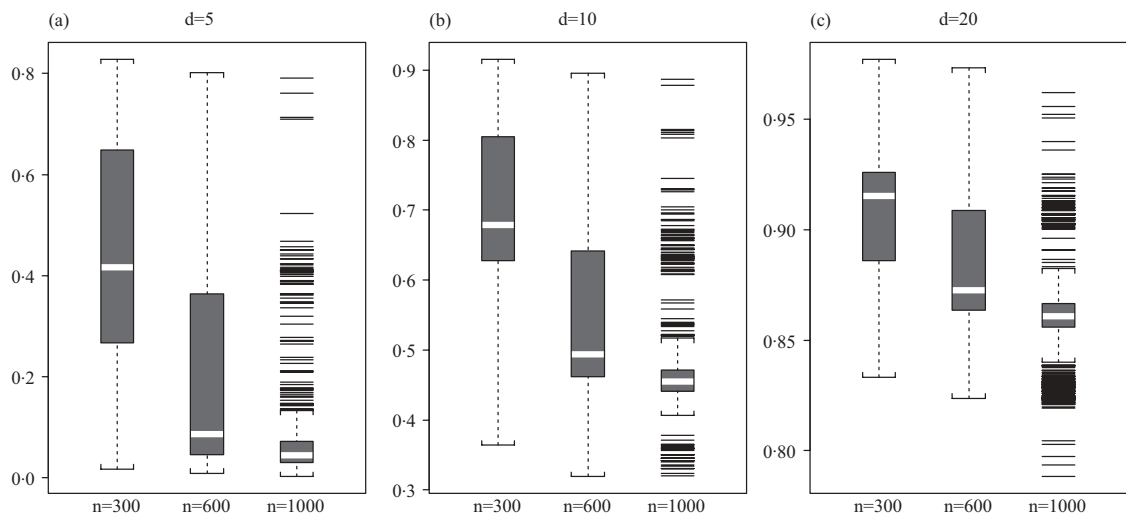


Fig. 1. Example 1. Boxplots of $D_1(A, \hat{A})$ for (a) $d = 5$, (b) $d = 10$, (c) $d = 20$.

model. The relative frequencies for \hat{r} taking different values are reported in Table 1, which shows that, when the sample size n increases, the estimation of r becomes more accurate. We used $L'_{m,p}$ given in (13) in our simulation, since both $L_{m,p}$ and $L^*_{m,p}$ produced substantially overestimated r -values when $d = 10$ and 20. Figure 1 presents the boxplots of the errors $D_1(A, \hat{A})$. As the sample size increases, $D_1(A, \hat{A})$ decreases. Furthermore, the errors increase when d increases.

Example 2. We use the same setting as in Example 1 but with X_{t1} , X_{t2} and X_{t3} replaced by

$$\begin{aligned} X_{t1} - 2t/n &= 0.8(X_{t-1,1} - 2t/n) + e_{t1}, \\ X_{t2} &= 3t/n, \\ X_{t3} &= X_{t-1,3} + (10/n)^{1/2}e_{t3} \quad \text{with} \quad X_{0,3} \sim N(0, 1). \end{aligned} \quad (20)$$

Thus X_{t1} is an AR(1) process with nonconstant mean, X_{t2} is a purely deterministic trend, and X_{t3} is a random walk. None of them is stationary. The relative frequencies for \hat{r} taking different values are reported in Table 2. The boxplots of the estimation errors $D_1(A, \hat{A})$ are depicted in Fig. 2. The general pattern observed in the stationary example, Example 1, remains. The quality of our

Table 2. Relative frequencies for \hat{r} taking different values in Example 2. The true value of r is 3

d	n	\hat{r}						
		0	1	2	3	4	5	≥ 6
5	300	0.000	0.000	0.255	0.743	0.002	0.000	
	600	0.000	0.000	0.083	0.907	0.010	0.000	
	1000	0.000	0.000	0.033	0.945	0.022	0.000	
10	300	0.000	0.000	0.283	0.695	0.022	0.000	0.000
	600	0.000	0.000	0.103	0.842	0.054	0.001	0.000
	1000	0.000	0.000	0.051	0.871	0.077	0.001	0.000
20	300	0.000	0.000	0.258	0.663	0.076	0.001	0.002
	600	0.000	0.000	0.035	0.673	0.278	0.012	0.002
	1000	0.000	0.000	0.099	0.733	0.162	0.006	0.000

