Automatica 65 (2016) 170-182

Contents lists available at ScienceDirect

Automatica

journal homepage: www.elsevier.com/locate/automatica

The Box–Jenkins Steiglitz–McBride algorithm*

Yucai Zhu^a, Håkan Hjalmarsson^b

^a State Key Laboratory of Industrial Control Technology, Department of Control Science and Engineering, Zhejiang University, Zheda Road 38, Hangzhou 310027, China

^b ACCESS Linnaeus Center, Electrical Engineering, KTH–Royal Institute of Technology, S-100 44 Stockholm, Sweden

ARTICLE INFO

Article history: Received 6 June 2013 Received in revised form 14 September 2015 Accepted 12 November 2015 Available online 21 December 2015

Keywords: System identification Steiglitz–McBride Box–Jenkins model High-order ARX-modeling

1. Introduction

In system identification, the prediction error method (PEM) is well established (Ljung, 1999). If correct model orders are used, a quadratic cost function yields consistent and asymptotically efficient estimates for both open-loop and closed-loop data when the noise is Gaussian. However, all these nice properties rely on the precondition that the global minimum of the cost function is found in the parameter estimation. With the exception of models with a linear regression structure such as ARX and FIR models, most model structures need (local) nonlinear numerical optimization routines for the parameter estimation. Thus, unless there are no non-global local minima, convergence to a global minimum cannot be guaranteed and the nice asymptotic properties may be lost. There are some asymptotic (in the sample-size) results for when there are no "false" minima: For ARMA-models this is always the case (Åström & Söderström, 1974). For Box-Jenkins models this holds when only one system pole is



An algorithm for identification of single-input single-output Box–Jenkins models is presented. It consists of four steps: firstly a high order ARX model is estimated; secondly, the input–output data is filtered with the inverse of the estimated disturbance model; thirdly, the filtered data is used in the Steiglitz–McBride method to recover the system dynamics; in the final step, the noise model is recovered by estimating an ARMA model from the residuals of the third step. The relationship to other identification methods, in particular the refined instrumental-variable method, are elaborated upon. A Monte Carlo simulation study with an oscillatory system is presented and these results are complemented with an industrial case study. The algorithm can easily be generalized to multi-input single-output models with common denominator. © 2015 Elsevier Ltd. All rights reserved.

estimated (Söderström, 1975a). For output error models this holds for general orders when the input is white (Söderström, 1975a). Necessary conditions on the input for this to hold are provided in Goodwin, Agüero, and Skelton (2003) (ARMAX-models), Zou and Heath (2009) (output-error models) and Eckhard, Bazanella, Rojas, and Hjalmarsson (2012). While experience is that the predictionerror method based on local non-linear optimization works well for many problems, there are also a range of examples when it gets stuck in local minima (Eckhard et al., 2012; Goodwin et al., 2003; Zou & Heath, 2009). It is also of general interest to develop alternative methods that can provide models that such non-linear optimization methods can use as initial estimates. As a consequence, a range of methods complementing PEM have been developed over time. Subspace identification (Van Overschee & De Moor, 1994; Verhaegen, 1994) and instrumental variable methods (Stoica & Söderström, 1983; Young, 1976) are two such families of methods.

The so-called Box–Jenkins model is a flexible and useful model structure (Box & Jenkins, 1970; Ljung, 1999). The model has compact rational descriptions for both the process model and the disturbance model. Zhu (Zhu, 2011) proposed an algorithm for Box–Jenkins model estimation and has outlined a proof of its convergence for open loop data. The idea was inspired by the analysis of the Steiglitz–McBride method in Stoica and Söderström (1981). We will continue this work in several aspects. We provide a detailed presentation of the algorithm and how it relates to other methods, in particular the Refined Instrumental Method (RIV) (Young, 2008). Theoretical justification of the algorithm is provided by way of results on convergence and asymptotic efficiency. Finally, we provide a simulation study complemented





코 IFA

automatica

[†] The work of Yucai Zhu is supported by 973 Program of China (No. 2012CB720500) and by National Science Foundation of China (No. 61273191). The work of Håkan Hjalmarsson was supported in part by the European Research Council under the advanced grant LEARN, contract 267381, and in part by the Swedish Research Council under contract 621-2009-4017. The material in this paper was partially presented at the 18th IFAC World Congress, August 28–September 2, 2011, Milan, Italy. This paper was recommended for publication in revised form by Associate Editor Wei Xing Zheng under the direction of Editor Torsten Söderström.

E-mail addresses: y.zhu@taijicontrol.com, yczhu@iipc.zju.edu.cn (Y. Zhu), hakan.hjalmarsson@ee.kth.se (H. Hjalmarsson).

with an industrial case study. The paper covers single-input single-output (SISO) models.

The outline of the paper is as follows. Section 2 introduces assumptions regarding the data generation mechanism, outlines the prediction-error method for Box–Jenkins models and reviews the Steiglitz–McBride method and ARX-model estimation. The new algorithm is presented in Section 3 where we also relate it to existing methods. Some of its asymptotic properties are presented in Section 4. Practical issues are briefly covered in Section 5. In Section 6, a simulation example and an industrial case study is used to demonstrate the accuracy and robustness of the algorithm. Section 7 contains conclusions.

Notation. q^{-1} is the time-shift operator: $q^{-1}x_t := x_{t-1}$. E[x] denotes the expectation of x and $E[x|\mathcal{F}]$ is conditional expectation with respect to the σ -algebra \mathcal{F} . Convergence with probability 1 is abbreviated w.p.1.

2. Preliminaries

2.1. Assumptions

We will make the following assumptions about the model and the true system.

Assumption 2.1 (*Model and True System*). The system has scalar input u_t , scalar output y_t and is subject to the scalar noise e_t . The model for the relationships between these signals is given by

$$y_t = G(q, \theta)u_t + v_t, \quad v_t = H(q, \gamma)e_t, \tag{1}$$

where $G(q, \theta)$ and $H(q, \gamma)$ are rational functions in q^{-1} :

$$G(q,\theta) = \frac{L(q,\theta)}{F(q,\theta)} = \frac{l_1 q^{-1} + \dots + l_{n_l} q^{-n_l}}{1 + f_1 q^{-1} + \dots + f_{n_f} q^{-n_f}},$$

$$H(q,\gamma) = \frac{C(q,\gamma)}{D(q,\gamma)} = \frac{1 + c_1 q^{-1} + \dots + c_{n_c} q^{-n_c}}{1 + d_1 q^{-1} + \dots + d_{n_d} q^{-n_d}}.$$

For ease of notation, we will in the following assume $n_l = n_f = n_c = n_d = m$ for some positive integer *m*. We will also drop the second argument in *L*, *F*, *C* and *D*. The model parameters are collected in $\theta = [f_1 \cdots f_m \ l_1 \cdots l_m]^T$, and $\gamma = [c_1 \cdots c_m \ d_1 \ \cdots d_m]^T$. We will assume that there is $\theta = \theta_0$ and $\gamma = \gamma_0$ such that (1)

We will assume that there is $\theta = \theta_0$ and $\gamma = \gamma_0$ such that (1) describes the true system. The polynomials $L^0(q) := L(q, \theta_0)$ and $F^0(q) := F(q, \theta_0)$ do not share common factors. The same is true for $C^0(q) := C(q, \gamma_0)$ and $D^0(q) := D(q, \gamma_0)$. It is further assumed that the transfer functions $G^0 := G(q, \theta_0)$ and $H^0 := H(q, \gamma_0)$ are stable, i.e. $F^0(z) = 0 \Rightarrow |z| \le 1$, $D^0(z) = 0 \Rightarrow |z| \le 1$, and that H^0 has a stable inverse.

The input $\{u_t\}$ will be assumed to be a realization of a stochastic process generated by a random sequence $\{w_t\}$. Let \mathcal{F}_{t-1} be the σ -algebra generated by $\{e_s, w_s, s \leq t - 1\}$. Then the following assumption is in force.

Assumption 2.2 (*Input*). The sequence $\{u_t\}$ is defined by

$$u_t = F_u(q)w_t,$$

where $F_u(q)$ is a stable and inversely stable finite dimensional filter, where $\{w_t\}$ is independent of $\{e_t\}$ satisfying

$$E[w_t | \mathcal{F}_{t-1}] = 0, \quad E[w_t^2 | \mathcal{F}_{t-1}] = 1, \quad |w_t| \le C, \ \forall t,$$

for some positive finite constant *C*.

Assumption 2.2 implies that the system is operating in openloop as the input and the noise are independent. Notice that the assumption that F_u is inversely stable implies that F_u cannot have any zeros on or outside the unit circle. It can be interpreted as the stable minimum phase spectral factor of the input spectrum. The noise satisfies the following assumption.

Assumption 2.3 (*Noise*). {*e*_{*t*}} is a stochastic process satisfying

$$\mathbb{E}[e_t|\mathcal{F}_{t-1}] = 0, \qquad \mathbb{E}[e_t^2|\mathcal{F}_{t-1}] = \sigma_o^2, \quad \mathbb{E}\left[|e_t|^{10}\right] \le C, \ \forall t.$$

2.2. The prediction error method

The prediction-error of the Box–Jenkins model (1) is given by

$$\varepsilon_t = \frac{D(q)}{C(q)} \left[y_t - \frac{L(q)}{F(q)} u_t \right].$$

In the prediction error method, using a quadratic loss function, the parameter estimates are determined by minimizing the loss function

$$V_N = \frac{1}{N} \sum_{t=1}^N \varepsilon_t^2, \tag{2}$$

where *N* is the number of data samples, with respect to the parameters in θ and γ .

It is well known (Ljung, 1999) that when the data set is informative and collected in open-loop, and the prediction-error method is applied to the model (1), the asymptotic covariance matrix of the parameter estimate $\hat{\theta}_N^{\text{PEM}}$ of the system parameters θ is under some mild technical assumptions given by

$$\lim_{N \to \infty} NE \left[(\hat{\theta}_N^{\text{PEM}} - \theta_o) (\hat{\theta}_N^{\text{PEM}} - \theta_o)^T \right] = \sigma_o^2 M_{CR}^{-1}, \tag{3}$$

where

$$M_{CR} = \frac{1}{2\pi\sigma_o^2} \int_{-\pi}^{\pi} \left[\frac{-\frac{G^0}{H^0F^0}}{\frac{1}{H^0F^0}} \Gamma_m \right] \left[\frac{-\frac{G^0}{H^0F^0}}{\frac{1}{H^0F^0}} \Gamma_m \right]^T |F_u|^2 d\omega, \qquad (4)$$

where in turn $\Gamma_m(q) = [q^{-1} \cdots q^{-m}]^T$. For brevity we have omitted the argument $e^{i\omega}$ in the integrand above; a practice we will employ hereafter. The positivity of M_{CR} follows from the assumption that the pairs $\{L^o, F^o\}$ and $\{C^o, D^o\}$ do not share common factors (identifiability) if it is assumed that the input spectrum $|F_u|^2$ is strictly positive on the unit circle (persistence of excitation), except possibly for a finite number of points, see Ljung (1999). When $\{e_t\}$ is Gaussian, M_{CR}/σ_o^2 is the per sample Fisher Information Matrix, and hence PEM with a quadratic cost function is asymptotically efficient, reaching the Cramér–Rao lower bound (Ljung, 1999).

Under the same assumptions it follows that the estimates of the parameters in the disturbance model C(q)/D(q) are asymptotically independent of the estimate of θ (Pierce, 1972).

2.3. The Steiglitz–McBride method

The Steiglitz–McBride (SM) method is an iterative method for estimating the system dynamics. Consider an output-error model, i.e. when C(q) = D(q) = 1, and that an estimate \hat{F}^k of F is available. Then (2) is given by

$$\begin{split} V_N &= \frac{1}{N} \sum_{t=1}^{N} \left| y_t - \frac{L(q)}{F(q)} u_t \right|^2 \\ &\approx \frac{1}{N} \sum_{t=1}^{N} \left| F(q) \frac{1}{\hat{F}^k(q)} y_t - L(q) \ \frac{1}{\hat{F}^k(q)} u_t \right|^2 \\ &= \frac{1}{N} \sum_{t=1}^{N} \left| F(q) y_t^k - L(q) u_t^k \right|^2, \end{split}$$

where u_t^k and y_t^k are u_t and y_t , respectively, prefiltered with $1/\hat{F}^k(q)$. The last expression is quadratic in the elements of *F* and *L* and can therefore be minimized explicitly, giving a new estimate \hat{F}^{k+1} of *F* as well as an estimate \hat{L}^{k+1} of *L*. Given an initial estimate \hat{F}^0 , this procedure of prefiltering the data and solving for new estimates can be iterated for $k = 0, 1, \ldots$.

