

Leading Impulse Response Identification via the Weighted Elastic Net Criterion

Giuseppe C. Calafiore, Carlo Novara, Michele Taragna

*Dipartimento di Automatica e Informatica, Politecnico di Torino,
Corso Duca degli Abruzzi 24, I-10129, Torino, Italy.*

Abstract

This paper deals with the problem of finding a low-complexity estimate of the impulse response of a linear time-invariant discrete-time dynamic system from noise-corrupted input-output data. To this purpose, we introduce an identification criterion formed by the average (over the input perturbations) of a standard prediction error cost, plus a weighted ℓ_1 regularization term which promotes sparse solutions. While it is well known that such criteria do provide solutions with many zeros, a critical issue in our identification context is *where* these zeros are located, since sensible low-order models should be zero in the tail of the impulse response. The flavor of the key results in this paper is that, under quite standard assumptions (such as i.i.d. input and noise sequences and system stability), the estimate of the impulse response resulting from the proposed criterion is indeed identically zero from a certain time index n_l (named the *leading order*) onwards, with arbitrarily high probability, for a sufficiently large data cardinality N . Numerical experiments are reported that support the theoretical results, and comparisons are made with some other state-of-the-art methodologies.

Key words: FIR identification, ℓ_1 regularization, Elastic Net, Lasso, Sparsity

1 Introduction

A large part of the literature on identification of linear time-invariant (LTI) dynamic systems follows a statistical approach (Ljung [1999a], Söderström and Stoika [1989]), where probabilistic assumptions are made, at least on the noise corrupting the measurements. The techniques available in this context may be classified in two main categories: parametric and nonparametric. Parametric techniques are mainly based on the prediction error methods (PEMs) or on the maximum likelihood approach, if Gaussian noise is assumed. The identified models belong to finite-dimensional spaces of given order, like FIR, ARX, ARMAX, OE, Laguerre, Kautz or orthonormal basis function models. In order to limit the model complexity and to avoid possible overfitting, a tradeoff between bias and variance is usually considered, and the model order selection is performed by optimizing some suitable cost function – such as the Akaike’s information criterion AIC (Akaike [1974]), the Rissanen’s Minimum Description Length MDL, or the Bayesian information criterion BIC (Rissanen [1978],

Schwarz [1978]) – and by applying some form of cross validation (CV), like hold-out or leave-one-out. Possible limits of these parametric methods have been pointed out in Pilonetto and De Nicolao [2010], Pilonetto et al. [2011], Chen et al. [2012], where it is shown that the sample properties of PEM approaches equipped with, e.g., AIC and CV, may be rather unsatisfactory and quite far from those predicted by standard (i.e., without model selection) statistical theory.

The nonparametric techniques aim to obtain the overall system’s impulse response as a suitable deconvolution of observed input-output data. In particular, very promising approaches have been recently developed, based on results coming from the machine learning field, see, e.g., Pilonetto et al. [2014] and the references therein. Rather than postulating finite-dimensional hypothesis spaces, the estimation problem is tackled in an infinite-dimensional space, and the intrinsic ill-posedness of the problem is circumvented by using suitable regularization methods. In particular, the system’s impulse response is modeled as a zero-mean Gaussian process, and the prior information is introduced by simply assigning a specific covariance, named *kernel* in the machine learning literature. This procedure can be interpreted as the counterpart of model order selection in the parametric PEM approach and, in some cases, it is shown to

Email addresses:

giuseppe.calafiore@polito.it (Giuseppe C. Calafiore),
carlo.novara@polito.it (Carlo Novara),
michele.taragna@polito.it (Michele Taragna).

be much more robust.

In the present paper, a novel nonparametric method is presented, whereby an estimate of the system's impulse response is obtained by minimizing a suitable cost function that directly takes into account the resulting model complexity. The aim is indeed to obtain a low-complexity model of the system, in the form of a reduced-order FIR (in this sense, the approach is not so far from parametric techniques). A key feature of the proposed approach, representing a relevant improvement over the state of the art, is that it allows for an effective model order selection, without using strong a-priori information on the true system. More specifically, we propose the use of an identification criterion which is a weighted combination of (a) a standard prediction error term, (b) an ℓ_2 regularization term, and (c) a weighted ℓ_1 penalty term which promotes sparse solutions; a full justification for such criterion is given in Section 3.2. This type of criterion corresponds to the so-called Elastic Net cost, which recently became popular in the machine learning community, see, e.g., Zou and Hastie [2005], De Mol et al. [2009]. Notice that, while it is well known that the use of ℓ_1 regularization leads to sparse solutions, sparsity alone is not a very interesting feature in our identification context. Indeed, reduced-order models are obtained only if the sparsity of the solution follows a specific pattern, whereby the zeros are all concentrated in the tail of the impulse response. Obtaining such a pattern is not obvious, nor a-priori granted by the ℓ_1 regularization. One of the key contributions of this paper is to prove that, under standard assumptions, the impulse response estimated via our Elastic-Net type of criterion has the property of being indeed nonzero only on the initial part of the impulse response (which we shall name the *leading response*), with arbitrarily high probability, if the number of data N is sufficiently large.

The present paper is organized as follows. In Section 2 the notation is set, and some preliminary results on a Chebyshev's type of convergence for random variables are stated. Section 3 describes the linear identification problem of interest, and contains the derivations of the Elastic Net cost. The main results on the recovery of the leading part of the impulse response are contained in Section 4. Section 5 illustrates a practical procedure for implementing the proposed identification scheme. Numerical experiments, including a comparative discussion with other identification methods, are given in Section 6. All proofs are contained in the Appendix.

2 Notation and preliminaries

2.1 Notation

For a vector $x \in \mathbb{R}^N$, we denote by $[x]_i$ the i -th entry of x , and we define its *support* as

$$\text{supp}(x) \doteq \{i \in \{1, \dots, N\} : [x]_i \neq 0\}.$$

The notation $\|x\|_p$ represents the standard ℓ_p norm of x , and $\|x\|_0$ denotes the cardinality of $\text{supp}(x)$, that is the number of nonzero entries of x .

For a matrix $X \in \mathbb{R}^{N, M}$ (with M possibly equal to ∞),

we denote by $[X]_{i,j}$ the entry of X in row i and column j . For $n \leq M$, we denote by $X_{\uparrow n} \in \mathbb{R}^{N, n}$ the sub-matrix formed by the first n columns of X , with $X_{\downarrow n} \in \mathbb{R}^{N, M-n}$ the sub-matrix formed by the columns of X of indices $n+1, \dots, M$, and with $X_{\#n}$ the $n \times n$ principal sub-matrix of X . The identity matrix is denoted by I , or by I_n , if we wish to specify its dimension. We denote by X^\dagger the Moore-Penrose pseudo-inverse of X ; if X has full column rank, then $X^\dagger = (X^\top X)^{-1} X^\top$.

If x is a random variable, then $\mathbb{E}\{x\}$ denotes the expected value of x , and $\text{var}\{x\}$ denotes its variance: $\text{var}\{x\} = \mathbb{E}\{(x - \mathbb{E}\{x\})^2\}$. \mathbb{P} denotes a probability measure on x . The symbol \rightsquigarrow implies almost sure convergence, and it is formally defined in Section 2.2.1.

2.2 Chebyshev's inequality for certain empirical means

Let x_i , $i = 1, \dots$, be a sequence of (not necessarily independent) random variables such that $\mathbb{E}\{x_i\} = \mu < \infty$ for all i , $\text{var}\{x_i\} = \sigma_i^2 \leq \bar{\sigma}^2 < \infty$ for all i , and $\mathbb{E}\{(x_i - \mu)(x_j - \mu)\} = 0$ for all $i \neq j$. For given $N \geq 1$, define the empirical mean

$$\hat{x}_N \doteq \frac{1}{N} \sum_{i=1}^N x_i.$$

Obviously, from linearity of the expectation, it holds that $\mathbb{E}\{\hat{x}_N\} = \mu$. Further, we have that

$$\begin{aligned} \sigma^2 \doteq \text{var}\{\hat{x}_N\} &= \mathbb{E}\{(\hat{x}_N - \mu)^2\} = \frac{1}{N^2} \mathbb{E}\left\{\left[\sum_{i=1}^N (x_i - \mu)\right]^2\right\} \\ &= \frac{1}{N^2} \left[\sum_{i=1}^N \mathbb{E}\{(x_i - \mu)^2\} + \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{E}\{(x_i - \mu)(x_j - \mu)\} \right] \\ &= \sum_{i=1}^N \sigma_i^2 / N^2 \leq \bar{\sigma}^2 / N, \end{aligned}$$

where the last passages follow from the fact that the x_i s are uncorrelated, and have first moment μ and variance $\sigma_i^2 \leq \bar{\sigma}^2$. Chebyshev's inequality applied to the random variable \hat{x}_N thus states that, for any $\eta > 0$,

$$\mathbb{P}\{|\hat{x}_N - \mu| \geq \eta\sigma\} \leq 1/\eta^2. \quad (1)$$

Since $\eta\sigma \leq \eta\bar{\sigma}/\sqrt{N}$, we have that $\mathbb{P}\{|\hat{x}_N - \mu| \geq \eta\bar{\sigma}/\sqrt{N}\} \leq \mathbb{P}\{|\hat{x}_N - \mu| \geq \eta\sigma\}$, whence, from (1), we obtain that $\mathbb{P}\{|\hat{x}_N - \mu| \geq \eta\bar{\sigma}/\sqrt{N}\} \leq 1/\eta^2$. Equivalently, we can state that, for any $\epsilon > 0$, it holds that

$$\mathbb{P}\{|\hat{x}_N - \mu| \geq \epsilon\} \leq \bar{\sigma}^2 / (N\epsilon^2).$$

We thus conclude that, for any given accuracy $\epsilon > 0$ and probability $\beta \in (0, 1)$, it holds that

$$\mathbb{P}\{|\hat{x}_N - \mu| \geq \epsilon\} \leq \beta, \quad \forall N \geq \lceil \bar{\sigma}^2 / (\beta\epsilon^2) \rceil.$$

Notice that (1) implies that $\mathbb{P}\{|\hat{x}_N - \mu| > \eta\sigma\} \leq 1/\eta^2$; hence, by considering the complementary event, it also holds that $\mathbb{P}\{|\hat{x}_N - \mu| \leq \eta\sigma\} \geq 1 - 1/\eta^2$, from which it follows that

$$\mathbb{P}\{|\hat{x}_N - \mu| \leq \epsilon\} \geq 1 - \bar{\sigma}^2 / (N\epsilon^2).$$

2.2.1 Meaning of the convergence symbol \rightsquigarrow

For a random variable z_N that depends on N and for a given real value \bar{z} , the notation $z_N \rightsquigarrow \bar{z}$ means that for any given $\epsilon > 0$ and $\beta \in (0, 1)$ there exists a *finite* integer $N_{\epsilon, \beta}$ such that

$$\mathbb{P}\{|z_N - \bar{z}| \geq \epsilon\} \leq \beta, \quad \forall N \geq N_{\epsilon, \beta}. \quad (2)$$

Notice that $z_N \rightsquigarrow \bar{z}$ implies that z_N converges to \bar{z} *almost surely* (that is, with probability one), as N tends

to infinity. However, we are specifically interested in the property in (2), that holds for possibly large, but finite, N .