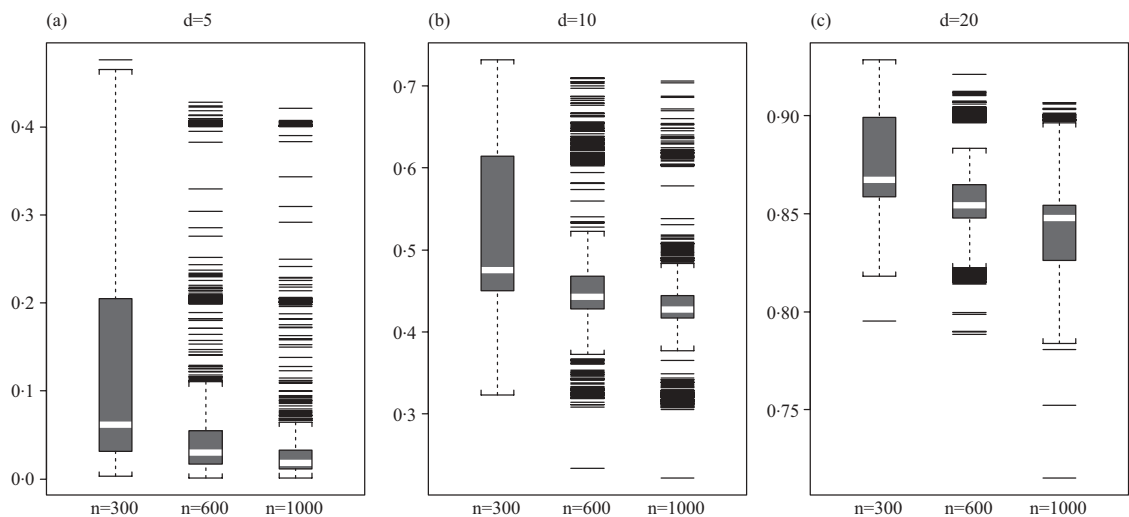


Fig. 2. Example 2. Boxplots of $D_1(A, \hat{A})$ for (a) $d = 5$, (b) $d = 10$, (c) $d = 20$.

estimation improves when sample sizes increase, because of the way in which the nonstationarity is specified in (20). For example, the sample $\{X_{12}, t = 1, \dots, n\}$ always consists of regular grid-points on the segment of the line $y = 3x$ between $(0, 0)$ and $(1, 3)$. Therefore, when n increases, we obtain more information from the same nonstationary system.

Our method rests on the simple fact that the quadratic forms of the sample cross-correlation function are close to zero along the directions perpendicular to the factor loading space, and are nonzero along the directions in the factor loading space; see Remarks 1(vi) and 2(ii). The departure from zero along the directions in the factor loading space in Example 2 is more pronounced than that in Example 1. This explains why the proposed method performs better in Example 2 than in Example 1, especially when $n = 300$ and 600 .

Example 3. Figure 3 displays the monthly temperatures in cities of Nanjing, Dongtai, Huoshan, Hefei, Shanghai, Anqing and Hangzhou in Eastern China in January 1954–December 1986. The sample size is $n = 396$ and $d = 7$. As expected, the data show strong periodic behaviour with period 12. We fitted the data with factor models (1). With $p = 12$, the estimated number of factors is $\hat{r} = 4$. We applied principal components analysis to the estimated factors. The variances of the