In the limit $N \rightarrow \infty$ and when the true system is also of outputerror type, the method is known to converge locally to the true parameter vector under standard persistence of excitation conditions, and also globally if the signal-to-noise ratio is sufficiently large (Stoica & Söderström, 1981). It is also known that the method is not asymptotically efficient. Furthermore, when the noise is colored, the true parameter vector θ_o is in general not a possible convergence point.

2.4. ARX modeling

Alternatively, the true system can be represented as

$$A^{o}(q)y_{t} = B^{o}(q)u_{t} + e_{t},$$
(5)

where

$$A^{o}(q) = \frac{1}{H^{o}(q)} = 1 + \sum_{k=1}^{\infty} a_{k}^{o} q^{-k},$$
$$B^{o}(q) = \frac{G^{o}(q)}{H^{o}(q)} = \sum_{k=1}^{\infty} b_{k}^{o} q^{-k}$$

are stable transfer functions (by Assumption 2.1). This suggests the use of an ARX-model for estimating the system transfer functions G^{0} and H^{0} . Therefore, let $\eta^{n} = \begin{bmatrix} a_{1} & \cdots & a_{n} & b_{1} & \cdots & b_{n} \end{bmatrix}^{T}$ and define

$$A(q, \eta^{n}) = 1 + \sum_{k=1}^{n} a_{k} q^{-k}, \qquad B(q, \eta^{n}) = \sum_{k=1}^{n} b_{k} q^{-k}.$$
 (6)

Consider the following ARX-model

 $A(q, \eta^n)y_t = B(q, \eta^n)u_t + e_t,$ which we can write as

$$y_t = (\varphi_t^n)^T \eta^n + e_t, \tag{7}$$

where

$$\varphi_t^n = \begin{bmatrix} -\Gamma_n^T(q) y_t & \Gamma_n^T(q) u_t \end{bmatrix}^T.$$
(8)

A regularized least-squares estimate of η^n in the model (7) is given by

$$\hat{\eta}_{N}^{n} = \begin{bmatrix} \hat{a}_{1}^{n,N} & \dots & \hat{a}_{n}^{n,N} & \hat{b}_{1}^{n,N} & \dots & \hat{b}_{n}^{n,N} \end{bmatrix}^{T} \\ \coloneqq \begin{bmatrix} R_{reg}^{n}(N) \end{bmatrix}^{-1} r^{n}(N),$$
(9)

where

$$r^{n}(N) = \frac{1}{N} \sum_{t=n+1}^{N} \varphi_{t}^{n} y_{t}, \qquad R^{n}(N) = \frac{1}{N} \sum_{t=n+1}^{N} \varphi_{t}^{n} (\varphi_{t}^{n})^{T},$$
$$R^{n}_{reg}(N) = \begin{cases} R^{n}(N) & \text{if } \|R^{n}(N)^{-1}\|_{2} < 2/\delta, \\ R^{n}(N) + \frac{\delta}{2} I_{2n}, & \text{otherwise,} \end{cases}$$

for some small $\delta > 0$. The reason for using the regularized $R_{reg}^n(N)$ rather than $R^n(N)$ as in standard least-squares, is that it facilitates the statistically analysis, see Ljung and Wahlberg (1992). Asymptotically (in the sample-size N), the first and second order properties of $\hat{\eta}_N^n$ do not depend on δ .

While ARX-models have the attractive property that the parameter estimate is given by the closed form expression (9),

it is in general not consistent when the underlying system is of Box–Jenkins type (1). This can, however, be remedied by allowing the model order n to depend on the sample size, i.e. n = n(N). For our theoretical results we will use the following assumption.

Assumption 2.4 (ARX Model Order). It holds that

$$n(N) \to \infty$$
 and $n(N)^{3+\delta}/N \to 0$, as $N \to \infty$,
for some $\delta > 0$.

Introduce the notation

$$\hat{\eta}_{N} \coloneqq \hat{\eta}_{N}^{n(N)},$$

$$\eta_{o}^{n} \coloneqq \begin{bmatrix} a_{1}^{o} & \dots & a_{n}^{o} & b_{1}^{o} & \dots & b_{n}^{o} \end{bmatrix}^{T},$$

$$\eta_{o} \coloneqq \begin{bmatrix} a_{1}^{o} & a_{2}^{o} & \dots & b_{1}^{o} & b_{2}^{o} & \dots \end{bmatrix}^{T}.$$
(10)

The asymptotic properties of $\hat{\eta}_N$ are established in Ljung and Wahlberg (1992). We will need the following result.

Lemma 2.1. Assume that Assumptions 2.1–2.4 hold. Then with probability 1,

$$\sup_{\omega} |A(e^{j\omega}, \hat{\eta}_N) - A^o(e^{j\omega})| = O(m(N)) \to 0, \quad N \to \infty,$$

where $m(N) = n(N)\sqrt{\log N/N}(1 + d(N)) + d(N)$, where

$$d(N) := \sum_{k=n(N)+1}^{\infty} |a_k^o| + |b_k^o| \le \tilde{C} \rho^{n(N)},$$
(11)

for some $\tilde{C} < \infty$, $\rho < 1$.

Proof. See Appendix A.

The lemma shows that by allowing the model order *n* to depend on the sample size *N* in a suitable manner, consistency of the inverse noise model estimate $A(e^{j\omega}, \hat{\eta}_N)$ follows. However, this comes at a cost, namely that the transfer function estimates are not asymptotically efficient (Ljung & Wahlberg, 1992).

3. The Box-Jenkins Steiglitz-McBride algorithm

In the previous section we have discussed two different ways to overcome the problem of local minima of PEM: the Steiglitz–McBride method and (high-order) ARX-model estimation. However, none of these two methods are asymptotically efficient, and furthermore the Steiglitz–McBride method is not even consistent unless the measurement noise is white, see Section 2.3. The idea in this paper is to combine the two methods to alleviate these problems.

3.1. Outline of algorithm

The algorithm comprises the following four steps:

- (1) Estimate an ARX model using the input-output data $\{u_t, y_t\}$, t = 1, 2, ..., N, using (9). Denote the corresponding A(q) and B(q) estimates by $\hat{A}(q)$ and $\hat{B}(q)$, respectively.
- (2) Filter the input and output signals using $\hat{A}(q)$, the inverse of the disturbance model obtained in Step 1:

$$y_t^f = \hat{A}(q)y_t, \qquad u_t^f = \hat{A}(q)u_t.$$

From (1) it follows that y_t^f and u_t^f are related by

$$y_t^f = \frac{L^0(q)}{F^0(q)} u_t^f + \hat{A}(q) \frac{C^0(q)}{D^0(q)} e_t.$$
 (12)

If $\hat{A}(q)$ is a good estimate of the inverse of the disturbance dynamics, it holds that the noise term in (12) approximately equals e_t so that

$$y_t^f pprox rac{L^o(q)}{F^o(q)} u_t^f + e_t,$$

which is the setting for which the Steiglitz–McBride method applies.

- (3) The third step is therefore to apply the Steiglitz–McBride method to the pre-filtered data y_t^f and u_t^f , yielding estimates $\hat{L}(q)$ and $\hat{F}(q)$ of $L^o(q)$ and $F^o(q)$, respectively.
- (4) In the last step a disturbance model $\hat{C}(q)/\hat{D}(q)$ is obtained by estimating an ARMA model of the output error residual \hat{v}_t given by

$$\hat{v}_t = y_t - \frac{\hat{L}(q)}{\hat{F}(q)} u_t.$$
(13)

Using PEM in this step, requires nonlinear optimization. Alternatively, an instrumental variable method can be used in this step (Young, 2006).

The above algorithm will be referred to as BJSM (Box–Jenkins Steiglitz–McBride). Note that the process model estimation in Step 3 will not be affected by the disturbance model estimation in Step 4.

3.2. A formal expression

We will now provide a formal expression for the estimate of θ in iteration k + 1 of the Steiglitz–McBride step of the BJSM algorithm. Let

$$\hat{\theta}_N^k = \begin{bmatrix} \hat{f}_1^{N,k} & \cdots & \hat{f}_m^{N,k} & \hat{l}_1^{N,k} & \cdots & \hat{l}_m^{N,k} \end{bmatrix}^T$$

denote the estimate at iteration k. For any signal x_t define

$$x_t(\eta^n,\theta) = \frac{A(q,\eta^n)}{F(q,\theta)} x_t, \qquad x_t(\eta_o,\theta) = \frac{A^o(q)}{F(q,\theta)} x_t.$$
(14)

The same definition applies to vector valued signals such as (8). Also define

$$e_t(\eta^n, \theta, F^o) := F^o(q)v_t(\eta^n, \theta) = \frac{F^o(q)}{F(q, \theta)} \frac{A(q, \eta^n)}{A^o(q)}e_t.$$

From (1) we have

 $F^{o}(q)y_{t} = L^{o}(q)u_{t} + F^{o}(q)v_{t},$

which implies

$$F^{o}(q)y_{t}(\hat{\eta}_{N},\hat{\theta}_{N}^{k}) = L^{o}(q)u_{t}(\hat{\eta}_{N},\hat{\theta}_{N}^{k}) + F^{o}(q)v_{t}(\hat{\eta}_{N},\hat{\theta}_{N}^{k})$$
$$= L^{o}(q)u_{t}(\hat{\eta}_{N},\hat{\theta}_{N}^{k}) + e_{t}(\hat{\eta}_{N},\hat{\theta}_{N}^{k},F^{o}),$$

which can be written in regression form as

$$\mathbf{y}_t(\hat{\eta}_N, \hat{\theta}_N^k) = [\varphi_t^m(\hat{\eta}_N, \hat{\theta}_N^k)]^T \theta_o + \mathbf{e}_t(\hat{\eta}_N, \hat{\theta}_N^k, F^o).$$
(15)

Given $\hat{\theta}_N^k$, $\hat{\theta}_N^{k+1}$ is now defined as the least squares estimate of θ_o in the linear regression (15), i.e.

$$\hat{\theta}_N^{k+1} = [R^m(N, \hat{\eta}_N, \hat{\theta}_N^k)]^{-1} r^m(N, \hat{\eta}_N, \hat{\theta}_N^k),$$
(16)
where

where

$$R^{m}(N, \eta^{n}, \theta) = \frac{1}{N} \sum_{t=m+1}^{N} \varphi_{t}^{m}(\eta^{n}, \theta) (\varphi_{t}(\eta^{n}, \theta))^{T},$$

$$r^{m}(N, \eta^{n}, \theta) = \frac{1}{N} \sum_{t=m+1}^{N} \varphi_{t}^{m}(\eta^{n}, \theta) y_{t}(\eta^{n}, \theta).$$

N.

3.3. Relation to other methods

Many researchers have proposed to use high order ARX models and then to apply model reduction; see, e.g., Söderström (1975b, Chapter 7) in Hsia (1977), Wahlberg (1989) and Zhu (1998). The first author has used high-order ARX-models in developing the so-called asymptotic (ASYM) method which has been successfully applied to many industrial processes; see Zhu (1998, 2009).

Furthermore, BJSM bears some resemblance to the refined instrumental variable (RIV) method and the closely related multistep algorithm presented in Stoica and Söderström (1983). Below we will point out similarities and differences. The RIV method was developed by Young and co-workers; see Jakeman and Young (1979) and Young (1976, 2008). For the Box–Jenkins model (1), RIV uses an iterative scheme as follows:

At iteration k, given the estimate $y_t = (\hat{L}^k(q)/\hat{F}^k(q))u_t + (\hat{C}^k(q)\hat{D}^k(q))\hat{\varepsilon}_t$, calculate the next estimate as follows

$$\hat{\theta}^{k+1} = \left[\frac{1}{N}\sum_{t=1}^{N}\hat{\varphi}_{t}(\hat{\theta}^{k})\varphi_{t}(\hat{\theta}^{k})^{T}\right]^{-1}\frac{1}{N}\sum_{t=1}^{N}\hat{\varphi}_{t}(\hat{\theta}^{k})y_{t}^{f,k}.$$
(17)

Here $y_t^{f,k} = \frac{\hat{D}^k(q)}{\hat{C}^k(q)\hat{F}^k(q)}y_t$ is the pre-filtered output,

$$\varphi_t(\hat{\theta}^k) = \frac{\hat{D}^k(q)}{\hat{C}^k(q)\hat{F}^k(q)} [-y_{t-1}, \dots, -y_{t-m}, u_{t-1}, \dots, u_{t-m}]^T$$

is the prefiltered input-output data vector and

$$\hat{\varphi}_{t}(\hat{\theta}^{k}) = \frac{\hat{D}^{k}(q)}{\hat{C}^{k}(q)\hat{F}^{k}(q)}[\hat{x}_{t-1}^{k}, \dots, \hat{x}_{t-n_{f}}^{k}, u_{t-1}^{k}, \dots, u_{t-n_{l}}^{k}]^{T}$$

is the prefiltered instrument vector, where $\hat{x}_k^k = (\hat{L}^k(q)/\hat{F}^k(q))u_t$ is the simulated model output at iteration k. Comparing RIV above, with BJSM in Section 3.2, we see that there is a similarity between the methods in that the methods are iterative, with the data prefilters of both methods being conceptually the same: they are the products of the disturbance model inverses and $1/\hat{F}^k(q)$. The two main differences are: (1) after the prefiltering, RIV is an instrumental variable method and BJSM is a least-squares method; (2) in RIV, the inverse disturbance model is re-estimated in each iteration and based on a low order (ARMA) model of the model residuals, while in the BJSM algorithm, the inverse disturbance model is the same in all iterations and based on a high-order (ARX) estimate.

