

# A software package for system identification in the behavioral setting

Ivan Markovsky

Department ELEC, Vrije Universiteit Brussel  
Pleinlaan 2, Building K, B-1050 Brussels, Belgium  
Fax: +32-2-629.28.50 , Email: imarkovs@vub.ac.be

## Abstract

An identification problem with no a priori separation of the variables into inputs and outputs and representation invariant approximation criterion is considered. The model class consists of linear time-invariant systems of bounded complexity and the approximation criterion is the minimum of a weighted 2-norm distance between the given time series and a time series that is consistent with the model. The problem is equivalent to and is solved as a mosaic-Hankel structured low-rank approximation problem. Software implementing the approach is developed and tested on benchmark problems. Additional nonstandard features of the software are specification of exact and missing variables and identification from multiple experiments.

**Keywords:** system identification; model reduction; behavioral approach; missing data; low-rank approximation; total least squares; reproducible research; DAISY.

## 1 Introduction

One of the main criticisms to the behavioral approach in system theory and control [24, 17] is that it is not supported by numerical algorithms and software. Despite its conceptual and theoretic advantages, the behavioral approach has produced only a few engineering tools; most notably the identification methods of [19] and [10]. Compared to the abundance of algorithms and software in the classical setting, this scarcity of tools accounts to a large extent for the slow acceptance of the behavioral approach by engineers.

This paper presents a flexible method for system identification in the behavioral setting, which is an outgrowth of the method in [10]. The method is implemented in a publicly available software package which is fast, easy to use, and has the following nonstandard features.

### 1. *Representation free problem formulation*

Classical identification problems are invariably defined in terms of representations of the system—transfer function, input/state/output, convolution kernel, *etc.* A problem can be equivalently posed in frequency domain as well as time domain, using input/state/output, convolution, or any other representation. Different representations have different advantages, so that all of them should be considered as options in solving the problem. The model representation should be viewed as an implementation detail of the particular solution method. A problem formulation that is representation free is defined in terms of the desired system behavior—set of trajectories of the system.

Related to the use of a representation invariant problem formulation, is the assumption in classical system identification problems that an input/output partitioning of the variables is given and fixed. The choice of inputs and outputs is not unique [25, Section VIII]. Different input/output partitions, generically, lead to the same optimal model's behavior [11]. In nongeneric cases, however, the optimal model is not compatible with a specified input/output partition. If such a partition is chosen, the identification problem is ill-posed. In practice, even a well-posed problem with fixed input/output partition can be ill-conditioned. Using the representation free setting and thus leaving the choice of the input-output partition to the identification method has the potential advantage of leading to methods that are more general conceptually (the input/output partition is identified together with the model), easier to analyze theoretically (solution always exists), and are computationally more robust (ill-conditioning is avoided).

## 2. Multiple experiments

Another nonstandard feature of the method presented in the paper is the use data of multiple measurement experiments. In static estimation problems, a measurement experiment yields a  $q$ -dimensional observation vector, so that the only way to increase the amount of data is to repeat the experiment. In dynamic problems, an experiment yields a  $q$ -dimensional vector time series. More data can be collected in this case also by increasing the length of the time series (number of samples in time). Our framework unifies the two approaches of data collection—repeated experiments and longer experiment. The former approach is of interest, for example, in the case of autonomous dynamical systems, where asymptotic analysis in time is not applicable.

The problem of using data of multiple experiments for identification of a single model is often addressed by heuristic methods. The rigorous method of [20] is based on concatenation of the separate time series into one time series, adding transient responses at the times of transition from one trajectory to another. This reduces the problem of identification from  $N$  time series to the problem of identification from a single time series *and* estimation of  $N$  transient responses—both being classical problems. This approach is implemented in all methods available in the System Identification Toolbox of Matlab [5].

The method presented in this paper uses a different approach. The time series parameterize separate block-Hankel matrices, with fixed row dimension but possibly different column dimensions. The block-Hankel matrices are appended to each other, forming a  $1 \times N$  block matrix with block-Hankel elements. Such a matrix is called a mosaic block-Hankel matrix. The identification problem of using data from multiple experiments is then equivalent to low-rank approximation of the mosaic block-Hankel matrix. Our approach is a natural generalization of the classical approach for solving static estimation problems where the blocks are the individual observations and the mosaic block-Hankel matrix is the unstructured  $q \times N$  matrix of the stacked next to each other observations.

## 3. Exact and missing data

Due to malfunctioning of measurement devices, communication channels, or storing devices, the available data may have missing measurements at some time instances. When the missing data span over  $\ell$  or more sequential time instances, where  $\ell$  is the lag of the system, the problem reduces to the problem of identification from two independent time series. Also for problems with a given input/output partitioning, missing output variables can be estimated by existing methods. More general patterns of missing data, *e.g.*, non sequential missing variables and missing variables in arbitrary combination of variables, require new methods. The approach used in the paper, treats missing data by minimizing a weighted approximation criterion with zero weights associated to the missing values.

The opposite of a missing data value is an exact data value. An example of how exact data occur in practice is data collection experiment for a system starting “at rest”, *i.e.*, zero initial conditions. This prior knowledge can be used in the identification problem by imposing equality constraints, setting the approximation of the exact values to the given data values. The approach of imposing exact values used in this paper, is by minimizing a weighted approximation criterion with infinite weights associated to the exact values. Thus, exact and missing values are treated uniformly as extreme cases of noise corrupted observations when the noise has zero and infinite variance, respectively.

## 4. Software implementation in a literature programming style

Scientific publications rarely include full implementation details of the methods described. A software implementation of the method (when available as open source) necessarily contains the full details about the method and is therefore the ultimate source of information. Even well documented code, however, may be difficult to read and may not correspond to a description in related scientific publication. The literature programming style [2] is a possible solution to this issue. A literate program includes the source code as an integral part of the document describing it. A literate program is designed to be read by a human being rather than a computer.

An extended version of this paper is a literate program for the described software package. The code is split into chunks and the chunks are presented in a logical sequence, which happens to be different from the sequence in which the chunks appear in the computer executable programs. The code chunks provide the implementation details for the algorithmic steps and, vice versa, the presentation in the paper serves as documentation of the

code. The particular literate programming tool used in this paper is `noweb` [18]. The computer executable code as well as the human readable text are automatically extracted from the literate program by `noweb`, so that the code and the paper coexist in common source files.

The approach for solving the identification problem, used in the paper, is weighted structured low-rank approximation [6]. Rank deficient mosaic-Hankel matrices are in a one-to-one correspondence with trajectories of linear time-invariant systems of bounded complexity. Therefore, data approximation by a linear time-invariant system of bounded complexity is a mosaic-Hankel low-rank approximation problem. The behavioral approach and the low-rank approximation setting are two sides of the same coin: they use representation free problem formulations, where the fundamental object of interest is a relation among variables rather than a map. The system theoretic notion of a linear time-invariant model of bounded complexity corresponds to the linear algebra notion of a rank deficient Hankel matrix. The linear algebra setting provides algorithms and software for addressing problems in the system theoretic setting.

Element-wise weights allow approximation with an emphasis on some variables and/or some time instances. Zero and infinite weights are allowed. A zero weight ignores the corresponding data point in the approximation and an infinite weight forces the approximation of the corresponding data point to be equal to the given one. This allows us to take into account missing and exact data.

Once the identification problem is cast as a mosaic-Hankel structured low-rank approximation problem, it is solved by the efficient methods of [22]. A practical implementation of the methods in C++, using optimization algorithms from the GNU scientific library [1], is available [8]. The identification function presented in this paper is a wrapper function to the structured low-rank approximation solver of [8]. The software is available from:

<http://homepages.vub.ac.be/~imarkovs/slra/software.html>

The paper is organized as follows: Section 2 sets the notation and states the considered identification problem. Section 3 presents the solution approach and gives details about the implementation of the software. Specific identification problems— $L_2$ -optimal model reduction, identification from data consisting of multiple trajectories of different lengths, and identification with missing data—are presented in Sections 5–7, respectively. In Section 8, the performance of the identification package is tested on real-life and simulated data from the database for system identification DAISY.

## 2 Problem formulation

### 2.1 Model class: linear time-invariant systems of bounded complexity

A discrete-time dynamical system  $\mathcal{B}$  is a collection of trajectories —  $q$ -variables time-series  $w : \mathbb{Z} \rightarrow \mathbb{R}^q$ . The class of finite dimensional linear time-invariant systems with  $q$  variables and at most  $m$  inputs is denoted by  $\mathcal{L}_m^q$ . A linear time-invariant system  $\mathcal{B} \in \mathcal{L}_m^q$  admits a representation by constant coefficients difference equation

$$\mathcal{B} = \mathcal{B}(R) := \{ w \mid R_0 w + R_1 \sigma w + \dots + R_\ell \sigma^\ell w = 0 \}, \quad (\text{DE})$$

where  $\sigma$  is the shift operator

$$(\sigma w)(t) = w(t+1).$$

*Note 1 (Kernel representation).* In the literature on the behavioral approach to systems and control, (DE) is called a kernel representation, because it can be written more compactly as the kernel

$$\ker(R(\sigma)) := \{ w \mid R(\sigma)w = 0 \}$$

of the operator  $R(\sigma)$ , where

$$R(z) := R_0 w + R_1 z + \dots + R_\ell z^\ell$$

is a polynomial matrix.

The minimal natural number  $\ell$ , for which there exists an  $\ell$ th order difference equation representation for  $\mathcal{B}$  is an important invariant of the system, called the *lag*. The number of inputs and the lag specify the complexity of the system in the sense that the dimension of the restriction of  $\mathcal{B}$  to the interval  $[1, T]$ , where  $T \geq \ell$ , is bounded by  $T\mathfrak{m} + \ell(q - \mathfrak{m})$ . The subset of  $\mathcal{L}_{\mathfrak{m}}^q$  with lag at most  $\ell$  is denoted by  $\mathcal{L}_{\mathfrak{m}, \ell}^q$ .