2.3 Lipschitz functions of random variables

If z_N is the empirical mean of N uncorrelated variables with common mean μ and variance bounded by $\bar{\sigma}^2$ then, from the discussion in Section 2.2, we conclude that indeed $z_N \rightsquigarrow \mu$ and, in particular, (2) holds for $N_{\epsilon, \beta} = \lceil \bar{\sigma}^2 / (\beta \epsilon^2) \rceil$. However, we shall use the convergence notation $z_N \rightsquigarrow \bar{z}$ also when \bar{z} is not necessarily the expected value of z_N , and/or when z_N is not necessarily an empirical mean. The following lemma holds.

Lemma 1 *For any fixed integer p , let y_1, \dots, y_p be (possibly correlated) scalar random variables that depend on N and such that $y_i \rightsquigarrow \bar{y}_i$, $i = 1, \dots, p$, for some given values $\bar{y}_1, \dots, \bar{y}_p$. Let f be a Lipschitz continuous function from \mathbb{R}^p into \mathbb{R} , such that $f(\bar{y}_1, \dots, \bar{y}_p)$ is finite. Then, it holds that $f(y_1, \dots, y_p) \rightsquigarrow f(\bar{y}_1, \dots, \bar{y}_p)$. ■*

Appendix A.1 contains a proof of Lemma 1.

3 Problem setup

3.1 A linear measurement model

We consider an identification experiment in which a discrete-time scalar input signal $\tilde{u}(k)$ enters an LTI dynamic system, which produces in response a scalar output signal $\tilde{y}(k)$. This output is acquired via noisy measurements over a time window $k = 1, \dots, N$, obtaining a sequence of output measurements $y(k) = \tilde{y}(k) + \delta_y(k)$, $k = 1, \dots, N$, where $\delta_y(k)$ is the measurement noise sequence. Since the unknown system is assumed to be LTI, there exists a linear relation between the output measurements and the unknown system's impulse response $h(i)$, $i = 1, \dots$. Assuming that the system is operating in steady state, this relation is given by the discrete-time convolution: for $k = 1, \dots, N$,

$$y(k) = \tilde{y}(k) + \delta_y(k) = \sum_{i=1}^{\infty} \tilde{u}(k - i + 1)h(i) + \delta_y(k). \quad (3)$$

Observe that, following a nonparametric approach, we do not assume to know in advance the order of the unknown system; therefore, in (3), all values $h(i)$ can be, a priori, nonzero. Letting

$$y \doteq \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix}; \quad \delta_y \doteq \begin{bmatrix} \delta_y(1) \\ \delta_y(2) \\ \vdots \\ \delta_y(N) \end{bmatrix}; \quad \tilde{u}_i \doteq \begin{bmatrix} \tilde{u}(2-i) \\ \tilde{u}(3-i) \\ \vdots \\ \tilde{u}(N+1-i) \end{bmatrix},$$

for $i = 1, 2, \dots$, we can write (3) in vector format as

$$y = \sum_{i=1}^{\infty} \tilde{u}_i h(i) + \delta_y. \quad (4)$$

For any integer $n \geq 0$, we define

$\tilde{U}_{\uparrow n} \doteq [\tilde{u}_1 \dots \tilde{u}_n] \in \mathbb{R}^{N, n}$, $h_{\uparrow n} \doteq [h(1) \dots h(n)]^T \in \mathbb{R}^n$, as well as the semi-infinite matrices and vectors

$$\tilde{U}_{\downarrow n} \doteq [\tilde{u}_{n+1} \tilde{u}_{n+2} \dots] \in \mathbb{R}^{N, \infty},$$

$$h_{\downarrow n} \doteq [h(n+1) h(n+2) \dots]^T \in \mathbb{R}^{\infty}.$$

Let now $q \leq N$ be a given integer: our goal is to estimate the first q elements of the impulse response h (i.e., to estimate $h_{\uparrow q} \in \mathbb{R}^q$), from N noisy output measurements. The value of q is fixed by the decision maker, based on

the available number of measurements N and on a priori knowledge. For instance, under a standard assumption of stability (see Assumption 2), since $h(i)$ decays exponentially, one may a priori assess that the response will be negligible for $i \geq q$, for some sufficiently large q . We can then rewrite (4) as

$$y = \tilde{U}_{\uparrow q} h_{\uparrow q} + \delta_y + y^{\text{ud}},$$

where

$$y^{\text{ud}} \doteq \tilde{U}_{\downarrow q} h_{\downarrow q}$$

represents the unmodelled dynamics due to the truncation of the impulse response to the q -th term. For simplifying the notation, we let from now on

$$\tilde{U} \doteq \tilde{U}_{\uparrow q},$$

which is an $N \times q$ Toeplitz matrix.

3.2 An Elastic Net identification criterion

The initial approach that we consider for identifying the unknown system's impulse response consists in finding an estimate of $h_{\uparrow q}$ that minimizes w.r.t. x the cost function

$$\frac{1}{\gamma} \|y - \tilde{U}x\|_2^2 + \|x\|_0, \quad (5)$$

where $\gamma > 0$ is a suitable tradeoff parameter. The first term in (5) is the standard prediction error, while the second term $\|x\|_0$ represents the cardinality of x , that is the number of nonzero entries in x . This term penalizes the complexity of the estimate, thus promoting solutions with a small number of nonzero entries. Note incidentally that, if $\delta_y(k)$ is a sequence of independent identically distributed (i.i.d.) Normal random variables with zero mean and known variance σ_y^2 then, for $\gamma = 2\sigma_y^2$, the above criterion coincides with the well-known Akaike's information criterion AIC. Other standard criteria, such as the BIC, can also be obtained for different values of γ .

3.2.1 Input uncertainty and averaged cost

In a realistic identification experiment, however, the input signal $\tilde{u}(k)$ that enters the unknown system is a possibly "perturbed" version of a nominal input signal $u(k)$ that the user intends to provide to the system. To model this situation, we assume that $\tilde{u}(k) = u(k) + \delta_u(k)$, where $u(k)$ is the nominal input signal, and $\delta_u(k)$ is an i.i.d. random noise sequence, which is assumed to have zero mean and variance σ_u^2 (setting $\sigma_u^2 = 0$ we recover the standard, no-input-noise, situation). Considering the time window $k = 1, \dots, N$, we have in matrix form that

$$\tilde{U} = U + \Delta, \quad (6)$$

where U is an $N \times q$ Toeplitz matrix containing the nominal input signal, and Δ is an $N \times q$ Toeplitz matrix containing the noise samples $\delta_u(k)$. Specifically, $U \doteq [u_1 \dots u_q]$, and $\Delta \doteq [\delta_1 \dots \delta_q]$, where for $i = 1, \dots, q$

$$u_i \doteq \begin{bmatrix} u(2-i) \\ u(3-i) \\ \vdots \\ u(N+1-i) \end{bmatrix}, \quad \delta_i \doteq \begin{bmatrix} \delta_u(2-i) \\ \delta_u(3-i) \\ \vdots \\ \delta_u(N+1-i) \end{bmatrix}.$$

We account for input uncertainty in the identification experiment by "averaging" the effect of this uncertainty

in the cost criterion (5). This leads to the following cost function:

$$\begin{aligned} J_0(x) &= \mathbb{E}_{\delta_u} \left\{ \frac{1}{\gamma} \|y - \tilde{U}x\|_2^2 + \|x\|_0 \right\} \\ &= \frac{1}{\gamma} \mathbb{E}_{\delta_u} \{ \|y - (U + \Delta)x\|_2^2 \} + \|x\|_0, \end{aligned} \quad (7)$$

where \mathbb{E}_{δ_u} denotes expectation w.r.t. the random sequence δ_u . Elaborating on the expression (7), we obtain

$$\begin{aligned} &\mathbb{E}_{\delta_u} \{ \|y - (U + \Delta)x\|_2^2 \} \\ &= \mathbb{E}_{\delta_u} \{ \|y - Ux\|_2^2 + \|\Delta x\|_2^2 - 2(y - Ux)^\top \Delta x \} \\ &= \|y - Ux\|_2^2 + \mathbb{E}_{\delta_u} \{ \|\Delta x\|_2^2 \} \\ &= \|y - Ux\|_2^2 + x^\top \mathbb{E}_{\delta_u} \{ \Delta^\top \Delta \} x, \end{aligned}$$

because $\mathbb{E}_{\delta_u} \{ \Delta \} = 0$. Since $\delta_u(k)$ is an i.i.d. sequence, and since Δ has Toeplitz structure, it is easy to verify that the off-diagonal terms in $\mathbb{E}_{\delta_u} \{ \Delta^\top \Delta \}$ are zero, while the diagonal terms are all equal to $N\sigma_u^2$. Therefore, it holds that $\mathbb{E}_{\delta_u} \{ \Delta^\top \Delta \} = N\sigma_u^2 I_q$, and the expected cost $J_0(x)$ is explicitly expressed as

$$J_0(x) = \frac{1}{\gamma} \|y - Ux\|_2^2 + \frac{N\sigma_u^2}{\gamma} \|x\|_2^2 + \|x\|_0. \quad (8)$$

Notice that this setting can be easily extended to wide-sense stationary input noise sequences $\delta_u(k)$, in which case the second term in the above expression takes the form $\frac{N}{\gamma} x^\top R_u x$, where R_u is the autocorrelation matrix of δ_u . For simplicity, however, we here focus on the basic case of an i.i.d. sequence, for which $R_u = \sigma_u^2 I_q$. Observe further that accounting for noise on the input signal results in the introduction of a Tikhonov-type regularization term in (8), a fact that has been previously observed in other contexts such as neural network training, see, e.g., Bishop [1995].

3.2.2 Normalizing the variables

We next rescale the variables in the cost (8) by normalizing the columns of the regression matrix. First, we rewrite $J_0(x)$ as

$$J_0(x) = \frac{1}{\gamma} \|b - \bar{A}x\|_2^2 + \|x\|_0, \quad (9)$$

where

$$b \doteq \begin{bmatrix} y \\ 0 \end{bmatrix}, \quad \bar{A} \doteq \begin{bmatrix} U \\ \sigma_u \sqrt{N} I_q \end{bmatrix}. \quad (10)$$

Second, we let $T \doteq \text{diag}(\|\bar{a}_1\|_2, \dots, \|\bar{a}_q\|_2)^{-1}$, where \bar{a}_i denotes the i -th column of \bar{A} , and perform the change of variable $\tilde{x} = T^{-1}x$, thus the right-hand side of (9) becomes

$$\tilde{J}_0(\tilde{x}) \doteq \frac{1}{\gamma} \|b - A\tilde{x}\|_2^2 + \|\tilde{x}\|_0, \quad (11)$$

where we defined $A \doteq \bar{A}T$, and we used the fact that $\|T\tilde{x}\|_0 = \|\tilde{x}\|_0$, since the cardinality of a vector does not depend on (nonzero) scalings of the entries of the vector. We observe that the columns a_1, \dots, a_q of A now have unit Euclidean norm. We let $\tilde{x}_0^* \doteq \arg \min \tilde{J}_0(\tilde{x})$, and $x_0^* \doteq \arg \min J_0(x)$, where it obviously holds that $x_0^* = T\tilde{x}_0^*$. These optimal solutions are hard to determine numerically in practice. However, we do not need to compute them, we only need them for theoretical purposes.

3.2.3 Weighted ℓ_1 relaxation of the cost function

We now introduce the following tractable relaxation of the cost (11):

$$\tilde{J}_1(\tilde{x}) \doteq \frac{1}{\gamma} \|b - A\tilde{x}\|_2^2 + \|W\tilde{x}\|_1. \quad (12)$$

where $W \doteq \text{diag}(w_1, \dots, w_q)$ is a suitable weighting matrix, with $\max_{k=1, \dots, q} w_k = 1$, $\min_{k=1, \dots, q} w_k > 0$. We shall henceforth assume that the weight sequence is non-decreasing: $w_1 \leq w_2 \leq \dots \leq w_q = 1$.