In the classification in Ljung (1999), Step 3 corresponds to the family of methods (7.110), where prefiltered prediction errors are correlated with past data, also hosting pseudo linear regression and instrumental variable methods. In BJSM, the prefilter is $\hat{A}(q)/\hat{F}^k$ and the correlation vector is $(\hat{A}(q)/\hat{F}^k)\varphi_t^m$. In RIV, the prefilter is $\hat{D}^k(q)/(\hat{C}^k\hat{F}^k)$ and the correlation vector is $\hat{\varphi}_t(\hat{\theta}^k)$.

4. Asymptotic properties

In this section we will study convergence and asymptotic variance of the BJSM method.

4.1. Convergence

We will follow the mean-value approach of Stoica and Söderström (1981) and analyze the limit case when N has reached infinity. As in Stoica and Söderström (1981), we begin the analysis by considering possible convergence points. For this we will first derive an equation that defines these points.

Suppose that the limit $\hat{\theta}_N := \lim_{k \to \infty} \hat{\theta}_N^k$ exists. Due to (16) it must hold for this limit that

$$\hat{\theta}_N = [R^m(N, \hat{\eta}_N, \hat{\theta}_N)]^{-1} r^m(N, \hat{\eta}_N, \hat{\theta}_N).$$
(18)

Due to (15),

$$y_t(\hat{\eta}_N, \hat{\theta}_N) = [\varphi_t^m(\hat{\eta}_N, \hat{\theta}_N)]^T \theta_o + e_t(\hat{\eta}_N, \hat{\theta}_N, F^o).$$
(19)

Using (18)–(19) gives

$$0 = r^{m}(N, \hat{\eta}_{N}, \hat{\theta}_{N}) - R^{m}(N, \hat{\eta}_{N}, \hat{\theta}_{N})\hat{\theta}_{N}.$$
(20)

Generalizing Theorem 1 in Stoica and Söderström (1981), our first result concern solutions to the equation

$$0 = r^{m}(N, \hat{\eta}_{N}, \theta) - R^{m}(N, \hat{\eta}_{N}, \theta)\theta, \qquad (21)$$

as the sample size tends to infinity.

Theorem 4.1. Let Assumptions 2.1–2.4 be in force. Then as $N \rightarrow \infty$, (21) converges with probability 1 to

$$0 = \mathbb{E}\left[\varphi_t^m(\eta_o, \theta) y_t(\eta_o, \theta)\right] - \tilde{R}(\theta)\theta,$$
(22)

where

$$\tilde{R}(\theta) = \mathbb{E}\left[\varphi_t^m(\eta_o, \theta)(\varphi_t^m(\eta_o, \theta))^T\right].$$
(23)

The unique solution to (22) is given by $\theta = \theta_0$.

Proof. See Appendix B.

In Stoica and Söderström (1981), the authors proceed with analyzing the behavior of the recursion

$$0 = \mathbb{E}\left[\varphi_t^m(\eta_o, \theta_k) y_t(\eta_o, \theta_k)\right] - \tilde{R}(\theta_k) \theta_{k+1}.$$
(24)

This equation corresponds to the Steiglitz–McBride iterations in the limit that the sample-size $N \rightarrow \infty$. The analysis in Stoica and Söderström (1981) provides local convergence (Theorem 2) and global convergence (Theorem 3). Both these results also apply to our setting since (24) is the same as the relation (5) analyzed in Stoica and Söderström (1981). The global convergence result requires a sufficiently high signal to noise ratio. Here we will in addition to these results present a new result where we trade this condition for a condition that the initial parameter estimate and θ_o both belong to a convex subset of the stability domain.

Following Stoica and Söderström (1981) we write (24) as

$$\begin{aligned} \theta_{k+1} &- \theta_o = \tilde{R}^{-1}(\theta_k) \\ & \mathsf{E}\left[\varphi_t^m(\eta_o, \theta_k) \left(y_t(\eta_o, \theta_k) - (\varphi_t^m(\eta_o, \theta_k))^T \theta_o\right)\right] \\ &= \tilde{R}^{-1}(\theta_k) \mathsf{E}\left[\varphi_t^m(\eta_o, \theta_k) \frac{F_o(q)}{F(q, \theta_k)} e_t\right] \\ &= \tilde{R}^{-1}(\theta_k) \mathsf{E}\left[\varphi_t^m(\eta_o, \theta_k) \left(-\frac{F(q, \theta_k) - F_o(q)}{F(q, \theta_k)} e_t\right)\right] \\ &= \tilde{R}^{-1}(\theta_k) D(\theta_k)(\theta_k - \theta_o), \end{aligned}$$
(25)

where

$$D(\theta) = \begin{bmatrix} E \begin{bmatrix} \frac{1}{F(q,\theta)} \Gamma_m(q) e_t & \frac{1}{F(q,\theta)} \Gamma_m^T(q) e_t \end{bmatrix} & \mathbf{0}_{m \times m} \\ \mathbf{0}_{m \times m} & \mathbf{0}_{m \times m} \end{bmatrix}.$$
 (26)

In (25) we have used (19) in the second equality, and that e_t is uncorrelated with $\varphi_t^m(\eta_o, \theta_k)$ in the third.

The system description (5) gives

$$\varphi_t^m(\eta_o,\theta) = \frac{A^o(q)}{F(q,\theta)}\varphi_t^m = \xi_t^u(\theta) + \xi_t^e(\theta),$$

where

$$\begin{split} \xi_t^u(\theta) &\coloneqq \frac{1}{F(q,\theta)} \begin{bmatrix} -B^o(q) \, \Gamma_m(q) \\ A^o(q) \, \Gamma_m(q) \end{bmatrix} u_t, \\ \xi_t^e(\theta) &\coloneqq \frac{1}{F(q,\theta)} \begin{bmatrix} \Gamma_m(q) \\ 0_{m \times 1} \end{bmatrix} e_t. \end{split}$$

Using that $\{u_t\}$ and $\{e_t\}$ are mutually independent, and Parseval's formula, gives

$$\tilde{R}(\theta) = M(\theta) + D(\theta),$$
(27)

where

$$M(\theta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \begin{bmatrix} -B^{0} \Gamma_{m} \\ A^{0} \Gamma_{m} \end{bmatrix} \begin{bmatrix} -B^{0} \Gamma_{m} \\ A^{0} \Gamma_{m} \end{bmatrix}^{*} \frac{|F_{u}|^{2}}{|F(\theta)|^{2}} d\omega.$$
(28)

We notice that $M(\theta) > 0$ whenever θ is in the stability domain for the coefficients of polynomials of degree *m*

$$\tilde{\delta} := \{\theta : F(z,\theta) = 0 \Rightarrow |z| < 1\} \subset \mathbb{R}^m.$$
⁽²⁹⁾

This follows from that we have assumed that L^o and F^o do not share common factors. The proof is omitted as it is almost identical to the proof of the well known result that the per sample information matrix for a Box–Jenkins model is positive definite if the input is informative.

When $\theta_k \in \hat{\mathscr{S}}$, we can thus write (25) as

$$\theta_{k+1} - \theta_o = (M(\theta_k) + D(\theta_k))^{-1} D(\theta_k) (\theta_k - \theta_o).$$
(30)

We now have the following theorem.

Theorem 4.2. Let F_u be defined in Assumption 2.2. Assume that there exist a compact and convex set $\mathscr{S} \subset \mathbb{R}^m$ which is a subset to the stability domain (29), i.e. $\mathscr{S} \subset \widetilde{\mathscr{S}}$, and that contains both θ_o and the initial point θ_1 . Then $\{\theta_k\}_{k=1}^{\infty}$, defined by (30), converges linearly to θ_o , i.e.

$$|\theta_{k+1} - \theta_o| \le \lambda^k |\theta_1 - \theta_o| \to 0, \quad k \to \infty, \tag{31}$$

for some $0 < \lambda < 1$.

Proof. Firstly, as discussed before Theorem 4.2, $M(\theta) > 0$ on \$ so the inverse of $M(\theta) + D(\theta)$ is well defined on this set. Now, since $1/|F(z, \theta)|$ is continuous on \$, the elements of $(M(\theta) + D(\theta))^{-1}D(\theta)$ are continuous on \$. Hence the maximum eigenvalue, $\lambda_{\max}(\theta)$ say, of $(M(\theta) + D(\theta))^{-1}D(\theta)$ is continuous on \$ (Horn & Johnson, 1985) and attains a maximum λ_{\max} on this compact set. Since

$$(M(\theta) + D(\theta))^{-1}D(\theta) = (I + M^{-1}(\theta)D(\theta))^{-1}M^{-1}(\theta)D(\theta) < I, \quad \forall \theta \in \mathcal{S},$$

it follows that $\lambda_{\max}(\theta) < 1 \quad \forall \theta \in \mathscr{S}$ and hence, using the compactness of \mathscr{S} , $\lambda_{\max} < 1$. This implies that (30) is a contraction and hence, since \mathscr{S} is convex, $\theta_k \in \mathscr{S} \Rightarrow \theta_{k+1} \in \mathscr{S}$. Furthermore, iterating (30) backwards *k* steps gives (31) with $\lambda =_{\max}$, which proves the result.

4.2. Asymptotic accuracy

Regarding the asymptotic accuracy, we have the following theorem.

Theorem 4.3. Let $\hat{\theta}_N$ be defined by (20). Assume that Assumptions 2.1–2.4 hold. Suppose that $\hat{\theta}_N \to \theta_o$, w.p.1 as $N \to \infty$. Then

$$\lim_{N \to \infty} N \mathbb{E} \left[(\hat{\theta}_N - \theta_o) (\hat{\theta}_N - \theta_o)^T \right] = \sigma_o^2 M_{CR}^{-1}.$$
(32)

Furthermore, $\sqrt{N}(\hat{\theta}_N - \theta_o) \sim AsN(0, \sigma_o^2 M_{CR}^{-1}).$

Proof. See Appendix C. ■

As M_{CR}/σ_o^2 is the Fisher information for Gaussian distributed noise, see Section 2.2, Theorem 4.3 shows that BJSM is asymptotically efficient for such noise.

We also remark that one reason for why Theorem 4.3 requires open-loop operation (see Assumption 2.2) is that then the estimates of system and disturbance dynamics become asymptotically uncorrelated. While we do not prove this, we conjecture that also the BJSM disturbance model estimate has the same asymptotic accuracy as the PEM.

174

4.3. Comparison with other methods

PEM. Even though restrictive, Theorem 4.2 indicates that there is a region in the parameter space where BJSM will converge. For PEM, as mentioned in the introduction, the cost function is only guaranteed to have no false local minima when only one system pole is estimated (Söderström, 1975a), and therefore unless this condition holds it cannot be guaranteed that this method ends up in the global minimum.

Comparing (32) with (3), we see that BJSM has the same asymptotic accuracy as PEM. In particular it is asymptotically efficient when the noise distribution is Gaussian.

