

Regularized System Identification: A Hierarchical Bayesian Approach [★]

Mohammad Khosravi,^{*} Andrea Iannelli,^{*} Mingzhou Yin,^{*}
Anilkumar Parsi,^{*} Roy S. Smith^{*}

^{*} *Automatic Control Laboratory, ETH, Zürich 8092, Switzerland*
(e-mail: {khosravm,iannelli,myin,aparsi,rsmith}@control.ee.ethz.ch).

Abstract: In this paper, the hierarchical Bayesian method for regularized system identification is introduced. To this end, a hyperprior distribution is considered for the regularization matrix and then, the impulse response and the regularization matrix are jointly estimated based on a maximum a posteriori (MAP) approach. Toward introducing a suitable hyperprior, we decompose the regularization matrix using Cholesky decomposition and reduce the estimation problem to the cone of upper triangular matrices with positive diagonal entries. Following this, the hyperprior is introduced on a designed sub-cone of this set. The method differs from the current trend in regularized system identification from various aspect, e.g., the estimation is performed by solving a single stage problem. The MAP estimation problem reduces to a multi-convex optimization problem and a sequential convex programming algorithm is introduced for solving this problem. Consequently, the proposed method is a computationally efficient strategy specially when the regularization matrix has a large size. The method is numerically verified on benchmark examples. Owing to the employed full Bayesian approach, the estimation method shows a satisfactory bias-variance trade-off.

Keywords: System identification, hierarchical Bayesian, sequential convex programming.

1. INTRODUCTION

System identification address the problem of estimating appropriate mathematical models of dynamical systems using a set of measured data and methods of statistics (Ljung, 1999). The importance of models for prediction and control of physical systems is well recognized, for example thermal energy system in the buildings, network of agents, robotics, biological systems among many others.

Including regularization in the identification problem has received an extensive attention. The underlying idea is to integrate the prior knowledge, which is available in addition to the measurement data, into in the model estimation problem. For the purpose of obtaining a system featuring low complexity in McMillan sense, various types of regularization such as the Hankel matrix rank, the nuclear norm and the atomic norm of system are introduced in the literature (Fazel et al., 2013; Mohan and Fazel, 2010; Smith, 2014; Shah et al., 2012). Additionally, starting from the seminal work of Pillonetto and De Nicolao (2010), the problem of studying Tikhonov-type regularizations for imposing constraints on the latent system, like the smoothness and stability of the impulse response, has received significant attention (Pillonetto et al., 2014). In this framework, the system identification problem is formulated as a regularized regression with a regularization term coming from the norm of a reproducing kernel Hilbert space

(RKHS) (Aronszajn, 1950) which penalizes the feasible solutions not agreeing with the prior knowledge. The main ingredient of these Hilbert spaces is the kernel function. Subsequently, designing a suitable regularization boils down to characterizing appropriate kernels. In the context of regularized system identification, the most common kernels are *tuned/correlated* (TC), *diagonal/correlated* (DC) and *stable spline* (SS) (Pillonetto et al., 2014). Moreover, many other kernels and also regularization matrices have been proposed following various approaches, e.g., inspired by machine learning and system theory (Chen, 2018; Khosravi et al., 2020), using the harmonic analysis of stochastic processes (Zorzi and Chiuso, 2018), based on filter design methods (Marconato et al., 2016). One can use various operation on given kernels or regularization matrices and introduce new ones (Hong et al., 2018; Chen et al., 2014, 2018). Once the type of kernel or regularization matrices is fixed, the characterizing hyperparameters are estimated using the given measurement data \mathcal{D} . In this regard, a number of approaches are introduced for hyperparameter estimation, e.g., empirical Bayes (EB) (Pillonetto et al., 2014), Stein unbiased risk estimator (SURE) (Hong et al., 2018) and cross-validation (CV) (Mu et al., 2018a,b). In summary, in the current trend of regularized system identification, first a prior or a regularization for the impulse response is considered and subsequently, the parameters of prior are tuned, and finally, the impulse response is estimated.

In this paper, the hierarchical Bayesian method for regularized system identification is investigated where a *joint estimation* for the impulse response and the regularization

[★] This research project is part of the Swiss Competence Center for Energy Research SCCER FEEB&D of the Swiss Innovation Agency Innosuisse, and supported by the Swiss National Science Foundation under grant no.: 200021.178890.

matrix is proposed based on a maximum a posteriori (MAP) approach. To this end, appropriate hyperprior which suitably determines the structure of regularization matrix is introduced, and subsequently the form of impulse response. Moreover, the hyperprior is designed so that the estimation method is numerical tractable. The regularization matrix is indeed decomposed using Cholesky decomposition and the estimation problem is reduced to a sub-cone of the space of upper triangular matrices with positive diagonal entries. The final MAP estimation problem has a multi-convex programming form which can be efficiently solved using an introduced sequential convex programming algorithm. The method is finally numerically verified on benchmark examples.

