

The role of Vector AutoRegressive Modeling in Predictor Based Subspace Identification

Alessandro Chiuso^{a,1}

^a*Dipartimento di Tecnica e Gestione dei Sistemi Industriali Università di Padova (sede di Vicenza), stradella San Nicola, 3 - 36100 Vicenza, Italy. e-mail: chiuso@dei.unipd.it*

Abstract

Subspace identification for closed loop systems has been recently studied by several authors. A class of new and consistent closed-loop subspace algorithms is based on identification of a predictor model, in a way similar as prediction error methods (PEM) do. Experimental evidence suggests that these methods have a behavior which is very close to PEM in certain examples. The asymptotical statistical properties of one of these methods have been studied recently allowing to show (i) its relation with CCA and (ii) that Cramér Rao lower bound is not reached in general. Very little however is known concerning their relative performance.

In this paper we shall discuss the link between these “predictor based” methods; to this purpose we exploit the role which Vector AutoRegressive with eXogenous inputs models play in all these algorithms. The results of this paper provide a unifying framework under which all these algorithms can be viewed; also the link with VARX modeling have important implications as to computational complexity is concerned, leading to very computationally attractive implementations.

We also hope that this framework, and in particular the relation with VARX modeling followed by model reduction will turn out to be useful in future developments of subspace identification, such as the quest for efficient procedures and the statistical analysis with finite-data.

Key words: Closed-loop Identification, Subspace Methods, Asymptotic Properties, Variance, Relative Efficiency.

1 Introduction

Subspace identification has attracted a lot of attention in the last two decades. It is also fair to say that the last few years have witnessed a renewed interest in this topic for essentially two reasons: first the introduction of new methods which have allowed subspace identification to be applied with closed loop data [12,11,40,25,16] and second a whole body of results on the asymptotic statistical properties of subspace methods which have allowed, on the one side, to asses accuracy of subspace estimators [4,2,13,9,24] and on the other to compare different methods [5,3,2,14,7,8].

Extension of subspace algorithms to closed loop operating conditions have required, at a certain stage, the in-

troductioin of two step procedures (see [25,44,31]) which were needed to eliminate undesired terms due to feedback. This was due to the lack of stochastic realization procedures indicating how the state space could be constructed in the presence of feedback. An overview of these realization procedures can be found in [12,11,16] and references therein. The reader is also referred to [38,37] for early contributions advocating for two-step procedures even for “open loop” identification. There, an “iterative” CCA which makes use of the Markov parameters estimated in a preliminary stage was proposed. In [38] it was also conjectured that the two step CCA could possibly lead to asymptotic efficiency².

Often the preliminary estimation has been performed using Vector AutoRegressive with eXogenous inputs (VARX) models. Some analysis regarding the role of VAR models in subspace identification was performed in [18] where it was shown that the CCA algorithm introduced in [30] is asymptotically equivalent (in the

¹ This work has been supported in part by the national project *New methods and algorithms for identification and adaptive control of technological systems* funded by MIUR. Part of this work has been presented at the 2006 IFAC SYSID Symposium held in Newcastle, Australia and at the 2006 IEEE CDC conference held in San Diego, USA.

² In this paper we shall always use the word efficient assuming Gaussian distributions of the data.

sense of having the same asymptotic distribution of the estimators) to a procedure which first estimates a long VAR model and then does balanced model reduction.

Also some preliminary recent work relating VARX models with subspace procedures can be found in [42,36].

In this paper we shall be concerned with a class of algorithms which we group under the name Predictor Based Subspace Identification. This terminology stems from the fact that these algorithms aim at identifying a “predictor model” in a way that reminds of prediction error methods (PEM). We refer the reader to the papers [16,17] for a thorough discussion of the basic issues. The algorithms that shall be discussed here are the SSARX algorithm by Jansson [25], the PBSID³ algorithm (introduced in [16] under the name “whitening filter”), its “optimized” version⁴ (PBSID_{opt} hereafter) introduced in [7] and the algorithm presented in [35] by Ljung and McKelvey.

We shall expand on two recent contributions (see [8] and [10]) and discuss the role of VARX models in subspace algorithms based on predictor identification; we shall show that the preliminary step based on VARX models, explicitly used in [25] is actually present, in a way or another, in all these algorithms. We also observe that the bank of predictors used in [35] to overcome problems due to feedback were constructed using VARX models. It turns out, as we shall see later in this paper, that the algorithm proposed in [35] is very much related to the PBSID_{opt} introduced in [7]. For this reason, even though PBSID_{opt} has been developed independently from [35] and actually derives from a theoretically sound optimization, we regard the paper [35] as a fundamental early contribution to closed-loop subspace identification.

In particular the main results of this paper can be enumerated as follows:

- (a) SSARX by Jansson [25], which requires a preliminary VARX modeling step, is asymptotically equivalent, in the sense of yielding the same asymptotic distribution of the estimators, to PBSID. Some preliminary results have appeared in [8].
- (b) The “optimally-weighted” projection step involved in PBSID_{opt} in [7] is actually equivalent (here in the sense of giving the *same numerical results*) to estimating a VARX model followed by the usual steps of subspace identification⁵. Some preliminary results can be found in [10].

³ Short for “Predictor Based Subspace IDentification”.

⁴ The word “optimized” refers to the fact that a projection step is replaced by an “optimally weighted” (Markov) estimator.

⁵ Even though some preliminary results along these lines have already been presented in [8], the author would like to thank an anonymous reviewer of the paper [7] which have

- (c) The algorithm presented by Ljung and McKelvey in [35] is equivalent to a weighted version of PBSID_{opt}

In our opinion the significance of these results with respect to the current state of the art in subspace identification can be described as follows:

- (a) One contribution of this paper, which can be seen as a natural continuation of previous works [3,5,14,7], is to provide a comparison between recently proposed methods, trying to obtain a more unified picture of subspace algorithms; we believe this is useful since subspace algorithms have grown rapidly in number in the last few years [35,25,41,16] with very little insight, if any, concerning their relative efficiency.
- (b) Second, by showing that PBSID_{opt} is numerically equivalent to estimating a VARX model followed by the usual steps of subspace identification, we provide a way to implement the PBSID_{opt} algorithm with a much lower computational complexity than originally discussed in [7]. We would like to remind that PBSID_{opt} share the advantages of PBSID in that it delivers consistent estimators with closed loop data while comparing favorably, in the sense of asymptotic variance, to CCA for the open loop case. In fact it was shown in [7], Theorem 5.3, that the asymptotic variance of PBSID_{opt} is less or equal than that of CCA for *any* choice of the input signal. The fact that CCA is known to be asymptotically efficient⁶ for time series identification (=no inputs) [3] and optimal for white inputs [5] strengthens the significance of our result.
- (c) Last but not least, the relation with VARX modeling followed by model reduction, together with the results in [18,2,5,7], might be very helpful, in the author’s opinion, in future developments of subspace identification. We also refer the reader to the paper [50] for a discussion and early references on the use of high order AR models for identification of Autoregressive Moving Average (ARMA) models. In particular all these results could provide:
 - i) suggestions on how subspace procedures could be modified such as to reach asymptotic efficiency; recall for instance that in the case of no observed inputs [50] derives an asymptotically efficient ARMA parameter estimation method based on AR modeling followed by model reduction.
 - ii) a tool to introduce structure in the identification problem (delays, inputs which do not affect certain outputs etc.) which might turn out to be very useful when handling systems with large numbers of input-output channels (see e.g. the

underlined the relevance of the comparison performed in this paper; part of the merit of this paper should also go to him.

⁶ When both the past and future horizon go to infinity with the number of data.

plenary lecture given by Yucai Zhu at the recent SYSID in Newcastle [53] and also [23]).

Concerning the relation of subspace methods with VARX modeling, recall that it was shown in [18] that, for time series identification (i.e. no inputs) VAR modeling followed by balanced model reduction is asymptotically equivalent to the CCA method, which is asymptotically efficient as shown in [3]. It has also been shown in [7] that PBSID (and therefore its “optimized” version) is asymptotically equivalent to CCA for time-series identification and when input signals are white.

Hence, at least for white inputs and time series identification PBSID “does model reduction right”. The situation is different when there are inputs, and they are colored. The PBSID_{opt} performs better than CCA but it is not clear whether it is efficient in general; note that in [38] it was conjectured that pre-estimation of certain Markov parameters might be a way to obtain efficient subspace procedures.

In [31] it is claimed that a procedure which is very much related to the SSARX algorithm might be efficient; however, in [31] it is claimed that efficiency is reached for both past and future horizon which go to infinity. This claim appears to be wrong since, on the contrary, depending upon the input characteristics, the asymptotic variance might increase or decrease as a function of the future horizon (see [7],[6]). Note also that in [6] the algorithm discussed in [31] is shown to be asymptotically equivalent to SSARX and hence, from the results of this paper, also to PBSID.

We believe existence of an efficient subspace procedure is worth investigating.

We warn the reader that this paper does not mean to provide an exhaustive coverage of the state of the art in subspace identification but rather an analysis of a specific class of algorithms as mentioned earlier in the introduction. Many algorithms are not discussed [26,28,27] or just mentioned in passing [42,36,41,44].

The structure of the paper is as follows. In Section 2 we state the problem and set up notation; Section 3 briefly recalls the algorithmic steps while Section 4 states the main results of this paper contained in Theorem 4.1, Theorem 4.2 and Proposition 4.6 together with some simulation results. Section 5 contains some conclusions. The most technical parts of the proofs are postponed to the Appendix.

2 Statement of the Problem and Notation

Let $\{\mathbf{z}(t)\}$, $t \in \mathbb{Z}$, $\mathbf{z} := [\mathbf{y}^\top \mathbf{u}^\top]^\top$, be a (weakly) stationary second-order ergodic stochastic process where $\mathbf{y}(t)$ and $\mathbf{u}(t)$ are respectively the output (p dimensional) and

input (m dimensional) signals of a linear stochastic system in innovation form

$$\begin{cases} \mathbf{x}(t+1) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{K}\mathbf{e}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) + \mathbf{e}(t) \end{cases} \quad t \geq t_0. \quad (2.1)$$

We allow for *feedback* from $\{\mathbf{y}(t)\}$ to $\{\mathbf{u}(t)\}$ [21], i.e. we consider “closed loop” identification. Without loss of generality we shall assume that the dimension n of the state vector $\mathbf{x}(t)$ is as small as possible, i.e. the representation (2.1) is minimal. For simplicity we assume that $D = 0$, i.e. there is no direct feedthrough. For future reference we define $\bar{A} := A - KC$. We shall assume that spectral density matrix of \mathbf{z} , $\Phi(z)$ is rational and bounded away from zero on the unit circle $z = e^{j\omega}$. Let μ_i denote the zeros of the spectral density matrix which are inside the closed unit disc. We define $\rho := \max(|\mu_i|)$. From the assumption $\Phi(e^{j\omega}) > cI > 0$ it follows that $\rho < 1$. Note in particular that $1 > \rho \geq \max(|\lambda_i(\bar{A})|)$ where $\lambda_i(\bar{A})$ is the i -th eigenvalue of \bar{A} .