The Steiglitz–McBride method. It is well known that the Steiglitz–McBride method is not asymptotically efficient. For output-error system and model, its asymptotic variance is given by (24) in Stoica and Söderström (1981):

$$P_{SM} \coloneqq \sigma_o^2 M^{-1}(\theta_o) (M(\theta_o) + D(\theta_o)) M^{-1}(\theta_o)$$

= $\sigma_o^2 M^{-1}(\theta_o) + \frac{1}{\sigma_o^2} \sigma_o^2 M^{-1}(\theta_o) D(\theta_o) \sigma_o^2 M^{-1}(\theta_o)$
 $\geq \sigma_o^2 M^{-1}(\theta_o).$

Theorem 4.3 also applies to the output-error case and shows that BJSM is superior to SM also in this case. This may appear paradoxical for the following reason: The main difference between BJSM and SM lies in the filtering that takes place in each iteration. In BJSM, all the signals used to compute the estimate at the next iteration are prefiltered with $A(q, \hat{\eta}_N)/F(q, \hat{\theta}_N^k)$, see (16), whereas SM uses $1/F(q, \hat{\theta}_N^k)$. Now, as the sample-size N grows the inverse noise model estimate $A(q, \hat{\eta}_N)$ will converge to the true inverse noise model, which in the output-error case is 1, i.e. the very one that is being used in SM. Thus, using a noisy estimate of the inverse noise model (BJSM) instead of the known true one (SM) gives better asymptotic accuracy. To understand the reason for this let us first observe that the reason for that SM is not asymptotically efficient can be traced to that the method does not take into account that the prefilter $1/F(q, \hat{\theta}_N^k)$ is a noisy estimate of the optimal prefilter $1/F(q, \theta_0)$. Very interestingly, in BJSM this uncertainty is accounted for by the error that is induced by the inverse noise model estimate $A(q, \hat{\eta}_N)$. This can be seen from the proof in Appendix C. The (normalized) estimation error in BJSM is given by (C.16) and consists of two terms $T_1(N)$ and $T_2(N)$ in (C.17)–(C.18). The second term $T_2(N)$ is due to the error in the (ARX) inverse noise model estimate and is what distinguishes the error in BJSM from that in SM. The analysis in Appendix C.4 shows that the size of this term corresponds exactly to the excess error that prevents SM from being asymptotically efficient. Moreover Appendix C.5 shows that $T_2(N)$ is perfectly negatively correlated with $T_1(N)$ and so this term, i.e. the error due to the ARX-model estimate, asymptotically cancels out the excess error that SM suffers from.

Instrumental variable methods.

For Box–Jenkins models, the multistep algorithm in Stoica and Söderström (1983) is limited to three steps: Firstly a IV-step to estimate L^o/F^o , secondly an ARMA noise model estimation step, and finally an optimal IV-step, i.e. (17), using the IV-estimate from the first step, and the noise model estimate from the second step, in the prefilter. Hence, as opposed to BJSM, for this method convergence is not an issue. For the multistep IV method, a critical issue for asymptotic efficiency of the estimate of the system dynamics L^o/F^o is that an efficient estimate of the ARMA noise model is obtained in the second step. As PEM is asymptotically efficient, and its cost function for ARMA models has no false local minima (Åström & Söderström, 1974), PEM can be used in this step (although for finite data there is always the risk of local minima). For BJSM, the asymptotic efficiency of the system dynamics L^o/F^o does not at all hinge on the properties of the estimate of the noise term $C^{o}(q)/D^{o}(q)$. In fact Step 4 in the BJSM algorithm can be omitted if no noise model is required.

For RIV there are no convergence results for Box–Jenkins models. However, practical experience, e.g. Young (2008), indicate good convergence properties. It follows from Stoica and Söderström (1983) that asymptotically, it is sufficient with the three steps outlined above for the multistep method. In Young (2008) it is suggested to use the IVARMA method (Young, 2006) in the ARMA noise model step. Simulations in Young (2008) indicate good performance of this method but asymptotic efficiency has not yet been established for the IVARMA method.

5. Practical aspects

5.1. ARX-model order selection

In the previous section we have shown that BJSM has some desirable asymptotic properties as the sample size grows, and when the order of the ARX-model grows with a rate satisfying Assumption 2.4. For a given data set, this theory does not help much in terms of how the order of the ARX-model should be selected. A guiding principle is that the order n should be sufficiently high that A(q) is able to model the inverse of the disturbance dynamics. By taking the order of the ARX-model which results in the Box–Jenkins model estimate having the smallest loss function (2), the ARX-model order can be optimized.

5.2. Initialization of model denominator

In the third step, the Steiglitz–McBride step, an initial estimate of $F^{o}(q)$ is required. When a convex subset of the stability set, as required by Theorem 4.2, is not known it is common, as in the SM method, to use F(q) = 1 as initial estimate.

5.3. Stopping criteria for Steiglitz-McBride iterations

In the Steiglitz–McBride step, the iterations can be terminated when some norm of the change of the parameter estimate is less than a pre-specified tolerance, or when a maximum number of iterations is reached.

5.4. Initialization of the ARMA estimate

The ARMA optimization routine can be initialized by a leastsquares estimate of *C* and *D* in the ARX-model $D(q)\hat{v}_t = C(q)\hat{\varepsilon}_t$, where \hat{v}_t is the output error residual in (13) and where $\hat{\varepsilon}_t$ is the residual of the ARX model.

6. Numerical studies on simulation and industrial data

Two methods are included in the study:

- (1) the prediction error method for Box–Jenkins model structures as implemented in the System Identification Toolbox of Matlab R2011b; this algorithm will be called PEM. Options *Tolerance* and *MaxIter* have been set to 1e - 4 and 50, respectively.
- (2) the BJSM where the ARX-model is estimated using the ARX command of the System Identification Toolbox of Matlab; the number of iterations in the Steiglitz–McBride step is fixed to 50 unless mentioned separately. The stopping criterion in Step (3) is that the relative change of the parameters in the *F*-polynomial is less than 0.0001, or that 50 iterations have been performed. In both examples, the Steiglitz–McBride iterations are initialized using F(q) = 1 as initial estimate.

The hardware used was a PC equipped with an Intel Core i3-2350M CPU running at 2.30 GHz, and having 4 GB RAM. Matlab



Fig. 1. The loss functions of the PEM method (x) and that of the BJSM method (o) for process (33), Example 1. The PEM method fails to converge in 21 simulations.

R2011b running under Windows 7 Home was used as software platform.

6.1. Example 1–a 4th order process with oscillation

The process is strongly oscillatory and given by

$$y_t = \frac{q^{-1} + 0.5q^{-2} - 2q^{-3} + q^{-4}}{1 - 1.5q^{-1} + 0.7q^{-2} + 0.3q^{-3} - 0.2q^{-4}}u_t + \frac{1 - 0.6q^{-1} + 0.4q^{-2}}{1 - 1.95q^{-1} + 0.9506q^{-2}}e_t.$$
(33)

The variance relative the variance of the noise free output is 20%. Models with correct orders are estimated for both methods; the order of the ARX model in the BJSM algorithm is 45. 200 simulations were performed with independent noise realizations.

For this process, BJSM seems to have converged to the global optimum in all 200 simulations as there are no big outliers in the loss function plot shown in Fig. 1, while from the same plot it is clear that PEM has failed to find the global optimum in 21 simulations as there are 21 outliers. When these 21 outliers are excluded, PEM and BJSM have almost identical accuracy. This is indicated by Fig. 1 in that the two methods have nearly the same loss functions, except for the outliers. The sum of the Mean-Squared Errors of the parameter estimates equals 0.0020 for both methods. The computation times are 6.894 s for PEM and 0.791 s for BISM.

6.2. Example 2-an industrial case study

This process is a crude unit at a European refinery which has been studied in Zhu (1998). The data set is obtained from an identification test of the main distillation column of the crude unit. The identified model was used in MPC control. Two open loop tests were carried out. The data set contains the inputs (MVs) and four outputs (CVs) of one test. Output 1 is the temperature difference of two trays; outputs 2–4 are product qualities measured by online analyzers. There are 7 inputs. Inputs 1–3 are temperature setpoints; input 4 is a flow setpoint; inputs 5 and 6 are flow ratio setpoints; input 7 is a flow measurement which is the measured disturbance in the MPC controller. The sampling time is 1 min. Interested researchers on system identification can contact the corresponding author to obtain the data.



Fig. 2. Inputs of one test of the main distillation column.



Fig. 3. Four outputs of the main distillation column.

The data were pre-processed using the following operations: (1) The means of all the input-output signals are removed. (2) The standard deviations of all the input-output signals are scaled to 1. (3) There are very large delays in the second output. This output signal is forward shifted by 80 samples, which means that 80 delays are removed (see Figs. 2 and 3).

6.2.1. Process delays

According to process knowledge there are considerable delays between process inputs and outputs. So different delays are used in model estimation using ARMAX models and the delays are determined so that the model simulation errors are minimal. To keep it simple the same delay is used for each output. The determined delays for the four outputs are: 4, 20, 20 and 20 samples.

Note that 80 delays have been removed for output 2, which means that the total delay for output 2 is 100 samples.

6.2.2. The BJSM algorithm

For BJSM, multi-input single-output (MISO) models are used. This means that for each output a common denominator polynomial is used for all the inputs. The high order ARX model is also of MISO type with diagonal denominator. The order of the high order ARX model is set to 100, and the number of iterations is fixed to 100.



Fig. 4. FOE of the BJSM models in Example 6.

Table 1

Identification results of the two methods using the crude unit data in Example 6. SO denotes orders that provided stable models, BO denotes the best model order according to FOE.

Method, output	SO	BO	Best FOE
PEM, y1	[1, 2]	2	0.1608
PEM, y2	-	-	-
PEM, y3	[1]	1	0.3147
PEM, y4	[1]	1	0.7846
BJSM, y1	[1:20]	17	0.1467
BJSM, y2	[1:20]	8	0.2196
BJSM, y3	[1:20]	18	0.0729
BJSM, y4	[1:20]	2	0.2296

For order selection, models from order 1 to order 20 (the same order for all polynomials in the model, including the noise model) are estimated and the final output error (FOE) criterion (Zhu, 2001) is used to determine the best order. Given a MISO process model in a prediction error structure

$$y_t = \hat{G}(q)u_t + \hat{H}(q)\hat{\varepsilon}_t, \tag{34}$$

where u_t is now a vector of input signals. The FOE is defined as

$$FOE = \frac{N+d}{N-d} \sum_{t=1}^{N} (y_t - \hat{G}(q)u_t)^2,$$
(35)

where *N* is the number of testing data and *d* is the number of model parameters. The test is evaluated on estimation data but notice that since it is not the prediction error criterion, FOE does not decrease monotonically with the model order. It is argued that this criterion is more control relevant than FPE (final prediction error) criterion; see Chapter 5 in Zhu (2001) for motivations and comparisons.

No numerical problem occurred during the estimation for BJSM and the FOE's are plotted in Fig. 4.

6.2.3. The PEM method

For MISO PEM models estimated using the MATLAB command bj, the obtained models were unstable for most of the model orders tried.¹ This happened even when the parameter 'Focus', was set to 'Stability'. A difference between PEM and BJSM is that different denominator polynomials are used for the different inputs in PEM, meaning that more parameters are estimated for the same order.

Table 1 summarizes the performance of the two methods.

7. Conclusions

An algorithm for identification of SISO Box–Jenkins models has been presented. It consists of ARX-model estimation followed by Steiglitz–McBride iterations, and an ARMA-estimation step for the noise model. Theoretical justification is provided in terms of results on convergence and asymptotic efficiency for data collected in open loop. These results show that the BJSM also has better accuracy than the Steiglitz–McBride method for outputerror models.

In the two numerical examples, the BJSM method proved more robust than PEM. However, being based on a (high-order) ARXmodeling step, BJSM may be expected to have limitations for short data records. For slowly decaying impulse responses, very high ARX-orders may be required, increasing computational burden. Before concluding, we remark that while derived for SISO BJ models, it is very easy to extend the method to MISO BJ models (we actually used this in Section 6.2).

Acknowledgments

The authors are very grateful for the extraordinary efforts of the three reviewers. Their extensive comments and suggestions have helped to significantly improve the quality of the paper. We also thank the Editor and the Associate Editor for their efforts in guiding this paper through the reviewing process.

Notation used in appendices

We will consider vector valued complex functions as row vectors and the inner product of two such functions $f, g : \mathbb{C} \to \mathbb{C}^{1 \times m}$ is defined as $\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(e^{j\omega})g^*(e^{j\omega})d\omega$ where g^* denotes the complex conjugate transpose of g. When f and g are matrix-valued functions, we will still use the notation $\langle f, g \rangle$ to denote $\frac{1}{2\pi} \int_{-\pi}^{\pi} f(e^{j\omega})g^*(e^{j\omega})d\omega$.

The \mathcal{L}_2 -norm of $f: \mathbb{C} \to \mathbb{C}^{n \times m}$ is given by $||f||_2 = \sqrt{\operatorname{Tr} \langle f, f \rangle}$. The space $\mathcal{L}_2^{n \times m}$ consists of all functions $f: \mathbb{C} \to \mathbb{C}^{n \times m}$ such that $||f||_2 < \infty$ and when n = 1, the notation is simplified to \mathcal{L}_2^m . The subspace of \mathcal{L}_2^m spanned by the rows of $\Psi \in \mathcal{L}_2^{n \times m}$ is denoted \mathcal{S}_{Ψ} .

