



# On asymptotic properties of hyperparameter estimators for kernel-based regularization methods<sup>☆</sup>

Biqiang Mu<sup>a</sup>, Tianshi Chen<sup>b,c,\*</sup>, Lennart Ljung<sup>a</sup>

<sup>a</sup> Division of Automatic Control, Department of Electrical Engineering, Linköping University, Linköping SE-58183, Sweden

<sup>b</sup> School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, China

<sup>c</sup> Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen, China

## ARTICLE INFO

### Article history:

Received 2 July 2017

Received in revised form 13 January 2018

Accepted 4 April 2018

### Keywords:

Kernel-based regularization

Empirical Bayes

Stein's unbiased risk estimator

Asymptotic analysis

## ABSTRACT

The kernel-based regularization method has two core issues: kernel design and hyperparameter estimation. In this paper, we focus on the second issue and study the properties of several hyperparameter estimators including the empirical Bayes (EB) estimator, two Stein's unbiased risk estimators (SURE) (one related to impulse response reconstruction and the other related to output prediction) and their corresponding Oracle counterparts, with an emphasis on the asymptotic properties of these hyperparameter estimators. To this goal, we first derive and then rewrite the first order optimality conditions of these hyperparameter estimators, leading to several insights on these hyperparameter estimators. Then we show that as the number of data goes to infinity, the two SUREs converge to the best hyperparameter minimizing the corresponding mean square error, respectively, while the more widely used EB estimator converges to another best hyperparameter minimizing the expectation of the EB estimation criterion. This indicates that the two SUREs are asymptotically optimal in the corresponding MSE senses but the EB estimator is not. Surprisingly, the convergence rate of two SUREs is slower than that of the EB estimator, and moreover, unlike the two SUREs, the EB estimator is independent of the convergence rate of  $\Phi^T \Phi / N$  to its limit, where  $\Phi$  is the regression matrix and  $N$  is the number of data. A Monte Carlo simulation is provided to demonstrate the theoretical results.

© 2018 Elsevier Ltd. All rights reserved.

## 1. Introduction

The kernel-based regularization methods (KRM) from machine learning and statistics were first introduced to the system identification community in Pilonetto and De Nicolao (2010) and then further developed in Chen, Andersen, Ljung, Chiuso, and Pilonetto (2014), Chen, Ohlsson, and Ljung (2012) and Pilonetto, Chiuso, and De Nicolao (2011). These methods attract increasing attention in the community and have become a complement to the classical maximum likelihood/prediction error methods (ML/PEM) (Chen et al., 2012; Ljung, Singh, & Chen, 2015; Pilonetto & Chiuso, 2015). In particular, KRM may have better average accuracy and robustness than ML/PEM when the data is short and/or has low signal-to-noise ratio (SNR).

There are two core issues for KRM: kernel design and hyperparameter estimation. The former is regarding how to parameterize the kernel matrix with a parameter vector, called hyperparameter, to embed the prior knowledge of the system to be identified, and the latter is regarding how to estimate the hyperparameter based on the data such that the resulting model estimator achieves a good bias–variance trade-off or equivalently, suitably balances the adherence to the data and the model complexity.

The kernel design plays a similar role as the model structure design for ML/PEM and determines the underlying model structure for KRM. In the past few years, many efforts have been spent on this issue and several kernels have been invented to embed various types of prior knowledge, e.g., Carli, Chen, and Ljung (2017), Chen (2018a), Chen et al. (2014), Chen et al. (2016), Chen et al. (2012), Chen and Pilonetto (2018), Dinuzzo (2015), Marconato, Schoukens, and Schoukens (2016), Pilonetto, Chen, Chiuso, Nicolao, and Ljung (2016), Pilonetto et al. (2011), Pilonetto and De Nicolao (2010) and Zorzi and Chiuso (2017). In particular, two systematic kernel design methods (one is from a machine learning perspective and the other one is from a system theory perspective) were developed in Chen (2018b) by embedding the corresponding type of prior knowledge.

<sup>☆</sup> The material in this paper was partially presented at the 20th World Congress of the International Federation of Automatic Control, July 9–14, 2017, Toulouse, France. This paper was recommended for publication in revised form by Associate Editor Thomas Bo Schön under the direction of Editor Torsten Söderström.

\* Corresponding author: Tianshi Chen, School of Science and Engineering and Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen, China.

E-mail addresses: [biqiang.mu@liu.se](mailto:biqiang.mu@liu.se) (B. Mu), [tschen@cuhk.edu.cn](mailto:tschen@cuhk.edu.cn) (T. Chen), [ljung@isy.liu.se](mailto:ljung@isy.liu.se) (L. Ljung).

The hyperparameter estimation plays a similar role as the model order selection in ML/PEM and its essence is to determine a suitable model complexity based on the data. As mentioned in the survey of KRM (Pillonetto, Dinuzzo, Chen, De Nicolao, & Ljung, 2014), many methods can be used for hyperparameter estimation, such as the cross-validation (CV), empirical Bayes (EB),  $C_p$  statistics and Stein's unbiased risk estimator (SURE) and etc. In contrast with the numerous results on kernel design, there are however few results on hyperparameter estimation except Aravkin, Burke, Chiuso, & Pillonetto (2012a, b, 2014), Chen et al. (2014) and Pillonetto & Chiuso (2015). In Aravkin et al. (2012a, b, 2014), two types of diagonal kernel matrices are considered. When  $\Phi^T \Phi / N$  is an identity matrix, where  $\Phi$  is the regression matrix and  $N$  is the number of data, the optimal hyperparameter estimate of the EB estimator has explicit form and is shown to be consistent in terms of the mean square error (MSE). When  $\Phi^T \Phi / N$  is not an identity matrix, the EB estimator is shown to asymptotically minimize a weighted MSE. In Chen et al. (2014), the EB with linear multiple kernel is shown to be a difference of convex programming problem and moreover, the optimal hyperparameter estimate is sparse. In Pillonetto and Chiuso (2015), the robustness of the EB estimator is analyzed.

In this paper, we study the properties of the EB estimator and two SUREs in Pillonetto and Chiuso (2015) with an emphasis on the asymptotic properties of these hyperparameter estimators. In particular, we are interested in the following questions: When the number of data goes to infinity,

- (1) what will be the best kernel matrix, or equivalently, the best value of the hyperparameter?
- (2) which estimator (method) shall be chosen such that the hyperparameter estimate tends to this best value in the given sense?
- (3) what will be the convergence rate of that the hyperparameter estimate tends to this best value? and what factors does this rate depend on?