No a priori separation of the variables  $w_1, \dots, w_q$  into inputs and output is made. However, the variables can always be partitioned into inputs  $u$  (free variables) and outputs  $y$  (dependent variables). A convenient way to represent an input/output partitioning is  $w = \Pi \begin{bmatrix} u \\ y \end{bmatrix}$ , where  $\Pi$  is a  $q \times q$  permutation matrix.

An input/output representation

$$\mathcal{B}(P, Q, \Pi) = \{w = \Pi(u, y) \mid Q(\sigma)u = P(\sigma)y\},$$

of the system  $\mathcal{B}(R)$  is obtained by partitioning the polynomial matrix  $R\Pi$  as

$$R\Pi = \begin{bmatrix} \mathfrak{m} & \mathfrak{p} \\ Q & -P \end{bmatrix}, \quad \text{with } P \text{ nonsingular.}$$

In addition, the system can be represented in the input/state/output form

$$\mathcal{B} = \mathcal{B}(A, B, C, D, \Pi) := \{w = \Pi(u, y) \mid \text{there is } x, \text{ such that } \sigma x = Ax + Bu, y = Cx + Du\}. \quad (\text{I/S/O})$$

If the permutation matrix  $\Pi$  is the identity matrix, then  $w = (u, y)$ , *i.e.*, the first  $\mathfrak{m}$  variables of  $w$  are inputs and the remaining  $\mathfrak{p}$  variables are outputs. In this case,  $\Pi$  is skipped in (I/S/O), *i.e.*,  $\mathcal{B}(A, B, C, D) = \mathcal{B}(A, B, C, D, I)$ .

The number of inputs  $\mathfrak{m}$ , the number of outputs  $\mathfrak{p} = q - \mathfrak{m}$ , and the minimal state dimension  $\mathfrak{n}$  of an input/state/output representation of  $\mathcal{B}$  are invariant of the representation and in particular of the input/output partitioning. The order  $\mathfrak{n}$  of a state-space representation of a linear time-invariant system  $\mathcal{B} = \mathcal{B}(R)$  with lag  $\ell$  and  $\mathfrak{p}$  outputs is  $\mathfrak{n} \leq \ell\mathfrak{p}$ . In the case when the block  $P_\ell \in \mathbb{R}^{\mathfrak{p} \times \mathfrak{p}}$  of  $R_\ell = \begin{bmatrix} Q_\ell & -P_\ell \end{bmatrix}$  is nonsingular,  $\mathfrak{n} = \ell\mathfrak{p}$  and  $w = (u, y)$  is a possible input/output partition, *i.e.*,  $\Pi$  can be chosen equal to  $I$ . This simplifying assumption is made in the rest of the paper. The class of systems with  $q$  variables and inputs, order, and lag bounded by, respectively  $\mathfrak{m}$ ,  $\mathfrak{n}$ , and  $\ell$  is denoted by  $\mathcal{L}_{\mathfrak{m}, \ell}^{q, \mathfrak{n}}$ .

## 2.2 Approximation criterion: distance between data and model

The *misfit* (lack of fit) between the data  $w_d$  and the model  $\mathcal{B}$  is measured by the orthogonal distance from  $w_d$  to  $\mathcal{B}$

$$M(w_d, \mathcal{B}) := \min_{\hat{w}^1, \dots, \hat{w}^N \in \mathcal{B}} \sqrt{\sum_{k=1}^N \|w_d^k - \hat{w}^k\|_2^2}. \quad (M)$$

Intuitively, the misfit shows how much the model  $\mathcal{B}$  fails to “explain” the data  $w_d$ .

Missing elements are marked by the symbol NaN and are excluded from the approximation criterion, *i.e.*,

$$\text{NaN} - \hat{w}_i(t) = 0, \quad \text{for all } \hat{w}_i(t) \in \mathbb{R}.$$

The optimal approximate modeling problem considered aims to find a system  $\hat{\mathcal{B}}$  in the model class  $\mathcal{L}_{\mathfrak{m}, \ell}^q$  that best fits the data according to the misfit criterion.

Given a set of time series

$$w_d = \{w_d^1, \dots, w_d^N\}, \quad \text{where } w_d^k = (w_d^k(1), \dots, w_d^k(T_k)), \quad w_d^k(t) \in \mathbb{R}^q,$$

and a complexity specification  $(\mathfrak{m}, \ell)$ , find the system

$$\hat{\mathcal{B}} := \arg \min_{\mathcal{B} \in \mathcal{L}_{\mathfrak{m}, \ell}^q} M(w_d, \mathcal{B}). \quad (\text{SYSID})$$

Special cases of (SYSID) are static data modeling ( $\ell = 0$ ) and output-only or autonomous system identification ( $\mathfrak{m} = 0$ ). The solution approach, described next, leads to an algorithm that covers these special cases. In addition,

1. elements of the given time series  $w_d$  can be specified as “missing” by passing the symbol NaN for their value;
2. elements of the given time series  $w_d$  can be specified as “exact”, in which case they appear unmodified in the approximation  $\hat{w}$ ;
3. the approximation  $\hat{w}$  can be constrained to be a trajectory of the model  $\mathcal{B}$ , generated under a priori fixed initial conditions  $w_{\text{ini}}$ , see [7], *i.e.*,

$$\begin{bmatrix} w_{\text{ini}} \\ \hat{w} \end{bmatrix} \in \mathcal{B}.$$

(Note that problem (SYSID) identifies the model without prior knowledge about the initial conditions, under which the data  $w_d$  is generated, *i.e.*,  $w_{\text{ini}}$  is a free variable.)

### 2.3 Stochastic interpretation

The system identification problem considered in the paper is defined as a deterministic approximation problem. As shown in [6, Section 3.1], however, it yields the maximum likelihood estimator in the errors-in-variables setting [21].

**Proposition 2** ([6, Proposition 6]). *Assume that the data  $w_d$  is generated in the errors-in-variables setting*

$$w_d = \bar{w} + \tilde{w},$$

where the true data  $\bar{w}$  is a trajectory of a true model  $\tilde{\mathcal{B}} \in \mathcal{L}_{m,\ell}^q$  and the measurement noise  $\tilde{w}$  is zero mean normally distributed with covariance matrix that is a multiple of the identity. Then the solution  $\hat{\mathcal{B}}^*$  of (SYSID) is a maximum likelihood estimator for the true model  $\tilde{\mathcal{B}}$ .

In [15, 3], it is proven that under additional mild assumptions the estimator  $\hat{\mathcal{B}}^*$  is consistent and the estimated parameters have asymptotically normal joint distribution. Therefore, asymptotic confidence regions can be obtained as a byproduct of the optimization algorithm (see [8, Section 5.1]). The statistical setting justifies the choice of the deterministic approximation criterion and provides a testbed for the method: the method works “well” (consistency) and is optimal (asymptotic efficiency) under specified conditions.

*Note 3* (Exact and missing data correspond to noisy data with, respectively, zero and infinite noise variance). In the classical regression model, an input-output partitioning of the variables is a priori given and the input variables (regressors) are assumed to be noise free. Classical regression is a special case of the errors-in-variables model when the input noise variance is zero. The opposite of a noise free variable is a missing variable, or equivalently a variable with infinite variance in the errors-in-variables setting.

## 3 Solution approach

For a given set of trajectories  $w_d$  and a natural number  $\ell$ , the mosaic-block-Hankel matrix (a  $1 \times N$  block matrix with block-Hankel blocks) is defined as

$$\mathcal{H}_{\ell+1}(w_d) := [\mathcal{H}_{\ell+1}(w_d^1) \quad \cdots \quad \mathcal{H}_{\ell+1}(w_d^N)], \quad \text{where} \quad \mathcal{H}_{\ell+1}(w_d^k) := \begin{bmatrix} w_d^k(1) & w_d^k(2) & \cdots & w_d^k(T-\ell) \\ w_d^k(2) & w_d^k(3) & \cdots & w_d^k(T-\ell+1) \\ \vdots & \vdots & & \vdots \\ w_d^k(\ell+1) & w_d^k(\ell+2) & \cdots & w_d^k(T) \end{bmatrix}.$$

The solution method is based on the following equivalence

$$(w_d(1), \dots, w_d(T_k - \ell)) \in \mathcal{B} \in \mathcal{L}_{m,\ell}^{q,n}, \quad \text{for } k = 1, \dots, N \quad \iff \quad \text{rank}(\mathcal{H}_{\ell+1}(w_d)) \leq (\ell+1)m + n,$$

*i.e.*, the time series  $w_d$ , possibly except for the last  $\ell$  samples, are exact trajectories of a linear time-invariant model  $\mathcal{B}$  with complexity bounded by  $(m, n)$  if and only if the mosaic-block-Hankel matrix  $\mathcal{H}_{\ell+1}(w_d)$  has rank bounded by

$$r = (\ell+1)m + n.$$

Therefore, the identification problem (SYSID) is equivalent to the structured low-rank approximation problem

$$\text{minimize over } \hat{w} \quad \|w_d - \hat{w}\|_2^2 \quad \text{subject to} \quad \text{rank}(\mathcal{H}_{\ell+1}(\hat{w})) \leq r. \quad (\text{SLRA})$$

*Note 4.* Under the assumption that  $P_\ell$  is full rank, the whole trajectories  $w_d^1, \dots, w_d^N$  are in the behavior of the model.