Notice that, expanding the squared norm in (12), we obtain the cost function \tilde{J}_1 in the form

$$\tilde{J}_1(\tilde{x}) = \frac{1}{\gamma} \|y - UT\tilde{x}\|_2^2 + \frac{N\sigma_u^2}{\gamma} \|T\tilde{x}\|_2^2 + \|W\tilde{x}\|_1, \quad (13)$$

which corresponds to the cost expressed in the original variable $x = T\tilde{x}$

$$J_1(x) = \frac{1}{\gamma} \|y - Ux\|_2^2 + \frac{N\sigma_u^2}{\gamma} \|x\|_2^2 + \|WT^{-1}x\|_1. \quad (14)$$

The cost function (13) is strongly convex, hence the optimal solution $\tilde{x}_1^* \doteq \arg \min \tilde{J}_1(\tilde{x})$ is unique and, equivalently, the minimization of (14) has a unique optimal solution $x_1^* = T\tilde{x}_1^*$. In the following section, we shall study the properties of x_1^* as an estimate of the impulse response $h_{\uparrow q}$. Note that only two parameters (γ and σ_u) have to be chosen to obtain this estimate. A systematic procedure is proposed in Section 5, allowing an effective choice of these parameters, based on the desired trade-off between model complexity and accuracy.

Remark 1 The cost criterion appearing in (13) is a particular version of the Lasso (see, e.g., Tibshirani [1996]), known as the Elastic Net (Zou and Hastie [2005]). The Elastic Net criterion includes an ℓ_2 regularization term which provides shrinkage and improves conditioning of the ℓ_2 -error cost (by guaranteeing strong convexity of the cost), as well as an ℓ_1 penalty term which promotes sparsity in the solution. Elastic Net-based methods are widely used in statistics and machine learning, see, e.g., De Mol et al. [2009], Hastie et al. [2009], and are amenable to very efficient large-scale solution algorithms (Friedman et al. [2010]). To the best of the authors' knowledge, this is the first work in which the Elastic Net criterion is used in the context of a system identification problem and the resulting sparsity pattern is rigorously analyzed.

4 Leading response recovery

This section contains the main results of the paper. First, we report a preliminary technical lemma (Lemma 2) stating that, under a certain condition, the minimizer x_1^* of (14) is supported on $\{1, \dots, n\}$, with $n \leq q$. Second, under some suitable assumptions on the input and noise signals, we show (Theorem 4) that if the unknown system is stable, then for a sufficiently large N and for a given $n \leq q$, there exist explicitly given γ values for which the support of x_1^* is contained in $\{1, \dots, n\}$, with any given high probability. This means that the estimated impulse response x_1^* is not only sparse but, with high probability, it is zero precisely on the tail of the system's impulse response $h_{\uparrow q}$. We next define the notions

of *leading response* and *leading support* of the system's impulse response, and show (Corollary 5) that if the unknown system is stable, then for a suitable γ and a sufficiently large N the support of x_1^* is contained in the leading support, with any given high probability; we call this property *leading response recovery* (LRR). Finally, we show (Corollary 6) that if the true unknown system is FIR then, for a sufficiently large N and for any $\gamma > 0$, the estimated impulse response x_1^* will be sparse, and of order no larger than the order of the true system, with high probability.

4.1 Preliminary results, assumptions and definitions

With the notation set in Section 3.2.2, for a given integer $n \leq q$, let $P_n \doteq A_{\uparrow n} A_{\uparrow n}^\dagger$ denote the orthogonal projector onto the span of $A_{\uparrow n}$, and define the *n-leading recovery coefficient* $\Upsilon_n(A) \doteq 1 - \max_{n < i \leq q} w_i^{-1} \|W_{\#n} A_{\uparrow n}^\dagger a_i\|_1$, where a_i is the i -th column of A . The following technical lemma, based on a result in Tropp [2006], holds.

Lemma 2 *Suppose that for some integer $n \leq q$ it holds that*

$$\|W^{-1} A^\top (b - P_n b)\|_\infty \leq \gamma \Upsilon_n(A)/2, \quad (15)$$

and let x_1^ be the minimizer of (14). Then, it holds that*

$$\text{supp}(x_1^*) \subseteq \{1, \dots, n\}. \quad \blacksquare$$

See Appendix A.2 for a proof of Lemma 2.

Let us now state the following working assumptions.

Assumption 1 *(on input and disturbance sequences)*

- (1) *The input $u(k)$ is an i.i.d. sequence with zero mean, bounded variance ν^2 and bounded 4-th order moment $\mathbb{E}\{u(k)^4\} = \overline{m}_4$.*
- (2) *The noise $\delta_y(k)$ is an i.i.d. sequence with zero mean and bounded variance σ_y^2 .*
- (3) *The input perturbation $\delta_u(k)$ is an i.i.d. sequence with zero mean and bounded variance σ_u^2 .*
- (4) *$u(k)$, $\delta_y(k)$, and $\delta_u(k)$ are mutually uncorrelated.*

Assumption 2 *(Stability)* *The unknown system's impulse response h is such that $|h(i)| \leq L\rho^{i-1}$, for $i = 1, 2, \dots$, for some given finite $L > 0$ and $\rho \in (0, 1)$.*

We next establish a preliminary lemma.

Lemma 3 *Under Assumption 1, for any pair of column vectors u_i and u_j it holds that*

$$\frac{1}{N} u_i^\top u_j \rightsquigarrow \begin{cases} \nu^2 & \text{if } i = j \\ 0 & \text{otherwise,} \end{cases} \quad (16)$$

where the notation \rightsquigarrow has the meaning specified in Section 2.2.1. Also, it holds that

$$\frac{1}{N} u_i^\top \delta_y \rightsquigarrow 0, \quad \forall i \quad (17)$$

$$\frac{1}{N} u_i^\top \delta_j \rightsquigarrow 0, \quad \forall i, j. \quad (18) \quad \blacksquare$$

Appendix A.3 contains a proof of Lemma 3.

We next define the notion of *leading order* of the system's impulse response, and the associated notions of *leading response* and *leading support*.

Definition 1 *Let Assumption 2 hold. We define the leading order, $n_l(N)$, of h as the largest integer $i \leq q$ such that*

$$L\rho^{i-1} \geq \frac{\sigma_y}{\nu} \times \frac{1}{\sqrt{N}}. \quad (19)$$

The leading response is $\{h(i), i = 1, \dots, n_l\}$ and the leading support is $\{1, \dots, n_l\}$.

Remark 2 We provide an intuitive interpretation of the definition in (19). The leading order is a value such that for time values larger than it the system's impulse response cannot essentially be discriminated from noise. Indeed, if a classical output error criterion would be used for estimating $h_{\uparrow q}$, then the covariance matrix of the estimated parameter would be of the form $\sigma_y^2 (U^\top U)^{-1}$, which tends to $\sigma_y^2 / (\nu^2 N) I_q$ as $N \rightarrow \infty$, see the proof of Theorem 4 for details. The standard error on the generic element $h(i)$ of the impulse response thus goes to zero as $1/\sqrt{N}$, where N is the number of measurements and the proportionality constant σ_y/ν is the noise-to-signal ratio. The leading order n_l is therefore defined as the time value after which the upper bound on $|h(i)|$ goes below the level $\eta = \frac{\sigma_y}{\nu} \frac{1}{\sqrt{N}}$, and hence $h(i)$ becomes essentially indistinguishable from noise, for all $i > n_l$; if this condition is not met for $i \leq q$, then we just set $n_l = q$. It is an immediate consequence of (19) that the leading order grows as the logarithm of N , until it saturates to q :

$$n_l(N) = \min \left(\left\lfloor \frac{\log(\nu L) + \frac{1}{2} \log N - \log(\sigma_y \rho)}{\log(\rho^{-1})} \right\rfloor, q \right).$$

4.2 Main results

We next establish the main results of this paper.

Theorem 4 *Let Assumptions 1 and 2 hold. Let $n \leq q$, $\kappa \doteq \nu / \sqrt{\nu^2 + \sigma_u^2}$, and*

$$\gamma = 2\mu w_n^{-1} L \rho^n \nu \kappa \times \sqrt{N}, \quad (20)$$

for some $\mu > 1$. Then, for any given $\beta \in (0, 1)$ there exists a finite integer N_β such that for any $N \geq N_\beta$ it holds that

$$\text{supp}(x_1^*) \subseteq \{1, \dots, n\}$$

with probability no smaller than $1 - \beta$, where x_1^ is the minimizer of (14).* \blacksquare

Appendix A.4 contains a proof of Theorem 4. The key point of this theorem is that if the tradeoff parameter γ is chosen proportional to \sqrt{N} then, with high probability and for a sufficiently large N , the minimization of (14) provides a solution which is not only sparse, but its sparsity pattern is identically zero on the tail of the impulse response, i.e., the estimated impulse response x_1^* is FIR of order at most n .

A consequence of Theorem 4 is stated in the following corollary: for a suitable constant value of γ , the minimizer x_1^* of (14) has its support contained in the leading support.

Corollary 5 *(Leading support recovery)*

Let Assumptions 1 and 2 hold. Let $\kappa \doteq \nu / \sqrt{\nu^2 + \sigma_u^2}$, and

$$\gamma > 2w_{n_l(N)}^{-1} \rho \sigma_y \kappa. \quad (21)$$

Then, for any given $\beta \in (0, 1)$ there exists a finite integer N_β such that for any $N \geq N_\beta$ it holds that

$$\text{supp}(x_1^*) \subseteq \{1, \dots, n_l(N)\}$$

with probability no smaller than $1 - \beta$, where x_1^* is the minimizer of (14), and $n_l(N)$ is the leading order of the unknown system's impulse response. ■

See Appendix A.5 for a proof of Corollary 5. Corollary 5 states that, under suitable conditions, an estimate of the impulse response based on the minimization of (14) is supported inside the leading support of the system, with high probability. The following corollary provides a similar result, for the case in which the true system is a-priori known to have finite impulse response (FIR).

Corollary 6 (FIR recovery) *Let Assumption 1 hold. Further, assume the “true,” unknown, system is FIR of order $n \leq q$, with n unknown. Then, for any $\gamma > 0$ and for any given $\beta \in (0, 1)$ there exists a finite integer N_β such that for any $N \geq N_\beta$ it holds that $\text{supp}(x_1^*) \subseteq \{1, \dots, n\}$ with probability no smaller than $1 - \beta$, where x_1^* is the minimizer of (14). ■*

See Appendix A.6 for a proof of Corollary 6. The key point of this corollary is that if the true system is known to be FIR, then the minimizer of (14) will tendentially recover the true order of the system, regardless of the value of $\gamma > 0$ (but, of course, the larger the value of γ , the sooner w.r.t. N the condition (15) will be satisfied).