In order to answer these questions, we employ the regularized least squares method for FIR model estimation in Chen et al. (2012). As a motivation, we first show that the regularized least squares estimate can have smaller MSE than the least squares estimate for any data length if the kernel matrix is chosen carefully. We then derive the first order optimality conditions of these hyperparameter estimators and their corresponding Oracle counterparts (relying on the true impulse response, see Section 3.2 for details). These first order optimality conditions are then rewritten in a way to better expose their relations, leading to several insights on these hyperparameter estimators. For instance, one insight is that for the Oracle estimators, for any data length, and without structure constraints on the kernel matrix, the optimal kernel matrices are same as the one in Chen et al. (2012) and equal to the outer product of the vector of the true impulse response and its transpose. Moreover, explicit solutions of the optimal hyperparameter estimate for two special cases are derived accordingly. Then we turn to the asymptotic analysis of these hyperparameter estimators. Regardless of the parameterization of the kernel matrix, we first show that the two SUREs actually converge to the best hyperparameter minimizing the corresponding MSE, respectively, as the number of data goes to infinity, while the more widely used EB estimator converges to the best hyperparameter minimizing the expectation of the EB estimation criterion. In general, these best hyperparameters are different from each other except for some special cases. This means that the two SUREs are asymptotically optimal in the corresponding MSE senses but the EB estimator is not. We then show that the convergence rate of two SUREs is slower than that of the EB estimator, and moreover, unlike the two SUREs, the EB

estimator is independent of the convergence rate of  $\Phi^T \Phi / N$  to its limit.

The remaining parts of the paper is organized as follows. In Section 2, we recap the regularized least squares method for FIR model estimation and introduce two types of MSE. In Section 3, we introduce six hyperparameter estimators, including the EB estimator, two SUREs, and their corresponding Oracle counterparts. In Section 4, we derive the first order optimality conditions of these hyperparameter estimators and put them in a form that clearly shows their relation, leading to several insights. In Section 5, we give the asymptotic analysis of these hyperparameter estimators, including the asymptotic convergence and the corresponding convergence rate. In Section 6, we illustrate our theoretical results with Monte Carlo simulations. Finally, we conclude this paper in Section 7. All proofs of the theoretical results are postponed to Appendix A.

## 2. Regularized least squares approach for FIR model estimation

### 2.1. Regularized least squares and two types of MSEs

Consider a single-input single-output linear discrete-time invariant, stable and causal system

$$y(t) = G_0(q)u(t) + v(t), \quad t = 1, \dots, N \quad (1)$$

where  $t$  is the time index,  $y(t)$ ,  $u(t)$ ,  $v(t)$  are the output, input and disturbance of the system at time  $t$ , respectively,  $G_0(q)$  is the rational transfer function of the system and  $q$  is the forward shift operator:  $qu(t) = u(t+1)$ . Assume that the input  $u(t)$  is known (deterministic) and the input–output data are collected at time instants  $t = 1, \dots, N$ , and moreover, the disturbance  $v(t)$  is a zero mean white noise with finite variance  $\sigma^2 > 0$ . The problem is to estimate a model for  $G_0(q)$  as well as possible based on the available data  $\{u(t-1), y(t)\}_{t=1}^N$ .

The transfer function  $G_0(q)$  can be written as

$$G_0(q) = \sum_{k=1}^{\infty} g_k^0 q^{-k} \quad (2)$$

where  $g_k^0$ ,  $k = 1, \dots, \infty$  form the impulse response of the system. Since the impulse response coefficients  $\{g_k^0\}$  of the stable rational transfer function  $G_0(q)$  decay exponentially, it is possible to truncate the infinite impulse response at a sufficiently high order, leading to the finite impulse response (FIR) model:

$$G(q) = \sum_{k=1}^n g_k q^{-k}, \quad \theta = [g_1, \dots, g_n]^T \in \mathbb{R}^n. \quad (3)$$

With the FIR model (3), system (1) is now written as

$$y(t) = \phi^T(t)\theta + v(t), \quad t = 1, \dots, N$$

where  $\phi(t) = [u(t-1), \dots, u(t-n)]^T \in \mathbb{R}^n$ , and its matrix–vector form is

$$Y = \Phi\theta + V, \quad \text{where} \quad (4)$$

$$Y = [y(1) y(2) \dots y(N)]^T$$

$$\Phi = [\phi(1) \phi(2) \dots \phi(N)]^T$$

$$V = [v(1) v(2) \dots v(N)]^T.$$

The well-known least squares (LS) estimator

$$\hat{\theta}^{\text{LS}} = \underset{\theta \in \mathbb{R}^n}{\operatorname{argmin}} \|Y - \Phi\theta\|^2 \quad (5a)$$

$$= (\Phi^T \Phi)^{-1} \Phi^T Y \quad (5b)$$

where  $\|\cdot\|$  is the Euclidean norm, is unbiased with respect to the FIR model (4) but may have large variance and mean square error

(MSE) (e.g., when the input is low-pass filtered white noise). The large variance can be mitigated if some bias is allowed and traded for smaller variance and smaller MSE.

One possible way to achieve this goal is to add a regularization term  $\sigma^2 \theta^T P^{-1} \theta$  in the LS criterion (5a), leading to the regularized least squares (RLS) estimator:

$$\hat{\theta}^R = \underset{\theta \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{Y} - \Phi \theta\|^2 + \sigma^2 \theta^T P^{-1} \theta \quad (6a)$$

$$= P \Phi^T (\Phi P \Phi^T + \sigma^2 I_N)^{-1} \mathbf{Y} \quad (6b)$$

where  $P$  is symmetric and positive semidefinite and is called the kernel matrix ( $\sigma^2 P^{-1}$  is often called the regularization matrix), and  $I_N$  is the  $N$ -dimensional identity matrix.

**Remark 1.** As is well known, the RLS estimator (6b) has a Bayesian interpretation. Specifically, assume that  $\theta$  and  $v(t)$  are independent and Gaussian distributed with

$$\theta \sim \mathcal{N}(0, P), \quad v(t) \sim \mathcal{N}(0, \sigma^2) \quad (7)$$

where  $P$  is the prior covariance matrix. Then  $\theta$  and  $\mathbf{Y}$  are jointly Gaussian distributed and moreover, the posterior distribution of  $\theta$  given  $\mathbf{Y}$  is

$$\theta | \mathbf{Y} \sim \mathcal{N}(\hat{\theta}^R, \hat{P}^R)$$

$$\hat{\theta}^R = P \Phi^T (\Phi P \Phi^T + \sigma^2 I_N)^{-1} \mathbf{Y}$$

$$\hat{P}^R = P - P \Phi^T (\Phi P \Phi^T + \sigma^2 I_N)^{-1} \Phi P.$$

Two types of MSE could be used to evaluate the performance of the RLS estimator (6b). The first one is the MSE related to the impulse response reconstruction, see e.g., Chen et al. (2012) and Pillonetto and Chiuso (2015),