2. NOTATIONS

The set of natural numbers, non-negative integers, integers, real numbers and non-negative real numbers are denoted here by \mathbb{N} , \mathbb{Z}_+ , \mathbb{Z} , \mathbb{R} and \mathbb{R}_+ , respectively. Let \mathbb{F} be either \mathbb{R} or \mathbb{R}_+ . Then, for any $n \in \mathbb{N}$, the space of n -dimensional vectors with entries in \mathbb{F} is denoted by \mathbb{F}^n and $\mathbf{h} \in \mathbb{F}^n$ is indicated by $\mathbf{h} = (h_k)_{k=1}^n$ to show explicitly in terms of its entries, i.e., h_k is the k^{th} entry of \mathbf{h} . The set of n -by- m matrices with entries in \mathbb{F} is denoted by $\mathbb{F}^{n \times m}$. The set of n -by- n real symmetric positive-definite matrices, n -by- n real symmetric semi-positive-definite matrices, n -by- n upper triangular matrices with non-negative diagonal entries and n -by- n upper triangular matrices with positive diagonal entries are denoted by \mathbb{S}_{++}^n , \mathbb{S}_+^n , \mathbb{U}_+^n and \mathbb{U}_{++}^n , respectively. The zero vector, the zero matrix, and the identity matrix are indicated with $\mathbf{0}$, $\mathbf{0}$, and \mathbb{I} , respectively. For any vector \mathbf{x} , the Euclidean norm of \mathbf{x} is shown by $\|\cdot\|$. The interior and closure of set \mathcal{S} is denoted by $\text{int}\mathcal{S}$ and $\text{cl}\mathcal{S}$, respectively. The convex cone of set \mathcal{X} is denoted by $\text{convcone}\mathcal{X}$. Given set \mathcal{X} , $\mathbf{1}_{\mathcal{X}}$ is a function which is one on \mathcal{X} and zero elsewhere. The space of real-valued signals defined over \mathbb{Z} is denoted by $\mathbb{R}^{\mathbb{Z}}$. The *forward shift operator*, denoted by \mathbf{q} , is an operator on the space of signals, $\mathbf{q} : \mathbb{R}^{\mathbb{Z}} \rightarrow \mathbb{R}^{\mathbb{Z}}$, defined by $(\mathbf{q}\mathbf{u})_t = u_{t+1}$, for any $t \in \mathbb{Z}$ and any $\mathbf{u} \in \mathbb{R}^{\mathbb{Z}}$. By $X \sim \mathcal{N}(\mu, \Sigma)$, we mean that X is a Gaussian random vector with mean μ and covariance Σ . The probability density functions are denoted by p . The *support* of distribution μ on the measurable space $(\mathcal{X}, \mathcal{A})$ is defined as $\text{supp}(\mu) = \{\mathbf{x} \in \mathcal{X} \mid \mu(\mathbf{x}) \neq 0\}$.

3. SYSTEM IDENTIFICATION

We consider a single-input-single-output causal and stable linear time-invariant (LTI) system described by transfer function $G(\mathbf{q})$. Let the input signal $\mathbf{u} \in \mathbb{R}^{\mathbb{Z}}$ be applied to the system and the output signal of system, denoted by $\mathbf{y} \in \mathbb{R}^{\mathbb{Z}}$, be measured. We assume that the measurement of the output is subject to additive white measurement noise. Accordingly, given that $\mathbf{u} = (u_t)_{t \in \mathbb{Z}}$ and $\mathbf{y} = (y_t)_{t \in \mathbb{Z}}$, at any measurement time instant t , we have

$$y_t = G(\mathbf{q})u_t + w_t, \quad (1)$$

where $\mathbf{w} := (w_t)_{t \in \mathbb{Z}} \in \mathbb{R}^{\mathbb{Z}}$ is the measurement noise, a zero mean white Gaussian signal with variance σ_w^2 . More precisely, \mathbf{w} is a sequence of independent and identically distributed (i.i.d.) random variables with distribution $\mathcal{N}(0, \sigma_w^2)$. Given a finite set of measurement data, it is desired to estimate the LTI system.

3.1 Maximum Likelihood and Prediction Error Method

In the classical system identification (Ljung, 1999), besides nonparametric approaches such as *empirical transfer function estimation* (ETFE), the system is represented in a parametric form and then, the parameters of the system are estimated using available measurement data. This procedure is done principally using a *maximum likelihood* (ML) approach. The system is assumed stable, one possible parametric representation of the system is in form of a finite impulse response (FIR). More precisely, the transfer function of the system $G(\mathbf{q})$ is approximated by

$$G(\mathbf{q}) = \sum_{k=0}^{n_g-1} g_k \mathbf{q}^{-k}, \quad (2)$$

where $\mathbf{g} := (g_k)_{k=0}^{n_g-1} \in \mathbb{R}^{n_g}$ is the corresponding FIR of the system. Now, let \mathcal{D} be a set of measurement data

$$\mathcal{D} = \{(u_t, y_t) \mid t = 0, 1, \dots, n_{\mathcal{D}} - 1\}. \quad (3)$$

Define vector φ_t as

$$\varphi_t = [u_t \ u_{t-1} \ \dots \ u_{t-n_g+1}]^{\top}, \quad (4)$$

for $t = 0, \dots, n_{\mathcal{D}} - 1$. Also, define vector \mathbf{y} , vector \mathbf{w} and matrix Φ respectively as

$$\mathbf{y} = \begin{bmatrix} y_0 \\ \vdots \\ y_{n_{\mathcal{D}}-1} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_{n_{\mathcal{D}}-1} \end{bmatrix}, \quad \Phi = \begin{bmatrix} \varphi_0^{\top} \\ \vdots \\ \varphi_{n_{\mathcal{D}}-1}^{\top} \end{bmatrix}. \quad (5)$$

From (1) and (2), one has

$$\mathbf{y} = \Phi \mathbf{g} + \mathbf{w}. \quad (6)$$

Since, for any t , we have $w_t \sim \mathcal{N}(0, \sigma_w^2)$, it holds $\mathbf{y} - \Phi \mathbf{g} \sim \mathcal{N}(0, \sigma_w^2 \mathbb{I})$. Therefore, it follows that

$$p(\mathbf{y} | \Phi, \mathbf{g}) = (2\pi\sigma_w^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma_w^2} \|\mathbf{y} - \Phi \mathbf{g}\|^2\right), \quad (7)$$

and the ML estimation of \mathbf{g} is defined as

$$\mathbf{g}^{\text{ML}} = \arg\max_{\mathbf{g} \in \mathbb{R}^{n_g}} p(\mathbf{y} | \Phi, \mathbf{g}). \quad (8)$$