5 Identification procedure

We next formalize a possible procedure illustrating how the proposed methodology can be used in a practical experimental setting. Suppose that a set of data $\{y(k), u(k)\}_{k=3-q}^N$ is available from a process of the form (3). Identification of the impulse response $h(i)$ is performed by minimizing the cost function (14). This operation requires the choice of two parameters (γ and σ_u). If σ_u and σ_y are known from some a-priori information on the noises affecting the system or can be reliably estimated, then γ can be chosen according to (21), where ρ can be estimated by means of the technique in Milanese et al. [2010] (see Section 6.1). If instead this information is not available, a systematic procedure for the choice of γ and σ_u is the following one:

- Take “reasonable” sets $\Gamma = \{\gamma^{(1)}, \gamma^{(2)}, \dots\}$ and $\Sigma_u = \{\sigma_u^{(1)}, \sigma_u^{(2)}, \dots\}$ for γ and σ_u values, respectively. If σ_u is known from some a-priori information on the noise affecting the input, then $\Sigma_u = \sigma_u$.
- Define y , U and T as shown in Section 3.
- Run the following algorithm:


```

      for  $i = 1 : \text{length}(\Sigma_u)$ 
        for  $j = 1 : \text{length}(\Gamma)$ 
           $\sigma_u = \sigma_u^{(i)}$ ;  $\gamma = \gamma^{(j)}$ ;
           $x^*(i, j) = \arg \min_x J_1(x)$ ;
           $E(i, j) = \|y - Ux^*(i, j)\|_2^2$ ;
           $C(i, j) = \|x^*(i, j)\|_0$ ;
        end
        plot( $C(i, :)$ ,  $E(i, :)$ )
      end
```

- The obtained plot shows how the model accuracy (measured by E) changes in function of its complexity (measured by C). Thus, γ and σ_u can be chosen according to the desired trade-off between model accuracy and complexity.

Choosing $\gamma^{(1)} > \gamma^{(2)} > \dots$ and using $x^*(i, j - 1)$ at the j th step as the initial condition for the optimization problem may significantly increase the speed of the algorithm. An example of application of this procedure is shown in Section 6.2 and, in particular, in Figure 3.

The weighting matrix W plays a relevant role in the model order selection, increasing the algorithm efficiency especially in situations where a low number of data is available. For simplicity, unitary weights w_i were here adopted in Section 6. Further research activity will be devoted to investigate how to automatically and optimally select these weights, in order to take into account possible priors on the unknown system.

6 Numerical examples

6.1 A simulated LTI system

For our first numerical test we considered a classical discrete-time LTI system proposed in Ljung [1999b]. This system is defined by the discrete-time transfer function

$$H(z) = \frac{z^3 + 0.5z^2}{z^4 - 2.2z^3 + 2.42z^2 - 1.87z + 0.7225}, \quad (22)$$

with sampling time 1 s. We assume that all necessary parameters (e.g., the noise variances and the impulse response's stability degree bounds) are known or have been estimated in advance by other means.

6.1.1 Experiments with a fixed number of data

Three i.i.d. input sequences with zero mean and variance $\nu^2 = 1$ were first generated. Each of these sequences was corrupted by an i.i.d. noise with zero mean and variance σ_u^2 , with $\sigma_u = 0.01$ for the first sequence, $\sigma_u = 0.03$ for the second one, and $\sigma_u = 0.05$ for the third one. These values correspond to noise-to-signal standard deviation ratios of 1%, 3%, and 5%, respectively.

For each noise-corrupted input sequence, the system (22) was simulated for 2000 s, assuming zero initial conditions. Note that the system reaches steady-state conditions after about 150 s. The resulting output sequence was corrupted by an i.i.d. noise with zero mean and variance σ_y^2 , with $\sigma_y = 0.1$ for the first sequence, $\sigma_y = 0.3$ for the second one, and $\sigma_y = 0.5$ for the third one. These values correspond to noise-to-signal standard deviation ratios of 1%, 3%, and 5%, respectively (the system static gain is about 10). Then, the last $N = 1000$ noise-corrupted output values were acquired. From these data, the following models of the unknown system impulse response were identified:

- Leading Response Recovery (LRR) model. This model was obtained minimizing the objective function (14). The parameters required for this minimization were taken as follows. The variances σ_u^2 and σ_y^2 were assumed known (or accurately estimated). The impulse

response bound parameters were estimated by means of the technique in Milanese et al. [2010], giving values $L = 6$ and $\rho = 0.93$ (note that only ρ is required by the LRR algorithm). Unitary weights w_i were adopted. The estimated length was taken as $q = 500$. The value of γ was chosen according to (21).

- Least Squares (LS) model. This model was identified using standard least squares, that is, by minimizing the objective function (14), with $\sigma_u^2 = 0$ and $T^{-1} = 0$.
- Tikhonov regularized Least Squares (TLS) model. This model was identified by minimizing the objective function (14), with $T^{-1} = 0$.

To validate the identified models, the following indices were computed:

- Best fit criterion:

$$\text{FIT} \doteq 100 \left(1 - \frac{\|y - \hat{y}\|_2}{\|y - \text{mean}(y)\|_2} \right)$$

where y is the measured system output vector and \hat{y} is the output vector simulated by the model. The FIT index was evaluated on $N_v = 2000$ validation data points (i.e., points not previously used for identification). Obviously, this index measures the model simulation accuracy: the closer it is to 100%, the more accurate the simulation is.

- Tail ℓ_0 quasi-norm:

$$\text{TN0} \doteq \|x_{\text{tail}}^*\|_0$$

where $x_{\text{tail}}^* \doteq [x^*(n_l+1) \cdots x^*(500)]^\top$ and x^* is the estimated model impulse response. This index is a measure of the model tail (the tail can be defined as the vector formed by the impulse response components with index $> n_l$). More precisely, it counts how many elements in the tail of the model impulse response are different from zero. Note that, for $\sigma_y = 0.1$, $n_l = 105$; for $\sigma_y = 0.3$, $n_l = 89$; for $\sigma_y = 0.5$, $n_l = 82$.

- Tail ℓ_1 norm:

$$\text{TN1} \doteq \|x_{\text{tail}}^*\|_1.$$

This index provides an indication on the average magnitude of the elements in the tail of the model impulse response.

A Monte Carlo simulation was then carried out, where the above identification-validation procedure was repeated for 100 trials. The averages $\overline{\text{FIT}}$, $\overline{\text{TN0}}$ and $\overline{\text{TN1}}$ of FIT, TN0 and TN1 obtained in this simulation are reported in Table 1. We observe that the three identification methods lead to very similar FIT values. However, the LRR models have a tail that is practically null (in average, about 4 non-null elements over about 280), even though the number of data used for identification is relatively low (1000 data). This fact shows that our identification algorithm is able to provide highly sparse models, without compromising their simulation accuracy. An even more important aspect is that sparsification does not occur for “random” indexes of the model impulse response but for large indexes, i.e., those indexes associated with the exponentially decaying tail of the impulse response.

It is important to remark that the LRR algorithm does not use the prior information in terms of L and ρ values

noise	model	$\overline{\text{FIT}}$	$\overline{\text{TN0}}$	$\overline{\text{TN1}}$
1%	LRR	98.6	6.0	0.012
	LS	98.6	315	1.40
	TLS	98.6	315	1.39
3%	LRR	95.9	4	0.019
	LS	96.0	267	3.29
	TLS	96.0	267	3.28
5%	LRR	93.3	3.3	0.025
	LS	93.4	246	4.97
	TLS	93.4	246	4.94

Table 1

Average indices obtained in the Monte Carlo simulation.

to impose strict constraints or weights on the samples of the leading response. The information on L and ρ is only used in the proof of Theorem 4 (see (A.7)) to derive a bound on the value of γ (see (21)).

It may be expected that using explicit constraints or weights based on L and ρ in the algorithm may lead to improvements in terms of model accuracy and/or complexity. To better investigate this aspect, we performed another Monte Carlo simulation, considering a 3% noise level, and applying standard constrained least squares and regularized Diagonal/Correlated kernel methods (the latter using the Matlab routine `impulseest.m`, see, e.g., Pilonetto et al. [2014]). Indeed, these methods use the L and ρ information (either known a priori or estimated from the data) to impose a desired exponential decay of the overall impulse response. The following index values were obtained with constrained least squares (CLS): $\overline{\text{FIT}} = 96.6$, $\overline{\text{TN0}} = 267$, $\overline{\text{TN1}} = 0.115$. The following index values were obtained with the regularized Diagonal/Correlated kernel method (DCK): $\overline{\text{FIT}} = 96.7$, $\overline{\text{TN0}} = 267$, $\overline{\text{TN1}} = 0.033$.

We can compare these results with those shown in Table 1. It can be noted that the CLS and DCK methods give slight improvements w.r.t. the other methods in terms of the FIT criterion, although the formers use a significantly stronger prior information. An interesting result of the CLS and DCK methods is that they lead to tails with very small (albeit nonzero) elements, giving a relevant reductions of the tail magnitude w.r.t. the LS and TLS methods, with $\overline{\text{TN1}}$ indexes not far from the one given by the LRR method. Nevertheless, the $\overline{\text{TN0}}$ values given by the LRR method are by far the lowest ones, showing that this method is the only one (among those considered) allowing effective and unsupervised model order selection.

6.1.2 Experiments with an increasing number of data

A “long” i.i.d. input sequence with zero mean and variance $\nu^2 = 1$ was generated and corrupted by an i.i.d. noise with zero mean and variance $\sigma_u^2 = 0.03^2$. The true system was then simulated using this input sequence, and the resulting output sequence was corrupted by an

i.i.d. noise with zero mean and variance $\sigma_y^2 = 0.3^2$. The data corresponding to the output values with time index $k = 1001, \dots, 1000 + N$ were selected, where $N = 500, \dots, 50000$. For each value of N , an LRR, an LS and a TLS model were identified from these data. The values of FIT, TN0 and TN1 obtained for these models are plotted as function of N in Figures 1 and 2. We can observe that the three identification methods lead to very similar FIT values, the LRR models giving slightly better results for low number of data. A key difference between the three techniques is that the LRR method is able to select the more appropriate impulse response components (i.e., the components with index in the interval $[1, n_l]$), forcing the others to vanish. After a certain value of N (about 32000), the tail of the LRR models is zero, confirming the theoretical result given in Corollary 5. Such an effective component selection is not guaranteed by the other two methods which, on the contrary, have tails with support cardinality (measured by the ℓ_0 quasi-norm) that grows with N .

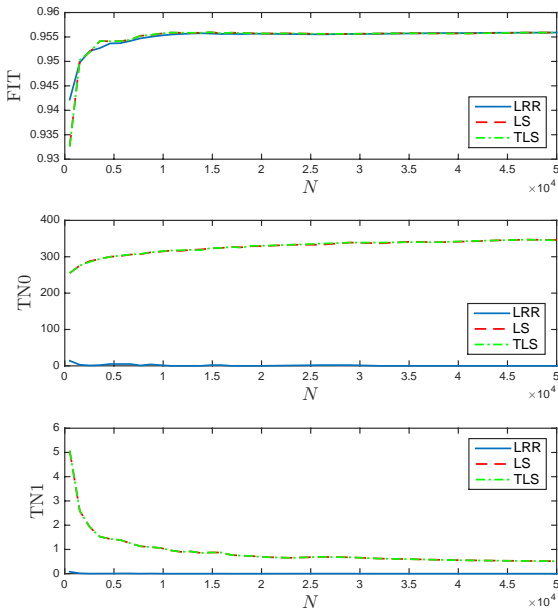


Fig. 1. Values of FIT, TN0 and TN1 for all models.

6.2 Experimental data from a flexible robot arm

The identification of poorly damped systems from experimental data is among the most challenging issues in many practical applications. For this reason, as second test we considered a system with a vibrating flexible robot arm described in Torfs et al. [1998], adopted as case study in various software packages (Kollár et al. [1994], Kollár [1994], National Instruments Corporation [2004-2006]). Data records from this process have been also analyzed in Pintelon and Schoukens [2012], Pilonetto et al. [2014]. The input is the driving torque and the output is the tangential acceleration of the tip of the robot arm.