$$\text{MSEg}(P) = E \|\hat{\theta}^R - \theta_0\|^2 \quad (8)$$

where  $\theta_0 = [g_1^0, \dots, g_n^0]^T$  with  $g_i^0, i = 1, \dots, n$ , defined in (2) and  $E(\cdot)$  is the mathematical expectation with respect to the noise distribution. The second one is the MSE related to output prediction, see e.g., Pillonetto and Chiuso (2015),

$$\text{MSEy}(P) = E \left[ \sum_{t=1}^N (\phi^T(t) \theta_0 + v^*(t) - \hat{y}(t))^2 \right] \quad (9)$$

where  $\hat{y}(t) = \phi^T(t) \hat{\theta}^R$  and  $v^*(t)$  is an independent copy of the noise  $v(t)$ . Interestingly, the two MSEs (8) and (9) are related with each other through

$$\text{MSEy}(P) = \text{Tr}(E(\hat{\theta}^R - \theta_0)(\hat{\theta}^R - \theta_0)^T \Phi^T \Phi) + N\sigma^2 \quad (10)$$

where  $\text{Tr}(\cdot)$  is the trace of a square matrix. Moreover, they have explicit expressions, which are given in the following proposition.

**Proposition 1.** For a given kernel matrix  $P$ , the two MSEs (8) and (9) take the following form

$$\begin{aligned} \text{MSEg}(P) &= \|P \Phi^T Q^{-1} \Phi \theta_0 - \theta_0\|^2 \\ &\quad + \sigma^2 \text{Tr}(P \Phi^T Q^{-1} Q^{-T} \Phi P^T) \end{aligned} \quad (11a)$$

$$\begin{aligned} \text{MSEy}(P) &= \|\Phi P \Phi^T Q^{-1} \Phi \theta_0 - \Phi \theta_0\|^2 + N\sigma^2 \\ &\quad + \sigma^2 \text{Tr}(\Phi P \Phi^T Q^{-1} Q^{-T} \Phi P^T \Phi^T) \end{aligned} \quad (11b)$$

where  $A^{-T}$  means  $(A^{-1})^T$  for a non-singular matrix  $A$  and

$$Q = \Phi P \Phi^T + \sigma^2 I_N. \quad (12)$$

## 2.2. RLS estimator can outperform LS estimator

It is interesting to investigate whether the RLS estimator (6b) with a suitable choice of the kernel matrix  $P$  can have smaller MSEs (8) and (9) than the LS estimator (5b). The answer is affirmative for MSEg (8) and for the ridge regression case, where  $P^{-1} = (\beta/\sigma^2)I_n$  with  $\beta > 0$ , Hoerl and Kennard (1970) and Theobald (1974). In what follows, we further show that this property also holds for more general  $P$  for MSEg (8) and MSEy (9).

**Proposition 2.** Consider the RLS estimator (6b) and the LS estimator (5b). Suppose that  $P^{-1} = \beta A/\sigma^2$ , where  $\beta > 0$  and  $A$  is symmetric and positive semidefinite. Then for the given  $A$ , there exists  $\beta > 0$  such that (6b) has a smaller MSEg (8) and MSEy (9) than (5b). Moreover, if  $A$  is positive definite, then (6b) has a smaller MSEg (8) and MSEy (9) than (5b) whenever  $0 < \beta < 2\sigma^2/(\theta_0^T A \theta_0)$ .

Proposition 2 shows that for any data length  $N$ , the RLS estimator (6b) can have smaller MSEg (8) and MSEy (9) than the LS estimator (5b) with a sufficiently small regularization “in any direction” and this merit motivates to further explore the potential of the RLS estimator (6b) by careful design of the kernel matrix  $P$ . It is worth to note the paper (Zorzi, 2017) also shows the Bayesian estimator is still optimal when a priori information is known with some uncertainty.

## 3. Design of kernel matrix and hyperparameter estimation

The regularization method has two core issues: kernel matrix design, namely parameterization of the kernel matrix by a parameter vector, called hyperparameter, and hyperparameter estimation.

### 3.1. Parameterization of kernel matrix

For efficient regularization, the symmetric and positive semidefinite kernel matrix  $P$  has to be chosen carefully. It is typically done by postulating a parameterized family of matrices

$$P(\eta), \quad \eta \in \Omega \subset \mathbb{R}^p \quad (13)$$

where  $\eta$  is called the *hyperparameter* and the feasible set  $\Omega$  of  $\eta$  is assumed to be compact. The choice of parameterization is a trade-off of the same kind as the choice of model class in identification: On one hand it should be a large and flexible class to allow as much benefits from regularization as possible. On the other hand, a large set requires larger dimensions of  $\eta$ , and the estimation of these comes with their own penalties (much in the spirit of the Akaike's criterion). Since  $P$  is the prior covariance of the true impulse response, the prior knowledge of the underlying system to be identified, e.g., exponential stability and smoothness, should be embedded in the parameterized matrix  $P(\eta)$ .

A popular way to achieve this goal is through a parameterized positive semidefinite kernel function. So far, several kernels have been invented, such as the stable spline (SS) kernel (Pillonetto & De Nicolao, 2010), the diagonal correlated (DC) kernel and the tuned-correlated (TC) kernel (Chen et al., 2012), which are defined as follows:

$$\begin{aligned} \text{SS} : P_{kj}(\eta) &= c \left( \frac{\alpha^{k+j+\max(k,j)}}{2} - \frac{\alpha^{3\max(k,j)}}{6} \right) \\ \eta &= [c, \alpha] \in \Omega = \{c \geq 0, 0 \leq \alpha \leq 1\}; \end{aligned} \quad (14a)$$

$$\begin{aligned} \text{DC} : P_{kj}(\eta) &= c \alpha^{(k+j)/2} \rho^{|j-k|} \\ \eta &= [c, \alpha, \rho] \in \Omega = \{c \geq 0, 0 \leq \alpha \leq 1, |\rho| \leq 1\}; \end{aligned} \quad (14b)$$

$$\begin{aligned} \text{TC} : P_{kj}(\eta) &= c \alpha^{\max(k,j)} \\ \eta &= [c, \alpha] \in \Omega = \{c \geq 0, 0 \leq \alpha \leq 1\}. \end{aligned} \quad (14c)$$

**Remark 2.** For nonlinearly parameterized kernels, e.g., the kernels in (14), the optimal hyperparameter should be located in the interior of the set  $\Omega$ . To justify this, we take the DC kernel as an example. For the case where either  $c = 0$  or  $\alpha = 0$ ,  $P(\eta) = 0$  and thus (6b) is trivially 0. For the case where  $\alpha = 1$ , this violates the stability of the system. For the case where  $|\rho| = 1$ , the coefficients of the impulse response are perfectly positive or negative correlated, but this is impossible for a stable system. In fact, more formal justification regarding this issue can be found in Pillonetto and Chiuso (2015, p. 115), which shows that the measure of the set containing all optimal estimates lying on the boundary of  $\Omega$  is zero and thus can be neglected when making almost sure convergence statement.