Since, we have

$$-\ln p(\mathbf{y} | \Phi, \mathbf{g}) = \frac{n}{2} \ln(2\pi\sigma_w^2) + \frac{1}{2\sigma_w^2} \|\mathbf{y} - \Phi \mathbf{g}\|^2, \quad (9)$$

it is easily deduced that

$$\mathbf{g}^{\text{ML}} = \arg\min_{\mathbf{g} \in \mathbb{R}^{n_g}} \|\mathbf{y} - \Phi \mathbf{g}\|^2. \quad (10)$$

From (10), one can see that the ML estimation, \mathbf{g}^{ML} , is obtained by solving a least square (LS) problem which is the minimization of the prediction error and the solution of prediction error method (PEM). Hence, one can also denote the ML estimation by \mathbf{g}^{LS} . One should note that when Φ is full column rank, then \mathbf{g}^{ML} or \mathbf{g}^{LS} equals to $(\Phi^{\top} \Phi)^{-1} \Phi^{\top} \mathbf{y}$.

Though the ML or the LS estimation is an unbiased method, it has a number of drawbacks. The main disadvantage of this approach is the requirement of having a large set of high quality data for obtaining an improved decent estimation \mathbf{g}^{ML} by solving problem (10). More precisely, the condition number of $\Phi^{\top} \Phi$ should not be large. Otherwise, the estimation can be non-unique or subject to high variance, specifically for large n_g . Moreover, choosing a suitable model order, here the length of impulse response \mathbf{g} , is a critical and difficult procedure. This is commonly performed by cross validation, model validation techniques or model complexity criteria such as Akaike information criterion (AIC) or Bayesian information criterion (BIC) (Ljung, 1999). These classical strategies are not desirably reliable (Pillonetto et al., 2014; Chen, 2018).

3.2 Maximum A Posteriori and Regularized Methods

In order to estimate the impulse response, one may use a *maximum a posteriori* (MAP) approach by introducing a prior for the impulse response. More precisely, we set a Gaussian prior for the impulse response as

$$\mathbf{g} \sim \mathcal{N}(0, \mathbf{R}^{-1}), \quad (11)$$

where \mathbf{R}^{-1} is a positive definite matrix. The covariance \mathbf{R}^{-1} can encode available prior knowledge and desired features for the latent impulse response, such as smoothness and stability. Once the prior is set, using Bayes rule, one has

$$p(\mathbf{g}|\Phi, \mathbf{y}) = \frac{p(\mathbf{y}|\Phi, \mathbf{g})p(\mathbf{g})}{p(\Phi, \mathbf{y})}, \quad (12)$$

and the MAP estimation of \mathbf{g} is derived by

$$\mathbf{g}^{\text{MAP}} = \underset{\mathbf{g} \in \mathbb{R}^{n_g}}{\operatorname{argmax}} p(\mathbf{g}|\Phi, \mathbf{y}) = \underset{\mathbf{g} \in \mathbb{R}^{n_g}}{\operatorname{argmax}} p(\mathbf{y}|\Phi, \mathbf{g})p(\mathbf{g}). \quad (13)$$

Since $\mathbf{y} - \Phi\mathbf{g} \sim \mathcal{N}(0, \sigma_w^2 \mathbb{I})$ and given prior, equations (7) and

$$p(\mathbf{g}) = (2\pi)^{-\frac{n_g}{2}} \det(\mathbf{R})^{\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{g}^T \mathbf{R} \mathbf{g}\right), \quad (14)$$

hold. Subsequently, it is concluded that

$$\mathbf{g}^{\text{MAP}} = \underset{\mathbf{g} \in \mathbb{R}^{n_g}}{\operatorname{argmax}} \exp\left(-\frac{1}{2\sigma_w^2}\|\mathbf{y} - \Phi\mathbf{g}\|^2 - \frac{1}{2}\mathbf{g}^T \mathbf{R} \mathbf{g}\right). \quad (15)$$

Due to the monotonicity of logarithmic function, it follows that

$$\begin{aligned} \mathbf{g}^{\text{MAP}} &= \underset{\mathbf{g} \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \frac{1}{2\sigma_w^2}\|\mathbf{y} - \Phi\mathbf{g}\|^2 + \frac{1}{2}\mathbf{g}^T \mathbf{R} \mathbf{g}, \\ &= (\Phi^T \Phi + \sigma_w^2 \mathbf{R})^{-1} \Phi^T \mathbf{y}. \end{aligned} \quad (16)$$

From (16), one can see that \mathbf{g}^{MAP} is the solution of a *regularized least square* (RLS) problem. The method is therefore called *regularized system identification* and one may alternatively denote this estimation by \mathbf{g}^{RLS} . For this reason, the covariance matrix \mathbf{R} is also called the *regularization matrix*.

The main role of regularization is to include the additionally available prior information in the estimation. In addition, it can also alleviate the issue of possible high variance. As a result, choosing a suitable regularization matrix \mathbf{R} has a significant impact on the estimation. This is commonly done in two steps. In the first step, a family of regularization matrices, denoted by \mathcal{R} , is chosen. Usually, each element of \mathbf{R}_η is a regularization matrix parameterized with a vector of scalars called *hyperparameters*, denoted by η , which belongs to a given set $\mathcal{E} \subset \mathbb{R}^{n_\eta}$. In other words, we have $\mathcal{R} = \{\mathbf{R}_\eta \mid \eta \in \mathcal{E}\} \subseteq \mathbb{S}_+^{n_g}$. In the literature (Pillonetto et al., 2014), the set of regularization matrices are generally defined based on the notion of the kernels,

$$\mathcal{R}_k = \{\mathbf{R}_\eta \in \mathbb{S}_+^{n_g} \mid \mathbf{R}_{\eta,ij} = \mathbf{k}_\eta(i-1, j-1), \quad 1 \leq i, j \leq n_g, \eta \in \mathcal{E}\}, \quad (17)$$