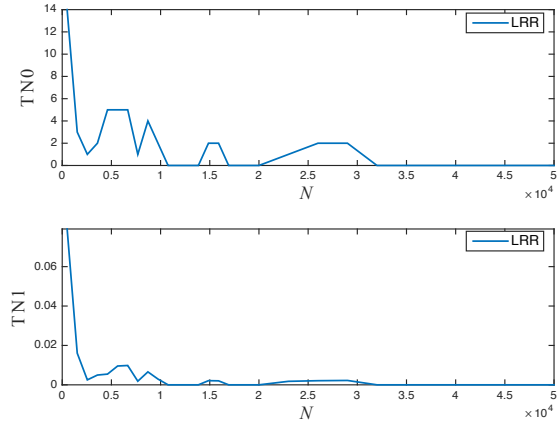


Fig. 2. Values of TN0 and TN1 for the LRR models.

Ten consecutive periods of the response to a multisine excitation signal were collected at a sampling frequency of 500 Hz, for a total of 40960 data points.

We have built models using different techniques: the Leading Response Recovery (LRR) and the regularized Diagonal/Correlated kernel (DCK) methods to obtain high-order FIR models, the standard Prediction Error Method (PEM) to estimate low-order state space models. Since the true system is unknown, the models cannot be evaluated by their fit to the actual system. Instead, we used the hold-out validation technique and measured how well the identified models can reproduce the output on validation portions of the data that were not used for estimation. We chose the estimation data to be the portion 1:7000 and the validation data to be the portion 10000:40960.

To identify the LRR models, the procedure described in Section 5 has been applied to suitably choose the values of σ_u and γ , taking $q = 7000$ as initial estimate length. No a-priori information was available about the input noise affecting the system, driven by an input signal u with sample variance $\nu^2 = 0.0298$. For this reason, three scenarios have been considered, with $\sigma_u = 0$, $\sigma_u = 0.02$ and $\sigma_u = 0.04$. These values correspond to noise-to-signal ratios of 0% (i.e., no-input-noise situation), 1.3%, and 5.3%, respectively. Then the LRR algorithm has run with values of γ in the range $[0.01, 1]$, using the MATLAB’s command `lasso` with optional input arguments `'RelTol', 4e-4, 'Standardize', false`. The results in terms of fitting error $\|y - Ux^*\|_2^2$ and complexity $\|x^*\|_0$ are shown in Figure 3 (lower error and higher complexity are achieved for lower values of γ). As expected, curves with lower σ_u dominate curves with higher σ_u (for any given γ , the solution obtained with lower σ_u has both lower error and complexity with respect to a solution obtained with higher σ_u). However, the choice of the actual curve to use depends on our confidence on the true value of σ_u , and underestimating this value may lead to worse-than-expected performance on validation data. Also, curves with higher σ_u show a flatter behavior after the “knee” for lower γ values.

We found a reasonable tradeoff for $\gamma = 0.2$, allowing a satisfactory fitting error (around 19.8 for $\sigma_u = 0$, which raises up to 21.1 in the worst-case $\sigma_u = 0.04$) with a small complexity (around 560, that raises up to 1220 for $\sigma_u = 0.04$). Alternatively, $\gamma = 0.1$ allows a lower error (around 14.8 for $\sigma_u = 0$, which raises up to 16.1 for $\sigma_u = 0.04$) with a still acceptable complexity (around 830, that raises up to 1740 for $\sigma_u = 0.04$).

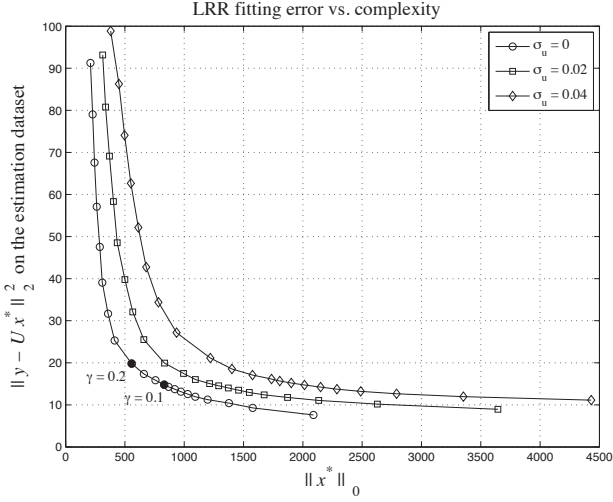


Fig. 3. LRR fitting error $\|y - Ux^*\|_2^2$ vs. complexity $\|x^*\|_0$.

To fairly compare the performances achieved by LRR and DCK methods, FIR models with the same given order have been identified. First, the LRR FIR model of order 2500 was identified with FIT value of 80.1%. Then, the DCK FIR model of order 2500 was estimated using regularized least squares, tuned by the marginalized likelihood method and with the unknown input data set to zero, via the MATLAB’s command `impulsest(data, 2500, 0, opt)` with option `opt` set as `opt.RegulKernel='dc'`; `opt.Advanced.AROrder=0`. The FIT for this DCK FIR model was 79.9%. For illustration, the FIT values of the LRR and DCK FIR models are shown in Figure 4 as horizontal lines. For comparison, we estimated n -th order state space PEM models without disturbance model for $n = 1, \dots, 30$ (via the MATLAB’s command `pem(data, n, 'dist', 'no')`) and calculated their FIT index to validation data. These FIT values are shown as function of n in Figure 4. The two best FITs were 78.9% and 78.6%, obtained for order $n = 21$ and $n = 18$, respectively, while the 5-th order model with FIT value of 69.6% could be a reasonable tradeoff between accuracy and complexity. In any case, any PEM fit is worst than those provided by the LRR and DCK FIR models.

One may observe that FIR models of order 2500 are quite large, but it is interesting to note that they can be easily reduced to low-order state space models by model order reduction methods, like balanced truncation, Hankel norm minimization and \mathcal{L}_2 reduction. For example, we applied the square root balanced truncation

method to the LRR FIR model, to obtain reduced state space models of order $n = 1, \dots, 30$ (via the MATLAB’s command `balancmr`), and we computed their FIT index on validation data. These FIT values are also shown as function of n in Figure 4. It can be observed that a reduced state-space model of order $n = 6$ provides a FIT of 74.2%, which is better than any PEM-estimated state space model of order $n = 1, \dots, 14$.

To discriminate the effects of the transient due to the mismatch between the initial states of the actual system and the identified models, the FIT index has been also computed by neglecting the initial 3000 samples of the validation data. The FIT values of the LRR and DCK FIR models of order 2500 raise to 83.4% and 83.6%, respectively; the FIT values of the PEM models of order $n = 5, 18, 21$ go to 71.2%, 83.2%, 83.7%, respectively; the FIT value of the reduced state-space model of order $n = 6$ increases to 76.2%. All these results are shown in Figure 5.

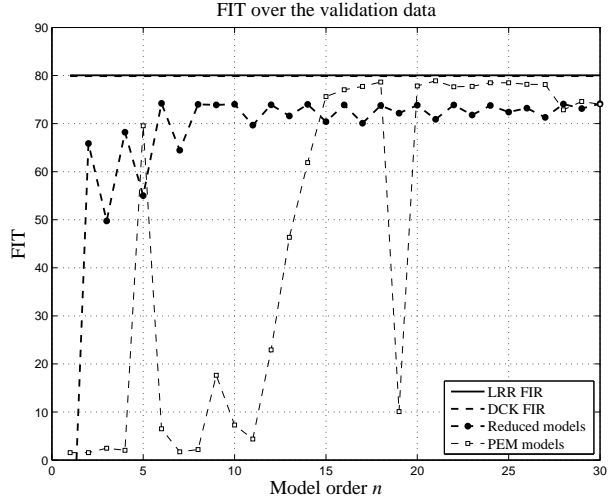


Fig. 4. Values of FIT for all models.

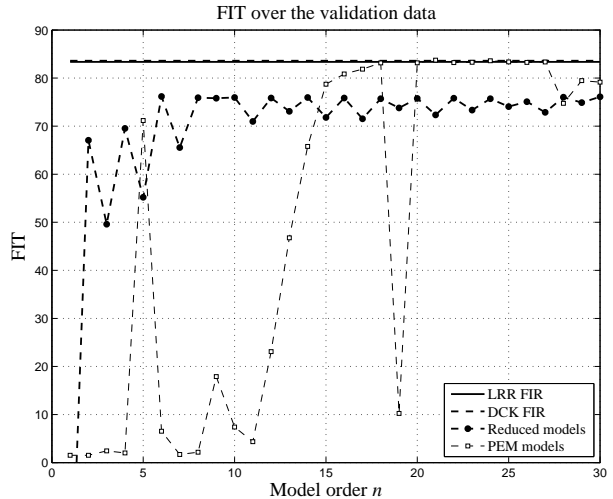


Fig. 5. Values of FIT, neglecting the starting 3000 samples.

It is worth to observe that, even if the fit performances of LRR and DCK FIR models are very close, the computational complexity of their corresponding algorithms is dramatically different. Referring to a workstation equipped with an Intel(R) Core(TM) i7-3770 CPU @ 3.40 GHz and with 16 GB of RAM, the overall CPU time used to estimate the LRR FIR model of order 2500 was around 30 seconds, while the computation of DCK FIR model of order 2500 required around 8700 seconds, i.e., 290 times more. For comparison, FIR models of order 3000 and 3500 were also identified using the same 7000 estimation data as before: the CPU times required by the LRR FIR models were around 30 and 36 seconds, respectively, while the CPU times required by the DCK FIR models were around 11000 and 23200 seconds, respectively. Widening the estimation data to the first 10000 samples, the identification of LRR FIR models of orders up to 5000 required no more than 75 seconds, thus showing that the approach proposed in this paper scales nicely with the problem dimensionality.

7 Conclusions

A novel method for the identification of low-complexity FIR models from experimental data is presented in this paper. The method is based on an Elastic Net criterion, which considers an identification cost defined as a weighted combination of a standard prediction error term, an ℓ_2 regularization term, and a weighted ℓ_1 penalty term. The main novelty of the method with respect to the state of the art is that it allows for an effective selection of the model order, while requiring only stability and standard statistical assumptions on the noises affecting the system; no additional information on the system impulse response behavior is needed. The effectiveness of the method has been tested through both extensive numerical simulations (considering two typical situations: one with a fixed number of data, and one with an arbitrarily large number of data) and real experimental data from a lightly damped mechanical system. In all situations, the method showed high numerical efficiency and satisfactory order selection capability and simulation accuracy.

Research activity is being devoted to developing a weighted version of the method proposed here. It is indeed expected that including suitable weights in the identification criterion may make the model order selection even more efficient, especially in situations where a low number of data is available.