### 3.2. Hyperparameter estimation

Once a parameterized family of the kernel matrix  $P(\eta)$  has been chosen, the task is to estimate, or “tune”, the hyperparameter  $\eta$  based on the data.

Several methods are suggested in the literature, see e.g., Section 14 of Pillonetto et al. (2014), including the empirical Bayes (EB) and SURE methods. The EB method uses the Bayesian interpretation in Remark 1. It follows from the assumption (7) that  $Y$  is Gaussian with mean zero and covariance matrix  $Q$ . As a result, it is possible to estimate the hyperparameter  $\eta$  by maximizing the (marginal) likelihood of  $Y$ , i.e.,

$$\text{EB} : \hat{\eta}_{\text{EB}} = \underset{\eta \in \Omega}{\operatorname{argmin}} \mathcal{F}_{\text{EB}}(P(\eta)) \quad (15a)$$

$$\mathcal{F}_{\text{EB}}(P) = Y^T Q^{-1} Y + \log \det(Q) \quad (15b)$$

where  $Q$  is defined in (12) and  $\det(\cdot)$  denotes the determinant of a square matrix. The SURE method first constructs a Stein's unbiased risk estimator (SURE) of the MSE and then estimates the hyperparameter by minimizing the constructed estimator. Two variants of the SURE method were considered in Pillonetto and Chiuso (2015), which construct the SUREs for  $\text{MSEg}(P)$  in (11a) and  $\text{MSEy}(P)$  in (11b), and are referred to as SUREg and SUREy, respectively:

$$\begin{aligned} \mathcal{F}_{\text{Sg}}(P) &= \|\hat{\theta}^{\text{LS}} - \hat{\theta}^{\text{R}}\|^2 + \sigma^2 \text{Tr}(2R^{-1} - (\Phi^T \Phi)^{-1}) \\ &= \sigma^4 Y^T Q^{-T} \Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y \\ &\quad + \sigma^2 \text{Tr}(2R^{-1} - (\Phi^T \Phi)^{-1}) \end{aligned} \quad (16a)$$

$$\begin{aligned} \mathcal{F}_{\text{Sy}}(P) &= \|Y - \Phi \hat{\theta}^{\text{R}}\|^2 + 2\sigma^2 \text{Tr}(\Phi P \Phi^T Q^{-1}) \\ &= \sigma^4 Y^T Q^{-T} Q^{-1} Y + 2\sigma^2 \text{Tr}(\Phi P \Phi^T Q^{-1}) \end{aligned} \quad (16b)$$

where

$$R = \Phi^T \Phi + \sigma^2 P^{-1}. \quad (17)$$

Then the hyperparameter  $\eta$  is estimated by minimizing the SUREg (16a) and SUREy (16b):

$$\text{SUREg} : \hat{\eta}_{\text{Sg}} = \underset{\eta \in \Omega}{\operatorname{argmin}} \mathcal{F}_{\text{Sg}}(P(\eta)) \quad (18a)$$

$$\text{SUREy} : \hat{\eta}_{\text{Sy}} = \underset{\eta \in \Omega}{\operatorname{argmin}} \mathcal{F}_{\text{Sy}}(P(\eta)). \quad (18b)$$

In the following sections, we will study the properties of the above three estimators EB, SUREg and SUREy. To set reference for these estimators, we introduce their corresponding Oracle counterparts that depend on the true impulse response  $\theta_0$ :

$$\begin{aligned} \text{MSEg} : \hat{\eta}_{\text{MSEg}} &= \underset{\eta \in \Omega}{\operatorname{argmin}} E[\mathcal{F}_{\text{Sg}}(P(\eta))] \\ &= \underset{\eta \in \Omega}{\operatorname{argmin}} \text{MSEg}(P(\eta)) \end{aligned} \quad (19a)$$

$$\text{MSEy} : \hat{\eta}_{\text{MSEy}} = \underset{\eta \in \Omega}{\operatorname{argmin}} E[\mathcal{F}_{\text{Sy}}(P(\eta))]$$

$$= \underset{\eta \in \Omega}{\operatorname{argmin}} \text{MSEy}(P(\eta)) \quad (19b)$$

$$\begin{aligned} \text{EEB} : \hat{\eta}_{\text{EEB}} &= \underset{\eta \in \Omega}{\operatorname{argmin}} E[\mathcal{F}_{\text{EB}}(P(\eta))] \\ &= \underset{\eta \in \Omega}{\operatorname{argmin}} \text{EEB}(P(\eta)) \end{aligned} \quad (19c)$$

where  $\text{MSEg}(P)$  and  $\text{MSEy}(P)$  are defined in (11a) and (11b), respectively, and

$$\text{EEB}(P) = \theta_0^T \Phi^T Q^{-1} \Phi \theta_0 + \sigma^2 \text{Tr}(Q^{-1}) + \log \det(Q). \quad (20)$$

The hyperparameter estimators (19a) and (19b) give the optimal hyperparameter estimates for any data length in the corresponding MSE sense and thus provide reference when evaluating the performance of hyperparameter estimators.

**Remark 3.** Among these hyperparameter estimators, only SUREg (16a) depends on  $(\Phi^T \Phi)^{-1}$ . When  $(\Phi^T \Phi)^{-1}$  is ill-conditioned, SUREg (16a) should be avoided for hyperparameter estimation. One may also note that  $(\Phi^T \Phi)^{-1}$  in the second term is independent of  $P$  and thus can actually be removed in the calculation.

**Remark 4.** It is interesting to note that the first terms of  $\mathcal{F}_{\text{Sg}}(P)$ ,  $\mathcal{F}_{\text{Sy}}(P)$ , and  $\mathcal{F}_{\text{EB}}(P)$  given in (16a), (16b), and (15b) contain the same factors  $Y$  and  $Q^{-1}$ . Moreover, similar to (10),  $\mathcal{F}_{\text{Sg}}(P)$  and  $\mathcal{F}_{\text{Sy}}(P)$  are related with each other through

$$\begin{aligned} \mathcal{F}_{\text{Sy}}(P) &= \text{Tr}\{[(\hat{\theta}^{\text{LS}} - \hat{\theta}^{\text{R}})(\hat{\theta}^{\text{LS}} - \hat{\theta}^{\text{R}})^T \\ &\quad + \sigma^2(2R^{-1} - (\Phi^T \Phi)^{-1})]\Phi^T \Phi\} \\ &\quad + \underbrace{Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y - Y^T Y - n\sigma^2}_{\text{independent of the kernel matrix } P}. \end{aligned} \quad (21)$$

In what follows, we will investigate the properties of the hyperparameter estimators EB, SUREg, and SUREy and their corresponding Oracle estimators EEB, MSEg and MSEy. Before proceeding to the details, we make, without loss of generality, the following assumption.