Problem (SLRA) is a classic problem that can be solved in different ways. The approach used in the paper is based on the kernel representation of the rank constraint

$$\text{rank}(\mathcal{H}_{\ell+1}(\hat{w})) \leq r \iff R\mathcal{H}_{\ell+1}(\hat{w}) = 0 \quad \text{and} \quad R^\top R = I_{(\ell+1)q-r}. \quad (\text{KER})$$

The matrix  $R$  in the right-hand-side of (KER) is related to the parameters  $R_0, R_1, \dots, R_\ell$  of the difference equation representation (DE) of the exact model for  $\hat{w}$  as follows:

$$R = [R_0 \quad R_1 \quad \dots \quad R_\ell], \quad \text{where} \quad R_i \in \mathbb{R}^{p \times q}.$$

Algorithms and software for mosaic-Hankel low-rank approximation are developed in [9, 22]. The software of [8] is used. The software presented in this paper is just an interface to the structured low-rank approximation solver for the purpose of linear time-invariant system identification, see Figure 1. On its own, the structured low-rank approximation solver computes a parameter  $\hat{R}$  of the difference equation representation of the optimal approximating system  $\hat{\mathcal{B}}$ . The system identification function converts the parameter  $\hat{R}$  to the parameters  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  of an input/state/output representation of  $\hat{\mathcal{B}}$  in order to facilitate the usage of the model by other analysis and synthesis tools. In addition, a system  $\mathcal{B}_{\text{ini}} \in \mathcal{L}_{m,\ell}^q$ , specified by parameters  $(A_{\text{ini}}, B_{\text{ini}}, C_{\text{ini}}, D_{\text{ini}})$  can be used as an initial approximation for the optimization algorithm. The parameters  $(A_{\text{ini}}, B_{\text{ini}}, C_{\text{ini}}, D_{\text{ini}})$  are converted to the parameter  $R_{\text{ini}}$ , used by the structured low-rank approximation solver. Details about these conversions are given in Appendix B. Although the identification function has as external interface the (I/S/O) representation of the model, the internal computations are done via the parameter  $R$  of the difference equation representation.

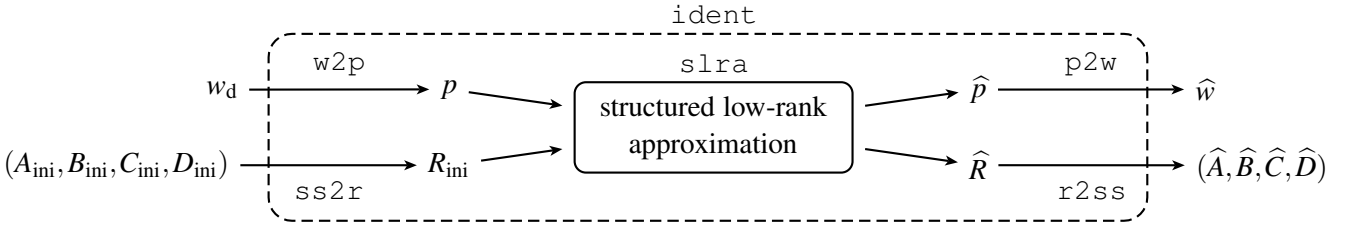


Figure 1: Implementation of the system identification function `ident`.

The function `ident` solves the approximate identification problem (SYSID), and `misfit` computes the misfit  $M(w_d, \mathcal{B})$ . They implement the following mappings (the input/output parameters are defined in Appendix A):

`ident`:  $(w_d, m, \ell) \mapsto (\hat{A}, \hat{B}, \hat{C}, \hat{D})$ , where  $(\hat{A}, \hat{B}, \hat{C}, \hat{D})$  define a (locally) optimal solution of (SYSID)

`(ident function definition)`  $\equiv$   
`function [sysh, info, wh, xini] = ident(w, m, ell, opt)`

`misfit`:  $(w_d, (A, B, C, D)) \mapsto (M, \hat{w})$ , where  $M$  is the misfit between  $\mathcal{B}(A, B, C, D)$  and  $w_d$ , and  $\hat{w}$  is the optimal approximation of  $w_d$  within  $\mathcal{B}(A, B, C, D)$  (the smoothed trajectory).

`(misfit function definition)`  $\equiv$   
`function [M, wh, xini] = misfit(w, sysh, opt)`

## 4 Identification without input/output partitioning

In this section, we show a simulation example that illustrates the independence of the representation free problem formulation (SYSID) of an a priori given input output partitioning of the variables. Noisy data  $(u_d, y_d)$  is generated from a single-input single-output second order system

`(permutation of the variables)`  $\equiv$   
`ell = 2; m = 1; p = 1; q = m + p; n = ell * p; T = 30; s = 0.1;`  
`sys0 = drss(n, p, m); u0 = rand(T, 1); y0 = lsim(sys0, u0);`  
`w0 = [u0 y0]; w = w0 + s * randn(T, q);`

and the identification method `ident` is applied, first on the data  $w_d = (u_d, y_d)$ :

```
<permutation of the variables>+≡
opt.solver = 'm';
[sysh1, info1] = ident(w, m, ell, opt);
```

and, second on the data  $w'_d = (y_d, u_d)$  with permuted variables:

```
<permutation of the variables>+≡
[sysh2, info2] = ident(fliplr(w), m, ell, opt);
```

The achieved misfit  $M(w_d, \hat{\mathcal{B}}) = \text{info1.M}$  and  $M(w'_d, \hat{\mathcal{B}}') = \text{info2.M}$  in the two cases is equal and the identified models  $\hat{\mathcal{B}} = \text{sys1}$  and  $\hat{\mathcal{B}}' = \text{sys2}$  are closely related. (If  $\mathcal{B}_1$  is invertible  $\mathcal{B}_2$  is equal to the inverse of  $\mathcal{B}_1$ .)

Although, identification from data with permuted variables leads to equivalent optimization problems, the behavior of optimization methods in these problems may differ significantly. This is illustrated on the plots of Figures 2 and 3. Figure 2 shows the convergence of the parameters  $R$  (identification using  $w_d$ ) and  $R'$  (identification using  $w'_d$ ) in the course of the optimization algorithm. Figure 3 shows the approximation errors  $M_k := M(w_d, \hat{\mathcal{B}}_k)$  and  $M'_k := M(w'_d, \hat{\mathcal{B}}'_k)$  as a function of the iteration step  $k$ .

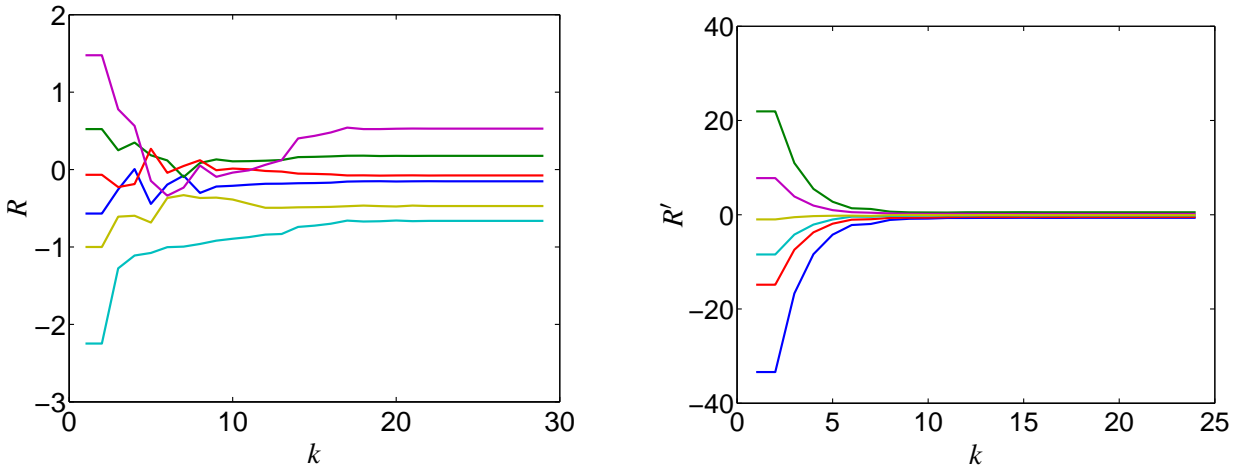


Figure 2: Convergence of the model parameters  $R$  and  $R'$  in the course of the optimization method.

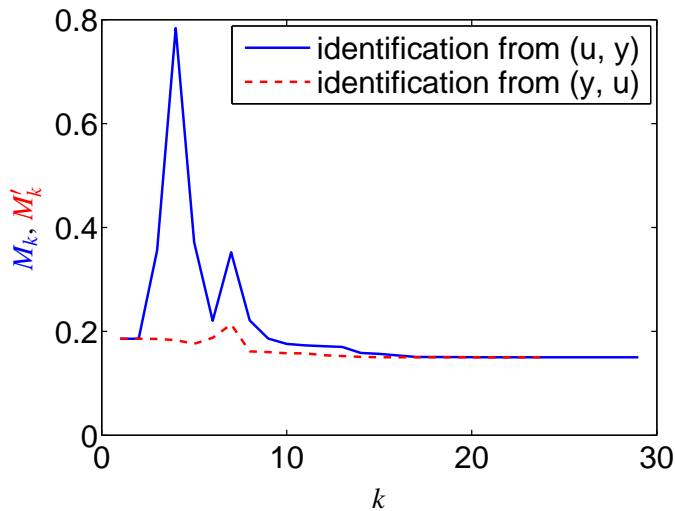


Figure 3: Misfits  $M_k := M(w_d, \hat{\mathcal{B}}_k)$  and  $M'_k := M(w'_d, \hat{\mathcal{B}}'_k)$  as a function of the iteration step  $k$ .

## 5 $L_2$ -optimal model reduction

In this section, the `ident` function is used for solution of a (finite horizon)  $L_2$ -optimal model reduction problem [6, Section 3.3]. Consider a discrete-time linear time-invariant system  $\mathcal{B}$  with an input/output partition  $w = (u, y)$  and let  $H(t) \in \mathbb{R}^{p \times m}$  be the  $t$ th sample of the impulse response of  $\mathcal{B}$ . The finite horizon- $T$   $L_2$  norm of  $\mathcal{B}$  (w.r.t. the given input/output partition  $w = (u, y)$ ) is defined as

$$\|\mathcal{B}\|_{2,T} := \|H\|_{2,T} := \sqrt{\sum_{t=0}^T \|H(t)\|_F^2},$$

where  $\|\cdot\|_F$  is the Frobenius norm.