where $\mathbf{k}_\eta : \mathbb{Z}_+ \times \mathbb{Z}_+$ is a *positive-definite kernel* and $\eta \in \mathcal{E}$ is the vector of hyperparameters (Pillonetto et al., 2014). The most common kernels are *tuned/correlated* (TC), *diagonal/correlated* (DC) and *stable spline* (SS) (Pillonetto et al., 2014). In the literature, many other kernels and also regularization matrices are designed in various approaches, e.g., inspired by machine learning and system theory (Chen, 2018), using the harmonic analysis of stochastic processes (Zorzi and Chiuso, 2018), based on

filter-design methods (Marconato et al., 2016). Given a set of regularization matrices \mathcal{R} , one can build a larger set of regularization matrices by considering the convex cone of \mathcal{R} . Also, Minkowski sum of sets of regularization matrices produces set of regularization matrices. Based on these ideas and similar ones, algorithms for designing the kernel and regularization have been developed (Hong et al., 2018; Chen et al., 2014, 2018). Once the set of regularization matrices is fixed, one should estimate the hyperparameters using the given data \mathcal{D} . In this regard, a number of approaches are introduced for hyperparameter estimation, e.g. empirical Bayes (EB) (Pillonetto et al., 2014), Stein unbiased risk estimator (SURE) (Hong et al., 2018), cross-validation (CV) and generalized cross-validation (GCV) (Mu et al., 2018a,b).

Compared to the ML or the LS estimation method, the regularized method have a satisfactory bias-variance trade-off and requires less amount of data. Moreover, the issue of model order selection is alleviated, specially when the regularization matrix is determined only by a few number of hyperparameters. However, when the number of hyperparameters is significantly large, the hyperparameters estimation methods are prone to high-variance or become computationally intractable. Motivated by this issue, we introduce a new method framed in a full Bayesian setting.

4. HIERARCHICAL BAYESIAN APPROACH

In this section, we introduce the *hierarchical Bayesian* method for joint estimation of the impulse response and the regularization matrix.

Let the noise model be the same as in Section 3. Moreover, given positive definite matrix \mathbf{R} , take the prior distribution of impulse response as

$$\mathbf{g} \sim \mathcal{N}(0, \mathbf{R}^{-1}), \quad (18)$$

Now, let $p_{\mathbf{R}}$ be a given probability distribution defined over $\mathbb{S}_+^{n_g}$, called the *hyperprior*, and \mathcal{R} be the support of $p_{\mathbf{R}}$, i.e.,

$$\mathcal{R} := \operatorname{supp}(p_{\mathbf{R}}) = \{\mathbf{R} \in \mathbb{S}_+^{n_g} \mid p_{\mathbf{R}}(\mathbf{R}) > 0\}. \quad (19)$$

Assume that almost surely the realizations of $p_{\mathbf{R}}$ are in $\mathbb{S}_{++}^{n_g}$. This ensures that the distribution (18) is well defined. From the Bayes rule, one has

$$p(\mathbf{g}, \mathbf{R}|\Phi, \mathbf{y}) = \frac{p(\mathbf{y}|\Phi, \mathbf{g})p(\mathbf{g}|\mathbf{R})p_{\mathbf{R}}(\mathbf{R})}{p(\Phi, \mathbf{y})}. \quad (20)$$

Consequently, the joint maximum a posteriori estimation of \mathbf{g} and \mathbf{R} is derived as

$$\begin{aligned} (\mathbf{g}^{\text{HB}}, \mathbf{R}^{\text{HB}}) &= \underset{\mathbf{g} \in \mathbb{R}^{n_g}, \mathbf{R} \in \mathcal{R}}{\operatorname{argmax}} p(\mathbf{g}, \mathbf{R}|\Phi, \mathbf{y}) \\ &= \underset{\mathbf{g} \in \mathbb{R}^{n_g}, \mathbf{R} \in \mathcal{R}}{\operatorname{argmax}} p(\mathbf{y}|\Phi, \mathbf{g})p(\mathbf{g}|\mathbf{R})p_{\mathbf{R}}(\mathbf{R}). \end{aligned}$$

From (7), (18), the prior distribution $p_{\mathbf{R}}$ and the monotonicity of logarithm function, it follows that

$$\begin{aligned} (\mathbf{g}^{\text{HB}}, \mathbf{R}^{\text{HB}}) &= \underset{\mathbf{g} \in \mathbb{R}^{n_g}, \mathbf{R} \in \mathcal{R}}{\operatorname{argmin}} \frac{1}{2\sigma_w^2}\|\mathbf{y} - \Phi\mathbf{g}\|^2 - \frac{1}{2} \ln \det \mathbf{R} \\ &\quad + \frac{1}{2}\mathbf{g}^T \mathbf{R} \mathbf{g} - \ln p_{\mathbf{R}}(\mathbf{R}). \end{aligned} \quad (21)$$

A suitable choice of $p_{\mathbf{R}}$ has two significant impacts: the correctness of the outcome of the estimation problem and also the numerical tractability of optimization problem (21). Regarding the former issue, one may suggest using a Wishart distribution with suitably chosen *scale matrix*

and *degree of freedom* (Wishart, 1928). However, the latter issue is still a case of concern since of the large number of tuning parameters in the definition of $p_{\mathbf{R}}$.