**Assumption 1.** The hyperparameter estimates  $\hat{\eta}_{\text{Sg}}$ ,  $\hat{\eta}_{\text{Sy}}$ ,  $\hat{\eta}_{\text{EB}}$ ,  $\hat{\eta}_{\text{MSEg}}$ ,  $\hat{\eta}_{\text{MSEy}}$  and  $\hat{\eta}_{\text{EEB}}$  are interior points of  $\Omega$ .

## 4. Properties of hyperparameter estimators: finite data case

In this section, focusing on the finite data case we first give the first order optimality conditions of the hyperparameter estimators and then we consider two special cases for which closed-form expressions of the optimal hyperparameter estimates are available.

### 4.1. First order optimality conditions

The hyperparameter estimates  $\hat{\eta}_{\text{Sg}}$ ,  $\hat{\eta}_{\text{Sy}}$ , and  $\hat{\eta}_{\text{EB}}$  in (18a), (18b), and (15a) should satisfy the first order optimality conditions if they are interior points of  $\Omega$ . For convenience, we let  $\mathcal{C}$  to denote one of the following estimation criteria  $\mathcal{F}_{\text{Sg}}$ ,  $\mathcal{F}_{\text{Sy}}$ ,  $\mathcal{F}_{\text{EB}}$ ,  $\text{MSEg}$ ,  $\text{MSEy}$  or  $\text{EEB}$ . Then the corresponding hyperparameter estimate is a root of the system of equations:

$$\frac{\partial \mathcal{C}(P(\eta))}{\partial \eta} = 0. \quad (22)$$

By the chain rule of compound functions, we have

$$\text{Tr}\left(\frac{\partial \mathcal{C}(P)}{\partial P} \left(\frac{\partial P(\eta)}{\partial \eta_i}\right)^T\right) = 0, \quad 1 \leq i \leq p \quad (23)$$

where, for convenience, when calculating  $\frac{\partial \mathcal{C}(P)}{\partial P}$  the symmetry of  $P$  is ignored according to Lemma B1, that is, the elements of  $P$  are treated independently. One merit of using  $\frac{\partial \mathcal{C}(P)}{\partial P}$  without imposing the structure information on  $P$  is that explicit expressions for the estimation criteria (16a), (16b), and (15b) are available.



**Proposition 3.** The first order partial derivatives of (16a), (16b), and (15b) with respect to  $P$  are, respectively,

$$\frac{\partial \mathcal{F}_{\text{Sg}}(P)}{\partial P} = -2\sigma^4 \Phi^T Q^{-T} \Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y Y^T Q^{-T} \Phi + 2\sigma^4 H^{-T} \bar{H}^{-T} \quad (24a)$$

$$\frac{\partial \mathcal{F}_{\text{Sy}}(P)}{\partial P} = -2\sigma^4 \Phi^T Q^{-T} Q^{-1} Y Y^T Q^{-T} \Phi + 2\sigma^4 \Phi^T Q^{-T} Q^{-T} \Phi \quad (24b)$$

$$\frac{\partial \mathcal{F}_{\text{EB}}(P)}{\partial P} = -\Phi^T Q^{-T} Y Y^T Q^{-T} \Phi + \Phi^T Q^{-T} \Phi \quad (24c)$$

where

$$H = P \Phi^T \Phi + \sigma^2 I_n, \quad \bar{H} = \Phi^T \Phi P + \sigma^2 I_n. \quad (25)$$

Similarly, the partial derivatives of  $\text{MSEg}(P)$ ,  $\text{MSEy}(P)$ , and  $\text{EEB}(P)$  with respect to  $P$  are also available.

**Proposition 4.** The first order partial derivatives of (11a), (11b), and (20) with respect to  $P$  are, respectively,

$$\frac{\partial \text{MSEg}(P)}{\partial P} = -2\sigma^4 H^{-T} H^{-1} \theta_0 \theta_0^T \Phi^T Q^{-T} \Phi + 2\sigma^4 H^{-T} H^{-1} P \Phi^T Q^{-T} \Phi \quad (26a)$$

$$\frac{\partial \text{MSEy}(P)}{\partial P} = -2\sigma^4 \Phi^T Q^{-T} Q^{-1} \Phi \theta_0 \theta_0^T \Phi^T Q^{-T} \Phi + 2\sigma^4 \Phi^T Q^{-T} Q^{-1} \Phi P \Phi^T Q^{-T} \Phi \quad (26b)$$

$$\frac{\partial \text{EEB}(P)}{\partial P} = -\Phi^T Q^{-T} \Phi \theta_0 \theta_0^T \Phi^T Q^{-T} \Phi + \Phi^T Q^{-T} \Phi P^T \Phi^T Q^{-T} \Phi \quad (26c)$$

where  $H$  is defined in (25).

In order to better expose the relation among the partial derivatives derived in Propositions 3 and 4, we define

$$S = P + \sigma^2 (\Phi^T \Phi)^{-1}. \quad (27)$$

With the use of (27) and the identities (B.11)–(B.13) in the Appendix, we rewrite the partial derivatives derived in Propositions 3 and 4 as follows.

**Corollary 1.** The partial derivatives derived in Propositions 3 and 4 can be rewritten as follows:

$$\frac{\partial \text{MSEg}(P)}{\partial P} = 2\sigma^4 S^{-T} (\Phi^T \Phi)^{-2} S^{-1} (P - \theta_0 \theta_0^T) S^{-T} \quad (28a)$$

$$\frac{\partial \mathcal{F}_{\text{Sg}}(P)}{\partial P} = 2\sigma^4 S^{-T} (\Phi^T \Phi)^{-2} S^{-1} (S - \hat{\theta}^{\text{LS}} (\hat{\theta}^{\text{LS}})^T) S^{-T} \quad (28b)$$

$$\frac{\partial \text{MSEy}(P)}{\partial P} = 2\sigma^4 S^{-T} (\Phi^T \Phi)^{-1} S^{-1} (P - \theta_0 \theta_0^T) S^{-T} \quad (28c)$$

$$\frac{\partial \mathcal{F}_{\text{Sy}}(P)}{\partial P} = 2\sigma^4 S^{-T} (\Phi^T \Phi)^{-1} S^{-1} (S - \hat{\theta}^{\text{LS}} (\hat{\theta}^{\text{LS}})^T) S^{-T} \quad (28d)$$

$$\frac{\partial \text{EEB}(P)}{\partial P} = S^{-T} (P^T - \theta_0 \theta_0^T) S^{-T} \quad (28e)$$

$$\frac{\partial \mathcal{F}_{\text{EB}}(P)}{\partial P} = S^{-T} (S^T - \hat{\theta}^{\text{LS}} (\hat{\theta}^{\text{LS}})^T) S^{-T}. \quad (28f)$$