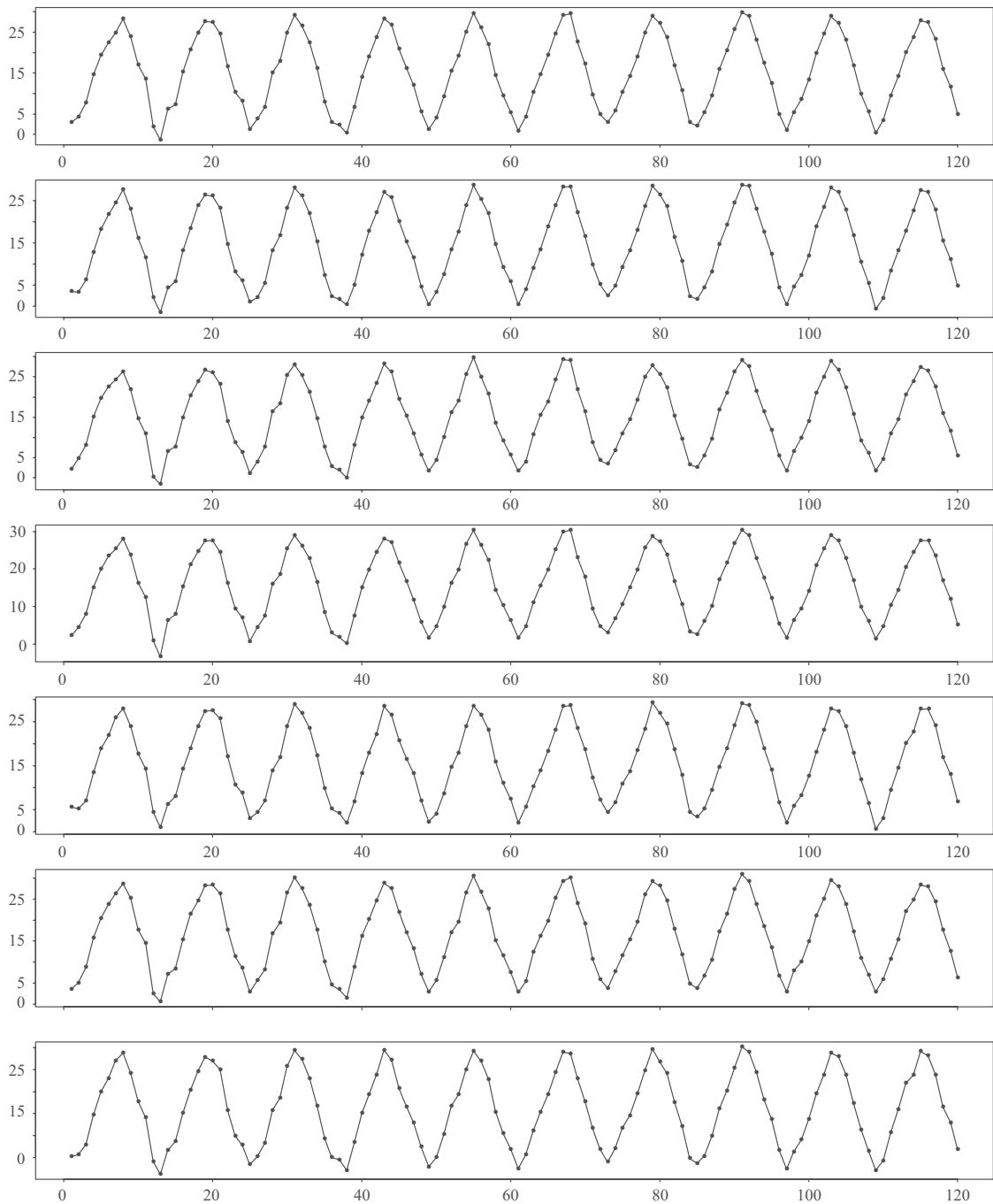


Fig. 3. Example 3. Time series plots of the monthly temperature in, from top to bottom, Nanjing, Dongtai, Huoshan, Hefei, Shanghai, Anqing and Hangzhou, corresponding to the first 10-year segments.

four principal component factors are, respectively, 542.08, 1.29, 0.07 and 0.06. The first factor accounts for over 99% of the total variation of the four factors, and 97.6% of the total variation of the original seven series. The first four principal component factors are plotted in Fig. 4, and their cross-correlation functions are displayed in Fig. 5. The periodic annual oscillation in the original

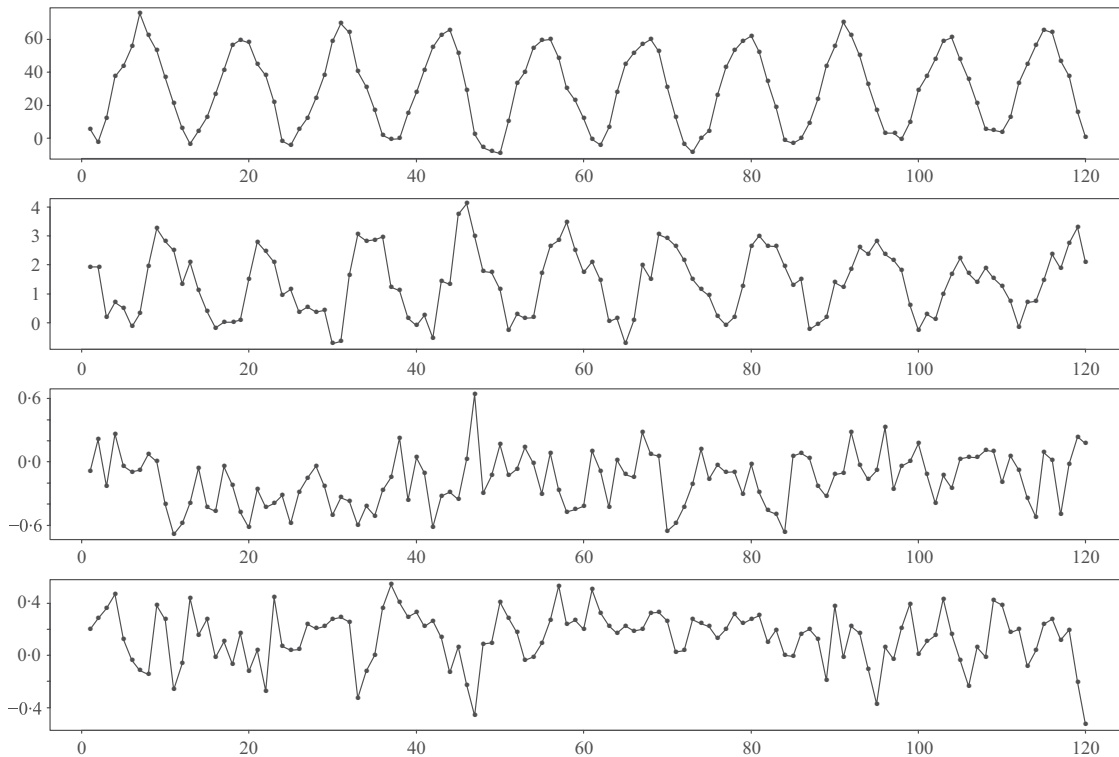


Fig. 4. Example 3. Time series plots of the first four principal component factors, corresponding to the first 10-year segments.