Tractability of problem (21) depends on the structure of $-\ln p_{\mathbf{R}}$ and therefore, it is desired to introduce a suitable structure for this distribution. The elements of \mathcal{R} are positive semi-definite matrices. Therefore, they have a Cholesky decomposition (Horn and Johnson, 2012). Motivated by this fact, we propose parameterization of the matrix factors in Cholesky decomposition. More precisely, let $n_{\eta} \in \mathbb{N}$ and $\mathbf{U}_1, \dots, \mathbf{U}_{n_{\eta}} \in \mathbb{U}_{++}^{n_{\eta}}$ be linearly independent matrices. Define set \mathcal{U} as the convex cone of these matrices, i.e., we have

$$\mathcal{U} = \left\{ \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i \mid \eta_i \geq 0, \quad \forall i = 1, \dots, n_{\eta} \right\}. \quad (22)$$

Denote by η as the vector defined by $[\eta_1, \dots, \eta_{n_{\eta}}]^{\top}$ in (22). Since $\mathbf{U}_1, \dots, \mathbf{U}_{n_{\eta}}$ are linearly independent, it can be seen from (22) that each element of \mathcal{U} uniquely determines a vector η in $\mathbb{R}_{+}^{n_{\eta}}$ and each $\eta \in \mathbb{R}_{+}^{n_{\eta}}$ uniquely identify an element of \mathcal{U} . Moreover, we know that $\{\mathbf{U}^{\top} \mathbf{U} \mid \mathbf{U} \in \mathcal{U}\} \setminus \{0\} \subseteq \mathbb{S}_{++}^{n_{\eta}}$. Let assume p_{η} be a probability distribution on $\mathbb{R}_{+}^{n_{\eta}}$, i.e., $\text{supp}(p_{\eta}) \subseteq \mathbb{R}_{+}^{n_{\eta}}$. Consequently, one can introduce a corresponding distribution $p_{\mathbf{R}}$ on $\mathcal{R} := \{\mathbf{U}^{\top} \mathbf{U} \mid \mathbf{U} \in \mathcal{U}\}$ as

$$p_{\mathbf{R}}(\mathbf{R}) = p_{\mathbf{R}}(\mathbf{U}^{\top} \mathbf{U}) = p_{\eta}(\eta), \quad (23)$$

where $\mathbf{R} = \mathbf{U}^{\top} \mathbf{U}$ and $\mathbf{U} = \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i$. Note that these decompositions are well-defined and unique due to properties of Cholesky decomposition and linear independency of $\mathbf{U}_1, \dots, \mathbf{U}_{n_{\eta}}$. From the Bayes rule, we can derive a MAP estimation as

$$\begin{aligned} (g^{\text{HB}}, \eta^{\text{HB}}) &= \underset{g \in \mathbb{R}^{n_{\mathbf{g}}}, \eta \in \mathbb{R}_{+}^{n_{\eta}}}{\text{argmax}} p(g, \eta \mid \Phi, \mathbf{y}), \\ &= \underset{g \in \mathbb{R}^{n_{\mathbf{g}}}, \eta \in \mathbb{R}_{+}^{n_{\eta}}}{\text{argmax}} p(\mathbf{y} \mid \Phi, g) p(g \mid \mathbf{U}) p_{\eta}(\eta). \end{aligned}$$

Define function $\tilde{J} : \mathbb{R}^{n_{\mathbf{g}}} \times \mathbb{R}_{+}^{n_{\eta}} \rightarrow \mathbb{R}$ as

$$\tilde{J}(g, \eta) = J(g, \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i). \quad (24)$$

Since, $\mathbf{U}_1, \dots, \mathbf{U}_{n_{\eta}}$ are upper triangular, one can see that

$$\ln \det \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i = \ln \prod_{j=1}^{n_{\eta}} \sum_{i=1}^{n_{\eta}} \eta_i U_{i,jj} = \sum_{j=1}^{n_{\eta}} \ln \sum_{i=1}^{n_{\eta}} \eta_i U_{i,jj},$$

where $U_{i,jj}$ is the j -th element on the diagonal of \mathbf{U}_i , for $i = 1, \dots, n_{\eta}$ and $j = 1, \dots, n_{\mathbf{g}}$. Therefore, we have

$$\begin{aligned} \tilde{J}(g, \eta) &= \frac{1}{2\sigma_w^2} \|\mathbf{y} - \Phi g\|^2 - \sum_{j=1}^{n_{\mathbf{g}}} \ln \sum_{i=1}^{n_{\eta}} \eta_i U_{i,jj} \\ &\quad + \frac{1}{2} \left\| \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i g \right\|^2 - \ln p_{\eta}(\eta). \end{aligned} \quad (25)$$

Subsequently, one can see that

$$(g^{\text{HB}}, \eta^{\text{HB}}) = \underset{g \in \mathbb{R}^{n_{\mathbf{g}}}, \eta \in \mathbb{R}_{+}^{n_{\eta}}}{\text{argmin}} \tilde{J}(g, \eta). \quad (26)$$

Potential candidates for p_{η} are distributions from family of exponential distributions like

$$p_{\eta}(\eta) = A(\lambda_1, \lambda_2) \exp(-\lambda_1 \sum_{i=1}^{n_{\eta}} \eta_i - \lambda_2 \sum_{i=1}^{n_{\eta}} \eta_i^2) \mathbf{1}_{\mathbb{R}_{+}^{n_{\eta}}}(\eta), \quad (27)$$

where $\lambda_2 \geq 0$ and $A(\lambda_1, \lambda_2)$ is the normalizing coefficient. Therefore, for $\eta \in \mathbb{R}_{+}^{n_{\eta}}$, it follows that

$$-\ln p_{\eta}(\eta) = -\ln A(\lambda_1, \lambda_2) - \lambda_1 \sum_{i=1}^{n_{\eta}} \eta_i - \lambda_2 \sum_{i=1}^{n_{\eta}} \eta_i^2. \quad (28)$$

Algorithm 1 Sequential Convex Programming for (29)