It follows from Corollary 1 that the difference between the partial derivatives of  $\mathcal{F}_{\text{Sg}}(P)$ ,  $\mathcal{F}_{\text{Sy}}(P)$ ,  $\mathcal{F}_{\text{EB}}(P)$  and that of their Oracle counterparts is that the factor  $S - \hat{\theta}^{\text{LS}} (\hat{\theta}^{\text{LS}})^T$  is replaced by  $P - \theta_0 \theta_0^T$ . Moreover, the difference between the partial derivative of  $\mathcal{F}_{\text{Sg}}(P)$  and that of  $\mathcal{F}_{\text{Sy}}(P)$  is that there is one extra factor  $(\Phi^T \Phi)^{-1}$ . The difference between the first order derivative of  $\mathcal{F}_{\text{Sy}}(P)$  and that of  $\mathcal{F}_{\text{EB}}(P)$  is that there is one extra factor  $2\sigma^4 (\Phi^T \Phi)^{-1} S^{-1} = 2\sigma^4 H^{-1}$ . The above relations extend to the partial derivatives of their Oracle counterparts.

**Remark 5.** It is important to note from Propositions 3 and 4 that only the first term of  $\frac{\partial \mathcal{F}_{\text{Sg}}(P)}{\partial P}$  depends on the possibly ill-conditioned  $(\Phi^T \Phi)^{-1}$ . With the use of  $S$  in (27), all partial derivatives of the hyperparameter estimators seemingly depend on the possibly ill-conditioned term  $(\Phi^T \Phi)^{-1}$ . However, it should be stressed that the partial derivatives derived in Corollary 1 are not intended for numerical calculation but for theoretical analysis and for better exposition of the relation among the partial derivatives derived in Propositions 3 and 4.

**Remark 6.** We keep the transpose notation for symmetric matrices in Propositions 3 and 4 and Corollary 1 because they are derived by using Lemma B1 without imposing the symmetric assumption on  $P$ . After we have derived  $\frac{\partial \mathcal{C}(P)}{\partial P}$  and made the symmetric assumption on  $P$ , the first optimality condition (23) can be written as

$$\text{Tr} \left( \frac{\partial \mathcal{C}(P)}{\partial P} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0, \quad 1 \leq i \leq p$$

where the transpose notation appearing in  $\frac{\partial \mathcal{C}(P)}{\partial P}$  can be dropped for symmetric matrices, e.g.  $S^T = S$ , and  $Q^{-T} Q^{-1}$  can be written as  $Q^{-2}$ .

Setting  $\frac{\partial \text{MSEg}(P)}{\partial P} = 0$ ,  $\frac{\partial \text{MSEy}(P)}{\partial P} = 0$ , and  $\frac{\partial \text{EEB}(P)}{\partial P} = 0$  in Corollary 1 leads to the next proposition.

**Proposition 5.** The optimal kernel matrix that minimizes  $\text{MSEg}(P)$ ,  $\text{MSEy}(P)$ , and  $\text{EEB}(P)$  without structure constraints on  $P$  is

$$P = \theta_0 \theta_0^T. \quad (29)$$

It was found in Chen et al. (2012) that (29) minimizes the MSE matrix  $E(\hat{\theta}^{\text{R}} - \theta_0)(\hat{\theta}^{\text{R}} - \theta_0)^T$  in the matrix sense. Here we further find that (29) is optimal for  $\text{MSEg}(P)$ ,  $\text{MSEy}(P)$  and  $\text{EEB}(P)$ , and for any data length  $N$ .

## 4.2. Two special cases

In general, there is no explicit expression of these hyperparameter estimators. However, there exist some special cases, for which it is possible to derive the explicit solution based on Corollary 1. In the following, we consider two special cases.

### 4.2.1. Ridge regression with $\Phi^T \Phi = N I_n$

Let  $P(\eta) = \eta I_n$  with  $\eta \geq 0$  and assume  $\Phi^T \Phi = N I_n$ . Then we have the following result.

**Proposition 6.** Consider  $P(\eta) = \eta I_n$  with  $\eta \geq 0$ . Further assume that  $\Phi^T \Phi = N I_n$ . Then we have

$$\hat{\eta}_{\text{Sg}} = \hat{\eta}_{\text{Sy}} = \hat{\eta}_{\text{EB}} = \max \left( 0, \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} - \frac{\sigma^2}{N} \right). \quad (30)$$

Moreover,

$$\hat{\eta}_{\text{MSEg}} = \hat{\eta}_{\text{MSEy}} = \hat{\eta}_{\text{EEB}} = \theta_0^T \theta_0 / n. \quad (31)$$

**Remark 7.** It is worth noting that the optimal hyperparameter  $\theta_0^T \theta_0 / n$  holds for any  $N$ . Moreover, one has

$$\text{MSEg}(\theta_0^T \theta_0 / n I_n) = \frac{n \sigma^2}{N + n \sigma^2 / (\theta_0^T \theta_0)} < \frac{n \sigma^2}{N}$$

where  $n \sigma^2 / N$  is the MSEg of the LS estimator (5b). This means that the ridge regression with  $P = \theta_0^T \theta_0 / n I_n$  has a smaller MSEg than the LS estimator (5b) when  $\Phi^T \Phi = N I_n$ . Finally, (30) is a consistent estimate of  $\theta_0^T \theta_0 / n$  if  $\hat{\theta}^{\text{LS}} \rightarrow \theta_0$  as  $N \rightarrow \infty$ .

#### 4.2.2. Diagonal kernel matrix with $\Phi^T \Phi = N I_n$

Let  $P(\eta)$  be a diagonal kernel matrix (in this case we have  $p = n$ ), i.e.,

$$P(\eta) = \text{diag}[\eta_1, \dots, \eta_n] \text{ with } \eta_i \geq 0, \quad 1 \leq i \leq n \quad (32)$$

where  $\eta_1, \dots, \eta_n$  are the main diagonal elements of the diagonal matrix  $\text{diag}[\eta_1, \dots, \eta_n]$ . Then under the assumption  $\Phi^T \Phi = N I_n$ , we have the following result.