The white noise process \mathbf{e} , the innovation of \mathbf{y} given the past of \mathbf{z} , is defined as the one step ahead prediction error of $\mathbf{y}(t)$ given the (strict) past of \mathbf{z} up to time t .

Given two sequences of (scalar) random variables \mathbf{x}_N and \mathbf{g}_N , we shall say that \mathbf{x}_N is $O_P(\mathbf{g}_N)$, which we shall write $\mathbf{x}_N = O_P(\mathbf{g}_N)$, if, $\forall \epsilon, \exists M$ s.t.

$$\sup_N P[|\mathbf{x}_N/\mathbf{g}_N| > M] < \epsilon$$

In particular if $\mathbf{x}_N = O_P(1)$ we say that \mathbf{x}_N is bounded in probability⁷; note that $\mathbf{x}_N = O_P(\mathbf{g}_N)$ means that $\mathbf{x}_N/\mathbf{g}_N$ is bounded in probability.

Similarly, $\mathbf{x}_N = o_P(\mathbf{g}_N)$ means that, $\forall \delta > 0$,

$$\lim_{N \rightarrow \infty} P[|\mathbf{x}_N/\mathbf{g}_N| > \delta] = 0$$

If both \mathbf{x}_N and \mathbf{g}_N are deterministic sequences, say x_N and g_N , then $x_N = o(g_N)$ has the usual meaning $\lim_{N \rightarrow \infty} x_N/g_N = 0$. The symbol $\dot{=}$ shall denote equality in probability up to $o_P(1/\sqrt{N})$ terms, which we shall call *asymptotic equivalence*. In fact, from standard results in asymptotic analysis (see for instance [19]) terms which are $o_P(1/\sqrt{N})$ can be neglected when studying the asymptotic distribution.

We shall use the notations $\underline{o}_P(\cdot)$, $\underline{O}_P(\cdot)$, $\underline{o}(\cdot)$ or $\underline{O}(\cdot)$ to denote random matrices (of suitable dimensions possibly depending on N) whose elements are respectively $o_P(\cdot)$, $O_P(\cdot)$, $o(\cdot)$ or $O(\cdot)$ uniformly. Uniformity is needed when

⁷ Sometimes this is stated saying that \mathbf{x}_N is “uniformly tight”. For instance every sequence of random variables converging in distribution is uniformly tight.

the matrices' sizes increase with N . In this paper uniformity shall be guaranteed by stationarity of the processes involved.

Note also, for future reference, that if $\mathbf{x}_N = O_P(1/\sqrt{N})$ and $\mathbf{y}_N = O_P(1/\sqrt{N})$, then $\mathbf{x}_N \mathbf{y}_N = o_P(1/\sqrt{N})$.

Our aim is to identify the system parameters (A, B, C, K) , or equivalently the transfer functions $F(z) = C(zI - A)^{-1}B$ and $G(z) = C(sI - A)^{-1}K + I$, starting from input-output data $\{y_s, u_s\}$, $s \in [t_0, T + N]$, generated by the system (2.1).

Throughout the paper the symbol t shall denote “present”, t_0 shall be the initial time from which data are collected, so that $t - t_0$ is the “past horizon”; T shall be a design parameter so that $T - t$ is the number of future lags used to form predictors, commonly known as the “future horizon”⁸, while N shall be the length of the finite tails⁹.

The analysis reported in this paper requires that both N , and $t - t_0$ go to infinity. We remind the reader that $t - t_0$ has to go to infinity at a certain rate depending on the number N of data available. Details can be found, for instance, in [5] where the following assumption is made:

Assumption 1 *The past horizon $t - t_0$ goes to infinity with N while satisfying:*

$$\begin{aligned} t - t_0 &\geq \frac{\log N^{-d/2}}{\log |\rho|}, \quad 1 < d < \infty \\ t - t_0 &= o(\log(N)^\alpha), \quad \alpha < \infty \end{aligned} \quad (2.2)$$

Under this assumption the effect of terms due to mishandling of the initial condition at time t_0 are $o_P(1/\sqrt{N})$ and therefore can be neglected. Moreover, (2.2) ensures that, when regressing onto past data and taking the limit as N goes to infinity, the computation of sample covariance matrices of increasing size (with $t - t_0$) does not pose any complication in the sense that their limit is well defined and equal to the population counterpart (see the discussion after Lemma 4 in [5]).

In order to simplify the analysis in this paper we shall keep the future horizon ν fixed (and finite). We warn the reader that, for instance, asymptotic efficiency of the CCA method for time series identification requires ν to grow with the sample size (see [2]). We believe however that our results are significant even under this assumption since the comparison holds for fixed yet arbitrarily large ν .

⁸ Respectively the number of block rows in the block Hankel data matrix containing the past and future data.

⁹ This is the parameter j in the notation of Van Overschee and De Moor [47] i.e. the number of columns in the block Hankel data matrices used in subspace identification.

According to Assumption 1 the total number of data is $T + N - t_0 + 1 = T - t + N + o(\log(N)^\alpha)$; therefore $\lim_{N \rightarrow \infty} \frac{T+N-t_0+1}{N} = 1$ (which implies, in particular, $O(N) = O(T + N - t_0 + 1)$). For this reason (and for convenience of notation) when dealing with asymptotic results, we shall refer to the length of the finite tails N rather to the total number of data $T + N - t_0 + 1$ (e.g. we shall use $o_P(1/\sqrt{N})$ and not $o_P(1/\sqrt{T + N - t_0 + 1})$).

Remark 2.1 [On the necessity of Assumption 1] *One may argue that having to deal with an “infinite” past horizon might not be very attractive. This condition, as discussed in detail in [16], is needed to ensure consistency in closed loop. There is however another important reason to keep this assumption. Essentially subspace algorithms are “covariance based” methods¹⁰; therefore, as discussed in [51, 39] for the MA/ARMA case, it is necessary to estimate an infinite number of covariances to obtain asymptotically efficient estimators (see also [50][Remark 4, pag. 289]). This requires that $t - t_0$ goes to infinity. Therefore, a necessary condition for a subspace algorithm to reach the Cramér-Rao lower bound is that $t - t_0$ goes to infinity, as required by Assumption 1. Of course this does not mean that an algorithm allowing $t - t_0 \rightarrow \infty$ will automatically be efficient.*

We shall use the standard notation of boldface (lowercase) letters to denote random variables. Lowercase letters denote sample values of a certain random variable. For example we shall denote with $\mathbf{z}(t)$ the random vector denoting the joint process and with z_t the sample value of $\mathbf{z}(t)$. We shall use capitals to denote the tail of length N . For instance $Z_t := [z_t \ z_{t+1}, \dots \ z_{t+N-1}]$. These are the block rows of the usual *block Hankel data matrices* which appear in subspace identification.

When dealing with tails of length different from N we shall add the number of columns as a superscript; for instance $Z_t^M := [z_t \ z_{t+1}, \dots \ z_{t+M-1}]$

For $-\infty \leq t_0 \leq \tau \leq t \leq T \leq +\infty$ we define the Hilbert space of scalar zero-mean random variables

$$\mathcal{Z}_{[\tau, t]} := \overline{\text{span}} \{ \mathbf{z}_k(s); k = 1, \dots, m + p, \tau \leq s < t \}$$

where the bar denotes closure in mean square, i.e. in the metric defined by the inner product $\langle \xi, \eta \rangle := \mathbb{E}\{\xi\eta\}$, the operator \mathbb{E} denoting mathematical expectation. Similar definitions hold for $\mathcal{Y}_{[\tau, t]}$ and $\mathcal{U}_{[\tau, t]}$.

When $\tau = -\infty$ we shall use the shorthands \mathcal{Z}_t^- for $\mathcal{Z}_{[-\infty, t]}$. The space generated by $\mathbf{z}(s)$, $-\infty < s < \infty$ shall be denoted with the symbol \mathcal{Z} . For convenience of notation we denote with $\nu := T - t$ the future horizon.

¹⁰ One may argue that there are “data based” methods, but this is just a computational aspect.

Given a subspace $\mathcal{C} \subseteq \mathcal{Z}$, we shall denote with $E[\mathbf{a} | \mathcal{C}]$ the orthogonal projection of the random variable \mathbf{a} onto \mathcal{C} ; in the Gaussian case the linear projection coincides with conditional expectation, i.e. $\mathbb{E}[\cdot | \mathcal{C}] = E[\cdot | \mathcal{C}]$. Let \mathbf{c} be a (finite) basis for \mathcal{C} . Using the notation $\Sigma_{\mathbf{ab}} := \mathbb{E}[\mathbf{ab}^\top]$ for the covariance matrix between the zero mean random vectors \mathbf{a} and \mathbf{b} , in the finite dimensional case $E[\mathbf{a} | \mathcal{C}]$ will be given by the usual formula

$$E[\mathbf{a} | \mathcal{C}] = \Sigma_{\mathbf{ac}} \Sigma_{\mathbf{cc}}^{-1} \mathbf{c}. \quad (2.3)$$

Defining also the projection errors $\tilde{\mathbf{a}} := \mathbf{a} - E[\mathbf{a} | \mathcal{C}]$ and $\tilde{\mathbf{b}} := \mathbf{b} - E[\mathbf{b} | \mathcal{C}]$, the symbol $\Sigma_{\mathbf{ab} | \mathcal{C}}$ will denote projection error covariance (conditional covariance in the Gaussian case) $\Sigma_{\mathbf{ab} | \mathcal{C}} := \Sigma_{\tilde{\mathbf{a}}\tilde{\mathbf{b}}} = \Sigma_{\mathbf{ab}} - \Sigma_{\mathbf{ac}} \Sigma_{\mathbf{cc}}^{-1} \Sigma_{\mathbf{cb}}$. Given two trivially intersecting subspaces $\mathcal{C} \subseteq \mathcal{Z}$, $\mathcal{B} \subseteq \mathcal{Z}$, $\mathcal{C} \cap \mathcal{B} = \{0\}$, $E_{\parallel \mathcal{B}}[\cdot | \mathcal{C}]$ shall denote the oblique projection onto \mathcal{C} along \mathcal{B} (see [20]) and can be computed by the formula:

$$E_{\parallel \mathcal{B}}[\mathbf{a} | \mathcal{C}] = \Sigma_{\mathbf{ac} | \mathcal{B}} \Sigma_{\mathbf{cc} | \mathcal{B}}^{-1} \mathbf{c}. \quad (2.4)$$

For column vectors formed by stacking past and/or future random variables we shall use the notation: $\mathbf{z}_{[t,s]} := \begin{bmatrix} \mathbf{z}^\top(t) & \mathbf{z}^\top(t+1) & \dots & \mathbf{z}^\top(s) \end{bmatrix}^\top$. Finite (block) Hankel data matrices will be denoted using capitals, i.e. $Z_{[t,s]} := \begin{bmatrix} Z_t^\top & Z_{t+1}^\top & \dots & Z_s^\top \end{bmatrix}^\top$.