- 1: **Input:** $\mathbf{y}, \Phi, \mathbf{U}_1, \dots, \mathbf{U}_{n_{\eta}}, \lambda_1, \lambda_2$ and initial $\eta^{(0)}$.
 - 2: $k \leftarrow 0$.
 - 3: **while** stopping condition is not met **do**
 - 4: $\mathbf{U}^{(k)} \leftarrow \sum_{i=1}^{n_{\eta}} \eta_i^{(k)} \mathbf{U}_i$.
 - 5: Compute

$$g^{(k)} = (\Phi^{\top} \Phi + \sigma_w^2 \mathbf{U}^{(k)\top} \mathbf{U}^{(k)})^{-1} \Phi^{\top} \mathbf{y}.$$
 - 6: Solve the following convex program:

$$\begin{aligned} \eta^{(k+1)} &= \underset{\eta \in \mathbb{R}_{+}^{n_{\eta}}}{\text{argmin}} \frac{1}{2} \left\| \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i g^{(k+1)} \right\|^2 \\ &\quad - \sum_{j=1}^{n_{\mathbf{g}}} \ln \sum_{i=1}^{n_{\eta}} \eta_i U_{i,jj} + \lambda_1 \sum_{i=1}^{n_{\eta}} \eta_i + \lambda_2 \sum_{i=1}^{n_{\eta}} \eta_i^2. \end{aligned}$$
 - 7: $k \leftarrow k + 1$.
 - 8: **end**
 - 9: **Output:** $(g^{\text{HB}}, \eta^{\text{HB}}, \mathbf{U}^{\text{HB}}, \mathbf{R}^{\text{HB}})$.
-

It is worth mentioning that for this choice of p_{η} , we have a lasso penalty function if $\lambda_1 > 0$ and $\lambda_2 = 0$, a ridge penalty function if $\lambda_1 = 0$ and $\lambda_2 > 0$, and an elastic-net penalty function if $\lambda_1 > 0$ and $\lambda_2 > 0$. Also, for particular choices of λ_1 and λ_2 , one has a *maximum entropy* prior on η (Dowson and Wragg, 1973). Moreover, since the distribution p_{η} is integrable, one should have $\lambda_2 \geq 0$, and therefore, the loss function is convex. When $\lambda_2 > 0$, the loss function is strongly convex.

It is noteworthy that (26) reduces to

$$\begin{aligned} (g^{\text{HB}}, \eta^{\text{HB}}) &= \underset{g \in \mathbb{R}^{n_{\mathbf{g}}}, \eta \in \mathbb{R}_{+}^{n_{\eta}}}{\text{argmin}} \frac{1}{2\sigma_w^2} \|\mathbf{y} - \Phi g\|^2 - \sum_{j=1}^{n_{\mathbf{g}}} \ln \sum_{i=1}^{n_{\eta}} \eta_i U_{i,jj} \\ &\quad + \frac{1}{2} \left\| \sum_{i=1}^{n_{\eta}} \eta_i \mathbf{U}_i g \right\|^2 + \lambda_1 \sum_{i=1}^{n_{\eta}} \eta_i + \lambda_2 \sum_{i=1}^{n_{\eta}} \eta_i^2. \end{aligned} \quad (29)$$

Theorem 4.1. i) Let $\eta^{(0)}$ be a non-zero vector in $\mathbb{R}_{+}^{n_{\eta}}$. Then, $\tilde{J}(\cdot, \mathbf{U}_0) : \mathbb{R}^{n_{\mathbf{g}}} \rightarrow \mathbb{R}$ is a strongly convex quadratic function and

$$\underset{g \in \mathbb{R}^{n_{\mathbf{g}}}}{\text{argmin}} J(g, \eta^{(0)}) = (\Phi^{\top} \Phi + \sigma_w^2 \mathbf{U}_0^{\top} \mathbf{U}_0)^{-1} \Phi^{\top} \mathbf{y}, \quad (30)$$

where $\mathbf{U}_0 = \sum_{i=1}^{n_{\eta}} \eta_i^{(0)} \mathbf{U}_i$.

ii) Let $p_{\mathbf{U}}$ be given as in (27). Then, for any $g_0 \in \mathbb{R}^{n_{\mathbf{g}}}$, the function $\tilde{J}(g_0, \cdot) : \mathbb{R}_{+}^{n_{\eta}} \rightarrow \mathbb{R}$ is a proper strictly convex function and the optimization problem

$$\min_{\eta \in \mathbb{R}_{+}^{n_{\eta}}} \tilde{J}(g_0, \eta) \quad (31)$$

is a convex optimization problem with a unique solution. Moreover, if $\lambda_2 > 0$, it is a strongly convex function.

Corollary 4.2. The optimization problem (29) is a bi-convex (multi-convex) programming.

Remark 4.3. Due to the Theorem 4.1, we can introduce a *sequential convex programming* approach for solving (29). The details of this procedure presented in Algorithm 1.

Remark 4.4. The stopping condition in Algorithm 1 can be a combination of maximum number of iterations, a cross-validation error and reaching to a preset minimum step length or reaching a predetermined minimum improvement of objective function.

5. NUMERICAL EXPERIMENT

In this section, the approach is numerically assessed, with main focus on Algorithm 1. To this end, it is required to introduce the settings of optimization, the test data-bank and the numerical setup.