data is predominately reflected by the fluctuation of the first factor, and to a much less extent by that of the second factor; see Fig. 4. Furthermore, no periodic pattern is present in the third and the fourth principal component factors. This suggests that the annual temperature oscillation over this area may be seen as driven by one, or at most two, ‘common factors’. The first two columns of the corresponding loading matrix \hat{A} are

$$\begin{pmatrix} 0.394 & 0.386 & 0.378 & 0.387 & 0.363 & 0.376 & 0.366 \\ -0.086 & 0.225 & -0.640 & -0.271 & 0.658 & -0.014 & 0.164 \end{pmatrix}^T,$$

which indicates that the first principal component factor is effectively the average temperature over the seven cities. The residuals $\hat{B}^T Y_t$ carry little dynamic information in the data; this is indicated by the cross-correlation functions, not shown. The sample mean and sample covariance matrix of e_t are, respectively,

$$\hat{\mu}_e = \begin{pmatrix} 3.41 \\ 2.32 \\ 4.39 \\ 4.30 \\ 3.40 \\ 4.91 \\ 4.77 \end{pmatrix}, \quad \hat{\Sigma}_e = \begin{pmatrix} 1.56 & & & & & & \\ 1.26 & 1.05 & & & & & \\ 1.71 & 1.34 & 1.91 & & & & \\ 1.90 & 1.49 & 2.10 & 2.33 & & & \\ 1.37 & 1.16 & 1.46 & 1.58 & 1.37 & & \\ 1.67 & 1.26 & 1.91 & 2.09 & 1.37 & 1.97 & \\ 1.41 & 1.14 & 1.58 & 1.67 & 1.39 & 1.56 & 1.53 \end{pmatrix}. \quad (21)$$

Figure 5 indicates that the first two factors are dominated by periodic components with period 12. We estimated those components simply by taking the averages of all values in each of the