Spaces generated by finite tails, i.e. spaces generated by the rows of finite block Hankel data matrices, will be denoted with the same symbol used for the matrix itself. Sample covariances will be denoted with the same symbol used for the corresponding random variables with a “hat” on top. For example, given finite sequences $A_t := [a_t, a_{t+1}, \dots, a_{t+N-1}]$ and $B_t := [b_t, b_{t+1}, \dots, b_{t+N-1}]$ we shall define the sample covariance matrix

$$\hat{\Sigma}_{\mathbf{ab}} := \frac{1}{N} \sum_{i=0}^{N-1} a_{t+i} b_{t+i}^\top.$$

Under our ergodic assumption $\lim_{N \rightarrow \infty} \hat{\Sigma}_{\mathbf{ab}} \stackrel{a.s.}{=} \Sigma_{\mathbf{ab}}$.

The orthogonal projection onto the row space of a matrix shall be denoted with the symbol \hat{E} ; for instance, given a matrix $C_t := [c_t, c_{t+1}, \dots, c_{t+N-1}]$, $\hat{E}[\cdot | C_t]$ will be the orthogonal projection onto the row space of the matrix C_t ; the symbol $\hat{E}[A_t | C_t]$ shall denote the orthogonal projection of the rows of the matrix A_t onto the row space of C_t , and is given by the formula

$$\hat{E}[A_t | C_t] = \hat{\Sigma}_{\mathbf{ac}} \hat{\Sigma}_{\mathbf{cc}}^{-1} C_t \quad (2.5)$$

As above, given a matrix C_t , we define the projection errors $\hat{A}_t := A_t - \hat{E}[A_t | C_t]$ and $\hat{B}_t := B_t - \hat{E}[B_t | C_t]$. The sample covariance (conditional sample covariance) of the

projection errors is denoted with the symbol $\hat{\Sigma}_{\mathbf{ab} | \mathcal{C}} := \hat{\Sigma}_{\tilde{\mathbf{a}}\tilde{\mathbf{b}}}$ and computed by the formula

$$\hat{\Sigma}_{\mathbf{ab} | \mathcal{C}} := \hat{\Sigma}_{\mathbf{ab}} - \hat{\Sigma}_{\mathbf{ac}} \hat{\Sigma}_{\mathbf{cc}}^{-1} \hat{\Sigma}_{\mathbf{cb}}.$$

We shall denote with $\hat{E}_{\parallel \mathcal{B}_t}[\cdot | C_t]$ the oblique projection along the space generated by the rows \mathcal{B}_t onto the space generated by the rows of C_t (provided they intersect only at zero). As above, the oblique projection can be computed using the formula:

$$\hat{E}_{\parallel \mathcal{B}_t}[A_t | C_t] = \hat{\Sigma}_{\mathbf{ac} | \mathcal{B}_t} \hat{\Sigma}_{\mathbf{cc} | \mathcal{B}_t}^{-1} C_t. \quad (2.6)$$

For future reference we also define the extended observability matrix

$$\bar{\Gamma}_\nu^\top := \begin{bmatrix} C^\top & \bar{A}^\top C^\top & (\bar{A}^\top)^2 C^\top & \dots & (\bar{A}^\top)^{\nu-1} C^\top \end{bmatrix}. \quad (2.7)$$

3 State Space Construction

It is well known [47,32,15] that identification using subspace methods can be seen as a two step procedure as follows:

- (a) Construct a basis \hat{X}_t for the state space via suitable projection operations on data sequences (block Hankel data matrices)
- (b) Given (coherent) bases for the state space at time t (\hat{X}_t) and $t+1$ (\hat{X}_{t+1}) solve

$$\begin{cases} \hat{X}_{t+1} \simeq A \hat{X}_t + B \hat{U}_t + K E_t \\ Y_t \simeq C \hat{X}_t + E_t \end{cases} \quad (3.1)$$

in the least squares sense.

Different subspace algorithms have different implementations of the first step while the second remains the same for virtually all algorithms¹¹. For this reason we compare algorithms on the basis of step (a). We shall identify procedures which are (asymptotically) equivalent, modulo change of basis, as the first step is concerned.

3.1 PBSID algorithm

The construction of the state space using this algorithm involves several oblique projections. The projection of

¹¹ In this paper we shall not be concerned with algorithms based on the so-called “shift invariance” method [2].

each (block) row Y_{t+h} , $h = 0, \dots, \nu$, can be seen as a long VARX model as follows

$$\begin{aligned}\hat{Y}_{t+h} &:= \hat{E} [Y_{t+h} | Z_{[t_0, t+h]}] = \\ &= \hat{\Psi}_{1,h} Z_{t+h-1} + \dots + \hat{\Psi}_{t+h-t_0, h} Z_{t_0}\end{aligned}\quad (3.2)$$

from which the oblique projections¹²

$$\begin{aligned}\hat{Y}_{t+h}^P &:= \hat{E}_{\|Z_{[t_0, t+h]}} [Y_{t+h} | Z_{[t_0, t]}] = \\ &= \sum_{i=h+1}^{t-t_0+h} \hat{\Psi}_{i,h} Z_{t+h-i} \simeq C\bar{A}^{h-1} X_t\end{aligned}\quad (3.3)$$

The last approximate equality has to be understood in the sense that, asymptotically in N ,

$$\hat{y}^P(t+h) := E_{\|Z_{[t_0, t+h]}} [\mathbf{y}(t+h) | Z_t^-] = C\bar{A}^{h-1} \mathbf{x}(t)\quad (3.4)$$

holds. Then one stacks all the predictors

$$\hat{Y}_{[t,T]}^P := \begin{bmatrix} \hat{Y}_t^P \\ \hat{Y}_{t+1}^P \\ \vdots \\ \hat{Y}_{T-1}^P \end{bmatrix} \simeq \bar{\Gamma}_\nu X_t.$$

From the Singular Value Decomposition

$$W_p^{-1} \hat{Y}_{[t,T]}^P = PDQ^\top = [P_n \tilde{P}_n] \begin{bmatrix} D_n & 0 \\ 0 & \tilde{D}_n \end{bmatrix} \begin{bmatrix} Q_n^\top & \tilde{Q}_n^\top \end{bmatrix}\quad (3.5)$$

where W_p is a weighting matrix which can be chosen appropriately, an estimate of the observability matrix $\bar{\Gamma}_\nu$ is obtained discarding the “less significant” singular values (i.e. pretending $\tilde{D}_n \simeq 0$) from

$$\hat{\Gamma}_\nu = W_p P_n D_n^{1/2}.$$

and consequently a basis for the state space

$$\begin{aligned}\hat{X}_t^{PBSID} &:= \hat{\Gamma}_\nu^{-L} \hat{Y}_{[t,T]}^P \\ \hat{X}_{t+1}^{PBSID} &:= \hat{\Gamma}_\nu^{-L} \hat{Y}_{[t+1,T]}^P\end{aligned}\quad (3.6)$$

where $\hat{\Gamma}_\nu^{-L}$ is the left inverse defined by

$$\hat{\Gamma}_\nu^{-L} := \left(\hat{\Gamma}_\nu^\top W_p^{-\top} W_p^{-1} \hat{\Gamma}_\nu \right)^{-1} \hat{\Gamma}_\nu^\top W_p^{-\top} W_p^{-1}.\quad (3.7)$$

¹² The superscript P reminds that the quantity has to do with the “predictor-based” algorithm.

3.2 SSARX Algorithm

The algorithm described in the previous section can be seen as a “geometric” version of the SSARX algorithm by Jansson [25]. Instead of computing the oblique projections (3.3), or, equivalently, instead of estimating $\nu+1$ long VARX models, Jansson estimates just one (long) VARX model

$$Y_T \simeq \hat{\Phi}_1 Z_{T-1} + \hat{\Phi}_2 Z_{T-2} + \dots + \hat{\Phi}_{T-t_0} Z_{t_0}\quad (3.8)$$

where without loss of generality we have taken the length of the VARX model equal to $T-t_0$; then the effect of the future inputs/outputs is removed using the estimated parameters $\hat{\Phi}_k$ as¹³:

$$\hat{Y}_{[t,T]}^S := \hat{E} \left[Y_{[t,T]} - \hat{H}_\nu^S Z_{[t,T]} | Z_{[t_0, t]} \right]\quad (3.9)$$

where

$$\hat{H}_\nu^S := \begin{bmatrix} 0 & 0 & \dots & 0 \\ \hat{\Phi}_1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ \hat{\Phi}_\nu & \dots & \hat{\Phi}_1 & 0 \end{bmatrix}$$

The remaining part is essentially the same as in the previous Section provided¹⁴ $\hat{Y}_{[t,T]}^S$ is substituted to $\hat{Y}_{[t,T]}^P$.

3.3 “Optimized” PBSID Algorithm

The optimized version of PBSID introduced in [7] (PBSID_{opt}) differs from the original PBSID algorithm in the computation of the predictors (3.2); in fact in the optimized algorithm the estimation of the predictors \hat{Y}_{t+h} is formulated as a weighted least squares problem as described in this Section.

Let us define $\mathcal{K} := [\bar{A}^{t-t_0-1} [K \ B] \ \bar{A}^{t-t_0-2} [K \ B] \ \dots \ [K \ B]]$. Recall that

$$\begin{aligned}Y_{t+h} &= C\bar{A}^h X_t + \\ &+ \sum_{i=1}^h C\bar{A}^{i-1} (KY_{t+h-i} + BU_{t+h-i}) + E_{t+h} \\ &= C\bar{A}^h \mathcal{K} Z_{[t_0, t]} + \\ &+ \sum_{i=1}^h C\bar{A}^{i-1} (KY_{t+h-i} + BU_{t+h-i}) + \\ &+ E_{t+h} + \underline{q}_P(1/\sqrt{N}) \\ &:= \Xi_h Z_{[t_0, t]} + \\ &+ \sum_{i=1}^h \Psi_{hi} Z_{t+h-i} + E_{t+h} + \underline{q}_P(1/\sqrt{N})\end{aligned}\quad (3.10)$$

¹³ The superscript S stands for “SSARX”.

¹⁴ Also a specific choice of W_p is done by Jansson. We leave this choice unspecified here since equivalence holds for every choice of W_p .

where the last equality defines the matrices Ξ_h and Ψ_{hi} . Stacking the data and using (3.10) (discarding $o_P(1/\sqrt{N})$ terms; this is a delicate matter see Appendix B in [7] for details) we obtain:

$$\begin{bmatrix} Y_t \\ Y_{t+1} \\ \vdots \\ Y_T \end{bmatrix} \doteq \begin{bmatrix} \Xi_0 \\ \Xi_1 \\ \vdots \\ \Xi_\nu \end{bmatrix} Z_{[t_0,t]} + \begin{bmatrix} 0 & 0 & \dots & 0 \\ \Psi_{11} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \Psi_{\nu\nu} & \dots & \Psi_{\nu 1} & 0 \end{bmatrix} Z_{[t,T]} + \begin{bmatrix} E_t \\ E_{t+1} \\ \vdots \\ E_T \end{bmatrix} \quad (3.11)$$

Observe that the lower triangular matrices in (3.11) are Toeplitz, since $\Psi_{ij} = C\bar{A}^{j-1}[K \ B]$, $\forall i, j$. The projection in (3.2) is equivalent to solving (3.11) ‘‘row by row’’; hence the Toeplitz structure is not preserved after estimation, i.e. $\hat{\Psi}_{ij} \neq \hat{\Psi}_{i'j}$, $i \neq i'$ almost surely.