Appendix A. Proof of Lemma 2.1

The result follows from Theorem D.1. Next, we verify the conditions of that theorem. Assumption 2.1 and the finite dimensionality of G^o and H^o implies that

$$\max(|a_k|, |b_k|) \le C\rho^{\kappa},\tag{A.1}$$

for some $C < \infty$ and $0 < \rho < 1$. This implies that Condition S1 in Appendix D holds. Furthermore, the bound (A.1) implies the inequality in (11) for some $\tilde{C} < \infty$. Assumption 2.3 clearly implies Condition S2 (for any $p \le 5$). Assumption 2.4 implies Conditions D1 and D3. Thus all conditions in Theorem D.1 have been verified and the result in the lemma follows from this theorem.

Appendix B. Proof of Theorem 4.1

We have

$$R^{m}(N, \hat{\eta}_{N}, \theta) = \frac{1}{N} \sum_{t=m+1}^{N} \varphi_{t}^{m}(\hat{\eta}_{N}, \theta)(\varphi_{t}(\hat{\eta}_{N}, \theta))^{T}$$
$$= \frac{1}{N} \sum_{t=m+1}^{N} \left(\frac{A(q, \hat{\eta}_{N})}{F(q, \theta)} \varphi_{t}^{m}\right) \left(\frac{A(q, \hat{\eta}_{N})}{F(q, \theta)} \varphi_{t}^{m}\right)^{T}$$
$$= R^{m}(N, \eta_{0}, \theta)$$

 $^{^{1}}$ As for BJSM the order was the same for all polynomials.

$$+\frac{1}{N}\sum_{t=m+1}^{N}\left(\frac{A(q,\,\hat{\eta}_{N})-A^{o}(q)}{F(q,\,\theta)}\varphi_{t}^{m}\right)\left(\frac{A^{o}(q)}{F(q,\,\theta)}\varphi_{t}^{m}\right)^{T}$$
$$+\frac{1}{N}\sum_{t=m+1}^{N}\left(\frac{A^{o}(q)}{F(q,\,\theta)}\varphi_{t}^{m}\right)\left(\frac{A(q,\,\hat{\eta}_{N})-A^{o}(q)}{F(q,\,\theta)}\varphi_{t}^{m}\right)^{T}.$$
(B.1)

Theorem 2B.1 in Ljung (1999) gives

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| R^{m}(N, \eta_{o}, \theta) - \mathbb{E} \left[\varphi_{t}^{m}(\eta_{o}, \theta) \left(\varphi_{t}^{m}(\eta_{o}, \theta) \right)^{T} \right] \right\|$$

$$\rightarrow 0, \quad \text{w.p. 1 as } N \rightarrow \infty.$$
(B.2)

Lemma 2.1 together with that both $\{e_t\}$ and $\{u_t\}$ are uniformly bounded (see Assumptions 2.2 and 2.3) imply

$$\|(A(q, \hat{\eta}_N) - A^o(q))\varphi_t^m\|_1 = O(m(N)), \tag{B.3}$$

with m(N) given in Lemma 2.1. Since $\{1/F(q, \theta) : \theta \in D_{\mathcal{M}}\}$ is uniformly stable, it follows that

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| \frac{1}{N} \sum_{t=m+1}^{N} \left(\frac{A(q, \hat{\eta}_{N}) - A^{o}(q)}{F(q, \theta)} \varphi_{t}^{m} \right) \left(\frac{A^{o}(q)}{F(q, \theta)} \varphi_{t}^{m} \right)^{T} \right\|$$

$$\rightarrow 0, \quad \text{w.p. 1 as } N \rightarrow \infty.$$
(B.4)

Now (B.1)-(B.4) gives

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| R^{m}(N, \hat{\eta}_{N}, \theta) - \mathbb{E} \left[\varphi_{t}^{m}(\eta_{o}, \theta) \left(\varphi_{t}^{m}(\eta_{o}, \theta) \right)^{T} \right] \right\|$$

$$\rightarrow 0, \quad \text{w.p. 1 as } N \rightarrow \infty.$$
(B.5)

Analogously it follows that

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| r^{m}(N, \hat{\eta}_{N}, \theta) - \mathbb{E} \left[\varphi_{t}^{m}(\eta_{o}, \theta) y_{t}(\eta_{o}, \theta) \right] \right\|$$

$$\to 0, \quad \text{w.p. 1 as } N \to \infty.$$
(B.6)

Combining (B.5)–(B.6) gives the first claim in the theorem.² Next, (22) is the same as (8) in Stoica and Söderström (1981). Theorem 1 in this reference then applies and since we assume the model order to be equal to the system order, there is according to this theorem only one solution.

Appendix C. Proof of Theorem 4.3

Denote the right-hand side of (20) by $Q(N, \hat{\eta}_N, \hat{\theta}_N)$. A first order Taylor expansion of $Q(N, \hat{\eta}_N, \hat{\theta}_N)$ around θ_o gives (partial derivatives are interpreted as row vectors)

$$0 = Q(N, \hat{\eta}_N, \theta_o) + \frac{\partial}{\partial \theta} Q(N, \hat{\eta}_N, \theta) \Big|_{\theta = \xi_N} (\hat{\theta}_N - \theta_o),$$
(C.1)

for some ξ_N between θ_o and $\hat{\theta}_N$. Thus

$$\sqrt{N}(\hat{\theta}_{N} - \theta_{o})$$

$$= -\left[\frac{\partial}{\partial\theta}Q(N, \hat{\eta}_{N}, \theta)\Big|_{\theta = \xi_{N}}\right]^{-1} \sqrt{N}Q(N, \hat{\eta}_{N}, \theta_{o}).$$
(C.2)

In the following sections we will analyze $\frac{\partial}{\partial \theta} Q(N, \hat{\eta}_N, \theta) \Big|_{\theta = \xi_N}$ and $\sqrt{N}Q(N, \hat{\eta}_N, \theta_o)$.

C.1.
$$\frac{\partial}{\partial \theta} Q(N, \hat{\eta}_N, \theta)|_{\theta = \xi_N}$$

We have

$$\frac{\partial}{\partial \theta} Q(N, \eta, \theta) = Q_1'(N, \eta, \theta) + Q_2'(N, \eta, \theta) + Q_3'(N, \eta, \theta), \quad (C.3)$$

where

$$\begin{aligned} Q_1'(N,\eta,\theta) &\coloneqq -\frac{1}{N} \sum_{t=m+1}^N \varphi_t^m(\eta,\theta) [\varphi_t^m(\eta,\theta)]^T, \\ Q_2'(N,\eta,\theta) &\coloneqq \frac{1}{N} \sum_{t=m+1}^N \frac{\partial}{\partial \theta} \varphi_t^m(\eta,\theta) (y_t(\eta,\theta) - [\varphi_t^m(\eta,\theta)]^T \theta), \\ Q_3'(N,\eta,\theta) &\coloneqq \frac{1}{N} \sum_{t=m+1}^N \varphi_t^m(\eta,\theta) \left(\frac{\partial}{\partial \theta} y_t(\eta,\theta) - \theta_o^T \frac{\partial}{\partial \theta} \varphi_t^m(\eta,\theta)\right). \end{aligned}$$

$$(C.4)$$

We observe that $Q'_1(N, \eta, \theta) = -R^m(N, \eta, \theta)$, and thus (B.1)–(B.5) in the proof of Theorem 4.1 imply

_ ...

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| Q_1'(N, \hat{\eta}_N, \theta) + \mathbb{E} \left[\varphi_t^m(\eta_o, \theta) \left(\varphi_t^m(\eta_o, \theta) \right)^T \right] \right\|$$

 $\to 0, \quad \text{w.p. 1 as } N \to \infty.$ (C.5)

Similar calculations give

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| Q_{2}'(N, \hat{\eta}_{N}, \theta) - E\left[\frac{\partial}{\partial \theta} \varphi_{t}^{m}(\eta_{o}, \theta)(y_{t}(\eta_{o}, \theta) - \varphi_{t}^{m}(\eta_{o}, \theta)) \right] \right\|$$

$$\rightarrow 0, \quad \text{w.p. 1 as } N \rightarrow \infty, \qquad (C.6)$$

$$\sup_{\theta \in D_{\mathcal{M}}} \left\| Q_{3}'(N, \hat{\eta}_{N}, \theta) - - E\left[\varphi_{t}^{m}(\eta_{o}, \theta) \left(\frac{\partial}{\partial \theta} y_{t}(\eta_{o}, \theta) - \theta^{T} \frac{\partial}{\partial \theta} \varphi_{t}^{m}(\eta_{o}, \theta) \right) \right] \right\|$$

$$\rightarrow 0, \quad \text{w.p. 1 as } N \rightarrow \infty.$$
 (C.7)

Since by assumption $\hat{\theta}_N \to \theta_o$ w.p.1 as $N \to \infty$, $\xi_N \to \theta_o$ w.p. 1. In view of this, it follows, c.f. (9A.29) in Ljung (1999), from (C.3), and (C.5)–(C.7) that as $N \to \infty$ (below $\frac{\partial}{\partial \theta} \varphi_t^m(\eta_o, \theta_o)$, $\frac{\partial}{\partial \theta} y_t(\eta_o, \theta_o)$ is shorthand for $\frac{\partial}{\partial \theta} \varphi_t^m(\eta_o, \theta)|_{\theta=\theta_o}$, $\frac{\partial}{\partial \theta} y_t(\eta_o, \theta)|_{\theta=\theta_o}$) w.p. 1

$$\frac{\partial}{\partial \theta} Q(N, \hat{\eta}_{N}, \theta) \Big|_{\theta = \xi_{N}} \rightarrow \\
- E \left[\varphi_{t}^{m}(\eta_{o}, \theta_{o}) \left(\varphi_{t}^{m}(\eta_{o}, \theta_{o}) \right)^{T} \right] \\
+ E \left[\frac{\partial}{\partial \theta} \varphi_{t}^{m}(\eta_{o}, \theta_{o}) (y_{t}(\eta_{o}, \theta_{o}) - (\varphi_{t}^{m}(\eta_{o}, \theta_{o}))^{T} \theta_{o}) \right] \\
+ E \left[\varphi_{t}^{m}(\eta_{o}, \theta_{o}) \left(\frac{\partial}{\partial \theta} y_{t}(\eta_{o}, \theta_{o}) - \theta_{o}^{T} \frac{\partial}{\partial \theta} \varphi_{t}^{m}(\eta_{o}, \theta_{o}) \right) \right]. \quad (C.8)$$

Observing that (15) implies

$$y_t(\eta_o, \theta_o) - [\varphi_t^m(\eta_o, \theta_o)]^T \theta_o = e_t,$$
(C.9)

and that e_t is independent of φ_t^m gives

$$E\left[\frac{\partial}{\partial\theta}\varphi_t^m(\eta_o,\theta_o) \left(y_t(\eta_o,\theta_o) - \left[\varphi_t^m(\eta_o,\theta_o)\right]^T\theta_o\right)\right]$$
$$= E\left[\frac{\partial}{\partial\theta}\varphi_t^m(\eta_o,\theta_o) e_t\right] = 0.$$
(C.10)

² Actually, we have proved a stronger result taking the supremum over θ .