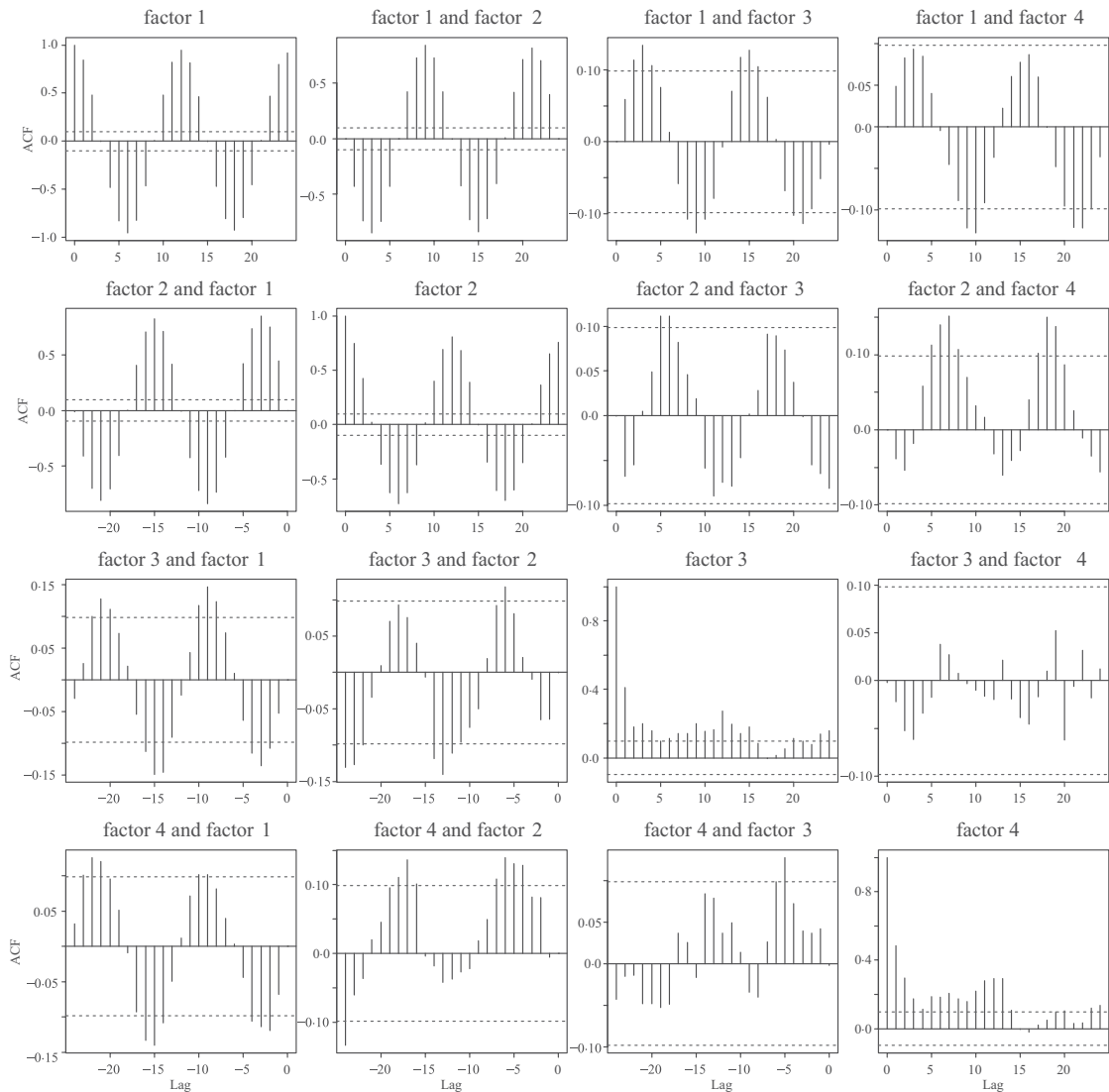


Fig. 5. Example 3. Sample cross-correlation functions of the first four principal component factors.

twelve months, leading to the estimated periodic components

$$\begin{aligned}
 (g_{1,1}, \dots, g_{12,1}) &= (-1.61, 1.33, 11.74, 28.06, 41.88, 54.51, 63.77, 62.14, 49.48, \\
 &\quad 33.74, 18.29, 3.50), \\
 (g_{1,2}, \dots, g_{12,2}) &= (1.67, 1.21, 0.47, 0.17, 0.41, 0.48, 1.37, 2.13, 2.98, 3.05, 2.78, 2.22) \quad (22)
 \end{aligned}$$

for, respectively, the first and the second factors. The cross-correlation functions of the four factors after removing the periodic components from the first two factors, not shown, indicate that the autocorrelation in each of those four series is not very strong, and that cross-correlations among those four series, at nonzero lag, are weak. We fitted a vector autoregressive model to those four series with the order 1 determined by the Akaike information criteria (AIC) (Brockwell

& Davis, 1991, p. 412) with the following estimated coefficients:

$$\hat{\varphi}_0 = \begin{pmatrix} 0.07 \\ -0.02 \\ -0.11 \\ 0.10 \end{pmatrix}, \quad \hat{\Phi}_1 = \begin{pmatrix} 0.27 & -0.31 & 0.72 & 0.40 \\ 0.01 & 0.36 & -0.04 & 0.04 \\ 0.00 & -0.01 & 0.42 & -0.02 \\ -0.00 & 0.03 & 0.03 & 0.48 \end{pmatrix}, \quad (23)$$

$$\hat{\Sigma}_u = \begin{pmatrix} 14.24 & & & \\ -0.17 & 0.23 & & \\ -0.02 & 0.03 & 0.05 & \\ 0.04 & 0.01 & -0.00 & 0.05 \end{pmatrix}. \quad (24)$$

The multivariate portmanteau tests, with lag $p = 12$, of both Li & Mcleod (1981) and Reinsel (1997, p. 149) for the residuals from the above fitted vector AR(1) model are nonsignificant at the 5% level. The univariate portmanteau test is nonsignificant at the 5% level for three out of the four component residual series, and is nonsignificant at the 1% level for the other component residual series. On the other hand, a vector AR(2) model was selected by the AIC for the 4-factor series with vector AR(1) as its closest competitor. In fact the AIC values are, respectively, 240.03, 0.11, 0.00, 6.38 and 18.76 for the autoregressive order 0, 1, 2, 3 and 4.

Overall, the fitted model for the monthly temperature vector Y_t is

$$Y_t = \hat{A}\xi_t + e_t,$$

where the white noise e_t has the mean and covariance matrix given in (21), and the 4×1 factor ξ_t follows the VAR(1) model

$$\xi_t - \alpha_t = \hat{\varphi}_0 + \hat{\Phi}_1(\xi_{t-1} - \alpha_{t-1}) + u_t,$$

in which the periodic component $\alpha_t^T = (g_{m(t),1}, g_{m(t),2}, 0, 0)$, $g_{t,i}$ is given in (22),

$$m(t) = \{k \mid 1 \leq k \leq 12 \text{ and } t = 12p + k \text{ for some integer } p \geq 0\},$$

and the white noise u_t has mean 0 and covariance matrix $\hat{\Sigma}_u$ given in (23).