This is equivalent to solving the least squares problem obtained vectorizing (3.11):

$$Y := \begin{bmatrix} \text{vec}(Y_t) \\ \text{vec}(Y_{t+1}) \\ \vdots \\ \text{vec}(Y_T) \end{bmatrix} \doteq S^P \Omega^P + \begin{bmatrix} \text{vec}(E_t) \\ \text{vec}(E_{t+1}) \\ \vdots \\ \text{vec}(E_T) \end{bmatrix} = S^P \Omega^P + E \quad (3.12)$$

where the matrix S^P has the form

$$S^P = \text{block diag} \left\{ (Z_{[t_0,t]}^\top \otimes I), \dots, (Z_{[t_0,T]}^\top \otimes I) \right\} \quad (3.13)$$

and Ω^P is given by

$$\Omega^P = \begin{bmatrix} \text{vec}^\top(\Xi_0) & \text{vec}^\top(\Xi_1) & \text{vec}^\top(\Psi_{11}) & \dots \\ \dots & \text{vec}^\top(\Xi_\nu) & \dots & \text{vec}^\top(\Psi_{\nu 1}) \end{bmatrix}^\top; \quad (3.14)$$

Finding an ‘‘optimal’’ solution $\hat{\Omega}^{P_{opt}}$ (Markov estimator) of

$$Y \doteq S^P \Omega^P + E, \quad (3.15)$$

where $o_P(1/\sqrt{N})$ terms have been neglected¹⁵, gives an estimator $\hat{\Omega}^{P_{opt}}$ of Ω^P which has the smallest asymptotic variance among all linear (asymptotically unbiased) estimators based on (3.12). Incidentally, this has allowed to show in [7] that this ‘‘optimized’’ version yields, asymptotically, a lower variance of the estimators of any system

invariant as compared to the standard PBSID and, more importantly, to the classical CCA algorithm [30,46].

To this purpose it is very useful to observe that the ‘‘noise term’’ E can be written in the form

$$E = L \text{vec}(E_t^{N+\nu}) \quad (3.16)$$

where L is a ‘‘selection’’¹⁶ matrix of size $pN\nu \times p(\nu + N)$. We refer the reader to the paper [7] for an explicit expression of L ; suffices it to remind that L has full column rank. We shall later use the specific structure of the column space of L and of its left kernel. Equation (3.16) shows that indeed E has a singular covariance matrix $R = \text{Var}\{E\} = L(I \otimes \Lambda)L^\top$.

In the paper [7] it is shown how (3.15) can be converted into a least squares problem with full rank noise covariance and equality constraints (see also [43,52,45,20]). Remarkably, as we shall see in the next Section, this is equivalent to estimating a long VARX model of length $t - t_0$, using data in the interval $[t_0, T + N - 1]$.

Using the estimator $\hat{\Omega}^{P_{opt}}$, the oblique projections \hat{Y}_{t+h}^P (3.3) can be substituted with $\hat{Y}_{t+h}^{P_{opt}} = \hat{\Xi}_h^{P_{opt}} Z_{[t_0,t]}$ in the SVD step (3.5); hence, defining $\hat{Y}_{[t,T]}^{P_{opt}} := \left[(\hat{Y}_t^{P_{opt}})^\top, (\hat{Y}_{t+1}^{P_{opt}})^\top, \dots, (\hat{Y}_{T-1}^{P_{opt}})^\top \right]^\top$, an estimator for the state shall be given by

$$\hat{X}_t^{P_{opt}} := \left(\hat{\Gamma}_\nu^{P_{opt}} \right)^{-L} \hat{Y}_{[t,T]}^{P_{opt}}. \quad (3.17)$$

where $\hat{\Gamma}_\nu^{P_{opt}}$ is the estimate obtained substituting $\hat{Y}_{t+k}^{P_{opt}}$ to \hat{Y}_{t+h}^P in (3.5).

Also the ‘‘shifted’’ oblique projections used for the computation of the state at time $t + 1$ (see (3.6)) can be substituted by

$$\hat{X}_{t+1}^{P_{opt}} := \left(\hat{\Gamma}_\nu^{P_{opt}} \right)^{-L} \begin{bmatrix} \hat{\Xi}_1^{P_{opt}} & \hat{\Psi}_{11}^{P_{opt}} \\ \hat{\Xi}_2^{P_{opt}} & \hat{\Psi}_{22}^{P_{opt}} \\ \vdots & \vdots \\ \hat{\Xi}_\nu^{P_{opt}} & \hat{\Psi}_{\nu\nu}^{P_{opt}} \end{bmatrix} Z_{[t_0,t+1]}. \quad (3.18)$$

Similarly an estimator of the innovation sequence E_t can be found by

$$\hat{E}_t^{P_{opt}} := Y_t - \hat{E} \left[\hat{Y}_t^{P_{opt}} | \hat{X}_t^{P_{opt}} \right] \quad (3.19)$$

¹⁵ See Appendix B in [7] for a rigorous discussion.

¹⁶ We call ‘‘selection matrix’’ a matrix formed with zeros and ones in which each row all entries are zero except for one.

4 Main Results

This section contains the main results of this paper; first we shall discuss the (asymptotic) equivalence between PBSID and SSARX and later we shall discuss how PBSID_{opt} can be implemented using VARX models. Last PBSID_{opt} is related to the algorithm of Ljung and McKelvey [35] which, in some sense, might be seen as a predecessor of all these methods.

For clarity of exposition we divide this task in three separate and self-contained subsections, each complemented with some simulation results.

4.1 Asymptotic Equivalence of PBSID and SSARX

The first main result of this paper can be summarized as follows:

Theorem 4.1 *Assume the past horizon $t-t_0$ grows with N according to Assumption 1. Denote with $\hat{\Theta}^P$ and $\hat{\Theta}^S$ the estimators of any system invariant Θ using respectively the PBSID algorithm and the SSARX algorithm. Then, under standard assumptions (see, e.g. [9,2]) on the innovation process \mathbf{e}*

$$\hat{\Theta}^P \doteq \hat{\Theta}^S \quad (4.1)$$

holds.

We first state the following technical lemma which shall be useful in the proof of this result:

Lemma 4.2 *Let the pair (\mathbf{y}, \mathbf{u}) satisfy the assumptions of Section 2. Assume also the coefficients of the following two VARX models*

$$\mathbf{y}(t) = \sum_{i=1}^{K_1} \alpha_i \mathbf{z}(t-i) + \mathbf{e}_{K_1}(t) \quad (4.2)$$

and

$$\mathbf{y}(t) = \sum_{i=1}^{K_2} \beta_i \mathbf{z}(t-i) + \mathbf{e}_{K_2}(t) \quad (4.3)$$

are estimated (in the least square sense) from data $\{y_s, u_s\}$ respectively in the intervals $s \in [t-K_1, t+N-1]$ and $s \in [t-K_2, t+N-1]$. Assume also that $K_1 \geq K_2 \geq K_{min}$ go to infinity with N while K_1, K_{min} satisfy¹⁷ Assumption 1. Then for any fixed and finite f

$$\hat{\alpha}_j \doteq \hat{\beta}_j \quad j = 1, \dots, f \quad (4.4)$$

The same holds if the parameters in (4.2) and (4.3) are estimated using data in the intervals $[t_1-K_1, t_1+N-1]$

¹⁷ Where the role of t_0 is played respectively by $t_0^1 := t-K_1$ and $t_0^{min} := t-K_{min}$.

and $[t_2-K_2, t_2+N-1]$ respectively as long as t_1-t_2 is fixed and finite.

Proof. The proof follows from equations (4.4) and (4.5) in [29], by letting $K_1 = h_{max}$, $K_2 = h_n$, $K_{min} = h_{min}$ and $l(h) = [0, \dots, 0, 1, 0, \dots, 0, \dots]^T$. Of course, in our case $P[K_2 \in [K_{min}, K_1]] = 1$. Note that the result in [29] is much stronger and holds for more general linear combinations $l(h)$ and for data dependent order selection rules K_2 s.t. $P[K_2 \in [K_{min}, K_1]] \rightarrow 1$ as $N \rightarrow \infty$. This is useful here since it makes it easy to extend our results also to the case in which the length of past and future horizons are estimated from data provided the conditions in [29] are still verified. However we shall not discuss this extension here. \square

Proof of Theorem 4.1. Our goal is essentially to show that $\hat{Y}_{[t,T]}^S$ and $\hat{Y}_{[t,T]}^P$ can be used interchangeably as far as asymptotic properties are concerned.

To this purpose, note that defining

$$\hat{H}_\nu^P := \begin{bmatrix} 0 & 0 & \dots & 0 \\ \hat{\Psi}_{1,1} & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ \hat{\Psi}_{\nu,\nu} & \dots & \hat{\Psi}_{1,\nu} & 0 \end{bmatrix}$$

$\hat{Y}_{[t,T]}^P$ can be rewritten as

$$\hat{Y}_{[t,T]}^P = \hat{E} \left[Y_{[t,T]} - \hat{H}_\nu^P Z_{[t,T]} \mid Z_{[t_0,t]} \right] \quad (4.5)$$

which has the same form as (3.9) provided \hat{H}_ν^P is substituted with \hat{H}_ν^S . Using this observation we can write

$$\hat{Y}_{[t,T]}^S - \hat{Y}_{[t,T]}^P = \left(\hat{H}_\nu^P - \hat{H}_\nu^S \right) \hat{E} \left[Z_{[t,T]} \mid Z_{[t_0,t]} \right] \quad (4.6)$$

It is obvious that, provided we can show that

$$\hat{H}_\nu^P \doteq \hat{H}_\nu^S, \quad (4.7)$$

using $\hat{Y}_{[t,T]}^P$ in lieu of $\hat{Y}_{[t,T]}^S$ does not change the asymptotic properties; in fact, under (4.7), also the difference $\hat{Y}_{[t,T]}^S - \hat{Y}_{[t,T]}^P$ will be $\underline{o}_P(1/\sqrt{N})$.

Inspecting the structure of the matrices \hat{H}_ν^S and \hat{H}_ν^P , it is rather simple to see that showing (4.7) is equivalent to prove that

$$\hat{\Phi}_i \doteq \hat{\Psi}_{i,h} \quad i = 1, \dots, h \quad h = 1, \dots, \nu \quad (4.8)$$

Hence the last part of the proof shall be concerned with (4.8).