**Problem 5** ( $L_2$ -optimal model reduction problem [6]). Given a linear time-invariant system  $\mathcal{B}_d \in \mathcal{L}_{m,\ell}^q$  and a complexity specification  $\ell_{\text{red}} < \ell$ , find an optimal approximation of  $\mathcal{B}_d$  with bounded complexity  $(m, \ell_{\text{red}})$

$$\hat{\mathcal{B}}^* := \arg \min_{\hat{\mathcal{B}}} \|\mathcal{B}_d - \hat{\mathcal{B}}\|_{2,T} \quad \text{subject to} \quad \hat{\mathcal{B}} \in \mathcal{L}_{m,\ell_{\text{red}}}^q.$$

The link between  $L_2$ -optimal model reduction and system identification is due to equivalence of an impulse response of  $\mathcal{B}$  and a set of responses of a related autonomous system  $\mathcal{B}_{\text{aut}}$ .

**Proposition 6.** *The shifted impulse response*

$$\sigma H = (H(1), H(2), \dots)$$

of  $\mathcal{B} = \mathcal{B}(A, B, C, D)$  is equal to the matrix  $[y_1 \ \dots \ y_m]$  of responses  $y_1, \dots, y_m$  of the autonomous system  $\mathcal{B}_{\text{aut}} = \mathcal{B}(A, C)$  to initial conditions  $b_1, \dots, b_m$ , where  $B = [b_1 \ \dots \ b_m]$ .

Proposition 6 and the fact that  $H(0)$  is equal to the feed through parameter  $D$  of the input/state/output representation  $\mathcal{B}(A, B, C, D)$  imply that the  $L_2$ -optimal model reduction problem is equivalent to an  $L_2$ -optimal identification problem for an autonomous linear time-invariant system from the set of  $m$  responses  $h_1, \dots, h_m$ .

A input/state/output representation of the reduced order model is obtained from the identified autonomous system  $\mathcal{B}(\hat{A}, \hat{C})$ , initial conditions  $X_{\text{ini}} = [\hat{x}_{\text{ini},1} \ \dots \ \hat{x}_{\text{ini},m}]$ , and  $H(0)$  as follows:

$$\mathcal{B}(\hat{A}, X_{\text{ini}}, \hat{C}, H(0)).$$

### A benchmark model reduction problem

As a test example, consider a mechanical system consisting of  $N$  point masses connected in a chain by ideal springs and ideal dampers. The first and the last masses are also connected to walls. Friction force, proportional to the speed, acts on all masses. The parameters of the system are the masses  $m_1, \dots, m_N$ , the spring and damper coefficients  $k_0, k_1, \dots, k_N$  and  $d_0, d_1, \dots, d_N$ , the length  $\delta$  of the springs, and the friction coefficient  $f$ , which are all nonnegative numbers. Setting the coefficients of the most left and most right springs and dampers to zero has the effect of detaching the chain of masses from the left and right walls.

Denoting by  $p_i$  the position of the  $i$ th point mass with respect to the left wall, the system dynamics is described by the system of differential equations

$$\begin{aligned} m_1 \ddot{p}_1 &= -(k_0 + k_1)p_1 + k_1 p_2 - (d_0 + d_1 + f)\dot{p}_1 + d_1 \dot{p}_2 + (k_0 - k_1)\delta, \\ m_2 \ddot{p}_2 &= k_1 p_1 - (k_1 + k_2)p_2 + k_2 p_3 + d_1 \dot{p}_1 - (d_1 + d_2 + f)\dot{p}_2 + d_2 \dot{p}_3 + (k_1 - k_2)\delta, \\ &\vdots \\ m_i \ddot{p}_i &= k_{i-1} p_{i-1} - (k_{i-1} + k_i)p_i + k_i p_{i+1} + d_{i-1} \dot{p}_{i-1} - (d_{i-1} + d_i + f)\dot{p}_i + d_i \dot{p}_{i+1} + (k_{i-1} - k_i)\delta, \\ &\vdots \\ m_N \ddot{p}_N &= k_{N-1} p_{N-1} - (k_{N-1} + k_N)p_N + k_N p_{N-1} - (d_{N-1} + d_N + f)\dot{p}_N + (k_{N-1} + k_N)\delta. \end{aligned} \tag{SYS}$$

The positions  $y$  of the masses with indexes in the set  $\mathcal{Y}$  are observed and external forces  $u$  are applied to the masses with indexes in the set  $\mathcal{U}$ . An input/state/output representation of the resulting model is given in Appendix C.



## Numerical results

The input data for the model reduction methods is the first  $T + 1$  samples of the impulse response  $H$  of the full order system and the lag of the reduced order model. The value of  $T$  is selected so that, after  $T$  steps, the impulse response has converged sufficiently close to zero. The  $L_2$ -optimal model reduction problem is solved by the `ident` function and the function `ar12` from the RARL2 package [12].

*Note 7* (Using  $H$  vs  $\mathcal{B}$  as input data). The `ar12` function accepts as an input either an input/state/output representation of the model or impulse response coefficients of the system. The latter option is used in the simulation examples.

*Note 8* (Initial approximation). The optimization based methods `ident` and `ar12` are initialized with the suboptimal approximation computed by Kung's method [4], implemented in the function `fc2ss` of the RARL2 package.

The approximate model  $\hat{\mathcal{B}}$  is evaluated by computing the finite horizon- $T$  relative  $L_2$  approximation error

$$e = \frac{\|\mathcal{B} - \hat{\mathcal{B}}\|_{2,T}}{\|\mathcal{B}\|_{2,T}}.$$

*Example 9.* In an example with model parameters

$$\begin{aligned} N = 10, \quad f = 0.1, \quad \delta = 1, \quad \mathcal{U} = 1, \quad \mathcal{Y} = 5, \quad \ell = 3, \\ [m_1 \ m_2 \ \dots \ m_{N-1} \ m_N] = [10 \ 9 \ \dots \ 2 \ 1], \\ [k_0 \ k_1 \ \dots \ k_{N-1} \ k_N] = [0.5 \ 9 \ \dots \ 1 \ 0.1], \\ [d_0 \ d_1 \ \dots \ d_{N-1} \ d_N] = 0.2 [0 \ 9 \ \dots \ 1 \ 1], \end{aligned}$$

the obtained results are:

'method'	'kung'	'ar12'	'ident'
'e'	[0.4380]	[0.4263]	[0.4263]
't'	[0.1035]	[2.3251]	[0.0579]

The high-order and reduced order systems' impulse and frequency responses are shown in Figure 4. The `ar12` and `ident` functions compute the same approximation, however, the computation time required by the `ident` function is much smaller. This is due to the efficient implementation of the structured low-rank approximation solver.

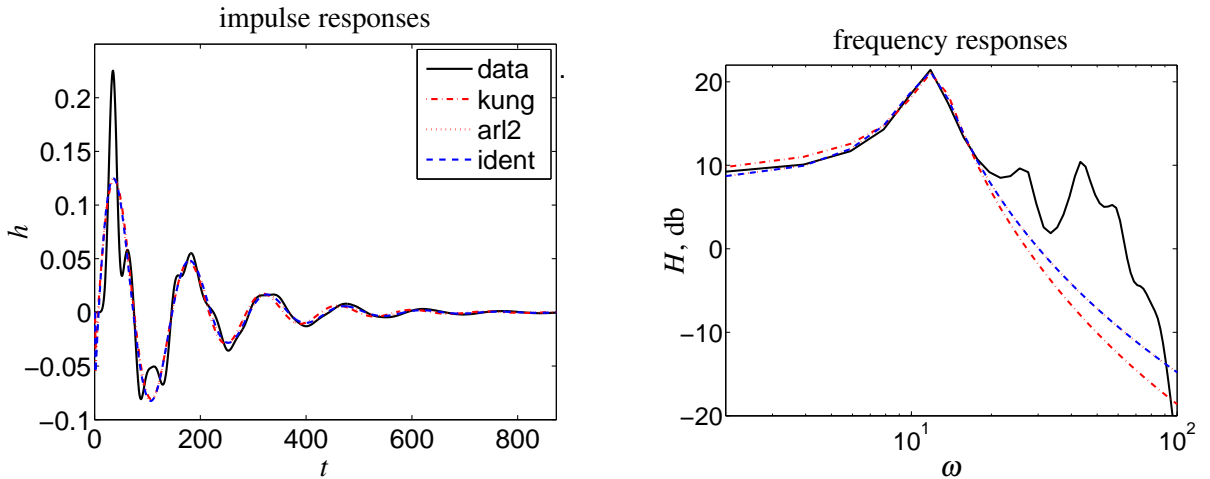


Figure 4: Comparison of the high order and reduced order models in example 9.

*Example 10.* In the setup of Example 9, detaching the system from the right wall ( $k_N = d_N = 0$ ) gives qualitatively similar results:

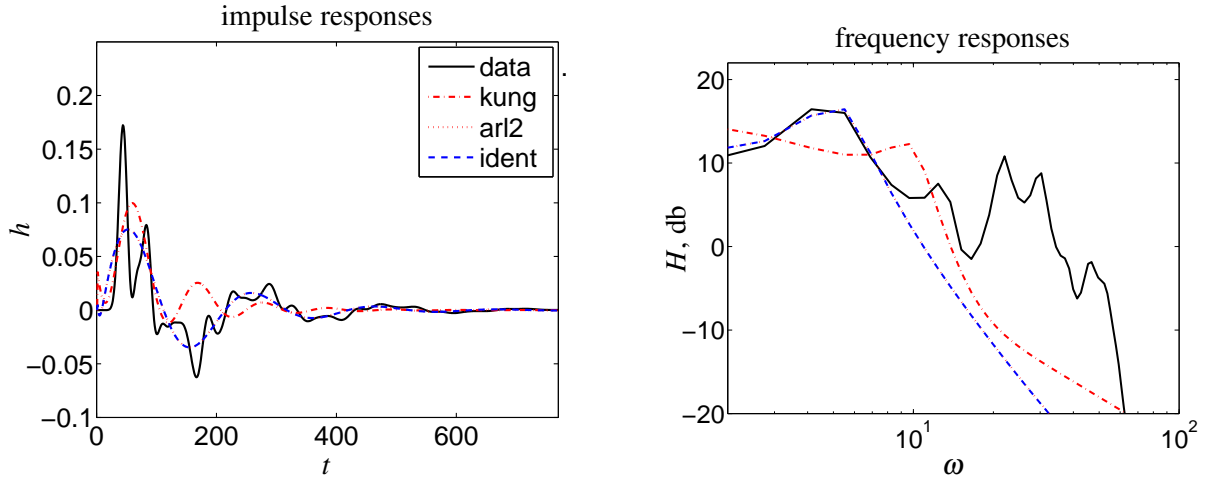


Figure 5: Comparison of the high order and reduced order models in example 10.