A Appendix

A.1 Proof of Lemma 1

From the hypothesis that $y_i \rightsquigarrow \bar{y}_i$, $i = 1, \dots, p$, applying the definition of symbol \rightsquigarrow , we have that for any $\tilde{\epsilon} > 0$ and $\tilde{\beta} \in (0, 1)$ there exists an integer $\tilde{N}_{\tilde{\epsilon}, \tilde{\beta}}$ such that

$$\mathbb{P}\{|y_i - \bar{y}_i| \leq \tilde{\epsilon}\} \geq 1 - \tilde{\beta}, \quad \forall N \geq \tilde{N}_{\tilde{\epsilon}, \tilde{\beta}}.$$

From Bonferroni's inequality we further have that the probability of the joint event $\{|y_i - \bar{y}_i| \leq \tilde{\epsilon}, i = 1, \dots, p\}$ is lower bounded as

$$\mathbb{P}\{|y_i - \bar{y}_i| \leq \tilde{\epsilon}, i = 1, \dots, p\} \geq 1 - p\tilde{\beta}, \quad \forall N \geq \tilde{N}_{\tilde{\epsilon}, \tilde{\beta}}.$$

Since $\mathbb{P}\{|y_i - \bar{y}_i| \leq \tilde{\epsilon}, i = 1, \dots, p\} = \mathbb{P}\{\|y - \bar{y}\|_\infty \leq \tilde{\epsilon}\}$, letting $\beta \doteq p\tilde{\beta}$, we write

$$\mathbb{P}\{\|y - \bar{y}\|_\infty \leq \tilde{\epsilon}\} \geq 1 - \beta, \quad \forall N \geq \tilde{N}_{\tilde{\epsilon}, \beta/p}.$$

Now, from the hypothesis that f is Lipschitz continuous, it follows that there exists a finite constant $C \geq 0$ such that

$$|f(y) - f(\bar{y})| \leq C\|y - \bar{y}\|_\infty.$$

Therefore, $\|y - \bar{y}\|_\infty \leq \tilde{\epsilon}$ implies that $|f(y) - f(\bar{y})| \leq \epsilon$, for $\epsilon \doteq C\tilde{\epsilon}$, whence

$$\mathbb{P}\{|f(y) - f(\bar{y})| \leq \epsilon\} \geq 1 - \beta, \quad \forall N \geq \tilde{N}_{\epsilon/C, \beta/p},$$

which proves that $f(y) \rightsquigarrow f(\bar{y})$. ■

A.2 Proof of Lemma 2

The claim is a direct consequence of the first point of Theorem 8 in Tropp [2006], where the index set Λ is $\{1, \dots, n\}$, and $\text{ERC}(\Lambda)$ in Tropp [2006] coincides with $\Upsilon_n(A)$. The symbol a_Λ used in Theorem 8 of Tropp [2006] corresponds to $P_n b$, that is the best ℓ_2 approximation of b using a linear combination of the first n columns of A . These first n columns have unit ℓ_2 norm and are indeed linearly independent, as requested by the hypotheses of Theorem 8 in Tropp [2006], due to the specific structure of $A = \bar{A}T$, where \bar{A} , shown in (10), has a multiple of the identity matrix I_q as a bottom block. ■

A.3 Proof of Lemma 3

Some parts of this result might possibly be derived as a particular case of Theorem 2.3 in Ljung [1999a]; we here report a full proof for the specific case of interest in the present work. For $i = 1, 2, \dots$, let us define $u_i = [u(2-i) u(3-i) \dots u(N+1-i)]^\top \in \mathbb{R}^N$. Then, for all i and j , we have that

$$\begin{aligned} \frac{1}{N} u_i^\top u_j &= \frac{1}{N} [u(2-i) u(3-i) \dots u(N+1-i)] \cdot \\ &\quad [u(2-j) u(3-j) \dots u(N+1-j)]^\top \\ &= \frac{1}{N} \sum_{k=2}^{N+1} u(k-i) u(k-j) \\ &= \frac{1}{N} \sum_{k=1}^N u(k+1-i) u(k+1-j). \end{aligned}$$

Consider first the case where $i = j$. Then

$$\frac{1}{N} u_i^\top u_i = \frac{1}{N} \sum_{k=1}^N u(k+1-i)^2$$

is the empirical mean of the elements of the sequence of length N of random variables $x_k = u(k+1-i)^2$, $k = 1, \dots, N$, such that, for all k , i and $l \neq k$:

$$\mathbb{E}\{x_k\} = \mathbb{E}\{u(k+1-i)^2\} = \text{var}\{u(k+1-i)\} = \nu^2 < \infty$$

$$\begin{aligned} \text{var}\{x_k\} &= \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})^2\} = \mathbb{E}\{(u(k+1-i)^2 - \nu^2)^2\} \\ &= \mathbb{E}\{u(k+1-i)^4 - 2\nu^2 u(k+1-i)^2 + \nu^4\} \\ &= \mathbb{E}\{u(k+1-i)^4\} - 2\nu^2 \mathbb{E}\{u(k+1-i)^2\} + \nu^4 \\ &= \overline{m_4} - \nu^4 < \infty \end{aligned}$$

$$\begin{aligned} &\mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} \\ &= \mathbb{E}\{(u(k+1-i)^2 - \nu^2)(u(l+1-i)^2 - \nu^2)\} \\ &= \mathbb{E}\{u(k+1-i)^2 u(l+1-i)^2 + \\ &\quad - \nu^2 [u(k+1-i)^2 + u(l+1-i)^2] + \nu^4\} \\ &= \mathbb{E}\{u(k+1-i)^2 u(l+1-i)^2\} + \\ &\quad - \nu^2 [\mathbb{E}\{u(k+1-i)^2\} + \mathbb{E}\{u(l+1-i)^2\}] + \nu^4 \\ &= \mathbb{E}\{u(k+1-i)^2\} \mathbb{E}\{u(l+1-i)^2\} - \nu^4 = \nu^2 \nu^2 - \nu^4 = 0 \end{aligned}$$

where the last derivation follows from the fact that x_k and x_l are mutually independent since the input $u(k)$ is an i.i.d. sequence. By applying the Chebyshev's inequality for sums of uncorrelated variables shown in Section 2.2, it holds that

$$\frac{1}{N} u_i^\top u_j \sim \mathbb{E}\{x_k\} = \nu^2, \quad \forall i.$$

Consider next the case where $i \neq j$. Then

$$\frac{1}{N} u_i^\top u_j = \frac{1}{N} \sum_{k=1}^N u(k+1-i)u(k+1-j)$$

is the empirical mean of the elements of the sequence of length N of random variables $x_k = u(k+1-i)u(k+1-j)$, $k = 1, \dots, N$, such that, for all $k, i, j = i + \tilde{i} \neq i$ and $l = k + \tilde{k} \neq k$, with $\tilde{i} \neq 0$ and $\tilde{k} \neq 0$:

$$\begin{aligned} \mathbb{E}\{x_k\} &= \mathbb{E}\{u(k+1-i)u(k+1-j)\} \\ &= \mathbb{E}\{u(k+1-i)\} \mathbb{E}\{u(k+1-j)\} = 0 \end{aligned}$$

$$\begin{aligned} \text{var}\{x_k\} &= \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})^2\} = \mathbb{E}\{(u(k+1-i)u(k+1-j))^2\} \\ &= \mathbb{E}\{u(k+1-i)^2 u(k+1-j)^2\} \\ &= \mathbb{E}\{u(k+1-i)^2\} \mathbb{E}\{u(k+1-j)^2\} = \nu^4 < \infty \end{aligned}$$

$$\begin{aligned} \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} &= \\ &= \mathbb{E}\{(u(k+1-i)u(k+1-j))(u(l+1-i)u(l+1-j))\} \\ &= \mathbb{E}\{u(k+1-i)u(k+1-i-\tilde{i})u(k+\tilde{k}+1-i)u(k+\tilde{k}+1-i-\tilde{i})\} \\ &= \mathbb{E}\{u(k+1-i)u(k+1-i-\tilde{i})u(k+1-i+\tilde{k})u(k+1-i-\tilde{i}+\tilde{k})\} \end{aligned}$$

if $\tilde{i} = \tilde{k}$, then:

$$\begin{aligned} \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} &= \\ &= \mathbb{E}\{u(k+1-i)^2 u(k+1-i-\tilde{i})u(k+1-i+\tilde{i})\} \\ &= \mathbb{E}\{u(k+1-i)^2\} \mathbb{E}\{u(k+1-i-\tilde{i})\} \mathbb{E}\{u(k+1-i+\tilde{i})\} = 0 \end{aligned}$$

otherwise, if $\tilde{i} \neq \tilde{k}$, then:

$$\begin{aligned} \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} &= \\ &= \mathbb{E}\{u(k+1-i)u(k+1-i-\tilde{i})u(k+1-i+\tilde{k})u(k+1-i-\tilde{i}+\tilde{k})\} \\ &= \mathbb{E}\{u(k+1-i)\}. \end{aligned}$$

$$\mathbb{E}\{u(k+1-i-\tilde{i})u(k+1-i+\tilde{k})u(k+1-i-\tilde{i}+\tilde{k})\} = 0$$

and this means that x_k and x_l are uncorrelated for all $k \neq l$. By applying the Chebyshev's inequality for sums of uncorrelated variables shown in Section 2.2, we obtain that

$$\frac{1}{N} u_i^\top u_j \sim \mathbb{E}\{x_k\} = 0, \text{ for all } i \neq j,$$

which proves (16).

We next prove (17). Since $\delta_y = [\delta_y(1) \delta_y(2) \dots \delta_y(N)]^\top \in \mathbb{R}^N$, then for $i = 1, 2, \dots$:

$$\begin{aligned} \frac{1}{N} u_i^\top \delta_y &= \frac{1}{N} [u(2-i) u(3-i) \dots u(N+1-i)] \cdot \\ &\quad [\delta_y(1) \delta_y(2) \dots \delta_y(N)]^\top \\ &= \frac{1}{N} \sum_{k=1}^N u(k+1-i) \delta_y(k) \end{aligned}$$

is the empirical mean of the elements of the sequence of length N of random variables $x_k = u(k+1-i)\delta_y(k)$, $k = 1, \dots, N$, such that, for all k, i and $l \neq k$:

$$\begin{aligned} \mathbb{E}\{x_k\} &= \mathbb{E}\{u(k+1-i)\delta_y(k)\} = \mathbb{E}\{u(k+1-i)\} \mathbb{E}\{\delta_y(k)\} = 0 \\ \text{var}\{x_k\} &= \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})^2\} \\ &= \mathbb{E}\{(u(k+1-i)\delta_y(k))^2\} = \mathbb{E}\{u(k+1-i)^2 \delta_y(k)^2\} \\ &= \mathbb{E}\{u(k+1-i)^2\} \mathbb{E}\{\delta_y(k)^2\} = \nu^2 \sigma_y^2 < \infty \end{aligned}$$

$$\begin{aligned} \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} &= \\ &= \mathbb{E}\{(u(k+1-i)\delta_y(k))(u(l+1-i)\delta_y(l))\} \\ &= \mathbb{E}\{u(k+1-i)u(l+1-i)\delta_y(k)\delta_y(l)\} \\ &= \mathbb{E}\{u(k+1-i)\} \mathbb{E}\{u(l+1-i)\} \mathbb{E}\{\delta_y(k)\} \mathbb{E}\{\delta_y(l)\} = 0. \end{aligned}$$

By applying the Chebyshev's inequality for sums of uncorrelated variables shown in Section 2.2, it thus holds that

$$\frac{1}{N} u_i^\top \delta_y \sim \mathbb{E}\{x_k\} = 0, \text{ for all } i.$$