5.1 Generics of Optimization Settings

As the first step, it is required to choose U_1, \dots, U_{n_η} . In this regard, we consider TC-kernels, denoted by $\mathbf{k}^{(\text{TC})} : \mathbb{Z}_+ \times \mathbb{Z}_+ \rightarrow \mathbb{R}_+$ and defined as

$$\mathbf{k}_\alpha^{(\text{DC})}(t_1, t_2) := \alpha^{\max(t_1, t_2)}, \quad \forall t_1, t_2 \in \mathbb{Z}_+, \quad (32)$$

where α is a real scalar in $(0, 1)$. Accordingly, one can define a regularization matrix $R_\alpha \in \mathbb{S}_{++}^{n_g}$ such that the entry of R_α^{-1} at location (t_1, t_2) is $\mathbf{k}_\alpha^{(\text{DC})}(t_1 - 1, t_2 - 1)$, for any $t_1, t_2 = 1, \dots, n_g$. Since R_α^{-1} is positive definite, it has a unique Cholesky decomposition (Marconato et al., 2016) as $R_\alpha^{-1} = U_\alpha^T U_\alpha$ where U_α is defined entry-wise by

$$U_{\alpha, st} = \begin{cases} \alpha^{\frac{1-s}{2}} (1 - \alpha)^{-\frac{1}{2}}, & \text{if } s = t, s < n_g, \\ \alpha^{\frac{1-s}{2}}, & \text{if } s = t = n_g, \\ -\alpha^{\frac{1-s}{2}} (1 - \alpha)^{-\frac{1}{2}}, & \text{if } s = t - 1, \\ 0, & \text{otherwise.} \end{cases} \quad (33)$$

Given $\alpha_1, \dots, \alpha_{n_\eta}$, we can set U_i as U_{α_i} , for $i = 1, \dots, n_\eta$. Regarding λ_1 and λ_2 , we consider two scenarios:

- $\lambda_1 = 0$, i.e., ℓ_2 or ridge penalty function, and,
- $\lambda_2 = 0$, i.e., ℓ_1 or lasso penalty function.

In each of these scenarios, the non-zero parameter can be tuned based on a cross-validation procedure. The initial η , i.e. $\eta^{(0)}$, can be chosen either randomly or by solving (31), where $\mathbf{g}^{(0)}$ is taken as the estimated impulse response obtained with another regularized identification method like TC-kernel regularization. The latter approach is used here. The final step for utilizing Algorithm 1 is specifying Φ and \mathbf{y} which are discussed in the following section.

5.2 Test Data-Bank

Define transfer functions $G(z)$ as

$$G(z) = \frac{z^3 + 0.5z^2}{z^4 - 2.2z^3 + 2.42z^2 - 1.87z + 0.7225}. \quad (34)$$

which are initially introduced in (Wahlberg and Ljung, 1986, Example 5.1) and also utilized in (Pillonetto and De Nicolao, 2010) as benchmark systems. To have a fair comparison, the systems are normalized by their \mathcal{H}_2 norm. The input signal consists of a pseudo-random binary signal (PRBS) of length $n_{\mathcal{D}} = 63$ for time instants $t = 0, \dots, n_{\mathcal{D}} - 1$. Moreover, it is assumed that the systems are initially at rest. At the output of each system, a white additive Gaussian noise with distribution $\mathcal{N}(0, \sigma_w^2)$ is added. Three values of σ_w^2 are considered, namely 0.001, 0.01, 0.1. For each of them 100 realizations of noise are generated and subsequently, 100 data set (3) are obtained by simulation.

5.3 Numerical Experiment Setup

In the numerical experiment, three impulse response estimation methods are tested:

- 1) regularized impulse response estimation using TC kernel,
- 2) hierarchical Bayes with ℓ_1 penalty function, and
- 3) hierarchical Bayes with ℓ_2 penalty function.

For the cases 2) and 3), we take 51 upper triangular matrices, $\{U_i \mid i = 1, \dots, 51\}$, defined according to (33). Here, $\alpha_1, \dots, \alpha_{51}$ are taken equidistantly from interval $[0.8, 0.995]$. The hyperparameters λ_1 and λ_2 are chosen based on a 85%-15% cross-validation scheme. The Algorithm 1 stops when either the number of iterations reaches 50 or when $\|(\mathbf{g}^{(k)}, \eta^{(k)}) - (\mathbf{g}^{(k-1)}, \eta^{(k-1)})\| \leq 10^{-4}$.

The measure of fit is defined (Ljung, 2012) as

$$\text{Fit} = 100 \times \left(1 - \frac{\|\mathbf{g} - \hat{\mathbf{g}}\|}{\|\mathbf{g} - \bar{\mathbf{g}}\|} \right), \quad (35)$$

where \mathbf{g} and $\hat{\mathbf{g}}$ are respectively the true and the estimated impulse responses, and $\bar{\mathbf{g}}$ is a vector of the same length as \mathbf{g} and with entries identically equal to $n_g^{-1} \sum_{k=0}^{n_g-1} g_k$. The results of the estimation is shown in Figure 1 and the statistics are shown in Table 1.

5.4 Discussion

From the numerical results, one can see that the algorithm shows an improved performance, especially when the variance of noise is not significantly high. Additionally, one should note that in all of the cases, the hierarchical Bayes method has smaller mean square error. The performance of ℓ_1 penalty is slightly better comparing to ℓ_2 penalty, however, more numerical analysis is required for providing a more solid argument.

It is worth noting that the optimization problem (29) is a non-convex program and thus can be hampered by the known issue of local minima. This revealed itself in the numerical results in form of outliers (Figure 1, red crosses). However, its effect can be alleviated by performing suitable cross-validation and invalidating the local minima. In fact, when the variance of noise is high, the cross validation does not have enough fidelity and the local minima issue becomes more severe in this case. One possible solution can be utilizing global optimization heuristics.