'method'	'kung'	'arl2'	'ident'
'e'	[0.8368]	[0.5950]	[0.5950]
't'	[0.0804]	[3.2495]	[0.0789]

The impulse and frequency responses are shown in Figure 5, right.

*Example 11.* In the setup of Example 9, detaching the system from both right and left walls ( $k_N = d_N = k_0 = d_0 = 0$ ) introduced a free motion, *i.e.*, an unstable mode. This gives qualitatively different result in the approximations (see Figure 6):

'method'	'kung'	'arl2'	'ident'
'e'	[0.6830]	[0.2244]	[0.0337]
't'	[0.0561]	[1.9533]	[0.0564]

The reason for the poor approximation of arl2 is that this method imposes stability of the approximation.

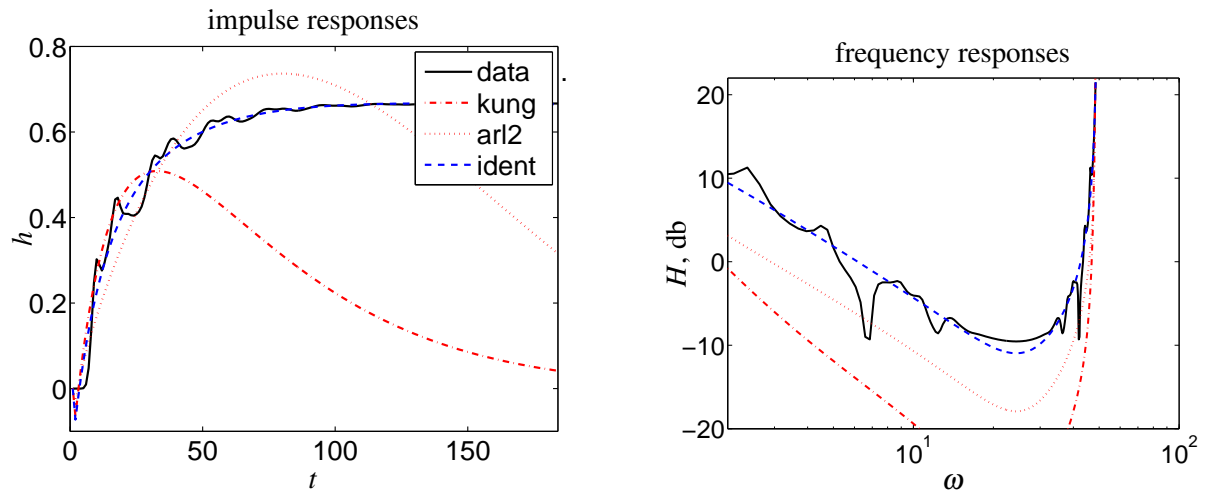


Figure 6: Comparison of the high order and reduced order models in example 10.

## 6 Identification from multiple trajectories of different experiments

In this section, the multiple experiments data feature of the identification method is tested. An alternative method to `ident`, used for comparison, is the function `pem` from the System Identification Toolbox of Matlab. The data is generated in the output-error setup and the methods are called with the corresponding options:

- `opt.exct = 1:m` for `ident` and
- `'dist', 'none'` for `pem`.

The true system is randomly generated with `drss` and the true data consists of random trajectories of the true system. The input is normally distributed with zero mean and identity covariance matrix. The data used for identification is a noise corrupted version of the true trajectory, where the output noise is zero mean white Gaussian. The simulation parameters are: number of inputs, specified by the variable `m`; number of outputs, specified by the variable `p`; lag of the system, specified by the variable `e11`; number of time series, specified by the variable `N`; number of samples of the time series, specified by the variable `T`; and a noise standard deviation, specified by the variable `NL`.

The identified models by the `ident` and `pem` functions are compared with respect to their relative output error

$$e = \sum_{i=1}^N \frac{\|y^i - \hat{y}^i\|_2^2}{\|y^i\|_2^2}$$

*Example 12.* In an example with parameters

*(Multiple experiments example 1)*≡

```
ex = 'mult-traj-ex1'; e11 = 3; p = 1; m = 1; T = [500 1000];
K = 5; NL = linspace(0, 0.1, K); test_mult_traj
```

the results obtained

<code>'nl'</code>	[	0]	[0.0250]	[0.0500]	[0.0750]	[0.1000]
<code>'e ident'</code>	[5.5929e-32]	[0.1137]	[0.3360]	[0.5445]	[0.6760]	
<code>'e pem'</code>	[3.4563e-32]	[0.1137]	[0.3360]	[0.5445]	[0.6760]	

show that the `ident` and `pem` methods achieve the same approximation error. (Difference of less than  $10^{-4}$  in the relative approximation errors can be attributed to exiting the optimization with the default convergence tolerance of  $10^{-5}$ .) Although, in this example `ident` and `pem` are functionally equivalent (compute the same answer), `ident` is on the average 28 times faster than `pem`.

*Example 13.* In the case of  $N = 50$  data sets of short ( $T = 20$ ) trajectories

*(Multiple experiments example 2)*≡

```
ex = 'mult-traj-ex2'; e11 = 2; p = 1; m = 1; T = 20 * ones(1, 50);
K = 5; NL = linspace(0, 0.1, K); test_mult_traj
```

<code>'nl'</code>	[	0]	[0.0250]	[0.0500]	[0.0750]	[0.1000]
<code>'e ident'</code>	[2.1716e-32]	[0.0015]	[0.0057]	[0.0130]	[0.0225]	
<code>'e pem'</code>	[1.9819e-32]	[0.0015]	[0.0057]	[0.0127]	[0.0225]	

the obtained results are again identical. In this case the `ident` function is on the average 59 times faster than `pem`. Note that the theoretical computational complexity of the structured low-rank approximation method, used by `ident`, scales linearly in  $N$ , see [22]. Empirical results suggest that the algorithms used in the implementation of the `pem` function scale worse with the number of time series.

## 7 System identification with missing data

This section shows examples of identification problems with missing data. The data  $w$  is a noisy  $T$  samples long random trajectory of a single-input single-output linear time-invariant system  $\bar{\mathcal{B}} = \mathcal{B}(\bar{R})$  with lag  $\ell = 2$ . Input, output, or both samples are missing at moments of time  $t \in \mathcal{T}_m$ . The true model parameters are

$$\bar{R}_0 = [0.3 \ 0.7], \quad \bar{R}_1 = [1 \ -1.4], \quad \bar{R}_2 = [0 \ 1]. \quad (*)$$

The approximation accuracy is measured by the angle

$$e = \angle(\bar{R}, \hat{R}) = \cos^{-1} \left( \frac{\bar{R}^\top \hat{R}}{\|\bar{R}\| \|\hat{R}\|} \right),$$

between the true  $\bar{R}$  and estimated  $\hat{R}$  parameter vectors.

An alternative method for solving the identification problem considered in this section is proposed in [14]. This method uses a frequency domain approach [16]. A Matlab implementation of the algorithm (called below `sysid`) was kindly provided by the authors and is used below for verification of the results obtained with `ident`.

The simulation parameters in the experiments are the number of samples  $T$ ; the set of missing values  $\mathcal{T}_m$ , specified by a variable `Tm`; and a noise standard deviation interval, specified by a vector `NL`. The reported results show the error  $e$  for the compared methods and for the different noise levels specified in `NL`. Three experiments are done for different distribution of the missing values: sequential, periodic, and random. Both input and output values are missing. A NaN value in the table of results indicates that the corresponding method fails in (or does not apply to) this case.

*Example 14* (Sequential missing data samples). In an example with sequential missing data,

*⟨Missing data example 1⟩*≡

```
ex = 'missing-data-ex1';
T = 100; K = 5; NL = linspace(0, 0.01, K); Tm = [30:70]; test_missing_data
```

the `ident` and `sysid` functions achieve comparable accuracy

'nl'	[ 0]	[0.0025]	[0.0050]	[0.0075]	[0.0100]
'ident'	[2e-08]	[0.0012]	[0.0078]	[0.0017]	[0.0067]
'sysid'	[ 0]	[0.0012]	[0.0078]	[0.0014]	[0.0065]

The `ident` function is 9 time faster than `sysid`.

*Example 15* (Periodic missing data samples). Similar results are obtained in the case of periodic missing data for small noise levels

*⟨Missing data example 2⟩*≡

```
ex = 'missing-data-ex2';
T = 100; K = 5; NL = linspace(0, 0.01, K); Tm = [30:3:70]; test_missing_data
```

'nl'	[ 0]	[0.0025]	[0.0050]	[0.0075]	[0.0100]
'ident'	[4e-07]	[0.0014]	[0.0058]	[0.0077]	[0.0068]
'sysid'	[ 0]	[0.0014]	[0.0073]	[0.0082]	[0.0077]

In this example, the `ident` function is 5 time faster than `sysid`.