Furthermore, with $W_m(q) = \begin{bmatrix} \Gamma_m^T(q) & \mathbf{0}_{1 \times m} \end{bmatrix}^T$,

$$\begin{aligned} \frac{\partial}{\partial \theta} y_{t}(\eta_{o}, \theta_{o}) &- \theta_{o}^{T} \frac{\partial}{\partial \theta} \varphi_{t}^{m}(\eta_{o}, \theta_{o}) = -\frac{1}{F^{o}(q)} \begin{pmatrix} W_{m}^{T}(q) y_{t}(\eta_{o}, \theta_{o}) \\ W_{m}^{T}(q) y_{t}(\eta_{o}, \theta_{o}) \\ & \\ & \\ - \theta_{o}^{T} \begin{bmatrix} -y_{t-1-1}(\eta_{o}, \theta_{o}) & \cdots & -y_{t-1-m}(\eta_{o}, \theta_{o}) & 0_{1 \times m} \\ \vdots & \ddots & \vdots & 0_{(m-2) \times m} \\ -y_{t-1-m}(\eta_{o}, \theta_{o}) & \cdots & -y_{t-1-2m}(\eta_{o}, \theta_{o}) & 0_{1 \times m} \\ \vdots & \ddots & \vdots & 0_{(m-2) \times m} \\ \vdots & \ddots & \vdots & 0_{(m-2) \times m} \\ -u_{t-1-m}(\eta_{o}, \theta_{o}) & \cdots & -u_{t-1-2m}(\eta_{o}, \theta_{o}) & 0_{1 \times m} \\ \end{bmatrix} \\ & = -\frac{1}{F^{o}(q)} W_{m}^{T}(q) e_{t}. \end{aligned}$$
(C.11)

From (C.10) and by noticing that

$$\varphi_t^m = \begin{bmatrix} -G^0 \Gamma_m \\ \Gamma_m \end{bmatrix} u_t + \begin{bmatrix} -\frac{1}{A^0} \Gamma_m \\ 0_{m \times 1} \end{bmatrix} e_t$$

with $\{u_t\}$ independent of $\{e_t\}$ (by assumption) it follows

$$E\left[\varphi_t^m(\eta_o, \theta_o) \left(\frac{\partial}{\partial \theta} y_t(\eta_o, \theta_o) - \theta_o^T \frac{\partial}{\partial \theta} \varphi_t^m(\eta_o, \theta_o)\right)\right]$$

= $-E\left[\frac{A^o(q)}{F^o(q)} \varphi_t^m \frac{1}{F^o(q)} W_m^T(q) e_t\right] = D(\theta_o),$ (C.12)

with $D(\theta)$ defined in (26). Inserting (23), (C.10) and (C.12) in (C.8), and using (27), gives

$$\lim_{N \to \infty} \frac{\partial}{\partial \theta} Q(N, \hat{\eta}_N, \theta) \Big|_{\theta = \xi_N} = -\tilde{R}(\theta_o) + D(\theta_o)$$
$$= -M(\theta_o) \quad \text{w.p. 1.}$$
(C.13)

C.2. $\sqrt{N}Q(N, \hat{\eta}_N, \theta_0)$

Returning to (C.2), we will now establish the asymptotic distribution and variance of $\sqrt{N}Q(N, \hat{\eta}_N, \theta_o)$. We have

$$Q(N, \hat{\eta}_{N}, \theta_{0}) = \frac{1}{N} \sum_{t=m+1}^{N} \varphi_{t}^{m}(\hat{\eta}_{N}, \theta_{0}) e_{t}(\hat{\eta}_{N}, \theta_{0}, F^{0})$$
$$= \frac{1}{N} \sum_{t=m+1}^{N} \frac{A(q, \hat{\eta}_{N})}{F^{0}(q)} \varphi_{t}^{m} \frac{A(q, \hat{\eta}_{N})}{A^{0}(q)} e_{t}.$$
(C.14)

It is straightforward to show that the mean-squared error between $\sqrt{N}Q(N, \hat{\eta}_N, \theta_o)$ and

$$\frac{1}{\sqrt{N}}\sum_{t=m+1}^{N}\frac{A^{o}(q)}{F^{o}(q)}\varphi_{t}^{m}\frac{A(q,\,\hat{\eta}_{N})}{A^{o}(q)}e_{t}$$
(C.15)

tends to zero as $N \rightarrow \infty$. Thus these two quantities have the same asymptotic distribution and the same asymptotic covariance. We will thus proceed and study the asymptotic properties of (C.15) instead of those of (C.14).

We split (C.15) into two terms which will be analyzed separately:

$$\frac{1}{\sqrt{N}}\sum_{t=m+1}^{N}\frac{A^{o}(q)}{F^{o}(q)}\varphi_{t}^{m}\frac{A(q,\,\hat{\eta}_{N})}{A^{o}(q)}e_{t}=T_{1}(N)+T_{2}(N), \tag{C.16}$$

where

$$T_{1}(N) := \frac{1}{\sqrt{N}} \sum_{t=m+1}^{N} \left(\frac{A^{o}(q)}{F^{o}(q)} \varphi_{t}^{m} \right) e_{t},$$
(C.17)
$$T_{2}(N) := \frac{1}{\sqrt{N}} \sum_{t=m+1}^{N} \left(\frac{A^{o}(q)}{F^{o}(q)} \varphi_{t}^{m} \right) \left(\frac{A(q, \hat{\eta}_{N}) - A^{o}(q)}{A^{o}(q)} e_{t} \right)$$
$$= \frac{1}{N} \sum_{t=m+1}^{N} \frac{A^{o}(q)}{F^{o}(q)} \varphi_{t}^{m} W_{n(N)}^{T}(q) \frac{1}{A^{o}(q)} e_{t} \sqrt{N} (\hat{\eta}_{N}^{n(N)} - \bar{\eta}^{N}).$$
(C.18)

It is also straightforward to show that the mean-squared error between $T_2(N)$ and

$$\tilde{T}_2(N) := Z^{n(N)} \sqrt{N} (\hat{\eta}_N^{n(N)} - \bar{\eta}^{n(N)}),$$

where

$$Z^{n} = \mathbb{E} \begin{bmatrix} \frac{A^{o}(q)}{F^{o}(q)} \varphi_{t}^{m} W_{n}^{T}(q) \frac{1}{A^{o}(q)} e_{t} \end{bmatrix}$$
$$= -\begin{bmatrix} \mathbb{E} \begin{bmatrix} \frac{1}{F^{o}} \Gamma_{m} e_{t} & \Gamma_{n}^{T} \frac{1}{A^{o}} e_{t} \end{bmatrix} & \mathbf{0}_{m \times n} \\ \mathbf{0}_{m \times n} & \mathbf{0}_{m \times n} \end{bmatrix},$$
(C.19)

tends to zero as $N \to \infty$ and we will analyze $\tilde{T}_2(N)$ rather than $T_2(N)$.

C.3. Asymptotic variance of $T_1(N)$

Since e_t is independent of φ_t it follows that

$$\lim_{N \to \infty} \mathbb{E}\left[T_1(N)T_1^T(N)\right] = \sigma_o^2 \tilde{R}(\theta_o).$$
(C.20)

C.4. Asymptotic variance of $\tilde{T}_2(N)$

In the proof of Lemma 2.1 Conditions S1, S2 (for $p \leq 5$), and D1 were verified under Assumptions 2.1–2.4. Furthermore D4 follows directly from Assumption 2.4. The bound in (11) implies that $\sqrt{N}d(N) \rightarrow 0$, $N \rightarrow \infty$ under Assumption 2.4, which establishes D2. This means that all conditions in Theorem D.2 are valid. Noticing that { $||Z^n||_2$ } is a bounded sequence, this theorem gives

$$\lim_{N \to \infty} \mathbb{E}\left[\tilde{T}_2(N)\tilde{T}_2^T(N)\right]$$
$$= \sigma_o^2 \lim_{N \to \infty} Z^{n(N)} \left[\bar{R}^{n(N)}\right]^{-1} (Z^{n(N)})^T, \qquad (C.21)$$

where $\bar{R}^n = E\left[\varphi_t^n(\varphi_t^n)^T\right]$ provided the right-hand side limit exists. This will be shown next.

We start the analysis of (C.21) by considering the inverse in the middle. Using the open loop assumption, we obtain

$$\bar{R}^{n} = \mathbb{E}\left[\begin{bmatrix}-G^{o}\Gamma_{n}u_{t}\\\Gamma_{n}u_{t}\end{bmatrix}\begin{bmatrix}-G^{o}\Gamma_{n}u_{t}\end{bmatrix}^{T}\right] + \mathbb{E}\left[\begin{bmatrix}\frac{1}{A^{o}}\Gamma_{n}e_{t}\\0_{n\times 1}\end{bmatrix}\begin{bmatrix}\frac{1}{A^{o}}\Gamma_{n}e_{t}\\0_{n\times 1}\end{bmatrix}^{T}\right] = \left\langle \begin{bmatrix}-G^{o}F_{u}\Gamma_{n} & \frac{\sigma_{o}}{A^{o}}\Gamma_{n}\\F_{u}\Gamma_{n} & 0_{n\times 1}\end{bmatrix}, \begin{bmatrix}-G^{o}F_{u}\Gamma_{n} & \frac{\sigma_{o}}{A^{o}}\Gamma_{n}\\F_{u}\Gamma_{n} & 0_{n\times 1}\end{bmatrix}\right\rangle.$$
(C.22)

For (C.19) we have

$$Z^{n} == - \left\langle \begin{bmatrix} 0_{m \times 1} & \frac{1}{F^{o}} \Gamma_{m} \\ 0_{m \times 1} & 0_{m \times 1} \end{bmatrix}, \begin{bmatrix} -G^{o} F_{u} \Gamma_{n} & \frac{\sigma_{o}}{A^{o}} \Gamma_{n} \\ F_{u} \Gamma_{n} & 0_{n \times 1} \end{bmatrix} \right\rangle.$$
(C.23)

Let now \mathscr{S}_n be the subspace in $\mathscr{L}_2^{1\times 2}$ spanned by the rows of

$$\begin{bmatrix} -G^{o}F_{u}\Gamma_{n} & \frac{\sigma_{o}}{A^{o}}\Gamma_{n} \\ F_{u}\Gamma_{n} & 0_{n\times 1} \end{bmatrix}.$$
 (C.24)

Then Lemma E.3, with $F_1 = -G^o$, $F_2 = \sigma_o/A^o$ (which has an exponentially stable inverse by Assumption 2.1), $F_3 = F_u$ and γ being an arbitrary row of $\begin{bmatrix} 0_{m \times 1} & \frac{1}{F^o} \Gamma_m \\ 0_{m \times 1} & 0_{m \times 1} \end{bmatrix}$, gives that

$$\|\gamma - P_{\delta_n}[\gamma]\|_2 \le \tilde{C}\tilde{\lambda}^n, \quad \text{for some } \tilde{C} < \infty, \ \tilde{\lambda} < 1.$$

But then Lemma E.2 immediately gives

$$\lim_{n \to \infty} Z^n \left[\bar{R}^n \right]^{-1} (Z^n)^T = \left\langle \begin{bmatrix} \mathbf{0}_{m \times 1} & \frac{1}{F^o} \Gamma_m \\ \mathbf{0}_{m \times 1} & \mathbf{0}_{m \times 1} \end{bmatrix}, \begin{bmatrix} \mathbf{0}_{m \times 1} & \frac{1}{F^o} \Gamma_m \\ \mathbf{0}_{m \times 1} & \mathbf{0}_{m \times 1} \end{bmatrix} \right\rangle$$
$$= D(\theta_o), \tag{C.25}$$

and hence, using (C.21), we have that as $n \to \infty$,

$$\lim_{N \to \infty} \mathbb{E}\left[\tilde{T}_2(N)\tilde{T}_2^T(N)\right] = \sigma_o^2 D(\theta_o).$$
(C.26)

C.5. Cross-correlation between $T_1(N)$ and $\tilde{T}_2(N)$

We now turn to the cross-correlation between $T_1(N)$ and $\tilde{T}_2(N)$. Firstly we observe that the mean-squared error between \tilde{T}_2 and

$$\bar{T}_2 := Z^{n(N)} \left[\bar{R}^{n(N)} \right]^{-1} \frac{1}{\sqrt{N}} \sum_{t=n(N)+1}^N \varphi_t^n e_t$$
(C.27)

tends to zero as $N \rightarrow \infty$. We will thus consider

$$E\left[\bar{T}_{2}(N)T_{1}^{T}(N)\right] = NZ^{n(N)}\left[\bar{R}^{n(N)}\right]^{-1}$$

$$E\left[\frac{1}{N}\sum_{t=n(N)+1}^{N}\varphi_{t}^{n}e_{t}\left(\frac{1}{N}\sum_{s=m+1}^{N}\frac{A^{o}}{F^{o}}\varphi_{s}e_{s}\right)^{T}\right].$$
(C.28)