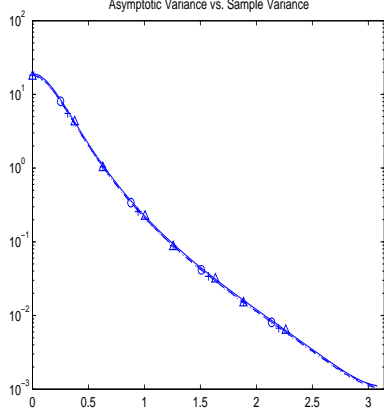


Fig. 1. Example 1. Sample Variance (Monte Carlo estimate) vs. normalized frequency ($\omega \in [0, \pi]$). Solid with triangles (\triangle): PEM. Dashed with crosses (+): PBSID. Dashed with circles (o): SSARX. Dotted with crosses (+): asymptotic variance for PBSID.

Let us fix for a moment $h = \bar{h}$. Showing that $\hat{\Phi}_i \doteq \hat{\Psi}_{i, \bar{h}}$ for $i = 1, \dots, \bar{h}$ amounts to prove that the estimators ($\hat{\Phi}_i$ and $\hat{\Psi}_{i, \bar{h}}$, $i = 1, \dots, \bar{h}$) of the first \bar{h} coefficients of two long VARX models satisfying

- the orders $T - t_0$ and $t - t_0 + \bar{h}$ differ of exactly $\nu - \bar{h}$ and both go to infinity at a rate specified by Assumption 1.
- the parameters are estimated essentially using the same data (essentially here means that there might be a finite number of data points which are used in one of the two and are not used in the other and vice versa)

are asymptotically equivalent. This result has been formalized in Lemma 4.2 above. Repeated application of Lemma 4.2 to the VARX regressions (3.2) and (3.8) allows indeed to prove (4.8) and hence (4.7), from which the statement of Theorem 4.1 follows. \square

We now report some simulation results concerning the equivalence of PBSID and SSARX. We consider two systems in innovation from (2.1) where the input given in closed loop

$$\mathbf{u}(t) = \mathbf{r}(t) - H_i(z)\mathbf{y}(t).$$

Example 1 is a first order ARMAX system with $A_1 = 0.7$, $B_1 = 1$, $K_1 = 1$, $C_1 = 1$, $D_1 = 0$, $\text{Var}\{\mathbf{e}_1\} = 1$, with a proportional controller $H_1(z) = 1.5$ and white reference signal $\mathbf{r}(t) = 5\mathbf{n}(t)$ where $\mathbf{n}(t)$ is zero mean unit variance white noise uncorrelated from $\mathbf{e}(t)$.

Example 2 is a second order ARMAX system with

$$\begin{aligned} A_2 &= \begin{bmatrix} 1.5 & -0.7 \\ 1 & 0 \end{bmatrix} & B_2 &= \begin{bmatrix} 1 \\ 0 \end{bmatrix} & K_2 &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ C_2 &= [1 \ 0] & D_2 &= 0 & \text{Var}\{\mathbf{e}_2\} &= 1 \end{aligned}$$

The reference signal is unit variance white noise uncorrelated with the innovation $\mathbf{e}(t)$ while the controller is a first order system of the form

$$H_2(z) = 0.2 \frac{0.1z - 0.5}{z - 0.5}.$$

We compare the Monte Carlo estimate (500 Monte Carlo runs) of the transfer function estimator ($\hat{F}(z) := \hat{C}(zI - \hat{A})^{-1}\hat{B}$) variance (normalized by N) of SSARX and PBSID algorithms. The parameters chosen in the three simulation reported respectively in figures 1,2 are summarized in table 4.1. The results of figure 1 refer to Example 1 while those in figure 2 to Example 2. SSARX and

#	$t - t_0$	$T - t = \nu$	N
Fig. 1	10	10	1000
Fig. 2	10	10	1000
Fig. 2	30	10	3000

Table 1
Parameters chosen in the implementation.

PBSID are indistinguishable as predicted by the theory in this paper in Example 1. As far as Example 2 is concerned, while equivalence does not hold for small $t - t_0$ and N , it does indeed hold (see figure 2, right plot) when increasing N and $t - t_0$.

4.2 Vector AutoRegressive implementation of the PBSID_{opt} method

The second main result of this paper shows that, indeed, the PBSID_{opt} can be efficiently implemented via VARX estimation. Even though VARX models were introduced also in previous contributions, among which [37,35,25], in our framework the VARX models pop up quite naturally from a theoretically sound “optimized” method. This consideration constitutes, in the author’s opinion, a starting point for future investigations.

Theorem 4.3 Consider the infinite VARX model

$$y_t = \sum_{i=1}^{\infty} \Phi_i z_{t-i} + e_t \quad (4.9)$$

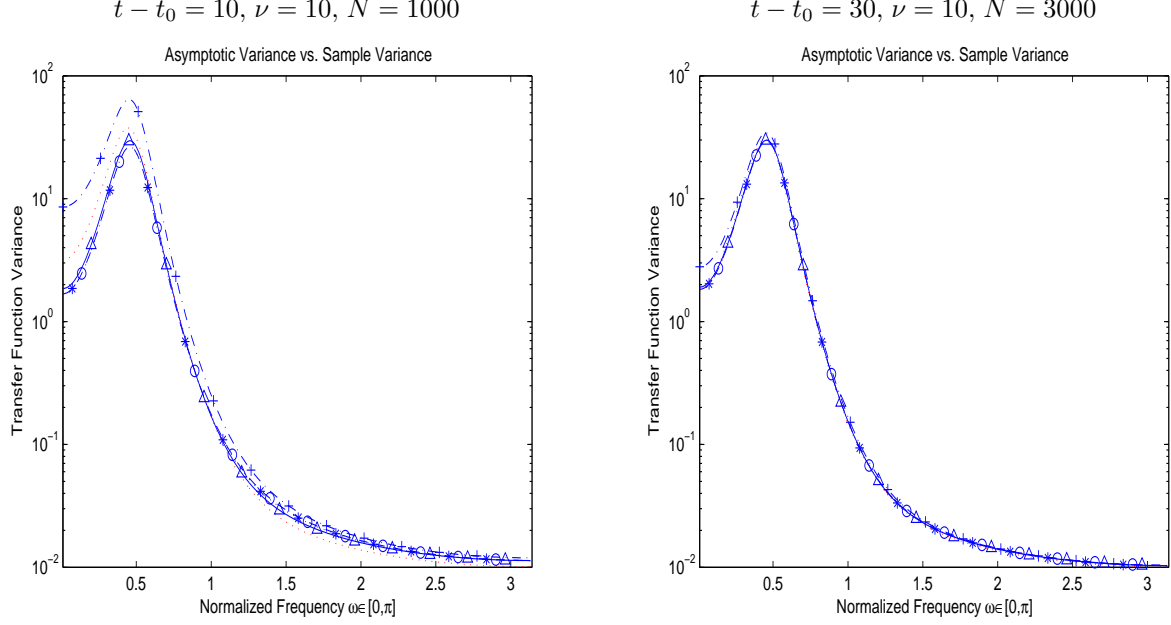


Fig. 2. Example 2. Left: “small” $t - t_0 = 10$. Right: “large” $t - t_0$. Variance (and its Monte Carlo estimate) vs. normalized frequency ($\omega \in [0, \pi]$). Solid with triangles (\triangle) PEM. Dashed with crosses (+) PBSID. Dashed with circles (\circ) PBSID_{opt} . Dotted with stars (*): SSARX. Dotted: asymptotic variance for PBSID.

and denote with $\hat{\Phi}_i$, $i = 1, \dots, t - t_0$, the estimators of the first $t - t_0$ coefficients in (4.9) obtained solving¹⁸

$$Y_t^{\nu+N} \simeq \sum_{i=1}^{t-t_0} \Phi_i Z_{t-i}^{\nu+N} \quad (4.10)$$

in the least squares sense.

The “optimally-weighted” solution to (3.15), i.e. the one that yields the least asymptotic variance of the estimators $\hat{\Omega}^{P_{opt}}$ among all linear, asymptotically unbiased estimators of Ω^P based on the regression (3.15), is equivalent to estimating the VARX model (4.10) in the sense that:

$$\begin{bmatrix} \hat{\Xi}_0^{P_{opt}} \\ \hat{\Xi}_1^{P_{opt}} \\ \vdots \\ \hat{\Xi}_\nu^{P_{opt}} \end{bmatrix} = \begin{bmatrix} \hat{\Phi}_{t-t_0} & \dots & \hat{\Phi}_{T-t_0} & \dots & \hat{\Phi}_1 \\ 0 & \hat{\Phi}_{t-t_0} & \dots & \dots & \hat{\Phi}_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & \hat{\Phi}_{t-t_0} & \dots & \hat{\Phi}_{\nu+1} \end{bmatrix} \quad (4.11)$$

and

$$\hat{\Psi}_{ji}^{P_{opt}} = \hat{\Phi}_i \quad (4.12)$$

¹⁸ Note that the estimators are function of N and $t - t_0$, which according to Assumption 1 grows with N . In order to streamline notation this dependence is not made explicit.

Proof. See the Appendix □

Remark 4.4 It is worth mentioning that, with the “optimally weighted” (Markov) estimator of the coefficients Ξ_i, Ψ_{ij} , the estimate of the lower triangular matrix in (3.11) is indeed Toeplitz (see eq. (4.12)). It is also interesting to note that the estimate of the VARX coefficients weighting the “far” past (i.e. Ψ_{ji} for $i > t - t_0$ in (3.2)) are set to zero by the “optimal” estimator (i.e. $\hat{\Psi}_{ji}^{P_{opt}} = 0$ for $i > t - t_0$). This is reasonable since, according to Assumption 1, for $i > t - t_0$ the Ψ_{ji} ’s go to zero faster than $1/\sqrt{N}$; on the contrary, estimating these coefficients would lead to errors which are of order $1/\sqrt{N}$ in probability. This also brings up the question of choosing the length of the past horizon $t - t_0$; the analysis of this paper gives, together with the results in [18], a more theoretically sound foundation to the (usually adopted) practice of determining $t - t_0$ using standard order selection criteria [33,22,45] for vector autoregressive models (see, e.g. [1]). The reader is also referred to the recent paper [29] which discusses automatic inference for infinite order autoregressions.

Using the result of Theorem 4.3 the PBSID_{opt} algorithm can be implemented as follows:

- (a) Estimate the VARX model (4.9) as described in (4.10); this may include estimation of the appropriate $t - t_0$ using standard criteria for VARX order estimation.

- (b) Use the estimated coefficients as described in formulas (4.11) and (4.12) to form the predictors

$$\hat{Y}_{t+h}^{P_{opt}} = \sum_{i=h+1}^{t-t_0+h} \hat{\Psi}_{hi}^{P_{opt}} Z_{t+h-i} = \sum_{i=h+1}^{t-t_0} \hat{\Phi}_i Z_{t+h-i}; \quad (4.13)$$

the state sequences $\hat{X}_t^{P_{opt}}$ and $\hat{X}_{t+1}^{P_{opt}}$ are then obtained as described in formulas (3.17) and (3.18).

This implementation has a much lower computational complexity w.r.t. the implementation described in [7] which involves solving the least squares problem (3.15) directly.