*Example 16* (Randomly distributed missing data samples). Finally, results for a simulation example with  $T = 1000$  data points from which 600 are randomly missing are shown:

*⟨Missing data example 3⟩*≡

```
ex = 'missing-data-ex3';
T = 1000; K = 5; NL = linspace(0, 0.01, K); Tm = sort(randperm(T, 600));
test_missing_data
```

'nl'	[ 0]	[0.0025]	[0.0050]	[0.0075]	[0.0100]
'ident'	[9e-05]	[0.0029]	[0.0087]	[0.0028]	[0.0123]
'sysid'	[ NaN]	[ NaN]	[ NaN]	[ NaN]	[ NaN]

**Discussion of the system identification results with missing data:** The problems solved by the `ident` and `sysid` functions are equivalent. In the stochastic setting of Section 2.3, they are errors-in-variables identification problems with missing data. Examples 14 and 15 numerically confirm that the results computed by the two methods coincide for small noise levels and a small number of missing values. Since the problems are nonconvex, however, for high noise levels the methods may converge to different locally optimal solutions.

In addition, there is an important algorithmic differences between `ident` and `sysid`, which explains the results of Example 16. In `ident`, the initial conditions are eliminated analytically at the level of the misfit computation (see Section 3). In `sysid`, the initial conditions are included in the parameter vector of the nonlinear least squares problem (see [14]). This results in a larger optimization problem, solved by `sysid` than the one solved by `ident`. Since in `sysid` a transient response is added for every missing data point, the number of initial conditions to-be-estimated is growing with the increase of number of missing data points. In contrast, the dimension of the nonlinear least

squares problem solved by `ident` is fixed. The difference in performance between `ident` and `sysid` due to the different way of dealing with the initial conditions becomes pronounced in problems with more missing values. This is illustrated in Example 16, where the problem is no longer solvable by `sysid`. Detailed statistical and numerical analysis of the identification methods in the case of missing data will be presented elsewhere.

## 8 Performance on real-life data

In this section, the performance of the `ident` function is tested on benchmark problems from the data base for system identification DAISY [13]. The data come from a number of applications in process industry, electrical, mechanical, and environmental engineering. A validation criterion is chosen that measures the predictive power of the model: how accurate the model can fit a part of the data that is not used for identification. Values for the identification methods’ parameters that correspond to this validation criterion are chosen and fixed for all data sets. Although “better” results may be obtained by preprocessing of the data and tuning “hyper parameters” (model structure and identification criterion) this is not done. Our tests reflect the view that the identification process should be as automatic as possible, *i.e.*, it should be done with as little human interaction as possible. The methods are applied on all data sets choosing only the model class (specified by a bound on the model complexity).

The considered data sets are listed in Table 1. References and details about the nature and origin of the data is given in [13]. The data is preprocessed by centering it. The model’s lag  $\ell$  is chosen manually for each data set from the complexity–accuracy trade-off curve (misfit as a function of the lag). For comparison, also the automatically selected lags  $\ell'$  by the `pem` function are shown in the table. (In the first two data sets, automatic order selection by `pem` is not possible due to a small number of samples.)

#	Data set name	$T$	$m$	$p$	$\ell$	$\ell'$
1	Data of a simulation of the western basin of Lake Erie	57	5	2	1	—
2	Data of ethane-ethylene distillation column	90	5	3	1	—
3	Heating system	801	1	1	2	1
4	Data from an industrial dryer (Cambridge Control Ltd)	867	3	3	1	2
5	Data of a laboratory setup acting like a hair dryer	1000	1	1	5	1
6	Data of the ball-and-beam setup in SISTA	1000	1	1	2	3
7	Wing flutter data	1024	1	1	5	2
8	Data from a flexible robot arm	1024	1	1	4	6
9	Data of a glass furnace (Philips)	1247	3	6	1	4
10	Heat flow density through a two layer wall	1680	2	1	2	1
11	Simulation data of a pH neutralization process	2001	2	1	6	2
12	Data of a CD-player arm	2048	2	2	1	1
13	Data from a test setup of an industrial winding process	2500	5	2	2	1
14	Liquid-saturated steam heat exchanger	4000	1	1	2	3
15	Data from an industrial evaporator	6305	3	3	1	1
16	Continuous stirred tank reactor	7500	1	2	1	2
17	Model of a steam generator at Abbott Power Plant	9600	4	4	1	1

Table 1: Examples from DAISY.  $T$ —number of data points,  $m$ —number of inputs,  $p$ —number of outputs,  $\ell$ —lag of the identified model.

The data  $w_d = (u_d, y_d)$  in all examples is split into identification and validation parts. The first 70% of the data, denoted  $w_{\text{idt}}$ , are used for identification, and the remaining 30% of the data, denoted  $w_{\text{val}}$ , are used for validation. A model  $\hat{\mathcal{B}}$  is identified from  $w_{\text{idt}}$  by an identification method and is validated on  $w_{\text{val}}$  by the validation criterion defined next. The model class is linear time-invariant systems with a bound  $\ell$  on the lag.

The validation criterion is the “simulation fit” computed by the function `compare` of the System Identification Toolbox. This choice is motivated by the fact that in this case the results obtained by the `ident` function can be validated, using the `pem` function. Correspondingly, the `ident` and `pem` functions are applied setting options for output error identification:

- `opt.exct = 1:m` for `ident` and
- `'dist'`, `'none'` for `pem`.

For a given time series  $w_d = (u_d, y_d)$  and a model  $\mathcal{B}$ , the approximation  $\hat{y}$  of  $y$  in  $\mathcal{B}$  is defined as follows:

$$\hat{y}((u_d, y_d), \mathcal{B}) := \min_{\hat{y}} \|y_d - \hat{y}\| \quad \text{subject to} \quad \text{col}(u_d, \hat{y}) \in \mathcal{B}.$$

(The optimization is carried over the initial conditions that generate  $\hat{y}$  from the given input  $u_d$ .) With this notation, the fit of  $w_d$  by  $\mathcal{B}$  is defined as

$$F(w_d, \mathcal{B}) := 100 \frac{\max(0, 1 - \|y_d - \hat{y}(w_d, \mathcal{B})\|)}{\|y_d - \sum_{t=1}^T y_d(t)/T\|}.$$

The fitting criterion  $F(w_{\text{val}}, \hat{\mathcal{B}})$  is compared for the models produced by `ident` and `pem`. The results are shown in Figure 7. In terms of accuracy, the `ident` and `pem` function show comparable performance. The differences in the results can be attributed to the nonconvexity of the optimization problem and the usage of different initial approximations (unstructured low-rank approximation for the `ident` function and a subspace identification method for the `pem` function). In terms of execution time, except for example 17, the `ident` function is faster than `pem`. It should be noted that in example 17, the accuracy achieved by `ident` is significantly higher on both identification and validation data than the one achieved by `pem`.

## 9 Conclusions

The paper presented a method for system identification in the behavioral setting, with the following salient features:

1. representation free problem formulation,
2. identification from multiple experiments,
3. specification of exact variables,
4. missing values in arbitrary variables and moments of time,
5. implementation in a literature programming style.

Application of the method for  $L_2$  optimal model reduction, identification from multiple data sets of different length, identification from data with missing values, and benchmark problems from the DAISY dataset was considered and illustrated on numerical examples. The developed software package was compared with state-of-the-art system identification packages, with respect to accuracy and computational efficiency. Despite of its generality and flexibility, the developed software is functionally equivalent to and computationally faster than existing alternatives.

## Acknowledgements

The structured low-rank approximation package [8], used in the implementation of the identification method, is a joint work with Konstantin Usevich. During the preparation of the manuscript, I had discussions with Peter Young, Rik Pintelon, Johan Schoukens, Lennart Ljung, and Martine Olivi on various aspects of the results presented in Sections 5–8.

The research leading to these results has received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007-2013) / ERC Grant agreement number 258581 “Structured low-rank approximation: Theory, algorithms, and applications”.