**Proposition 7.** Consider  $P(\eta)$  in (32) and assume that  $\Phi^T \Phi = N I_n$ . Then we have

$$\begin{aligned} \hat{\eta}_{\text{Sg}} = \hat{\eta}_{\text{Sy}} = \hat{\eta}_{\text{EB}} = & [\max\{0, \hat{g}_1^2 - \sigma^2/N\}, \\ & \dots, \max\{0, \hat{g}_n^2 - \sigma^2/N\}]^T \end{aligned} \quad (33)$$

where  $\hat{g}_i$  is the  $i$ th element of the LS estimate (5b),  $i = 1, \dots, n$ . Moreover,

$$\hat{\eta}_{\text{MSEg}} = \hat{\eta}_{\text{MSEy}} = \hat{\eta}_{\text{EEB}} = [(g_1^0)^2, \dots, (g_n^0)^2]^T. \quad (34)$$

**Remark 8.** In the papers Aravkin et al. (2012b, 2014), the linear model (4) but with a slightly different setting is considered, where the parameter  $\theta$  is partitioned into  $m$  sub-vectors  $\theta = [\theta^{(1)T}, \dots, \theta^{(m)T}]^T$  and the dimension of  $\theta^{(i)}$  is  $n_i$  so that  $n = \sum_{i=1}^m n_i$ . In addition, the prior distribution of  $\theta^{(i)}$  is set to be  $\mathcal{N}(0, \eta_i I_{n_i})$  and  $\eta_i$  is an independent and identically distributed exponential random variable with probability density  $p_\gamma(\eta_i) = \gamma \exp(-\gamma \eta_i) \chi(\eta_i)$  where  $\gamma$  is a positive scalar and  $\chi(t) = 1$  for  $t \geq 0$  and 0 otherwise. Under the setting given above, the solution maximizing the marginal posterior of  $\eta$  given the data and the optimal solution of the MSEg are derived in Aravkin et al. (2012b, 2014) when  $\Phi^T \Phi = N I_n$ . When  $m = 1, n_1 = n, \gamma = 0$ , their estimates become (30) and (31), respectively. When  $n = m, n_i = 1$  for  $i = 1, \dots, n$ , and  $\gamma = 0$ , their solutions become (33) and (34), respectively. In contrast, we study here the SUREg, SUREy, MSEy, and EEB estimators besides the EB and MSEg estimators and find their solutions are the same under the simplified setting. Clearly,  $\max\{0, \hat{g}_i^2 - \sigma^2/N\}$  is a consistent estimator of  $(g_i^0)^2, i = 1, \dots, n$ .

### 5. Properties of hyperparameter estimators: when the data length goes to infinity

In this section, we investigate the asymptotic properties of these hyperparameter estimators. For this purpose, it is useful to first consider the asymptotic property of the partial derivatives derived in Corollary 1. Noting the finding of Corollary 1 and that  $S - \hat{\theta}^{\text{LS}}(\hat{\theta}^{\text{LS}})^T$  converges to  $P - \theta_0 \theta_0^T$  under proper conditions, we can derive the following proposition.

**Proposition 8.** Consider the partial derivatives derived in Corollary 1. Assume that  $P$  is nonsingular and  $\Phi^T \Phi/N \rightarrow \Sigma$  almost surely as  $N \rightarrow \infty$ , where  $\Sigma$  is positive definite. Then we have as  $N \rightarrow \infty$

$$N^2 \frac{\partial \text{MSEg}(P)}{\partial P} \rightarrow 2\sigma^4 P^{-T} \Sigma^{-2} P^{-1} (P - \theta_0 \theta_0^T) P^{-T} \quad (35a)$$

$$N^2 \frac{\partial \mathcal{F}_{\text{Sg}}(P)}{\partial P} \rightarrow 2\sigma^4 P^{-T} \Sigma^{-2} P^{-1} (P - \theta_0 \theta_0^T) P^{-T} \quad (35b)$$

$$N \frac{\partial \text{MSEy}(P)}{\partial P} \rightarrow 2\sigma^4 P^{-T} \Sigma^{-1} P^{-1} (P - \theta_0 \theta_0^T) P^{-T} \quad (35c)$$

$$N \frac{\partial \mathcal{F}_{\text{Sy}}(P)}{\partial P} \rightarrow 2\sigma^4 P^{-T} \Sigma^{-1} P^{-1} (P - \theta_0 \theta_0^T) P^{-T} \quad (35d)$$

$$\frac{\partial \text{EEB}(P)}{\partial P} \rightarrow P^{-T} (P^T - \theta_0 \theta_0^T) P^{-T} \quad (35e)$$

$$\frac{\partial \mathcal{F}_{\text{EB}}(P)}{\partial P} \rightarrow P^{-T} (P^T - \theta_0 \theta_0^T) P^{-T} \quad (35f)$$

almost surely.

Proposition 8 shows that the three pairs,  $N^2 \frac{\partial \text{MSEg}(P)}{\partial P}$  and  $N^2 \frac{\partial \mathcal{F}_{\text{Sg}}(P)}{\partial P}$ ,  $N \frac{\partial \text{MSEy}(P)}{\partial P}$  and  $N \frac{\partial \mathcal{F}_{\text{Sy}}(P)}{\partial P}$ , and  $\frac{\partial \text{EEB}(P)}{\partial P}$  and  $\frac{\partial \mathcal{F}_{\text{EB}}(P)}{\partial P}$ , have respectively the same limit as  $N$  goes to  $\infty$ . This observation motivates to explore if this property also holds for the estimation criteria of these hyperparameter estimators. The answer is affirmative and we have the following result.

**Proposition 9.** Consider the hyperparameter estimation criteria SUREg (16a), SUREy (16b), and EB (15b), and their corresponding Oracle counterparts MSEg (11a), MSEy (11b), and EEB (20). Assume that  $P$  is nonsingular and  $\Phi^T \Phi/N \rightarrow \Sigma$  almost surely as  $N \rightarrow \infty$ , where  $\Sigma$  is positive definite. Then we have as  $N \rightarrow \infty$

$$N^2 (\text{MSEg}(P) - \sigma^2 \text{Tr}((\Phi^T \Phi)^{-1})) \rightarrow W_g(P, \Sigma, \theta_0) \quad (36a)$$

$$N^2 (\mathcal{F}_{\text{Sg}}(P) - \sigma^2 \text{Tr}((\Phi^T \Phi)^{-1})) \rightarrow W_g(P, \Sigma, \theta_0) \quad (36b)$$

$$N (\text{MSEy}(P) - (n + N) \sigma^2) \rightarrow W_y(P, \Sigma, \theta_0) \quad (36c)$$

$$N (\mathcal{F}_{\text{Sy}}(P) + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y - Y^T Y - 2n \sigma^2) \rightarrow W_y(P, \Sigma, \theta_0) \quad (36d)$$

$$\begin{aligned} \text{EEB}(P) - (N - n) \\ - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \rightarrow W_b(P, \theta_0) \end{aligned} \quad (36e)$$

$$\begin{aligned} \mathcal{F}_{\text{EB}}(P) + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y / \sigma^2 - Y^T Y / \sigma^2 \\ - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \rightarrow W_b(P, \theta_0) \end{aligned} \quad (36f)$$