ACKNOWLEDGEMENT

This project was partially supported by a grant from the U.K. Engineering and Physical Science Research Council. Jiazhu Pan was partially supported by the National Natural Science Foundation of China and the National Basic Research Program of China. The authors thank Professor Valdimir Spokoiny for helpful discussion, Mr Da Huang for making available the temperature data analysed in Example 3. Thanks also go to Professor D. M. Titterton and two referees for their helpful comments and suggestions.

APPENDIX

Proofs

We first introduce two lemmas.

LEMMA A1. (i) *It holds for any $H_1, H_2, H_3 \in \mathcal{H}$ that*

$$D(H_1, H_3) \leq D(H_1, H_2) + D(H_2, H_3).$$

(ii) *For any $H_1, H_2 \in \mathcal{H}$, $\Psi(H_1) = \Psi(H_2)$ and $\Psi_n(H_1) = \Psi_n(H_2)$ provided that $D(H_1, H_2) = 0$.*

Proof. (i) For any symmetric matrices M_1 , M_2 and M_3 , it follows from the standard triangle inequality for the matrix norm $\|\cdot\|$ that $\|M_1 - M_3\| \leq \|M_1 - M_2\| + \|M_2 - M_3\|$; that is

$$\{\text{tr}(M_1^2 + M_3^2 - 2M_1M_3)\}^{1/2} \leq \{\text{tr}(M_1^2 + M_2^2 - 2M_1M_2)\}^{1/2} + \{\text{tr}(M_2^2 + M_3^2 - 2M_2M_3)\}^{1/2}. \quad (\text{A1})$$

Let $M_1 = H_1H_1^T$, $M_2 = H_2H_2^T$ and $M_3 = H_3H_3^T$. Since now $\text{tr}(M_i^2) = \text{tr}(M_i) = d - r$ for $i = 1, 2, 3$, the inequality required follows from (A1) and (16) directly.

(ii) Under the condition $D(H_1, H_2) = 0$, $H_1H_1^T = H_2H_2^T$ as it is the projection matrix onto the linear space $\mathcal{M}(H_1) = \mathcal{M}(H_2)$. Now

$$\|H_1^T \Sigma_k H_1\|^2 = \text{tr}\{(H_1^T \Sigma_k H_1)^T H_1^T \Sigma_k H_1\} = \text{tr}(\Sigma_k^T H_1 H_1^T \Sigma_k H_1 H_1^T) = \|H_2^T \Sigma_k H_2\|^2.$$

Hence $\Psi(H_1) = \Psi(H_2)$. The equality for Ψ_n may be proved in the same manner. \square

LEMMA A2. Let $\{Y_t\}$ be a φ -mixing process and let $E(S_k) \rightarrow \Sigma_k$. Suppose that Y_t can be represented as $Y_t = U_t + V_t$, where U_t and V_t are uncorrelated for each t , $\sup_{t \geq 1} E\|U_t\|^h < \infty$ for some constant $h > 2$, and

$$\frac{1}{n} \sum_{t=1}^n V_t \rightarrow c \text{ in probability}, \quad \frac{1}{n} \sum_{t=1}^n E V_t \rightarrow c, \quad (\text{A2})$$

where c is a constant vector. It holds that

- (i) $S_k \rightarrow \Sigma_k$ in probability,
- (ii) $S_k \rightarrow \Sigma_k$ almost surely provided that the mixing coefficients satisfy the condition

$$\varphi(m) = \begin{cases} O\left(\frac{(m-b)}{(2b-2)} - \delta\right), & \text{if } 1 < b < 2, \\ O\left(\frac{(m-3)}{(b-\delta)}\right), & \text{if } b \geq 2, \end{cases} \quad (\text{A3})$$

where $\delta > 0$ is a constant, and the convergence in condition (A2) is also almost sure.

Proof. Assertion (i) follows from the law of large numbers for φ -mixing processes; see Theorem 8.1.1 of Lin & Lu (1997). Applying the result of Chen & Wu (1989) to the sequences $\{U_t\}$ and $\{U_t U_{t-i}^T\}$, and using the almost sure version of Condition (A2), we may obtain (ii). \square

Proof of Theorem 1. Applying the Cauchy–Schwarz inequality to the matrix norm, we have

$$\begin{aligned} |\Psi_n(H) - \Psi(H)| &\leq \sum_{k=1}^p \left| \|H^T S_k H\|^2 - \|H^T \Sigma_k H\|^2 \right| \\ &\leq \sum_{k=1}^p \|H^T (S_k - \Sigma_k) H\| (\|H^T S_k H\| + \|H^T \Sigma_k H\|) \\ &\leq \|H\|^4 \sum_{k=1}^p \|S_k - \Sigma_k\| (\|S_k\| + \|\Sigma_k\|). \end{aligned}$$