Finally, we prove (18). For $j = 1, 2, \dots$, let us define $\delta_j = [\delta_u(2-j) \delta_u(3-j) \dots \delta_u(N+1-j)]^\top \in \mathbb{R}^N$. Then, $\forall i, j$:

$$\begin{aligned} \frac{1}{N} u_i^\top \delta_j &= \frac{1}{N} [u(2-i) u(3-i) \dots u(N+1-i)] \cdot \\ &\quad [\delta_u(2-j) \delta_u(3-j) \dots \delta_u(N+1-j)]^\top \end{aligned}$$

$$\begin{aligned} &= \frac{1}{N} \sum_{k=2}^{N+1} u(k-i) \delta_u(k-j) \\ &= \frac{1}{N} \sum_{k=1}^N u(k+1-i) \delta_u(k+1-j) \end{aligned}$$

is the empirical mean of the elements of the sequence of length N of random variables $x_k = u(k+1-i)\delta_u(k+1-j)$, $k = 1, \dots, N$, such that, for all k, i, j and $l \neq k$:

$$\begin{aligned} \mathbb{E}\{x_k\} &= \mathbb{E}\{u(k+1-i)\delta_u(k+1-j)\} \\ &= \mathbb{E}\{u(k+1-i)\} \mathbb{E}\{\delta_u(k+1-j)\} = 0 \end{aligned}$$

$$\begin{aligned} \text{var}\{x_k\} &= \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})^2\} = \mathbb{E}\{(u(k+1-i)\delta_u(k+1-j))^2\} \\ &= \mathbb{E}\{u(k+1-i)^2 \delta_u(k+1-j)^2\} \\ &= \mathbb{E}\{u(k+1-i)^2\} \mathbb{E}\{\delta_u(k+1-j)^2\} = \nu^2 \sigma_u^2 < \infty \end{aligned}$$

$$\begin{aligned} \mathbb{E}\{(x_k - \mathbb{E}\{x_k\})(x_l - \mathbb{E}\{x_l\})\} &= \\ &= \mathbb{E}\{(u(k+1-i)\delta_u(k+1-j))(u(l+1-i)\delta_u(l+1-j))\} \\ &= \mathbb{E}\{u(k+1-i)u(l+1-i)\delta_u(k+1-j)\delta_u(l+1-j)\} \\ &= \mathbb{E}\{u(k+1-i)\} \mathbb{E}\{u(l+1-i)\} \cdot \\ &\quad \mathbb{E}\{\delta_u(k+1-j)\} \mathbb{E}\{\delta_u(l+1-j)\} = 0. \end{aligned}$$

By applying the Chebyshev's inequality for sums of uncorrelated variables shown in Section 2.2, it holds that

$$\frac{1}{N} u_i^\top \delta_j \sim \mathbb{E}\{x_k\} = 0, \text{ for all } i \text{ and } j. \quad \blacksquare$$

A.4 Proof of Theorem 4

A.4.1 Preliminaries

For any integer $n \leq q$, let $h^n \doteq [h(1) \dots h(n) 0 \dots 0]^\top \in \mathbb{R}^q$ denote the n -leading truncation of $h_{\uparrow q} \in \mathbb{R}^q$, let

$$h_{\downarrow n} = [h(n+1) h(n+2) \dots]^\top,$$

and let, for $i = 1, 2, \dots$,

$$\tilde{u}_i \doteq \begin{bmatrix} \tilde{u}(2-i) \\ \tilde{u}(3-i) \\ \vdots \\ \tilde{u}(N+1-i) \end{bmatrix}, u_i \doteq \begin{bmatrix} u(2-i) \\ u(3-i) \\ \vdots \\ u(N+1-i) \end{bmatrix}, \delta_i \doteq \begin{bmatrix} \delta_u(2-i) \\ \delta_u(3-i) \\ \vdots \\ \delta_u(N+1-i) \end{bmatrix}.$$

For any integer $k \geq 1$, let $\Delta_{\uparrow k} \doteq [\delta_1 \dots \delta_k]$, $\Delta_{\downarrow k} \doteq [\delta_{k+1} \dots]$ and define $\Delta_{\downarrow 0} = [\delta_1 \delta_2 \dots] \in \mathbb{R}^{N, \infty}$. Considering the expression in (4), and splitting the summation at n , we can write

$$y = \tilde{U}_{\uparrow n} h_{\uparrow n} + \delta_y + \tilde{U}_{\downarrow n} h_{\downarrow n}.$$

Further, using (6), we have

$$y = U_{\uparrow n} h_{\uparrow n} + (\Delta_{\uparrow n} h_{\uparrow n} + \delta_y + U_{\downarrow n} h_{\downarrow n} + \Delta_{\downarrow n} h_{\downarrow n}).$$

Since $U_{\uparrow n} h_{\uparrow n} = U_{\uparrow q} h^n \doteq U h^n$, we can write

$$y = U h^n + e_0,$$

where

$$e_0 \doteq U_{\downarrow n} h_{\downarrow n} + \Delta_{\downarrow 0} h + \delta_y,$$

being $h \doteq [h(1) h(2) \dots]^\top$. Then, using the notation in (10), we have that

$$b = \bar{A} h^n + e,$$

where

$$e \doteq \begin{bmatrix} e_0 \\ -\sigma_u \sqrt{N} h^n \end{bmatrix}, \quad (\text{A.1})$$

and, by the change of variable $\tilde{h}^n = T^{-1} h^n$,

$$b = A \tilde{h}^n + e.$$

Since $A = \bar{A} T$, where T is diagonal, we can write

$$A_{\uparrow n} = \bar{A}_{\uparrow n} T_{\#n},$$

where $T_{\#n}$ is the $n \times n$ principal submatrix of T . Therefore,

$$\begin{aligned} A_{\uparrow n}^\dagger &= (A_{\uparrow n}^\top A_{\uparrow n})^{-1} A_{\uparrow n}^\top = T_{\#n}^{-1} \bar{A}_{\uparrow n}^\dagger \\ &= T_{\#n}^{-1} \left(U_{\uparrow n}^\top U_{\uparrow n} + N \sigma_u^2 I_n \right)^{-1} \left[U_{\uparrow n}^\top \sigma_u \sqrt{N} I_{\uparrow n}^\top \right], \end{aligned}$$

where $I_{\uparrow n}$ is the submatrix formed by the first n columns of the identity matrix I_q .

The orthogonal projector P_n onto the span of $A_{\uparrow n}$ is given by $P_n = A_{\uparrow n} A_{\uparrow n}^\dagger = \bar{A}_{\uparrow n} \bar{A}_{\uparrow n}^\dagger$

$$= \begin{bmatrix} U_{\uparrow n} \\ \sigma_u \sqrt{N} I_{\uparrow n} \end{bmatrix} \left(U_{\uparrow n}^\top U_{\uparrow n} + N \sigma_u^2 I_n \right)^{-1} \left[U_{\uparrow n}^\top \sigma_u \sqrt{N} I_{\uparrow n}^\top \right].$$

For any given vector b , the best ℓ_2 approximation of b using the columns in $A_{\uparrow n}$ is given by $b_n = P_n b$, where, by the Projection theorem, $b_n \perp (b - b_n)$. The corresponding optimal coefficient vector is $x_n = A_{\uparrow n}^\dagger b = A_{\uparrow n}^\dagger b_n$.

For a column a_i of A , $i = 1, \dots, q$, we have that

$$\begin{aligned} A_{\uparrow n}^\dagger a_i &= t_i A_{\uparrow n}^\dagger \bar{a}_i \\ &= t_i T_{\#n}^{-1} \left(U_{\uparrow n}^\top U_{\uparrow n} / N + \sigma_u^2 I_n \right)^{-1} \left(U_{\uparrow n}^\top u_i / N + \sigma_u^2 I_{\uparrow n}^\top \zeta_i \right), \end{aligned}$$

where ζ_i is the i -th column of the identity matrix I_q , and $t_i \doteq [T]_{i,i} = \|\bar{a}_i\|_2^{-1} = (u_i^\top u_i + N \sigma_u^2)^{-1/2}$.

We shall next examine the condition in (15).

A.4.2 The large N sparsity pattern

From Lemma 3, we have that $U_{\uparrow n}^\top U_{\uparrow n} / N \rightsquigarrow \nu^2 I_n$, and $U_{\uparrow n}^\top u_i / N \rightsquigarrow 0_{n,1}$, if $i > n$. Moreover, $t_i T_{\#n}^{-1} \rightsquigarrow I_n$, and $I_{\uparrow n}^\top \zeta_i = \begin{bmatrix} I_n & 0_{n,q-n} \end{bmatrix} \zeta_i = 0_{n,1}$, if $i > n$. Therefore, considering the scalar-valued function $\|W_{\#n} A_{\uparrow n}^\dagger a_i\|_1$, which is Lipschitz continuous w.r.t. the entries of $U_{\uparrow n}^\top U_{\uparrow n} / N$ and $U_{\uparrow n}^\top u_i / N$, and applying Lemma 1, we obtain that, for $i > n$,

$$\begin{aligned} &\left\| W_{\#n} A_{\uparrow n}^\dagger a_i \right\|_1 \rightsquigarrow \\ &\left\| W_{\#n} t_i T_{\#n}^{-1} \frac{\sigma_u^2}{\nu^2 + \sigma_u^2} I_{\uparrow n}^\top \zeta_i \right\|_1 = \left\| W_{\#n} I_n \frac{\sigma_u^2}{\nu^2 + \sigma_u^2} 0_{n,1} \right\|_1 = 0. \end{aligned}$$

Hence it holds that

$$\Upsilon_n(A) = 1 - \max_{i > n} w_i^{-1} \left\| W_{\#n} A_{\uparrow n}^\dagger a_i \right\|_1 \rightsquigarrow 1 - 0 = 1. \quad (\text{A.2})$$

Let us now consider the left-hand side in the condition (15). Using the fact that $b = \bar{A} h^n + e$, with e given in (A.1), we have

$$\begin{aligned} W^{-1} A^\top (b - P_n b) &= W^{-1} T \bar{A}^\top (\bar{A} h^n + e - P_n \bar{A} h^n - P_n e) \\ &= W^{-1} T (\bar{A}^\top \bar{A} h^n + \bar{A}^\top e - \bar{A}^\top P_n \bar{A} h^n - \bar{A}^\top P_n e). \end{aligned} \quad (\text{A.3})$$

Defining $\tilde{T} \doteq T \sqrt{N}$ and dividing (A.3) by \sqrt{N} , we obtain

$$\begin{aligned} &\frac{1}{\sqrt{N}} W^{-1} A^\top (b - P_n b) = \\ &= W^{-1} \tilde{T} (\bar{A}^\top \bar{A} h^n + \bar{A}^\top e - \bar{A}^\top P_n \bar{A} h^n - \bar{A}^\top P_n e) / N. \end{aligned} \quad (\text{A.4})$$

Now we evaluate

$$\begin{aligned} \bar{A}^\top \bar{A} / N &= U^\top U / N + \sigma_u^2 I_q \\ \bar{A}^\top e / N &= U^\top U_{\downarrow n} h_{\downarrow n} / N + U^\top \Delta_{\downarrow 0} h / N + U^\top \delta_y / N - \sigma_u^2 h^n \\ \bar{A}^\top P_n \bar{A} h^n / N &= (U^\top U / N + \sigma_u^2 I_q) h^n \\ \bar{A}^\top P_n e / N &= (U^\top U_{\uparrow n} / N + \sigma_u^2 I_{\uparrow n}) (U_{\uparrow n}^\top U_{\uparrow n} / N + \sigma_u^2 I_n)^{-1} \\ &\quad \{ [U_{\uparrow n}^\top U_{\downarrow n} h_{\downarrow n} + U_{\uparrow n}^\top (\Delta_{\downarrow 0} h + \delta_y)] / N - \sigma_u^2 I_{\uparrow n}^\top h^n \} \end{aligned}$$

and observe that

$$\begin{aligned} U^\top U / N &\rightsquigarrow \nu^2 I_q \\ U^\top U_{\downarrow n} h_{\downarrow n} / N &\rightsquigarrow \nu^2 (h^q - h^n) \\ U^\top U_{\uparrow n} / N &\rightsquigarrow \begin{bmatrix} \nu^2 I_n \\ 0 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} U^\top \Delta_{\downarrow 0} h / N &\rightsquigarrow 0 \\ U^\top \delta_y / N &\rightsquigarrow 0 \\ U_{\uparrow n}^\top U_{\downarrow n} h_{\downarrow n} / N &\rightsquigarrow 0 \end{aligned}$$

hence

$$\begin{aligned} \bar{A}^\top \bar{A} / N &\rightsquigarrow (\nu^2 + \sigma_u^2) I_q \\ \bar{A}^\top e / N &\rightsquigarrow \nu^2 (h^q - h^n) - \sigma_u^2 h^n \\ \bar{A}^\top P_n \bar{A} h^n / N &\rightsquigarrow (\nu^2 + \sigma_u^2) h^n \\ \bar{A}^\top P_n e / N &\rightsquigarrow -\sigma_u^2 h^n. \end{aligned}$$