But

$$NE\left[\frac{1}{N}\sum_{t=n(N)+1}^{N}\varphi_{t}^{n}e_{t}\left(\frac{1}{N}\sum_{s=m+1}^{N}\frac{A^{o}}{F^{o}}\varphi_{s}e_{s}\right)^{T}\right]$$
$$=\sigma_{o}^{2}E\left[\begin{bmatrix}-G^{o}\Gamma_{n}\\\Gamma_{n}\end{bmatrix}u_{t}\begin{bmatrix}-\frac{A^{o}}{F^{o}}G^{o}\Gamma_{m}\\\frac{A^{o}}{F^{o}}\Gamma_{m}\end{bmatrix}^{T}u_{t}\right]$$
$$+\sigma_{o}^{2}E\left[\begin{bmatrix}-\frac{1}{A^{o}}\Gamma_{n}\\0_{n\times 1}\end{bmatrix}e_{t}\begin{bmatrix}-\frac{1}{F^{o}}G^{o}\Gamma_{m}\\0_{m\times 1}\end{bmatrix}^{T}e_{t}\right]$$
$$=\sigma_{o}^{2}X^{n}-\sigma_{o}^{2}(Z^{n})^{T},$$
(C.29)

where

$$X^{n} = \mathbb{E}\left[\begin{bmatrix}-G^{o}\Gamma_{n}u_{t}\\\Gamma_{n}u_{t}\end{bmatrix}\begin{bmatrix}-\frac{A^{o}}{F^{o}}G^{o}\Gamma_{m}u_{t}\\\frac{A^{o}}{F^{o}}\Gamma_{m}u_{t}\end{bmatrix}^{T}\right]$$
$$= \left\langle \begin{bmatrix}-F_{u}G^{o}\Gamma_{n}\\F_{u}\Gamma_{n}\end{bmatrix}, \begin{bmatrix}-\frac{A^{o}}{F^{o}}F_{u}G^{o}\Gamma_{m}\\\frac{A^{o}}{F^{o}}F_{u}\Gamma_{m}\end{bmatrix}\right\rangle$$
$$= \left\langle \begin{bmatrix}-F_{u}G^{o}\Gamma_{n}&\frac{1}{A^{o}}\Gamma_{n}\\F_{u}\Gamma_{n}&0_{n\times1}\end{bmatrix}, \begin{bmatrix}-\frac{A^{o}}{F^{o}}F_{u}G^{o}\Gamma_{m}&0_{m\times1}\\\frac{A^{o}}{F^{o}}F_{u}\Gamma_{m}&0_{m\times1}\end{bmatrix}\right\rangle, \quad (C.30)$$

so that (C.27) can be written as

$$E\left[\bar{T}_{2}(N)T_{1}^{T}(N)\right] = \sigma_{o}^{2}Z^{n(N)}\left[\bar{R}^{n(N)}\right]^{-1}(X^{n(N)} - (Z^{n(N)})^{T}).$$
(C.31)

As in Appendix C.4 denote by \mathscr{S}_n the span of the rows of (C.24). Then Lemma E.3, with $F_1 = -G^o$, $F_2 = \sigma_o/A^o$ (which has an exponentially stable inverse by Assumption 2.1), $F_3 = F_u$ and γ being an arbitrary row of

$$\begin{bmatrix} 0_{m\times 1} & \frac{1}{F^o}\Gamma_m \\ 0_{m\times 1} & 0_{m\times 1} \end{bmatrix},$$

gives that

$$\|\gamma - P_{\delta_n}[\gamma]\|_2 \leq \tilde{C}\tilde{\lambda}^n$$
, for some $\tilde{C} < \infty$, $\tilde{\lambda} < 1$.

But then Lemma E.2 immediately gives

$$\lim_{n \to \infty} Z^{n} \left[\bar{R}^{n} \right]^{-1} X^{n}$$

$$= \left\langle \begin{bmatrix} 0_{m \times 1} & \frac{1}{F^{o}} \Gamma_{m} \\ 0_{m \times 1} & 0_{m \times 1} \end{bmatrix}, \begin{bmatrix} -\frac{A^{o}}{F^{o}} F_{u} G^{o} \Gamma_{m} & 0_{m \times 1} \\ \frac{A^{o}}{F^{o}} F_{u} \Gamma_{m} & 0_{m \times 1} \end{bmatrix} \right\rangle$$

$$= 0_{2m \times 2m}.$$
(C.32)

Using (C.26) and (C.32) in (C.31) gives

$$\lim_{N \to \infty} \mathbb{E}\left[\bar{T}_2(N)T_1^T(N)\right] = -\sigma_o^2 D(\theta_o).$$
(C.33)

C.6. Asymptotic normality

Consider

$$\begin{bmatrix} T_1(N) \\ \bar{T}_2(N) \end{bmatrix} = \frac{1}{\sqrt{N}} \sum_{t=m+1}^N \zeta_t(N) e_t,$$
(C.34)

where

$$\zeta_{t}(N) = \begin{bmatrix} \frac{A^{o}(q)}{F^{o}(q)}\varphi_{t}^{m} \\ Z^{n(N)} \left[\bar{R}^{n(N)}\right]^{-1}\varphi_{t}^{n(N)} \end{bmatrix}.$$
(C.35)

As already observed at the beginning of Appendix C.4, all conditions of Theorem D.3 are satisfied. Asymptotic normality of (C.34) now follows as in the proof of Theorem D.3 (the reader is referred to Ljung and Wahlberg (1992) for details). The asymptotic covariance matrix is obtained from (C.20), (C.26) and (C.33). In summary

$$\begin{bmatrix} T_1(N) \\ \bar{T}_2(N) \end{bmatrix} \sim AsN\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma_o^2 \begin{bmatrix} \tilde{R}(\theta_o) & -D(\theta_o) \\ -D(\theta_o) & D(\theta_o) \end{bmatrix} \right).$$
(C.36)

C.7. Summing up

Using that (C.15) is equivalent to (C.14), the split (C.16), and (C.20), (C.26) and (C.33), we have that

$$\lim_{N \to \infty} NE \left[Q(N, \hat{\eta}_N, \theta_o) Q^T(N, \hat{\eta}_N, \theta_o) \right] = \sigma_o^2$$

($\tilde{R}(\theta_o) - D(\theta_o) - D(\theta_o) + D(\theta_o)$) = $\sigma_o^2 M(\theta_o)$. (C.37)
Using (C.13) and (C.36) in (C.2) gives

$$\lim_{N \to \infty} NE\left[(\hat{\theta}_N - \theta_o) (\hat{\theta}_N - \theta_o)^T \right] = \sigma_o^2 M^{-1}(\theta_o).$$
(C.38)

which in view of that $M(\theta_o) = M_{CR}$, cf. (4) with (28), proves (32). The asymptotic normality of $\sqrt{N}(\hat{\theta}_N - \theta_o)$ follows from (C.2) and (C.36), together with that $\sqrt{NQ}(N, \hat{\eta}_N, \theta_o)$ has the same asymptotic distribution as $T_1(N) + \bar{T}_2(N)$.

Appendix D. Results from Ljung and Wahlberg (1992)

We first state some conditions that will be used.

Condition S1. The filters A^o and B^o satisfy $\sum_{k=0}^{\infty} k^{\frac{1}{2}} |a_k^o| < \infty$, $\sum_{k=0}^{\infty} k^{\frac{1}{2}} |b_k^o| < \infty$. The term b_0^o is always assumed to be 0. **Condition S2.** $\{e_t\}$ is a stochastic process such that

$$\mathbf{E}[e_t|\mathcal{F}_{t-1}] = \mathbf{0}, \qquad \mathbf{E}[e_t^2|\mathcal{F}_{t-1}] \qquad = \sigma_o^2, \quad \mathbf{E}[e_t^{2p}] \le C < \infty,$$

for some *p* to be specified. Here \mathcal{F}_{t-1} is the σ -algebra generated by $\{e_s, u_s, s \leq t-1\}$.

Condition D1. $n(N) \to \infty, N \to \infty$.

Condition D3. $n^2(N) \log(N)/N \to 0, N \to \infty$.

Condition U1. We need the following lemma.

Lemma D.1. Assumption 2.2 implies Condition U1 in Ljung and Wahlberg (1992).

Proof. Theorems 4.1 and 4.2 in Ljung and Wahlberg (1992) give that Assumption 2.2 implies that $\{u_t\}$ is f_N -quasi-stationary with $f_N = (\log N/N)^{1/2}$. Furthermore, $\{u_t\}$ is uniformly bounded and since F_u does not have zeros on the unit circle its spectral density is bounded from below by some $\delta > 0$.

The next result is a part of Theorem 3.1 in Ljung and Wahlberg (1992).

Theorem D.1. Assume Conditions S1, S2 (with $p = 2 + \delta$, $\delta > 0$), D1 and D3 to hold, as well as Assumption 2.2. Let the estimate $A(e^{i\omega}, \hat{\eta}_N)$ be defined by (6), (9), and (10). Then with probability 1,

$$\sup_{\omega} |A(e^{j\omega}, \hat{\eta}_N) - A^o(e^{j\omega})| = O(m(N)),$$
(D.1)

where $m(N) = n(N)\sqrt{\log N/N}(1 + d(N)) + d(N)$, where

$$d(N) := \sum_{k=n(N)+1}^{\infty} |a_k^o| + |b_k^o|.$$
(D.2)

The next result is Theorem 7.1 in Ljung and Wahlberg (1992).

Theorem D.2. Assume Conditions S1, S2 (with $p \ge 5$), D1, D2 and D4 to hold, as well as Assumption 2.2. Then

$$\left\| \mathbb{E} \left[N(\hat{\eta}_N - \bar{\eta}^{n(N)})(\hat{\eta}_N - \bar{\eta}^{n(N)})^T \right] - \sigma_o^2 \left[\bar{R}^n \right]^{-1} \right\|_2 \to 0,$$

as $N \to \infty$.

The final result from Ljung and Wahlberg (1992) that will be used is Theorem 7.3.

Theorem D.3. Assume Conditions S1, S2 (with $p \ge 4$), D1, D2 and D4 to hold, as well as Assumption 2.2. Let $\{\gamma^n\}$ be a sequence of deterministic $2m \times 2m$ matrices and assume that $\{\|\gamma^n\|_2\}$ is a bounded sequence. Then

$$\sqrt{N}\Upsilon^{n}(N)(\hat{\eta}_{N}-\bar{\eta}^{n(N)})\sim AsN(0,Q),$$

where $Q=\lim_{n\to\infty}\Upsilon^{n}\left[\bar{R}^{n}\right]^{-1}(\Upsilon^{n})^{T}.$

Appendix E. Orthogonal projections

We will make use of the following lemma, which is Lemma II.3 in Hjalmarsson and Mårtensson (2011).

Lemma E.1. Let $\gamma \in \mathcal{L}_2^{q \times m}$ and $\Psi \in \mathcal{L}_2^{n \times m}$. Then the orthogonal projection of the rows of γ on \mathcal{S}_{Ψ} (the subspace of \mathcal{L}_2^m spanned by the rows of Ψ) is given by

$$\operatorname{Proj}_{\mathcal{S}_{\mathcal{W}}}\{\gamma\} = \langle \gamma, \Psi \rangle \langle \Psi, \Psi \rangle^{\dagger} \Psi, \tag{E.1}$$

where *A*[†] is the pseudo-inverse of *A*. Furthermore,

$$\langle \gamma, \Psi \rangle \langle \Psi, \Psi \rangle^{\dagger} \langle \Psi, \gamma \rangle = \langle \operatorname{Proj}_{\delta_{\Psi}} \{\gamma\}, \operatorname{Proj}_{\delta_{\Psi}} \{\gamma\} \rangle.$$
(E.2)

Lemma E.2. Let $\gamma_i \in \mathcal{L}_2^{q_i \times m}$, i = 1, 2, and $\Psi_n \in \mathcal{L}_2^{n \times m}$, $n = 1, 2, \ldots$ Furthermore, let \mathscr{S}_n denote the subspace of \mathcal{L}_2^m spanned by the rows of Ψ_n . Suppose that

$$\|\gamma_i - P_{\delta_n}[\gamma]\|_2 \to 0 \quad \text{as } n \to \infty, \tag{E.3}$$

for i = 1 or/and i = 2. Then

$$\lim_{n \to \infty} \langle \gamma_1, \Psi_n \rangle \ \langle \Psi_n, \Psi_n \rangle^{\dagger} \ \langle \Psi_n, \gamma_2 \rangle = \langle \gamma_1, \gamma_2 \rangle. \tag{E.4}$$

Proof. A slight extension of Lemma E.1 gives

$$\langle \gamma_1, \Psi_n \rangle \langle \Psi_n, \Psi_n \rangle^{\dagger} \langle \Psi_n, \gamma_2 \rangle = \langle \operatorname{Proj}_{\delta_n} \{ \gamma_1 \}, \operatorname{Proj}_{\delta_n} \{ \gamma_2 \} \rangle.$$
(E.5)

Using that $\langle P_{s_n}[\gamma_i], \gamma_j - P_{s_n}\gamma_j \rangle = 0$ for i = 1, j = 2 and vice versa, gives

$$0 \leq \langle \gamma_1, \gamma_2 \rangle - \langle \operatorname{Proj}_{\delta_n} \{\gamma_1\}, \operatorname{Proj}_{\delta_n} \{\gamma_2\} \rangle = \langle \gamma_1 - \operatorname{Proj}_{\delta_n} \{\gamma_1\}, \gamma_2 - \operatorname{Proj}_{\delta_n} \{\gamma_2\} \rangle.$$
(E.6)

Now consider the absolute value of the *kl*th element, $1 \le k \le q_1$, $1 \le l \le q_2$, of the $q_1 \times q_2$ matrix in the right-most expression of (E.6). Applying Cauchy–Schwarz inequality and using (E.3) together with that $\|\gamma_i - \operatorname{Proj}_{s_n}\{\gamma_i\}\|_2$ is bounded since $\gamma_i \in \mathcal{L}_2^{q_i \times m}$ gives that this element converges to zero as $n \to \infty$. Combining this with (E.5) now proves the lemma since *k* and *l* are arbitrary.