In fact, step a) above involves the estimation of a VARX model of length $t - t_0$ (which, according to Assumption 1, is $O(\log(N))$); solving (4.10) has complexity $O(N(\log N)^2)$ (see [20] pag. 248). The order and state estimation (step b) above) can be performed on the “squared” version of the matrix $\hat{Y}_{[t,T]}^{P_{opt}}$. This second step is common to all subspace algorithms. Instead step a) has the same “order” of complexity than, e.g., CCA and PBSID; however both these algorithms essentially estimate ν long VARX models, increasing the complexity of the first step roughly by a factor ν .

Hence the implementation described above of the $PBSID_{opt}$ compares favorably to a variety of subspace procedures (among which PBSID or CCA) as far computational complexity is concerned while, according to Theorem 5.3 in [7], yielding lower asymptotic variance than CCA. We remind also that the $PBSID_{opt}$ algorithm works (i.e. is consistent) regardless of the presence of feedback.

These considerations make the algorithm described above a strong alternative to standard used methods for a variety of reasons, among which computational complexity and asymptotic statistical properties (it is consistent also in closed loop and gives lower variance than the original PBSID and CCA).

Remark 4.5 [PBSID_{opt} vs. SSARX]

The main differences between the $PBSID_{opt}$ and SSARX algorithms are as follows: (i) the length of the ARX model estimated is (in general) different for the two methods; in particular SSARX uses order larger than ν (but in [25] it is just required that the order be “high” to ensure consistency), possibly chosen according to infinite-order ARX models selection rules [29]; only the first ν coefficients are then used; instead $PBSID_{opt}$ the order is exactly $t - t_0$; this results in the $PBSID_{opt}$ filling with zeros the Toeplitz matrix used to construct the bank of predictors (see equation (4.11)); (ii) the SSARX methods projects the “corrected future” to form $\hat{Y}_{[t,T]}^S$ (see eq. (3.9)) while the $PBSID_{opt}$ uses directly the estimated coefficients from the VARX modeling step to form the bank

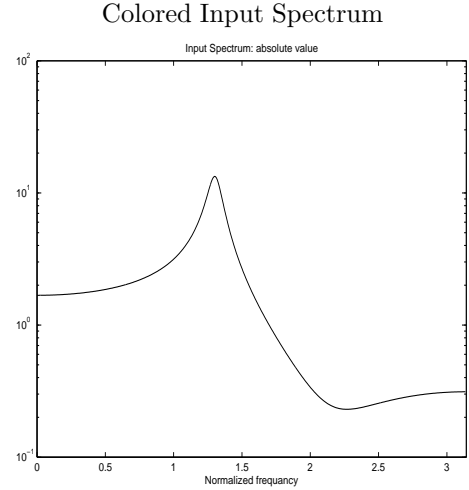


Fig. 3. Colored input spectrum: absolute value.

of predictors $\hat{Y}_{t+k}^{P_{opt}}$ (see eq. 3.17 and (4.11)). This makes $PBSID_{opt}$ even more advantageous from the computational point of view, since it does not require computing the projection (3.9).

We consider the following examples, frequently used in the literature of subspace identification, to illustrate the result.

The first is an “open loop” experiment which contains all the essential features of the “optimized” method i.e.: (i) it is not efficient (it does not reach Cramér Rao) and (ii) it gives (strictly) lower asymptotic variance than CCA. Of course this example is performed in “open loop” to allow the comparison with CCA. In this example the original PBSID and the “optimized” version have the same asymptotic behavior.

We consider the first order ARMAX model

$$\mathbf{y}(t) - 0.5\mathbf{y}(t-1) = \mathbf{u}(t-1) + \mathbf{e}(t) + 0.5\mathbf{e}(t-1)$$

The input is unit variance white noise passed through the filter $H_u(z)$

$$H_u(z) = \frac{z^2 + 0.8z + 0.55}{z^2 - 0.5z + 0.9};$$

the input spectrum is plotted in Figure 3.

We report in figure 4 results concerning the asymptotic variance and the sample variance estimated over 100 Monte Carlo runs multiplied by the number $N = 1000$ of data points used in each experiment of the deterministic transfer function $F(z) = \frac{1}{z-0.5}$

As a second example we consider a fifth order (marginally stable) system in state space form (2.1) where (see

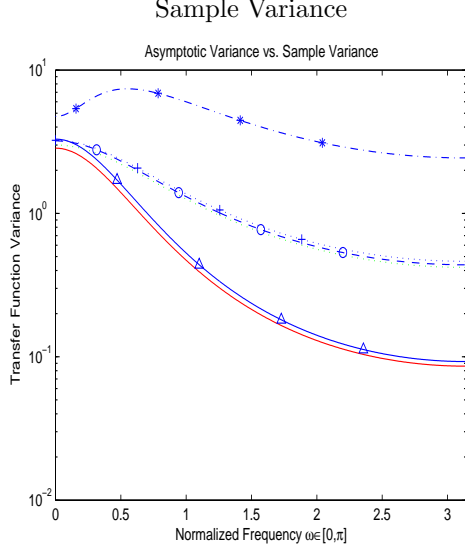


Fig. 4. Asymptotic Variance (and its Monte Carlo estimate) vs. normalized frequency ($\omega \in [0, \pi]$) (*ARMAX of order 1*). Solid with triangles (\triangle) PEM, dashed-dotted with stars ($*$): CCA, dotted with crosses ($+$): (PBSID), dashed with circles (o): PBSID_{opt} ; dotted: asymptotic variance for PBSID, solid: Cramér Rao lower bound.

[49,48]):

$$A = \begin{bmatrix} 4.40 & 1 & 0 & 0 & 0 \\ -8.09 & 0 & 1 & 0 & 0 \\ 7.83 & 0 & 0 & 1 & 0 \\ -4 & 0 & 0 & 0 & 1 \\ 0.86 & 0 & 0 & 0 & 0 \end{bmatrix} \quad C = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} 0.00098 & 0.01299 & 0.01859 & 0.0033 & -0.00002 \end{bmatrix}^\top$$

$$K = \begin{bmatrix} 2.3 & -6.64 & 7.515 & -4.0146 & 0.86336 \end{bmatrix}^\top \quad D = 0$$

and $\mathbf{e}(t)$ is unit variance white noise. The input \mathbf{u} is generated in closed loop by $\mathbf{u}(t) = 5\mathbf{r}(t) - H(z)\mathbf{y}(t)$; the reference signal $\mathbf{r}(t)$ is unit variance white noise; $H(z)$ is given by:

$$H(z) = \frac{0.63 - 2.083z^{-1} + 2.8222z^{-2} - 1.865z^{-3} + 0.4978z^{-4}}{1 - 2.65z^{-1} + 3.11z^{-2} - 1.75z^{-3} + 0.39z^{-4}}$$

We have also used $t - t_0 = 30$, $\nu = 10$ and $N = 2000$.

In this example PBSID_{opt} outperforms PBSID in the low frequency band and performs slightly better than the innovation estimation method (IEM hereafter) [41] in the high frequency band (see figure 4). The algorithm by Ljung and McKelvey, which as shown in the next section is a weighted version of PBSID_{opt} , performs worse than PBSID_{opt} and IEM [41].

Sample Variance (100 Monte Carlo runs)

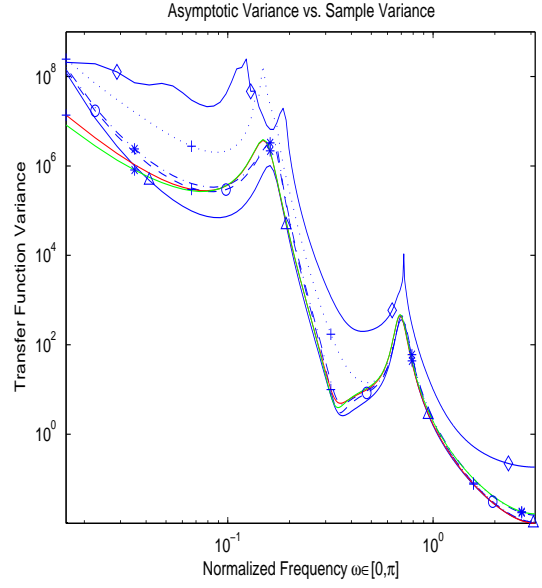


Fig. 5. Asymptotic variance (and its Monte Carlo estimate) Dashed-dotted with circles (o): PBSID_{opt} . Dashed-dotted with stars ($*$): IEM [41]. Dotted with crosses ($+$): PBSID. Solid with diamonds (\diamond): Ljung-McKelvey [35]. Dotted with triangles (\triangle): PEM. Solid with stars (green) ($*$): asymptotic variance for IEM. Dotted with crosses (red) ($+$): asymptotic variance for PBSID.

It has been checked that the original algorithm presented in [7] and its alternative implementation presented in this paper give indeed the same result. In particular conditions (4.11) and (4.12) have been verified to hold for the estimated coefficients of the PBSID_{opt} described in [7].

4.3 Relation with the method by Ljung and McKelvey

In this Section we shall briefly discuss the relation of PBSID_{opt} with the algorithm presented in [35] by Ljung and McKelvey. We shall not enter into a detailed description of the algorithm for which we refer the reader to the original paper; our description of the algorithm follows the Matlab code provided in [34].

Suffices here to say that the first step is to construct a matrix $\hat{Y}_{[t,T]}^{LK}$ formed with predictors from which a basis of the state space is extracted; $\hat{Y}_{[t,T]}^{LK}$ shall play the same role as $\hat{Y}_{[t,T]}^{Pop}$ and $\hat{Y}_{[t,T]}^S$ in PBSID_{opt} and SSARX respectively.

The main result can be stated as follows:

Proposition 4.6 *Assume the model orders n_a and n_b in (4.14) are chosen according to $n_a = n_b = t - t_0$ and the VARX coefficients $H_i := [H_{y,i} \ H_{u,i}]$ in (4.14) are estimated letting $\hat{H}_i := \hat{\Phi}_i$ where $\hat{\Phi}_i$ are the least squares solution of (4.10). Then the algorithm proposed in Ljung and McKelvey [35] is a weighted version of PBSID_{opt} , in the sense that the two*

state construction steps differ only for the choice of a (row) weighting matrix W_{LK} , as made precise by formula (4.20).

Proof. Consider the VARX model

$$\hat{Y}_{t|t-1} = \sum_{i=1}^{n_a} \hat{H}_{y,i} Y_{t-i} + \sum_{i=1}^{n_b} \hat{H}_{u,i} U_{t-i} \quad (4.14)$$

Essentially the algorithm in [35] construct the state space using a bank of predictors¹⁹

$$\hat{Y}_{[t,T]}^{LK} := \left[\hat{Y}_{t|t-1}^\top \hat{Y}_{t+1|t-1}^\top \cdots \hat{Y}_{T|t-1}^\top \right]^\top \quad (4.15)$$

where $\hat{Y}_{t+k|t-1}$ is computed recursively as

$$\begin{aligned} \hat{Y}_{t+k|t-1} &:= \sum_{i=1}^k \hat{H}_{y,i} \hat{Y}_{t+k-i|t-1} + \sum_{i=k+1}^{n_a} \hat{H}_{y,i} Y_{t+k-i} \\ &\quad + \sum_{i=k+1}^{n_b} \hat{H}_{u,i} U_{t+k-i} \end{aligned} \quad (4.16)$$

The remaining steps (i.e. state construction and estimation of A, B, C, K) follow the same lines as described in the previous Sections.