Note that $\|H\|^2 = d - r$ for any $H \in \mathcal{H}$, $\|S_k - \Sigma_k\| \rightarrow 0$ in probability, which is implied by Condition 1, and $\|S_k\| + \|\Sigma_k\| = O_P(1)$. Hence, in probability,

$$\sup_{H \in \mathcal{H}_D} |\Psi_n(H) - \Psi(H)| \rightarrow 0. \quad (\text{A4})$$

LEMMA A1. (i) ensures that (\mathcal{H}_D, D) is a well-defined metric space which is complete. Lemma A1(ii) guarantees that $\Psi_n(\cdot)$ is a well-defined stochastic process indexed by $H \in \mathcal{H}_D$, and $\Psi(\cdot)$ is a well-defined function on the metric space (\mathcal{H}_D, D) . It follows from the argmax theorem, see Theorem 3.2.2 and Corollary 3.2.3 of van der Vaart & Wellner (1996), that $D(\hat{B}, B) \rightarrow 0$ in probability.

To show the convergence with probability 1, note that the convergence in (A4) is with probability 1 provided that $S_k \rightarrow \Sigma_k$ with probability 1. Suppose by contradiction that there exists a δ such that $\text{pr}\{\limsup_{n \rightarrow \infty} D(\hat{B}, B_0) > \delta\} > 0$. Let $\mathcal{H}'_D = \mathcal{H}_D \cap \{B : D(B, B_0) \geq \delta\}$. Then \mathcal{H}'_D is a compact subset of \mathcal{H}_D . Note that, if $\sup_{H \in \mathcal{H}_D} |\Psi_n(H) - \Psi(H)| \rightarrow 0$ almost surely, then there exists a set of sample points Ω' satisfying $\Omega' \subset \{\limsup_{n \rightarrow \infty} D(\hat{B}, B_0) > \delta\}$ and $\text{pr}(\Omega') > 0$ such that, for each $\omega \in \Omega'$, one can find a subsequence $\{\hat{B}_{n_k}(\omega)\} \subset \mathcal{H}'_D$ with $\hat{B}_{n_k}(\omega) \rightarrow B \in \mathcal{H}'_D$. Then, by the definition of \hat{B} ,

$$\Psi(B) = \lim_{k \rightarrow \infty} \Psi_{n_k}(\hat{B}_{n_k}(\omega)) \leq \lim_{k \rightarrow \infty} \Psi(B_0) = \Psi(B_0)$$

holds for $\omega \in \Omega'$ and with positive probability. This is a contradiction to Condition 2. Therefore, it must hold that $D(\hat{B}, B_0) \rightarrow 0$ with probability 1. \square

Proof of Theorem 2. Denote by $s_{(i,j),k}$ and $\sigma_{(i,j),k}$, respectively, the (i, j) th elements of S_k and Σ_k . By the Central Limit Theorem for φ -mixing processes, see Lin & Lu (1997) and an unpublished London School of Economics dissertation paper by J. Davidson, it holds that $n^{1/2}(s_{(i,j),k} - E s_{(i,j),k}) \rightarrow N_{(i,j),k}$ in distribution, where $N_{(i,j),k}$ denotes a Gaussian random variable, $i, j = 1, \dots, d$. Hence, $\|n^{1/2}(S_k - E S_k)\| = O_P(1)$. It holds now that

$$\begin{aligned} & \sup_{H \in \mathcal{H}_D} n^{1/2} |\Psi_n(H) - \Psi(H)| \\ & \leq \sup_{H \in \mathcal{H}_D} n^{1/2} \sum_{k=1}^p \|H^T S_k H\|^2 - \|H^T \Sigma_k H\|^2 \\ & \leq \sup_{H \in \mathcal{H}_D} \sum_{k=1}^p \|H^T n^{1/2}(S_k - E S_k) H\| (\|H^T S_k H\| + \|H^T \Sigma_k H\|) \\ & \quad + \sup_{H \in \mathcal{H}_D} \sum_{k=1}^p \|H^T \{n^{1/2}(E S_k - \Sigma_k)\} H\| (\|H^T S_k H\| + \|H^T \Sigma_k H\|) \\ & \leq p \sup_{H \in \mathcal{H}_D, 1 \leq k \leq p} \|H^T n^{1/2}(S_k - E S_k) H\| (\|H^T S_k H\| + \|H^T \Sigma_k H\|) \\ & \quad + p \sup_{H \in \mathcal{H}_D, 1 \leq k \leq p} \|H^T \{n^{1/2}(E S_k - \Sigma_k)\} H\| (\|H^T S_k H\| + \|H^T \Sigma_k H\|) \\ & \leq p(d-r)^4 \left\{ \sup_{1 \leq k \leq p} \|n^{1/2}(S_k - E S_k)\| (\|S_k\| + \|\Sigma_k\|) \right. \\ & \quad \left. + \sup_{1 \leq k \leq p} \|n^{1/2}(E S_k - \Sigma_k)\| (\|S_k\| + \|\Sigma_k\|) \right\} = O_P(1). \end{aligned} \tag{A5}$$