6. CONCLUSION

A new hierarchical Bayesian method is proposed in this paper. This essentially consists of a maximum a posteriori (MAP) approach for joint estimation of impulse response and regularization matrix. The crucial step is choosing an appropriate hyperprior. To this end, the Cholesky decomposition is used for decomposing the regularization matrix and then the prior is introduced on a suitable sub-cone of the cone of upper triangular matrices with positive diagonal entries. By doing so, one can specify the prior simply on the positive orthant of the space of hyperparameters. Utilizing suitable distributions from the exponential family, the estimation is reduced to a tractable optimization problem with multi-convex structure. Efficient sequential convex programming algorithms are proposed to find optimal solutions. Due to the non-convexity, reaching global optimal solution cannot be guaranteed. The approach is verified numerically on benchmark examples and compared, as an example, with the results

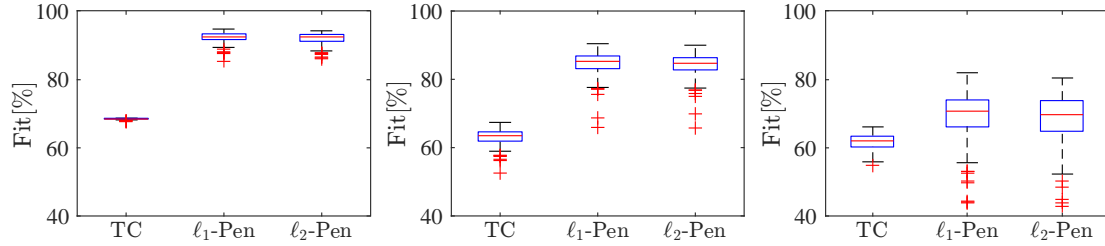


Figure 1. Comparison of fitting performance under different noise levels of $\sigma_w^2 = 0.001$ (left), $\sigma_w^2 = 0.01$ (middle) and $\sigma_w^2 = 0.1$ (right), for system $G(z)$ given in (34).

of TC-kernel. The numerical results show the significant performance of the proposed method.

REFERENCES

- Aronszajn, N. (1950). Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3), 337–404.
- Chen, T. (2018). On kernel design for regularized LTI system identification. *Automatica*, 90, 109–122.
- Chen, T., Andersen, M.S., Ljung, L., Chiuso, A., and Pillonetto, G. (2014). System identification via sparse multiple kernel-based regularization using sequential convex optimization techniques. *IEEE Transactions on Automatic Control*, 59(11), 2933–2945.
- Chen, T., Andersen, M.S., Mu, B., Yin, F., Ljung, L., and Qin, S.J. (2018). Regularized LTI system identification with multiple regularization matrix. *Ifac-papersonline*, 51(15), 180–185.
- Dowson, D. and Wragg, A. (1973). Maximum-entropy distributions having prescribed first and second moments (corresp.). *IEEE Transactions on Information Theory*, 19(5), 689–693.
- Fazel, M., Pong, T.K., Sun, D., and Tseng, P. (2013). Hankel matrix rank minimization with applications to system identification and realization. *SIAM Journal on Matrix Analysis and Applications*, 34(3), 946–977.
- Hong, S., Mu, B., Yin, F., Andersen, M.S., and Chen, T. (2018). Multiple kernel based regularized system identification with SURE hyper-parameter estimator. *IFAC-papersonline*, 51(15), 13–18.
- Horn, R.A. and Johnson, C.R. (2012). *Matrix analysis*. Cambridge university press.
- Khosravi, M., Yin, M., Iannelli, A., Parsi, A., and Smith, R.S. (2020). Low-complexity identification by sparse hyperparameter estimation. *IFAC-papersonline*.
- Ljung, L. (1999). *System identification: theory for the user*. Prentice Hall.
- Ljung, L. (2012). *System identification toolbox: for use with matlab*.
- Marconato, A., Schoukens, M., and Schoukens, J. (2016). Filter-based regularisation for impulse response modelling. *IET Control Theory & Applications*, 11(2), 194–204.
- Mohan, K. and Fazel, M. (2010). Reweighted nuclear norm minimization with application to system identification. In *Proceedings of the 2010 American Control Conference*, 2953–2959. IEEE.
- Mu, B., Chen, T., and Ljung, L. (2018a). Asymptotic properties of generalized cross validation estimators for regularized system identification. *IFAC-PapersOnLine*, 51(15), 203–208.
- Mu, B., Chen, T., and Ljung, L. (2018b). Asymptotic properties of hyperparameter estimators by using cross-validations for regularized system identification. In *2018 IEEE Conference on Decision and Control (CDC)*, 644–649. IEEE.
- Pillonetto, G. and De Nicolao, G. (2010). A new kernel-based approach for linear system identification. *Automatica*, 46(1), 81–93.
- Pillonetto, G., Dinuzzo, F., Chen, T., De Nicolao, G., and Ljung, L. (2014). Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3), 657–682.
- Shah, P., Bhaskar, B.N., Tang, G., and Recht, B. (2012). Linear system identification via atomic norm regularization. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, 6265–6270. IEEE.
- Smith, R.S. (2014). Frequency domain subspace identification using nuclear norm minimization and hankel matrix realizations. *IEEE Transactions on Automatic Control*, 59(11), 2886–2896.
- Wahlberg, B. and Ljung, L. (1986). Design variables for bias distribution in transfer function estimation. *IEEE Transactions on Automatic Control*, 31(2), 134–144.
- Wishart, J. (1928). The generalised product moment distribution in samples from a normal multivariate population. *Biometrika*, 32–52.
- Zorzi, M. and Chiuso, A. (2018). The harmonic analysis of kernel functions. *Automatica*, 94, 125–137.

$\sigma_w^2 = 0.001$	TC	ℓ_1	ℓ_2
Bias ² [$\times 10^4$]	916	16.2	12.7
Var [$\times 10^4$]	9.34	40.3	49.1
MSE [$\times 10^4$]	925	56.5	61.8
$\sigma_w^2 = 0.01$	TC	ℓ_1	ℓ_2
Bias ² [$\times 10^3$]	120.8	3.33	2.88
Var [$\times 10^3$]	6.68	20.1	22.0
MSE [$\times 10^3$]	127	23.4	24.9
$\sigma_w^2 = 0.1$	TC	ℓ_1	ℓ_2
Bias ² [$\times 10^3$]	114	11.9	11.9
Var [$\times 10^3$]	20.9	85.3	94.2
MSE [$\times 10^3$]	135	97.2	106

Table 1. The table provides the statistics of estimation performance. One can see that the proposed method shows better bias-variance trade-off.