In order to make clear the link between the ‘‘predictor’’ used in $PBSID_{opt}$ and $\hat{Y}_{[t,T]}^{LK}$, we rewrite (4.16) as follows:

$$\begin{aligned} \hat{Y}_{t+k|t-1} - \sum_{i=1}^k \hat{H}_{y,i} \hat{Y}_{t+k-i|t-1} &= \sum_{i=k+1}^{n_a} \hat{H}_{y,i} Y_{t+k-i} + \\ &\quad + \sum_{i=k+1}^{n_b} \hat{H}_{u,i} U_{t+k-i} \end{aligned} \quad (4.17)$$

Using the assumption that $n_a = n_b = t - t_0$, letting $\hat{H}_i := [\hat{H}_{y,i} \hat{H}_{u,i}]$ and defining

$$W_{LK} := \begin{bmatrix} I & 0 & \cdots & 0 \\ -\hat{H}_{y,1} & I & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ -\hat{H}_{y,\nu} & \cdots & -\hat{H}_{y,1} & I \end{bmatrix}, \quad (4.18)$$

equation (4.17) can be rewritten in matrix form as follows:

$$W_{LK} \hat{Y}_{[t,T]}^{LK} = \begin{bmatrix} H_{t-t_0} & \cdots & \cdots & \cdots & \cdots & H_1 \\ 0 & H_{t-t_0} & \cdots & \cdots & \cdots & H_2 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & H_{t-t_0} & \cdots & H_{\nu+1} \end{bmatrix} Z_{[t_0,t]}; \quad (4.19)$$

from the assumption that the VARX model (4.14) has been estimated using the same data as (4.9) it also follows that $\hat{H}_i = \hat{\Phi}_i$. Therefore, using (4.19) the stacked predictors in

(4.15) can be rewritten as

$$\begin{aligned} \hat{Y}_{[t,T]}^{LK} &= W_{LK}^{-1} \begin{bmatrix} \hat{\Phi}_{t-t_0} & \cdots & \cdots & \cdots & \cdots & \hat{\Phi}_1 \\ 0 & \hat{\Phi}_{t-t_0} & \cdots & \cdots & \cdots & \hat{\Phi}_2 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \hat{\Phi}_{t-t_0} & \cdots & \hat{\Phi}_{\nu+1} \end{bmatrix} Z_{[t_0,t]} \\ &= W_{LK}^{-1} \hat{Y}_{[t,T]}^{P_{opt}} \end{aligned} \quad (4.20)$$

where (4.13) has been used in the last equality. \square

It is remarkable that the bank of predictors used in the original paper [34] is indeed equivalent to a weighted version of the bank of predictors used in $PBSID_{opt}$. It should not be surprising that the original algorithm in [34,35] does not perform as well as $PBSID_{opt}$ (see figure 4); in fact it is well known that the (row) weighting does affect the asymptotic statistical properties of the estimators using the ‘‘state sequence’’ approach (see [2]). Note that in [42] it was conjectured (even though not proved) that an algorithm named HOARX is equivalent (asymptotically) to the algorithm by Ljung and McKelvey. As shown in this section (see, in particular, the right hand side of (4.19) and formula (8) in [42]) this is not exactly true; instead they differ up to a row weighting. To be precise also some Markov parameters (those weighting the ‘‘far past’’) are set to zero (see (4.19)). Using the techniques in [29] it would be possible to see that this does not make any difference asymptotically as long as ν remains fixed (or bounded above).

5 Conclusions

In this paper we have discussed several subspace algorithms based on predictor model identification. It is shown that all these algorithms can be formulated as VARX estimation followed by model reduction.

In particular it has been shown that

- (a) SSARX [25] and PBSID [16] are asymptotically equivalent;
- (b) $PBSID_{opt}$ [7] is exactly equivalent (i.e. give the same numerical results on finite data) to estimating a suitable VARX model followed by the usual steps of subspace identification (i.e. state estimation via SVD followed by estimation of the system matrices)
- (c) The algorithm by Ljung-McKelvey [34] is a weighted version of $PBSID_{opt}$.

Experimental results (on simulated data) are included which support the theoretical derivations. The results of this paper, together with the comparison performed in [7], indicate that $PBSID_{opt}$ should be considered as one of the most appealing in this class of algorithms for the following reasons:

- (a) it is consistent under closed-loop operating conditions

¹⁹ This implementation has been taken from the Matlab code reported in [34].

- (b) it performs no worse than CCA (regardless of the choice of input) for open loop data and better than SSARX/PBSID with feedback
- (c) it is prone to a very simple and computationally attractive implementation via VARX modeling.

The simulation results reported in this paper seem to support these statements. Even though PBSID_{opt} can be verified to be asymptotically efficient in a number of examples²⁰, it is not so in general.

It was conjectured in [31] that an algorithm which is essentially equivalent to PBSID is (asymptotically) efficient for large ν (actually for $t-t_0 = \nu \rightarrow \infty$). Instead, as clearly seen in figure 4 PBSID is not efficient for large ν . We have verified that indeed the performance does not change increasing ν . Instead PBSID is nearly efficient for this example with $\nu = 1$ (see [7], figure 4 and [6].)

There is certainly much work to be done; in particular not completely clear is, at the moment, the relation of these methods with the new algorithms introduced in [42,36] and with the IEM of [41].

Also the question of finite-data behavior is certainly of interest and deserves, in our opinion, further investigation.

6 Acknowledgments

The author would like to thank Manfred Deistler, Lennart Ljung, Giorgio Picci and Joe Qin for fruitful discussions on the subject. Also an anonymous reviewer of the paper [7] is gratefully acknowledged for comments and suggestions contained in his report.

Appendix A: Proofs

Proof of Theorem 4.3. The proof makes use of the fine structure of the matrix L . Let us denote with L_I a matrix which columns span the image of L and with L_K a matrix spanning the left kernel of L , so that $[L_I L_K]$ is a full rank square ($pN(\nu+1) \times pN(\nu+1)$) matrix. The least squares problem (3.15) can be transformed into the equivalent form

$$\begin{bmatrix} L_I^\top \\ L_K^\top \end{bmatrix} Y = \begin{bmatrix} L_I^\top \\ L_K^\top \end{bmatrix} S^P \Omega^P + \begin{bmatrix} L_I^\top \\ L_K^\top \end{bmatrix} E \quad (\text{A.1})$$

Note that, by construction, $L_I^\top L$ has full rank and therefore $L_I^\top E$ has full rank covariance. Similarly $L_K^\top L = 0$ and hence $L_K^\top E = 0$ (in the mean square sense).

In this way the least squares problem (3.15) with singular noise covariance (3.15) is transformed into a least squares problem with full rank noise covariance (the ‘‘top’’ part of (A.1)) and equality constraints (the ‘‘bottom’’ part of (A.1)).

²⁰ This might require that also ν grows with N . However the analysis in this paper deals only with the case of fixed ν .

It is easy to show that L_I can be chosen to be a selection matrix so that $L_I^\top Y = \text{vec}(Y_t^{N+\nu}) := Y_I$. For future reference observe that $L_I^\top E = \text{vec}(E_t^{N+\nu}) := E_I$ so that (A.1) can be rewritten as

$$\begin{aligned} Y_I &= L_I^\top S^P \Omega^P + E_I \\ \text{s.t. } L_K^\top Y &= L_K^\top S^P \Omega^P \end{aligned} \quad (\text{A.2})$$

Let us introduce the pair of indexes (j, \bar{j}) such that $\nu \geq \bar{j} > j \geq 0$ and define $\delta := \bar{j} - j$. Then it is easy to see that there exist matrices $L_K(j, \bar{j}, l)$ so that

$$\begin{aligned} 0 = L_K^\top(j, \bar{j}, l) Y &= \Xi_j \begin{bmatrix} z_{t_0+\delta+l} \\ z_{t_0+\delta+1+l} \\ \vdots \\ z_{t+\delta-1+l} \end{bmatrix} - \Xi_{\bar{j}} \begin{bmatrix} z_{t_0+l} \\ z_{t_0+1+l} \\ \vdots \\ z_{t-1+l} \end{bmatrix} + \\ &+ \sum_{k=1}^j \Psi_{jk} z_{t+\delta-1+k+l} - \sum_{k=1}^{\bar{j}} \Psi_{\bar{j}k} z_{t-1+k+l} \end{aligned} \quad (\text{A.3})$$

where, for each pair (j, \bar{j}) with $\bar{j} > j \in [0, \nu-1]$, l ranges in the interval $[0, N-\delta-1]$.

It is possible to extract exactly $N\nu - (\nu+1)$ independent constraints (recall that L_K has rank $p(N\nu - (\nu+1))$) of the form (A.3) by letting $j \in [0, \nu-1]$, $\bar{j} = j' := j+1$, (so that $\delta = 1$) and $l \in [0, N-2]$. With these choices the constraints (A.3) can be written in the form:

$$\begin{aligned} \left[\Xi_j \Psi_{j1} \dots \Psi_{jj} \right] Z_{[t_0+\delta, t+\delta+j-1]}^{N-1} &= \\ = \left[\Xi_{j'} \Psi_{j'1} \dots \Psi_{j'j'} \right] Z_{[t_0, t+j'-1]}^{N-1} \end{aligned} \quad (\text{A.4})$$

Recalling that $\delta = j' - j = 1$, and defining 0_δ to be the zero matrix of size $p \times \delta(p+m)$, we can rewrite (A.4) in the form

$$\begin{aligned} \left[0_1 \Xi_j \Psi_{j1} \dots \Psi_{jj} \right] Z_{[t_0, t+j'-1]}^{N-1} &= \\ \left[\Xi_{j'} \Psi_{j'1} \dots \Psi_{j'j'} \right] Z_{[t_0, t+j'-1]}^{N-1} \end{aligned} \quad (\text{A.5})$$

From the assumption that the joint spectrum is coercive, it follows that, for N large enough (i.e. and hence $N-1$ large), the matrix $Z_{[t_0, t+j'-1]}^{N-1}$ is of full row rank for all possible choices of j ; therefore (A.5) is equivalent to the ‘‘dual’’ equation for the coefficients:

$$\left[0_1 \Xi_j \Psi_{j1} \dots \Psi_{jj} \right] = \left[\Xi_{j'} \Psi_{j'1} \dots \Psi_{j'j'} \right]$$