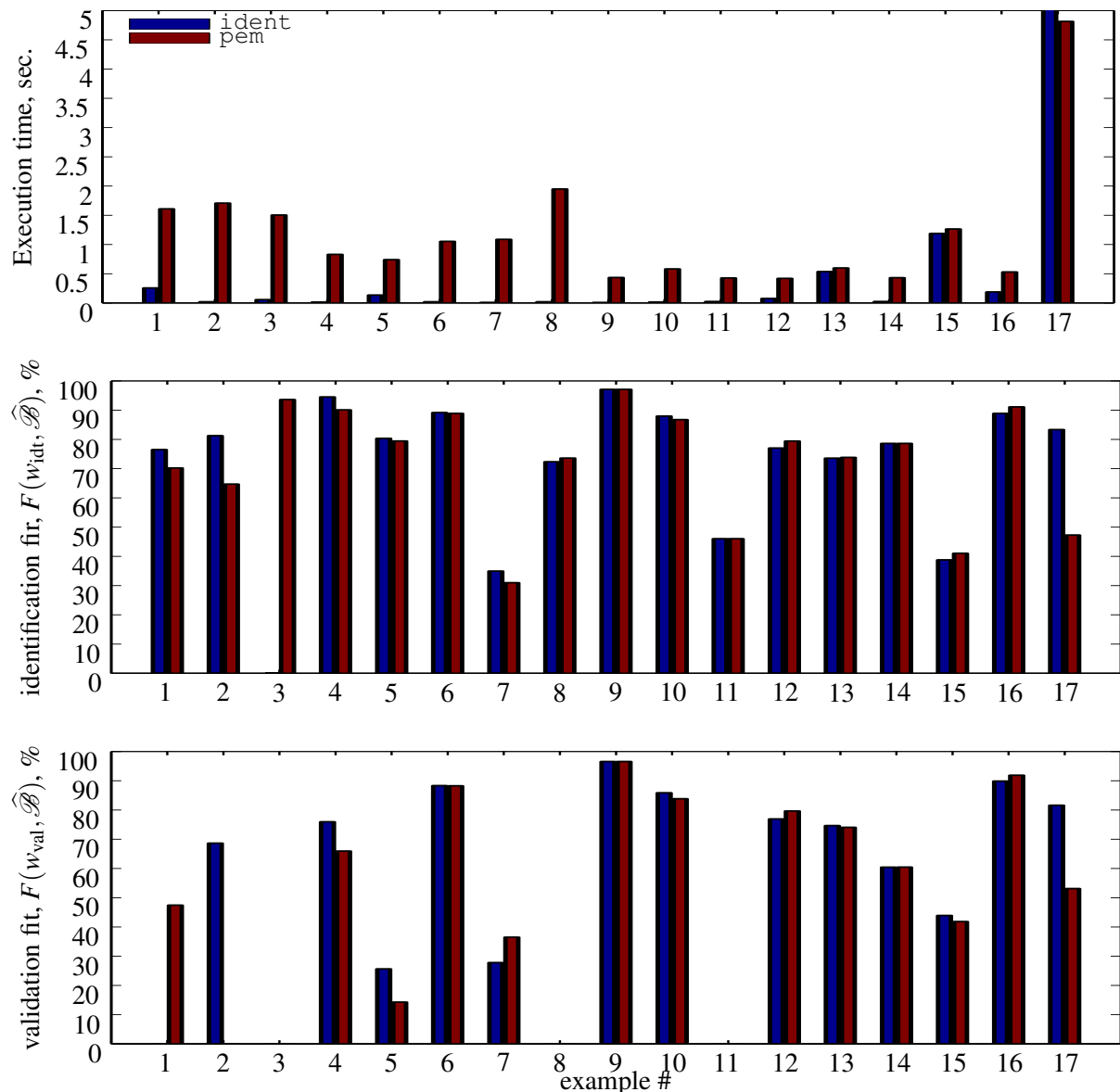


Figure 7: Results on all data sets by splitting of the data into first 70% for identification and remaining 30% for validation.

## References

- [1] M. Galassi et al. GNU scientific library reference manual. <http://www.gnu.org/software/gsl/>.
- [2] D. Knuth. *Literate programming*. Cambridge University Press, 1992.
- [3] A. Kukush, I. Markovsky, and S. Van Huffel. Consistency of the structured total least squares estimator in a multivariate errors-in-variables model. *J. Statist. Plann. Inference*, 133(2):315–358, 2005.
- [4] S. Kung. A new identification method and model reduction algorithm via singular value decomposition. In *Proc. 12th Asilomar Conf. Circuits, Systems, Computers*, pages 705–714, Pacific Grove, 1978.
- [5] L. Ljung. *System Identification Toolbox: User’s guide*. The MathWorks.
- [6] I. Markovsky. Structured low-rank approximation and its applications. *Automatica*, 44(4):891–909, 2008.
- [7] I. Markovsky and P. Rapisarda. Data-driven simulation and control. *Int. J. Control*, 81(12):1946–1959, 2008.

- [8] I. Markovsky and K. Usevich. Software for weighted structured low-rank approximation. Technical Report 339974, Univ. of Southampton, <http://eprints.soton.ac.uk/339974>, 2012.
- [9] I. Markovsky, S. Van Huffel, and R. Pintelon. Block-Toeplitz/Hankel structured total least squares. *SIAM J. Matrix Anal. Appl.*, 26(4):1083–1099, 2005.
- [10] I. Markovsky, J. C. Willems, S. Van Huffel, and B. De Moor. Software for approximate linear system identification. In *Proc. 44th Conf. on Decision and Control*, pages 1559–1564, Seville, Spain, 2005.
- [11] I. Markovsky, J. C. Willems, S. Van Huffel, B. De Moor, and R. Pintelon. Application of structured total least squares for system identification and model reduction. *IEEE Trans. Automat. Control*, 50(10):1490–1500, 2005.
- [12] Jean-Paul Marmorat and Martine Olivi. RARL2 software (realizations and rational approximation in  $L_2$  norm). <http://www-sop.inria.fr/apics/RARL2>.
- [13] B. De Moor, P. De Gerssem, B. De Schutter, and W. Favoreel. DAISY: A database for identification of systems. *Journal A*, 38(3):4–5, 1997. Available from <http://homes.esat.kuleuven.be/~smc/daisy/>.
- [14] R. Pintelon and J. Schoukens. Frequency domain system identification with missing data. *IEEE Trans. Automat. Control*, 45(2):364–369, February 2000.
- [15] R. Pintelon and J. Schoukens. *System Identification: A Frequency Domain Approach*. IEEE Press, Piscataway, NJ, 2001.
- [16] R. Pintelon and J. Schoukens. *System Identification: A Frequency Domain Approach*. IEEE Press, Piscataway, NJ, second edition, 2012.
- [17] J. Polderman and J. C. Willems. *Introduction to mathematical systems theory*. Springer-Verlag, New York, 1998.
- [18] N. Ramsey. Literate programming simplified. *IEEE Software*, 11:97–105, 1994.
- [19] B. Roorda. *Global Total Least Squares—a method for the construction of open approximate models from vector time series*. PhD thesis, Tinbergen Institute, 1995.
- [20] J. Schoukens, G. Vandersteen, Y. Rolain, and R. Pintelon. Frequency response function measurements using concatenated subrecords with arbitrary length. *IEEE Trans. on Instr. and Measurement*, 61(10):2682–2688, 2012.
- [21] T. Söderström. Errors-in-variables methods in system identification. *Automatica*, 43:939–958, 2007.
- [22] K. Usevich and I. Markovsky. Variable projection for affinely structured low-rank approximation in weighted 2-norm. *J. Comput. Appl. Math.*, 2013. Available from <http://arxiv.org/abs/1211.3938>.
- [23] P. Van Overschee and B. De Moor. *Subspace identification for linear systems: Theory, implementation, applications*. Kluwer, Boston, 1996.
- [24] J. C. Willems. From time series to linear system—Part I. Finite dimensional linear time invariant systems, Part II. Exact modelling, Part III. Approximate modelling. *Automatica*, 22, 23:561–580, 675–694, 87–115, 1986, 1987.
- [25] J. C. Willems. Paradigms and puzzles in the theory of dynamical systems. *IEEE Trans. Automat. Control*, 36(3):259–294, 1991.



## A Input/output parameters of `ident` and `misfit`

- `w` is the given set of time series  $w_d$  — a real Matlab array of dimension  $T \times q \times N$ , where  $T$  is the number of samples,  $q$  is the number of variables, and  $N$  is the number of time series. In case of multiple experiments of different duration, `w` should be specified as a cell array with  $N$  cells, each one of which is a  $T_i \times q$  matrix containing the  $i$ th time series, *i.e.*,

$$w(\mathbf{t}, :, \mathbf{k}) = (w^k(t))^\top \quad \text{or} \quad w\{\mathbf{k}\}(\mathbf{t}, :) = (w^k(t))^\top.$$

- `(m, ell)` is the complexity specification (input dimension and lag).
- `opt` is a optional argument specifying exact variables, exact initial conditions, and options for the optimization solver, used by the `ident` function. The options are passed to the functions `ident` and `misfit` as fields of a structure.
  - `'exct'` (default value `[]`) —  $q$ -dimensional vector or  $N$ -dimensional cell array with  $q$ -dimensional vector elements, specifying the indices of the exact variables.
  - `'wini'` — specifies exact initial conditions. If `wini = 0`, exact zero initial conditions are specified, *i.e.*,  $\text{col}(0, \hat{w}^k) \in \hat{\mathcal{B}}$ . More generally, `wini = wini` is an  $\ell$  samples long trajectory (specified by an  $\ell \times q \times N$  array or a  $N$ -dimensional cell array of  $\ell \times q$  matrices), defining initial conditions for the time series  $\hat{w}$ , *i.e.*,  $\text{col}(w_{\text{ini}}^k, \hat{w}^k) \in \hat{\mathcal{B}}$ .
  - `'sys0'` — initial approximation: an input/state/output representation of a system, given as an `ss` object, with `m` inputs, `p := q - m` outputs, and order `n := lp`. Default value is computed by the `slra` function, using unstructured low-rank approximation.
  - Arguments allowing the user to specify different optimization algorithms (`'solver'` and `'method'`), control the displayed information (`'disp'`), and change the convergence criteria are described in the structured low-rank approximation user's manual [8].
- `sysh` is an input/state/output representation of the identified or validated system  $\hat{\mathcal{B}}$ , given by an `ss` object.
- `info` is a structure, containing exit information from the structured low-rank approximation solver: `info.M` is the misfit  $M(w_d, \hat{\mathcal{B}})$ , `info.time` is the execution time, and `info.iter` is the number of iterations.
- `M` is the misfit  $M(w_d, \hat{\mathcal{B}})$ .
- `wh` is the optimal approximating time series.
- `xini` is a matrix whose columns are the initial condition, under which  $\hat{w}^k, k = 1, \dots, N$  are obtained.

## B Implementation

### B.1 Main functions

Figure 1 visualizes the data processing done when the `ident` function is called. A number of auxiliary functions, explained and defined next, are used. The main one is the structured low-rank approximation solver `slra`, which computes a locally optimal solution of (SLRA). Technically, `slra` is a mex-file calling a C++ solver. For our purposes the `slra` function is a black box with inputs the structure parameter vector  $p$ , the structure specification  $\mathcal{S}$ , upper bound for the rank  $r$ , and (optionally) an initial value of the kernel parameter  $R_{\text{ini}}$ ; and outputs the structure parameter vector  $\hat{p}$  of a locally optimal approximation, and the corresponding kernel parameter  $\hat{R}$ .

The computation done in the `ident` function is

1. transformation of the user defined data  $w_d, m, \ell$ , exact variables, and (optionally) initial model into input parameters for the `slra` function; and

2. transformation of the `slra` function's solution  $(\hat{p}, \hat{R})$  into state-space representation of the optimal approximate model  $\hat{\mathcal{B}} = \mathcal{B}(\hat{A}, \hat{B}, \hat{C}, \hat{D})$ , optimal approximation  $\hat{w}$  of the data  $w_d$ , and corresponding initial conditions  $x_{\text{ini}}$ .