Lemma E.3. Let \mathscr{S}_n be the subspace of \mathscr{L}_2^2 spanned by the rows of

$$\begin{bmatrix} F_1 F_3 \Gamma_n & F_2 \Gamma_n \\ F_3 \Gamma_n & 0 \end{bmatrix},$$

where $\Gamma_m(q) = \begin{bmatrix} q^{-1} & \cdots & q^{-m} \end{bmatrix}^T$, $F_i(q) = \sum_{k=0}^{\infty} f_k^i q^{-k}$. Suppose that F_i , i = 1, 2, 3 are exponentially stable, i.e.

 $|f_k^i| \leq C\lambda^k$, for some $C < \infty$, $\lambda < 1$,

and that there is a causal exponentially stable inverse \tilde{F}_2 for F_2 , i.e. $\tilde{F}_2(q)F_2(q) = 1$ where

$$\tilde{F}_2(q) = \sum_{k=0}^{\infty} \tilde{f}_k^2 q^{-k}, \qquad |\tilde{f}_k^2| < C\lambda^k.$$

Let $\gamma(q) = \begin{bmatrix} 0 & \sum_{k=1}^{\infty} d_k q^{-k} \end{bmatrix}$ be exponentially stable. Then $\|\gamma - P_{\delta_n}[\gamma]\|_2 \leq \tilde{C} \tilde{\lambda}^n$, for some $\tilde{C} < \infty$, $\tilde{\lambda} < 1$.

Proof. We will construct an explicit approximation to γ that belongs to \mathscr{S}_n . First we obtain an approximation to $\begin{bmatrix} 0 & z^{-k}F_2(z) \end{bmatrix}$ by (below k = 1, ..., n)

$$b_{n,k}(z) = \left[z^{-k} F_3(z) \sum_{l=n-k+1}^{\infty} f_l^1 z^{-l} \quad z^{-k} F_2(z) \right] \in \mathscr{S}_n$$

We will now use $\{b_{n,k}\}_{k=1}^n$ to approximate γ . To this end consider the expansion $\tilde{F}_2(z)\gamma(z) = \begin{bmatrix} 0 & \sum_{l=1}^{\infty} \alpha_l z^{-l} \end{bmatrix}$, where $|\alpha_l| \leq C_{\alpha} \lambda_{\alpha}^l$ for some $C_{\alpha} < \infty$ and $\lambda_{\alpha} < 1$ since both \tilde{F}_2 and γ are assumed to be exponentially stable. As approximation of γ we take

$$\hat{\gamma}_n(z) := \sum_{k=1}^n \alpha_l b_{n,k}(z) \\ \left(\approx \sum_{k=1}^\infty \alpha_k z^{-k} \begin{bmatrix} 0 & F_2(z) \end{bmatrix} \approx \gamma(z) \tilde{F}_2(z) F_2(z) = \gamma(z) \right)$$

which according to its construction belongs to δ_n . Hence

$$\|\gamma - P_{\delta_n}[\gamma]\|_2 \le \|\gamma - \hat{\gamma}_n\|_2$$
(E.7)

since $P_{\delta_n}[\gamma]$ has the smallest approximation error of γ of all functions in δ_n .

We introduce the notation $\gamma - \hat{\gamma}_n = [\delta_{n,1} \quad \delta_{n,2}]$. Using the exponentially decaying bounds on $|f_l^1|$ and $|\alpha_k|$,

$$\begin{split} \|\delta_{n,1}\|_{2} &= \left\|F_{3}(z)\sum_{k=1}^{n}\alpha_{k}\sum_{l=n-k+1}^{\infty}f_{l}^{1}z^{-l}\right\|_{2} \\ &\leq \|F_{3}\|_{2}\left\|\sum_{k=1}^{n}\alpha_{k}\sum_{l=n-k+1}^{\infty}f_{l}^{1}z^{-l}\right\|_{2} \\ &= \|F_{3}\|_{2}\left\|\sum_{l=1}^{n}f_{l}^{1}\left(\sum_{k=n-l+1}^{n}\alpha_{k}\right)z^{-l} + \sum_{l=n+1}^{\infty}f_{l}^{1}\left(\sum_{k=1}^{n}\alpha_{k}\right)z^{-l}\right\|_{2} \\ &= \|F_{3}\|_{2}\sqrt{\sum_{l=1}^{n}|f_{l}^{1}|^{2}\left(\sum_{k=n-l+1}^{n}\alpha_{k}\right)^{2} + \sum_{l=n+1}^{\infty}|f_{l}^{1}|^{2}\left(\sum_{k=1}^{n}\alpha_{k}\right)^{2}} \\ &\leq C_{1}\lambda_{1}^{n}, \quad C_{1} < \infty, \ 0 < \lambda_{1} < 1. \end{split}$$
(E.8)

For $\delta_{n,2}$ we have

$$\begin{aligned} \left\| \delta_{n,2} \right\|_{2} &= \left\| \gamma(z) - \sum_{k=1}^{n} \alpha_{k} z^{-k} F_{2}(z) \right\|_{2} \\ &= \left\| F_{2}(z) \left(\tilde{F}_{2}(z) \gamma(z) - \sum_{k=1}^{n} \alpha_{k} z^{-k} \right) \right\|_{2} \\ &\leq \|F_{2}\|_{2} \left\| \sum_{k=n+1}^{\infty} \alpha_{k} z^{-k} \right\|_{2} \leq C_{4} \lambda_{4}^{n}, \end{aligned}$$
(E.9)

for some $C_4 < \infty$ and $\lambda_4 < 1$, since F_2 and $\tilde{F}_2 \gamma$ are exponentially stable. Combining (E.7) and (E.8)–(E.9) gives

$$\|\gamma - P_{s_n}[\gamma]\|_2 \le \|\gamma - \hat{\gamma}_n\|_2 = \|\delta_{n,1}\|_2 + \|\delta_{n,2}\|_2 < \tilde{C}\tilde{\lambda}^n,$$

for some $\tilde{C} < \infty$ and $\tilde{\lambda} < 1$. This concludes the proof.

References

- Åström, K., & Söderström, T. (1974). Uniqueness of the maximum likelihood estimates of the parameters of an arma model. *IEEE Transactions on Automatic Control*, AC-19, 769–773.
- Box, G., & Jenkins, G. (1970). Time series analysis, forecasting and control. Oakland, California: Holden-Day.
- Eckhard, D., Bazanella, A., Rojas, C., & Hjalmarsson, H. (2012). On the convergence of the prediction error method to its global minimum. In 16th IFAC symposium on system identification, Brussels, Belgium (pp. 698–703).
- Goodwin, G., Agüero, J., & Skelton, R. (2003). Conditions for local convergence of maximum likelihood estimation for ARMAX models. In Proc. 13th IFAC symposium on system identification, Rotterdam, The Netherlands (pp. 797–802).
- Hjalmarsson, H., & Mårtensson, J. (2011). A geometric approach to variance analysis in system identification. IEEE Transactions on Automatic Control, 56, 983–997.
- Horn, R., & Johnson, R. (1985). Matrix analysis. Cambridge, UK: Cambridge University Press.

- Hsia, T. (1977). Identification: least squares methods. Lexington, MA, USA: Lexington Books.
- Jakeman, A., & Young, P. (1979). Refined instrumental variable methods of timeseries analysis: Part II, multivariable systems. *International Journal of Control*, 29, 621–644.
- Ljung, L. (1999). System identification: theory for the user (2nd ed.). Englewood Cliffs, NJ: Prentice-Hall.
- Ljung, L., & Wahlberg, B. (1992). Asymptotic properties of the least-squares method for estimating transfer functions and disturbance spectra. Advances in Applied Probability, 24, 412–440.
- Pierce, D. (1972). Least squares estimation in dynamic-disturbance time series models. *Biometrika*, 59(1), 73–78.
- Söderström, T. (1975a). On the uniqueness of maximum likelihood identification. Automatica, 11, 193–197.
- Söderström, T. (1975b). Tests of pole-zero cancellation in estimated models. Automatica, 11, 537–541.
- Stoica, P., & Söderström, T. (1981). The Steiglitz–Mcbride identification algorithm revisited: Convergence analysis and accuracy aspects. *IEEE Transactions on Automatic Control*, 26(3), 712–717.
- Stoica, P., & Söderström, T. (1983). Optimal instrumental variable estimation and approximate implementations. *IEEE Transactions on Automatic Control*, 28(7), 757–772.
- Van Overschee, P., & De Moor, B. (1994). N4SID–Subspace algorithms for the identification of combined deterministic stochastic-systems. *Automatica*, 30(1), 75–93.
- Verhaegen, M. (1994). Identification of the deterministic part of MIMO state-space models given in innovations form from input-output data. *Automatica*, 30(1), 61–74.
- Wahlberg, B. (1989). Model reduction of high-order estimated models: The asymptotic ML approach. International Journal of Control, 49, 169–192.
- Young, P. (1976). Some observations on instrumental variable methods of timeseries analysis. International Journal of Control, 23, 593–612.
- Young, P. (2006). An instrumental variable approach to ARMA model identification and estimation. In 14th IFAC symposium on system identification, Newcastle, Australia.
- Young, P. (2008). The refined instrumental variable method: Unified estimation of discrete and continuous-time transfer function models. *Journal Européen des* Systèmes Automatisés, 42, 149–179.
- Zhu, Y. (1998). Multivariable process identification for MPC: the asymptotic method and its applications. *Journal of Process Control*, 8(2), 101–115.
- Zhu, Y. (2001). Multivariable system identification for process control. Oxford: Elsevier Science Ltd..
- Zhu, Y. (2009). System identification for process control: recent experience and method and its applications. *International Journal of Modelling, Identification and Control*, 6(2), 89–103.
- Zhu, Y. (2011). A Box-Jenkins method that is asymptotically globally convergent for open loop data. In Proceedings of 17th IFAC world congress, Milan, Italy (pp. 9047–9051).
- Zou, Y., & Heath, W. (2009). Conditions for attaining global minimum in maximum likelihood system identification. In Proc. 15th IFAC symposium on system identification, Saint-Malo, France (pp. 1110–1115).



Yucai Zhu graduated from Xi-an Jiaotong University in 1982. He received his Master degree and Ph.D. in 1985 and 1990, both from Eindhoven University of Technology. During 1990 and 1993, he worked at IPCOS B.V. as a Control Engineer and he was also a co-founder of the company. He worked at Setpoint Inc. and later AspenTech from 1993 to 1996. He founded Tai-Ji Control BV in 1996. During 1998 and 2011, he has worked as a part-time Researcher at Eindhoven University of Technology. He has developed the so-called ASYM method of identification and several process control software packages. He has done consulting

work for major companies including BP, Dow, Statoil, ExxonMobil, Saudi Aramco in the field of model predictive control (MPC). Since 2011, he joined Zhejiang University as a full professor. At CPC 2012, he received the "Slowest publication award". His research interests are system identification, industrial MPC and related applications.



Håkan Hjalmarsson (M'98, SM'11, F'13) was born in 1962. He received the M.S. degree in Electrical Engineering in 1988, and the Licentiate degree and the Ph.D. degree in Automatic Control in 1990 and 1993, respectively, all from Linköping University, Sweden. He has held visiting research positions at California Institute of Technology, Louvain University and at the University of Newcastle, Australia. He has served as an Associate Editor for Automatica (1996–2001), and IEEE Transactions on Automatic Control (2005–2007) and been Guest Editor for European Journal of Control and Control Engineering Practice. He is Professor

at the School of Electrical Engineering, KTH, Stockholm, Sweden. He is an IEEE Fellow and Chair of the IFAC Coordinating Committee CC1 Systems and Signals. In 2001 he received the KTH award for outstanding contribution to undergraduate education. His research interests include system identification, signal processing, control and estimation in communication networks and automated tuning of controllers.