By Condition 3, (A5) and the definitions of B and \hat{B} , we have that

$$\begin{aligned} 0 & \leq \Psi_n(B) - \Psi_n(\hat{B}) \\ & = \Psi(B) - \Psi(\hat{B}) + O_P(n^{-1/2}) \leq -a\{D(\hat{B}, B)\}^c + O_P(n^{-1/2}). \end{aligned}$$

Now let $n \rightarrow \infty$ in the above expression. It must hold that $D(\hat{B}, B) = O_P(n^{-1/(2c)})$. \square

REFERENCES

- AHN, S. K. (1997). Inference of vector autoregressive models with cointegration and scalar components. *J. Am. Statist. Assoc.* **93**, 350–6.
- ANDERSON, T. W. (2002). Canonical correlation analysis and reduced rank regression in autoregressive models. *Ann. Statist.* **30**, 1134–54.

- BAI, J. & NG, S. (2002). Determining the number of factors in approximate factor models. *Econometrica* **70**, 191–222.
- BRILLINGER, D. R. (1981). *Time Series Data Analysis and Theory*, extended ed. San Francisco: Holden-Day.
- BOX, G. & TIAO, G. (1977). A canonical analysis of multiple time series. *Biometrika* **64**, 355–65.
- BROCKWELL, J. P. & DAVIS, R. A. (1991). *Time Series Theory and Methods*, 2nd ed. New York: Springer.
- CHEN, X. R. & WU, Y. H. (1989). Strong law for a mixing sequence. *Acta Math. Appl. Sinica* **5**, 367–71.
- ENGLE, R. & WATSON, M. (1981). A one-factor multivariate time series model of metropolitan wage rates. *J. Am. Statist. Assoc.* **76**, 774–81.
- FORNI, M., HALLIN, M., LIPPI, M. & REICHLIN, L. (2000). The generalized dynamic factor model: identification and estimation. *Rev. Econ. Statist.* **82**, 540–54.
- FRANCO, C., ROY, R. & ZAKOÏAN, J.-M. (2005). Diagnostic checking in ARMA models with uncorrelated errors. *J. Am. Statist. Assoc.* **100**, 532–44.
- GEWEKE, J. (1977). The dynamic factor analysis of economic time series models. In *Latent Variables in Socio-Economic Models*, Ed. D. J. Aigner and A. S. Goldberger, pp. 365–83 Amsterdam: North-Holland.
- GEWEKE, J. & SINGLETON, K. (1981). Maximum likelihood confirmatory factor analysis of economic time series. *Int. Econ. Rev.* **22**, 37–54.
- LI, W. K. & MCLEOD, A. I. (1981). Distribution of the residuals autocorrelations in multivariate ARMA time series models. *J. R. Statist. Soc. B*, **43**, 231–9.
- LIN, Z. & LU, C. (1997). *Limit Theory for Mixing Dependent Random Variables*. New York: Kluwer.
- LJUNG, G. M. & BOX, G. E. P. (1978). On a measure of lack of fit in time series models. *Biometrika* **65**, 297–303.
- PEÑA, D. & BOX, G. E. P. (1987). Identifying a simplifying structure in time series. *J. Am. Statist. Assoc.* **82**, 836–43.
- PEÑA, D. & PONCELA, P. (2006). Nonstationary dynamic factor analysis. *J. Statist. Plan. Infer.* **136**, 1237–57.
- PRESS, W. H., TEUKOLSKY, S. A., VETTERLING, W. T. & FLANNERY, B. P. (1992). *Numerical Recipes in C*. Cambridge: Cambridge University Press.
- PRIESTLEY, M. B., SUBBA RAO, T. & TONG, H. (1974). Applications of principal component analysis and factor analysis in the identification of multivariate systems. *IEEE Trans. Auto. Contr.* **19**, 703–4.
- REINSEL, G. C. (1997). *Elements of Multivariate Time Series Analysis*, 2nd ed. New York: Springer.
- REINSEL, G. C. & VELU, R. P. (1998). *Multivariate Reduced Rank Regression*. New York: Springer.
- STOCK, J. H. & WATSON, M. W. (2002). Forecasting using principal components from a large number of predictors. *J. Am. Statist. Assoc.* **97**, 1167–79.
- TIAO, G. C. & TSAY, R. S. (1989). Model specification in multivariate time series (Discussion). *J. R. Statist. Soc. B* **51**, 157–213.
- VAN DER VAART, A. W. & WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. New York: Springer.

[Received September 2006. Revised September 2007]