Obtaining a input/state/output representation of the identified model is possible in two different ways: 1) using the  $\hat{R}$  parameter and 2) using the trajectory  $\hat{w}$ . The first option is a realization problem (see Sections B.2) and the second one is a deterministic identification problem [23]. Both transformations are classic problems, for which solutions exists.

The misfit  $M(w_d, \mathcal{B})$  between the data  $w_d$  and the model  $\mathcal{B}$  is the cost function of the approximate identification problem. Its fast computation is a key element of the optimization method used by the `slra` function. A convenient way to access the misfit computation is to call the `ident` function with the data  $w_d$ , initial approximation corresponding to the model  $\mathcal{B}$ , and with specification of zero iterations for the optimization solver. Then the misfit is returned in the field `M` of the output parameter `info`.

## B.2 $R \mapsto$ input/state/output representation

The transformation from kernel to input/state/output representation is done using the standard observer canonical form, defined in the following proposition.

**Proposition 17** ([11, Section IV.A]). *Consider a kernel representation  $\mathcal{B}(R)$  of a linear time-invariant system and let  $R =: [Q \quad -P]$ , with*

$$R_i = [Q_i \quad -P_i], \quad \text{where } Q_i \in \mathbb{R}^{p \times m} \text{ and } P_i \in \mathbb{R}^{p \times p}.$$

*Assuming that  $P_\ell$  is nonsingular, a minimal state-space representation  $\mathcal{B}(A, B, C, D)$  of the system  $\mathcal{B}(R)$ , i.e.,*

$$\mathcal{B}(A, B, C, D) = \mathcal{B}([Q \quad P]),$$

*is given by*

$$A = \begin{bmatrix} 0 & \cdots & 0 & -P_\ell^{-1}P_0 \\ I_p & & & -P_\ell^{-1}P_1 \\ & \ddots & & \\ & & I_p & -P_\ell^{-1}P_{\ell-1} \end{bmatrix}, \quad B = \begin{bmatrix} P_\ell^{-1}(Q_0 - P_0P_\ell^{-1}Q_\ell) \\ P_\ell^{-1}(Q_1 - P_1P_\ell^{-1}Q_\ell) \\ \vdots \\ P_\ell^{-1}(Q_{\ell-1} - P_{\ell-1}P_\ell^{-1}Q_\ell) \end{bmatrix}, \quad C = [0 \quad \cdots \quad 0 \quad I_p], \quad D = P_\ell^{-1}Q_\ell.$$

## B.3 Estimation of the initial state

Given an input/state/output representation  $\mathcal{B}(A, B, C, D)$  and a trajectory  $w = \text{col}(u, y)$  of that system, the corresponding initial state  $x_{\text{ini}}$  is computed from the system of linear equations

$$y - y_{\text{forced}} = \mathcal{O}_L(A, C)x_{\text{ini}}, \tag{x_{\text{ini}}}$$

where  $y_{\text{forced}}$  is the output of the system to input  $u$  and zero initial conditions and

$$\mathcal{O}_L(A, C) := \text{col}(C, CA, \dots, CA^{L-1})$$

is the extended observability matrix.

In order to determine uniquely  $x_{\text{ini}}$  from  $(x_{\text{ini}})$ , it is sufficient to choose the parameter  $L$  equal to the lag  $\ell$ . For  $L > \ell$ , existence of solution for  $(x_{\text{ini}})$  can be used as a test for ‘‘exactness’’ of the trajectory  $(w(1), \dots, w(L))$  with respect to the model  $\mathcal{B}(A, B, C, D)$ , i.e., a test for  $(w(1), \dots, w(L)) \in \mathcal{B}(A, B, C, D)$ .

## B.4 Input/state/output representation $\mapsto R$

The transformation from input/state/output to kernel representation is done using the following proposition

**Proposition 18.** *Consider an input/state/output representation  $\mathcal{B}(A, B, C, D)$  of a linear time-invariant system with order that is a multiple of the number of outputs. We have,*

$$\mathcal{B}(A, B, C, D) = \mathcal{B}([Q \quad -P]),$$

where

$$P := (\mathcal{O}_{\ell+1}^\perp(A, C))^\top \quad \text{and} \quad Q := P \begin{bmatrix} D & 0 & 0 & \dots & 0 \\ CB & D & 0 & \ddots & \vdots \\ CAB & CB & D & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & 0 \\ CA^{\ell-2}B & \dots & CAB & CB & D \end{bmatrix}. \quad ((A, B, C, D) \mapsto (P, Q))$$

(For a matrix  $M$ ,  $M^\perp$  is a full column rank matrix of maximal dimension, such that  $M^\perp M = 0$ .)

*Proof.* The columns of  $\mathcal{O}_{\ell+1}(A, C)$  are free trajectories of the system  $\mathcal{B}(A, B, C, D)$ , so that they are annihilated by the  $P$  parameter

$$P \mathcal{O}_{\ell+1}(A, C) = 0.$$

Moreover, by the assumption that the order is a multiple of the number of outputs, we have that

$$\dim(\text{leftker}(\mathcal{O}_{\ell+1}^\perp(A, C))) = p.$$

Therefore, the rows of  $P$  form a basis for the left kernel of  $\mathcal{O}_{\ell+1}^\perp(A, C)$ . This proves that

$$P := (\mathcal{O}_{\ell+1}^\perp(A, C))^\top.$$

Consider the sequence of the first  $\ell + 1$  samples of the impulse response of  $\mathcal{B}(A, B, C, D)$ , padded with  $\ell$  zeros:

$$Y := (\underbrace{0, \dots, 0}_\ell, D, \underbrace{CB, CAB, \dots, CA^{\ell-1}B}_\ell),$$

and the corresponding input sequence

$$U := (\underbrace{0, \dots, 0}_\ell, I_m, \underbrace{0, 0, \dots, 0}_\ell).$$

Since  $(U, Y)$  is a matrix valued trajectory of the system  $\mathcal{B}(\begin{bmatrix} Q & -P \end{bmatrix})$ , we have that

$$Q \mathcal{H}_{\ell+1}(U) = P \mathcal{H}_{\ell+1}(Y) \quad \implies \quad Q = P \mathcal{H}_{\ell+1}(Y) \begin{bmatrix} I_p \\ \vdots \\ I_p \end{bmatrix}. \quad (*)$$

The second equation in  $((A, B, C, D) \mapsto (P, Q))$  follows from (\*).  $\square$

## C Input/state/output representation of the model reduction benchmark

Let  $p := \text{col}(p_1, \dots, p_N)$  and define the state vector by

$$x = \text{col}\left(p, \frac{d}{dt}p, x_{N+1}\right),$$

where  $x_{N+L}$  is a constant,  $x_{N+1}(t) = x_{N+1}(0)$ , for all  $t$ . The order of the state space representation is  $n = 2N + 1$ . Using (SYS), we obtain an input/state/output  $\mathcal{B}(A, B, C, D)$  representation with parameters

$$A = \begin{bmatrix} 0 & I_N & 0 \\ A_{21} & A_{22} & A_{23} \\ 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{(2N+1) \times (2N+1)}, \quad B = \begin{bmatrix} 0 \\ E_{\mathcal{Y}} \\ 0 \end{bmatrix}, \quad C = [E_{\mathcal{Y}}^\top \quad 0 \quad 0], \quad D = 0, \quad x_{2N+1}(0) = 1,$$

where

$$A_{21} = \begin{bmatrix} -\frac{k_0+k_1}{m_1} & \frac{k_2}{m_1} & & & & \\ \frac{k_1}{m_2} & -\frac{k_1+k_2}{m_2} & \frac{k_3}{m_2} & & & \\ & \ddots & \ddots & \ddots & & \\ & & \frac{k_{N-2}}{m_{N-2}} & -\frac{k_{N-2}+k_{N-1}}{m_{N-2}} & \frac{k_N}{m_{N-2}} & \\ & & & \frac{k_{N-1}}{m_N} & -\frac{k_{N-1}+k_N}{m_N} & \end{bmatrix}$$

and

$$A_{22} = \begin{bmatrix} -\frac{d_0+d_1+f}{m_1} & \frac{d_2}{m_1} & & & & \\ \frac{d_1}{m_2} & -\frac{d_1+d_2+f}{m_2} & \frac{d_3}{m_2} & & & \\ & \ddots & \ddots & \ddots & & \\ & & \frac{d_{N-2}}{m_{N-2}} & -\frac{d_{N-2}+d_{N-1}+f}{m_{N-2}} & \frac{d_N}{m_{N-2}} & \\ & & & \frac{d_{N-1}}{m_N} & -\frac{d_{N-1}+d_N+f}{m_N} & \end{bmatrix}, \quad A_{23} = \begin{bmatrix} \frac{k_0-k_1}{m_1} \\ \frac{k_1-k_2}{m_2} \\ \vdots \\ \frac{k_{N-2}-k_{N-1}}{m_{N-1}} \\ \frac{k_{N-1}+k_N}{m_N} \end{bmatrix} \delta.$$

Define the unit vector  $e_i$ , as the  $i$ th column of the  $n \times n$  identity matrix. Let  $p$  be the number of elements of the set  $\mathcal{Y}$  and  $m$  be the number of elements of the set  $\mathcal{U}$ . The matrix  $E_{\mathcal{Y}}$  is the  $n \times m$  matrix with  $i$ th column  $1/m_i \cdot e_i$  and, similarly,  $E_{\mathcal{U}}$  is the  $n \times p$  matrix with  $i$ th column  $e_i$ .

The order of the system  $\mathcal{B}(A, B, C, D)$  is  $2N$ , *i.e.*, the state space representation is nonminimal. (It has an uncontrollable mode  $x_{N+1} = \text{const.}$ )