As mentioned above this should hold for each pair (j, j') , $\nu > j \geq 0$; this is equivalent to the following constraints on the estimated coefficients:

$$\begin{aligned} \begin{bmatrix} \hat{\Xi}_1^{P_{opt}} & \hat{\Psi}_{11}^{P_{opt}} \end{bmatrix} &= \begin{bmatrix} 0_1 & \hat{\Xi}_0^{P_{opt}} \end{bmatrix} \\ \begin{bmatrix} \hat{\Xi}_2^{P_{opt}} & \hat{\Psi}_{22}^{P_{opt}} & \hat{\Psi}_{21}^{P_{opt}} \end{bmatrix} &= \begin{bmatrix} 0_2 & \hat{\Xi}_0^{P_{opt}} \end{bmatrix} \\ \vdots & \vdots \\ \begin{bmatrix} \hat{\Xi}_\nu^{P_{opt}} & \hat{\Psi}_{\nu\nu}^{P_{opt}} & \dots & \hat{\Psi}_{\nu 1}^{P_{opt}} \end{bmatrix} &= \begin{bmatrix} 0_\nu & \hat{\Xi}_0^{P_{opt}} \end{bmatrix} \end{aligned} \quad (\text{A.6})$$

For convenience, let us define $\Xi_0 = \begin{bmatrix} \Phi_{t-t_0} & \dots & \Phi_2 & \Phi_1 \end{bmatrix}$. Using the constraints above, some extra algebra will show that (A.2) can be written in the form

$$Y_I = \begin{bmatrix} \Phi_{t-t_0} & \dots & \Phi_1 & 0 & \dots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \dots & 0 & \Phi_{t-t_0} & \dots & \Phi_1 \end{bmatrix} \text{vec} \left(Z_{t_0}^{N+\nu} \right) + E_I. \quad (\text{A.7})$$

A more compact expression of (A.7) is:

$$Y_t^{N+\nu} = \begin{bmatrix} \Phi_{t-t_0} & \dots & \Phi_2 & \Phi_1 \end{bmatrix} Z_{(t_0,t)}^{N+\nu} + E_t^{N+\nu} \quad (\text{A.8})$$

The “optimal” (Markov) solution to (A.7) is obtained by pre-whitening the residual vector E_I , which can be obtained pre-multiplying by $(I \otimes \Lambda^{-1/2})$ both sides of (A.7) or, equivalently, pre-multiplying both sides of (A.8) by $\Lambda^{-1/2}$. It is a simple calculation to check that, indeed, solving in the least squares sense

$$\Lambda^{-1/2} Y_t^{N+\nu} \simeq \Lambda^{-1/2} \begin{bmatrix} \Phi_{t-t_0} & \dots & \Phi_2 & \Phi_1 \end{bmatrix} Z_{(t_0,t)}^{N+\nu}$$

is equivalent to solving (4.10); this implies that $\hat{\Xi}_0^{P_{opt}} = \begin{bmatrix} \hat{\Phi}_{t-t_0} & \dots & \hat{\Phi}_2 & \hat{\Phi}_1 \end{bmatrix}$; conditions (4.11) and (4.12) can then be obtained using the constraints in (A.6), which completes the proof. \square

References

- [1] D. Bauer, *Order estimation for subspace methods*, *Automatica* **37** (2001), 1561–1573.
- [2] ———, *Asymptotic properties of subspace estimators*, *Automatica* **41** (2005), 359–376.
- [3] ———, *Comparing the CCA subspace method to pseudo maximum likelihood methods in the case of no exogenous inputs*, *Journal of Time Series Analysis* **26** (2005), 631–668.
- [4] D. Bauer and M. Jansson, *Analysis of the asymptotic properties of the MOESP type of subspace algorithms*, *Automatica* **36** (2000), 497–509.
- [5] D. Bauer and L. Ljung, *Some facts about the choice of the weighting matrices in Larimore type of subspace algorithm*, *Automatica* **38** (2002), 763–773.
- [6] A. Chiuso, *Some insights on the choice of the future horizon in CCA-type subspace algorithms*, Submitted to ACC 2007, available at <http://www.dei.unipd.it/~chiuso>.
- [7] ———, *On the relation between CCA and predictor based subspace identification*, Submitted to IEEE Trans. on Aut. Control (2005), available at <http://www.dei.unipd.it/~chiuso>.
- [8] ———, *Asymptotic equivalence of certain closed-loop subspace identification methods*, Proc. of SYSID 2006 (Newcastle, Australia), 2006.
- [9] ———, *Asymptotic variance of closed-loop subspace identification algorithms*, IEEE Trans. on Aut. Control **51** (2006), no. 8, 1299–1314.
- [10] ———, *The role of Vector AutoRegressive modeling in subspace identification*, Proc. of CDC 2006 (San Diego (USA)), Dec. 2006.
- [11] A. Chiuso and G. Picci, *Constructing the state of random processes with feedback*, Proc. of the IFAC Int. Symposium on System Identification (SYSID) (Rotterdam), August 2003.
- [12] ———, *Geometry of oblique splitting, minimality and hankel operators*, Lect. Notes in Control and Information Sciences, no. 286, pp. 85–124, Springer, 2003.
- [13] ———, *The asymptotic variance of subspace estimates*, *Journal of Econometrics* **118** (2004), no. 1-2, 257–291.
- [14] ———, *Asymptotic variance of subspace methods by data orthogonalization and model decoupling: A comparative analysis*, *Automatica* **40** (2004), no. 10, 1705–1717.
- [15] ———, *On the ill-conditioning of subspace identification with inputs*, *Automatica* **40** (2004), no. 4, 575–589.
- [16] ———, *Consistency analysis of some closed-loop subspace identification methods*, *Automatica* **41** (2005), no. 3, 377–391.
- [17] ———, *Prediction error vs. subspace methods in closed-loop identification*, Proc. of the 16th IFAC World Congress (Prague), July 2005.
- [18] A. Dahlén and W. Scherrer, *The relation of CCA subspace method to a balanced reduction of an autoregressive model*, *Journal of Econometrics* **118** (2004), no. 1-2, 293–312.
- [19] T. Ferguson, *A course in large sample theory*, Chapman and Hall, 1996.
- [20] G.H. Golub and C.R. Van Loan, *Matrix computation*, 2nd ed. ed., The Johns Hopkins Univ. Press., 1989.
- [21] C.W.J. Granger, *Economic processes involving feedback*, *Information and Control* **6** (1963), 28–48.
- [22] E.J. Hannan and M. Deistler, *The statistical theory of linear systems*, Wiley, 1988.
- [23] Z. Hong, M. Harmse, J. Guiver, and W. Canney, *Subspace identification in industrial MPC applications - a review of recent progress and industrial experience (i)*, Proc. of SYSID 2006 (Newcastle, Australia), 2006.
- [24] M. Jansson, *Asymptotic variance analysis of subspace identification methods*, Proceedings of SYSID2000 (S. Barbara Ca.), 2000.
- [25] ———, *Subspace identification and ARX modeling*, Proceedings of SYSID 2003 (Rotterdam), 2003.
- [26] ———, *A new subspace identification method for open and closed loop data*, Proc. of the 16th IFAC World Congress (Prague), 2005.
- [27] T. Katayama, H. Kawauchi, and G. Picci, *Subspace identification of closed loop systems by orthogonal decomposition*, *Automatica* **41** (2005), 863–872.

- [28] T. Katayama, H. Tanaka, and T. Enomoto, *A simple subspace identification method of closed-loop systems using orthogonal decomposition*, Proc. of the 16th IFAC World Congress (Prague), 2005.
- [29] G.M. Kuersteiner, *Automatic inference for infinite order vector autoregressions*, *Econometric Theory* **21** (2005), 85–115.
- [30] W.E. Larimore, *System identification, reduced-order filtering and modeling via canonical variate analysis*, Proc. American Control Conference, 1983, pp. 445–451.
- [31] ———, *Large sample efficiency for ADAPTX subspace identification with unknown feedback*, Proc. of IFAC DYCOPS'04 (Boston, MA, USA), 2004.
- [32] A. Lindquist and G. Picci, *Canonical correlation analysis, approximate covariance extension and identification of stationary time series*, *Automatica* **32** (1996), 709–733.
- [33] L. Ljung, *System identification, theory for the user*, Prentice Hall, 1997.
- [34] L. Ljung and T. McKelvey, *Subspace identification from closed loop data*, Tech. Report LiTH-ISY-R-1752, 1995.
- [35] ———, *Subspace identification from closed loop data*, *Signal Processing* **52** (1996), no. 2, 209–216.
- [36] K. Onodera, G. Emoto, and S.J. Qin, *A new subspace identification method for closed loop systems*, Proceedings of SYSID 2006 (Newcastle, Australia), 2006.
- [37] K. Peternell, *Subspace methods for subspace identification*, Ph.D. thesis, Technical University of Vienna, 1995.
- [38] K. Peternell, W. Scherrer, and M. Deistler, *Statistical analysis of novel subspace identification methods*, *Signal Processing* **52** (1996), 161–178.
- [39] B. Porat and B. Friedlander, *Asymptotic accuracy of ARMA parameter estimation methods based on sample covariances*, Proceedings of the 7th IFAC SYSID (York), 1985.
- [40] S.J. Qin and L. Ljung, *Closed-loop subspace identification with innovation estimation*, Proceedings of SYSID 2003 (Rotterdam), 2003.
- [41] ———, *Parallel QR implementation of subspace identification with parsimonious models*, Proceedings of SYSID 2003 (Rotterdam), 2003.
- [42] ———, *On the role of future horizon in closed-loop subspace identification*, Proceedings of SYSID 2006 (Newcastle, Australia), 2006.
- [43] C.R. Rao, *Representations of the best linear unbiased estimators in the Gauss-Markov model with a singular dispersion matrix*, *J. Multivariate Anal.* **3** (1973), 276–292.
- [44] F. Shi and J.F. MacGregor, *A framework for subspace identification*, Proc. of IEEE ACC (Arlington, VA), 2001.
- [45] T. Söderström and P. Stoica, *System identification*, Prentice-Hall, 1989.
- [46] P. Van Overschee and B. De Moor, *Subspace algorithms for the stochastic identification problem*, *Automatica* **29** (1993), 649–660.
- [47] ———, *N4SID: Subspace algorithms for the identification of combined deterministic–stochastic systems*, *Automatica* **30** (1994), 75–93.
- [48] ———, *Closed loop subspace systems identification*, Proceedings of 36th IEEE Conference on Decision & Control (San Diego, CA), 1997, pp. 1848–1853.
- [49] M. Verhaegen, *Application of a subspace model identification technique to identify lti systems operating in closed-loop*, *Automatica* **29** (1993), 1027–1040.
- [50] B. Wahlberg, *Estimation of autoregressive moving-average models via high order autoregressive approximations*, *Journal of Time Series Analysis* **10** (1989), no. 3, 283–299.
- [51] A.M. Walker, *Large-sample estimation of parameters for moving average models*, *Biometrika* **48** (1961), no. 3/4, 343–357.
- [52] H.J. Werner and C. Yapar, *On inequality constrained generalized least squares selections in the general possibly singular Gauss-Markov model: A projector theoretical approach*, *Linear Algebra and its Applications* **237/238** (1996), no. 1–3, 359–393, Special issue honoring Calyampudi Radhakrishna Rao.
- [53] Y. Zhu, *System identification for process control: Recent experiences and outlook*, Plenary Lecture delivered at SYSID 2006 (Newcastle, Australia), 2006.