Substituting in (A.4) we obtain that

$$\frac{1}{\sqrt{N}} W^{-1} A^\top (b - P_n b) \rightsquigarrow \nu^2 W^{-1} \tilde{T} (h^q - h^n). \quad (\text{A.5})$$

Finally, observe that for the i -th diagonal element t_i of T it holds that (by Lemma 1)

$$t_i^2 = \frac{1}{\|\bar{a}_i\|_2^2} = \frac{1}{\|u_i\|_2^2 + \sigma_u^2 N} \rightsquigarrow \frac{1}{N(\nu^2 + \sigma_u^2)}$$

and thus, for the i -th diagonal element \tilde{t}_i of \tilde{T} , we have

$$\tilde{t}_i^2 \rightsquigarrow \frac{1}{\nu^2 + \sigma_u^2}.$$

Therefore, from (A.5), we obtain that

$$\frac{1}{\sqrt{N}} [W^{-1} A^\top (b - P_n b)]_i \rightsquigarrow z_i \doteq \begin{cases} 0, & \text{for } i = 1, \dots, n; \\ w_i^{-1} \nu \kappa h(i), & \\ & \text{for } i = n + 1, \dots, q. \end{cases}$$

where $\kappa \doteq \nu / \sqrt{\nu^2 + \sigma_u^2}$.

From the definition of the symbol \rightsquigarrow , the above expression implies that for any given $\epsilon_1 > 0$ and $\beta_1 \in (0, 1)$ there exists an integer N_1 such that, for any $N \geq N_1$, it results

$$\mathbb{P} \left\{ \left| \frac{1}{\sqrt{N}} [W^{-1} A^\top (b - P_n b)]_i - |z_i| \right| \geq \epsilon_1 \right\} \leq \beta_1. \quad (\text{A.6})$$

Further, under the Assumption 2 that $|h(i)| \leq L \rho^{i-1}$ and since the weight sequence is assumed to be nondecreasing, we have that

$$|z_i| \leq w_n^{-1} \nu \kappa L \rho^n, \quad \forall i = 1, \dots, q. \quad (\text{A.7})$$

Since, for all $i = 1, \dots, q$,

$$\begin{aligned} &\left| \frac{1}{\sqrt{N}} [W^{-1} A^\top (b - P_n b)]_i - |z_i| \right| \geq \\ &\geq \frac{1}{\sqrt{N}} |[W^{-1} A^\top (b - P_n b)]_i| - |z_i| \\ &\geq \frac{1}{\sqrt{N}} |[W^{-1} A^\top (b - P_n b)]_i| - w_n^{-1} \nu \kappa L \rho^n, \end{aligned}$$

from (A.6) it follows that, for any $N \geq N_1$,

$$\mathbb{P} \left\{ \frac{1}{\sqrt{N}} |[W^{-1} A^\top (b - P_n b)]_i| - w_n^{-1} \nu \kappa L \rho^n \geq \epsilon_1 \right\} \leq \beta_1;$$

hence, from Bonferroni's inequality, for any $N \geq N_1$ we have

$$\mathbb{P} \left\{ \frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty > w_n^{-1} \nu \kappa L \rho^n + \epsilon_1 \right\} \leq q \beta_1.$$

Taking the complementary event, for any $N \geq N_1$ it results

$$\mathbb{P} \left\{ \frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty \leq w_n^{-1} \nu \kappa L \rho^n + \epsilon_1 \right\} \geq 1 - q \beta_1. \quad (\text{A.8})$$

Similarly, from (A.2) it follows that for any given $\epsilon_2 > 0$ and $\beta_2 \in (0, 1)$ there exists an integer N_2 such that

$$\mathbb{P} \{ |\Upsilon_n(A) - 1| \leq \epsilon_2 \} = \mathbb{P} \{ 1 - \Upsilon_n(A) \leq \epsilon_2 \} \geq 1 - \beta_2, \quad \forall N \geq N_2; \quad (\text{A.9})$$

thus

$$\mathbb{P} \left\{ \frac{\gamma}{2\sqrt{N}} (1 - \epsilon_2) \leq \frac{\gamma}{2\sqrt{N}} \Upsilon_n(A) \right\} \geq 1 - \beta_2, \quad \forall N \geq N_2. \quad (\text{A.10})$$

Considering the joint events in (A.8) and (A.10), we have from Bonferroni's inequality that

$$\begin{aligned} &\left\{ \frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty \leq w_n^{-1} \nu \kappa L \rho^n + \epsilon_1 \right\} \cap \\ &\quad \left\{ \frac{\gamma}{2\sqrt{N}} (1 - \epsilon_2) \leq \frac{\gamma}{2\sqrt{N}} \Upsilon_n(A) \right\} \end{aligned}$$

holds with probability no smaller than $1 - \beta$, for any $N \geq N_\beta \doteq \max(N_1, N_2)$, with $\beta = q\beta_1 + \beta_2$. Next, observe that if it holds that

$$w_n^{-1} \nu \kappa L \rho^n + \epsilon_1 \leq \frac{\gamma}{2\sqrt{N}} (1 - \epsilon_2), \quad (\text{A.11})$$

then we may conclude with confidence at least $1 - \beta$ that

$$\frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty \leq \frac{\gamma}{2\sqrt{N}} \Upsilon_n(A). \quad (\text{A.12})$$

Suppose that condition (20) holds, thus $\gamma = 2\mu w_n^{-1} L \rho^n \nu \kappa \sqrt{N}$ for some $\mu > 1$; substituting this expression into (A.11), we obtain the condition

$$\epsilon_1 + \mu w_n^{-1} L \rho^n \nu \kappa \epsilon_2 \leq (\mu - 1) w_n^{-1} L \rho^n \nu \kappa.$$

Since $\mu > 1$, and since ϵ_1, ϵ_2 can be chosen arbitrarily, this condition is satisfied for a sufficiently small choice of ϵ_1, ϵ_2 . Therefore, condition (A.11) is satisfied, and hence (A.12) is satisfied with probability no smaller than $1 - \beta$. The statement then follows from Lemma 2. ■

A.5 Proof of Corollary 5

We apply Theorem 4 with n being equal to the leading order $n_l(N)$ of the system. Since (19) holds for $i = n_l(N)$, substituting this expression into (20) we have the condition

$$\gamma \geq 2\mu w_{n_l(N)}^{-1} \rho \sigma_y \kappa,$$

for some $\mu > 1$, which is equivalent to (21). The claim then follows by applying Theorem 4. ■

A.6 Proof of Corollary 6

We follow the same reasoning as in Section A.4 up to (A.5). Then, we observe that since h is FIR of order n , then $h_{\downarrow n}$ is identically zero, hence from (A.5) it follows that

$$\frac{1}{\sqrt{N}} W^{-1} A^\top (b - P_n b) \sim 0,$$

which means that for any given $\epsilon_1 > 0$ and $\beta_1 \in (0, 1)$ there exists an integer N_1 such that

$$\mathbb{P} \left\{ \frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty \leq \epsilon_1 \right\} \geq 1 - q\beta_1, \quad \forall N \geq N_1.$$

Following a reasoning similar to the one in (A.9)–(A.12), we claim that if

$$\epsilon_1 \leq \frac{\gamma}{2\sqrt{N}} (1 - \epsilon_2), \quad (\text{A.13})$$

then we may conclude with confidence at least $1 - \beta$ that

$$\frac{1}{\sqrt{N}} \|W^{-1} A^\top (b - P_n b)\|_\infty < \frac{\gamma}{2\sqrt{N}} \Upsilon_n(A). \quad (\text{A.14})$$

But, since $\gamma > 0$, condition (A.13) can always be satisfied for some ϵ_1, ϵ_2 , and hence (A.14) holds with probability at least $1 - \beta$. The claim then follows from Lemma 2. ■

References

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, AC-19(6):716–723, 1974.
- C.M. Bishop. Training with noise is equivalent to Tikhonov regularization. *J. Neural Computation*, 7(1):108–116, 1995.
- T. Chen, H. Ohlsson, and L. Ljung. On the estimation of transfer functions, regularizations and Gaussian processes – revisited. *Automatica*, 48(8):1525–1535, 2012.
- C. De Mol, E. De Vito, and L. Rosasco. Elastic-net regularization in learning theory. *J. of Complexity*, 25(2):201–230, 2009.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *J. of Statistical Software*, 33(1):1–22, 2010.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Series in Statistics. Springer, New York, second edition, 2009.
- I. Kollár. *Frequency Domain System Identification Toolbox User’s Guide*. The MathWorks, Inc., Natick, MA, 1994.
- I. Kollár, R. Pintelon, and J. Schoukens. Frequency domain system identification toolbox for Matlab: a complex application example. In *Proc. of IFAC SYSID’94*, pages 23–28, vol. 4, Copenhagen, Denmark, 1994.
- L. Ljung. *System Identification: Theory for the User*. Prentice-Hall, Englewood Cliffs, second ed., 1999a.
- L. Ljung. Model validation and model error modeling. In B. Wittenmark and A. Rantzer, editors, *The Åström Symposium on Control*, pages 15–42, Lund, Sweden, Aug. 1999b. Studentlitteratur.
- M. Milanese, F. Ruiz, and M. Taragna. Direct data-driven filter design for uncertain LTI systems with bounded noise. *Automatica*, 46(11):1773–1784, 2010.
- National Instruments Corporation. *LabVIEW System Identification Toolkit User Manual*. Austin, TX, 2004–2006.
- G. Pillonetto and G. De Nicolao. A new kernel-based approach for linear system identification. *Automatica*, 46(1):81–93, 2010.
- G. Pillonetto, A. Chiuso, and G. De Nicolao. Prediction error identification of linear systems: a nonparametric Gaussian regression approach. *Automatica*, 47(2):291–305, 2011.
- G. Pillonetto, F. Dinuzzo, T. Chen, G. De Nicolao, and L. Ljung. Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3):657–682, 2014.
- R. Pintelon and J. Schoukens. *System identification: a frequency domain approach*. John Wiley & Sons, second edition, 2012.
- J. Rissanen. Modelling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- T. Söderström and P. Stoika. *System Identification*. Prentice-Hall, 1989.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal Statist. Soc. B*, 58(1):267–288, 1996.
- D. E. Torfs, R. Vuerinckx, J. Swevers, and J. Schoukens. Comparison of two feedforward design methods aiming at accurate trajectory tracking of the end point of a flexible robot arm. *IEEE Transactions on Control Systems Technology*, 6(1):2–14, January 1998.
- J. A. Tropp. Just relax: convex programming methods for identifying sparse signals in noise. *IEEE Transactions on Information Theory*, 52(3):1030–1051, 2006.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *J. Royal Statist. Soc. B*, 67(2):301–320, 2005.