almost surely, where

$$\begin{aligned} W_g(P, \Sigma, \theta_0) = & \sigma^4 \theta_0^T P^{-1} \Sigma^{-2} P^{-1} \theta_0 \\ & - 2\sigma^4 \text{Tr}(\Sigma^{-1} P^{-1} \Sigma^{-1}) \end{aligned} \quad (37a)$$

$$\begin{aligned} W_y(P, \Sigma, \theta_0) = & \sigma^4 \theta_0^T P^{-1} \Sigma^{-1} P^{-1} \theta_0 \\ & - 2\sigma^4 \text{Tr}(\Sigma^{-1} P^{-1}) \end{aligned} \quad (37b)$$

$$W_b(P, \theta_0) = \theta_0^T P^{-1} \theta_0 + \log \det(P). \quad (37c)$$

**Remark 9.** For these hyperparameter estimation criteria,  $W_g(P, \Sigma, \theta_0)$ ,  $W_y(P, \Sigma, \theta_0)$  and  $W_b(P, \theta_0)$  contain all information about the asymptotic benefits of regularization: how it depends on any kernel matrix  $P$ , any true impulse response vector  $\theta_0$  and any stationary properties of the input covariance matrix  $\Sigma$ .

Proposition 9 enables us to derive asymptotic properties of these hyperparameter estimators for any parameterization  $P(\eta)$  of the kernel matrix. Moreover, it also implies that the estimators  $\hat{\eta}_{\text{Sg}}$ ,  $\hat{\eta}_{\text{Sy}}$ , and  $\hat{\eta}_{\text{EB}}$  possibly share the same limits with their corresponding Oracle counterparts  $\hat{\eta}_{\text{MSEg}}$ ,  $\hat{\eta}_{\text{MSEy}}$ , and  $\hat{\eta}_{\text{EEB}}$ , respectively.

To state the result, we need an extra assumption. It is worth to note that the limit functions  $W_g(P(\eta), \Sigma, \theta_0)$ ,  $W_y(P(\eta), \Sigma, \theta_0)$  and  $W_b(P(\eta), \theta_0)$  may not have a unique global minimum, respectively. In this case, the analysis of how minimizing elements of a sequence of functions  $M_N(\eta)$  converge to the minimizing element of the limit function  $\lim M_N(\eta)$ , i.e.,

$$\lim \arg \min M_N(\eta) = \arg \min \lim M_N(\eta)' \quad (38)$$

where  $M_N(\eta)$  denotes any function on the left hand side of “ $\rightarrow$ ” in (36), follows the same idea as for prediction error identification methods, see, e.g. Lemma 8.2 and Theorem 8.2 in Ljung (1999). Accordingly, it is useful in this context to let “arg min” denote the set of minimizing arguments in case where  $W_g(P(\eta), \Sigma, \theta_0)$ ,  $W_y(P(\eta), \Sigma, \theta_0)$  and  $W_b(P(\eta), \theta_0)$  do not have a unique global minimum, respectively.:

$$\arg \min_{\eta \in \Omega} M(\eta) = \{\eta | \eta \in \Omega, M(\eta) = \min_{\eta' \in \Omega} M(\eta')\} \quad (39)$$

where  $M(\eta)$  could be any one of  $W_g(P(\eta), \Sigma, \theta_0)$ ,  $W_y(P(\eta), \Sigma, \theta_0)$  and  $W_b(P(\eta), \theta_0)$ .

Now we define

$$\eta_g^* = \arg \min_{\eta \in \Omega} W_g(P(\eta), \Sigma, \theta_0) \quad (40)$$

$$\eta_y^* = \arg \min_{\eta \in \Omega} W_y(P(\eta), \Sigma, \theta_0) \quad (41)$$

$$\eta_B^* = \arg \min_{\eta \in \Omega} W_B(P(\eta), \theta_0). \quad (42)$$

**Remark 10.** The optimal hyperparameter  $\eta_B^*$  has the following interpretation: under the assumption  $\theta_0 \sim \mathcal{N}(0, P(\eta))$ ,  $\eta_B^*$  maximizes the value of the probability density function of  $\theta_0$ .

The following assumption is also needed.

**Assumption 2.** The sets  $\eta_g^*$ ,  $\eta_y^*$  and  $\eta_B^*$  consist of interior points of  $\Omega$ , and are discrete, i.e., made up of only isolated points, respectively.

Then we have the following theorem.

**Theorem 1.** Assume that  $P(\eta)$  is any parameterization of the kernel matrix such that  $P(\eta)$  is positive definite and moreover,  $\Phi^T \Phi / N \rightarrow \Sigma$  almost surely as  $N \rightarrow \infty$ , where  $\Sigma$  is positive definite. Then we have as  $N \rightarrow \infty$

$$\hat{\eta}_{\text{MSEg}} \rightarrow \eta_g^*, \quad \hat{\eta}_{\text{Sg}} \rightarrow \eta_g^* \quad (43a)$$

$$\hat{\eta}_{\text{MSEy}} \rightarrow \eta_y^*, \quad \hat{\eta}_{\text{Sy}} \rightarrow \eta_y^* \quad (43b)$$

$$\hat{\eta}_{\text{EEB}} \rightarrow \eta_B^*, \quad \hat{\eta}_{\text{EB}} \rightarrow \eta_B^* \quad (43c)$$

almost surely. Moreover,  $\eta_g^*$ ,  $\eta_y^*$ , and  $\eta_B^*$  are roots of the system of equations, respectively,  $i = 1, \dots, p$ :

$$\text{Tr} \left( P(\eta)^{-1} \Sigma^{-2} P(\eta)^{-1} (P(\eta) - \theta_0 \theta_0^T) P(\eta)^{-1} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0$$

$$\text{Tr} \left( P(\eta)^{-1} \Sigma^{-1} P(\eta)^{-1} (P(\eta) - \theta_0 \theta_0^T) P(\eta)^{-1} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0$$

$$\text{Tr} \left( P(\eta)^{-1} (P(\eta) - \theta_0 \theta_0^T) P(\eta)^{-1} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0.$$

The Oracle estimators  $\hat{\eta}_{\text{MSEg}}$  and  $\hat{\eta}_{\text{MSEy}}$  are optimal for any data length  $N$  in the average sense if we are concerned with the ability to reproduce the true impulse response and predict the future outputs of the system respectively, while the SUREg  $\hat{\eta}_{\text{Sg}}$  and the SUREy  $\hat{\eta}_{\text{Sy}}$  are not optimal in general. Surprisingly, a nice property of  $\hat{\eta}_{\text{Sg}}$  and  $\hat{\eta}_{\text{Sy}}$  is that they converge to the best possible hyperparameter  $\eta_g^*$  and  $\eta_y^*$ , respectively, for any chosen parameterized kernel matrix  $P(\eta)$ . It is so to speak that the two SURE methods are “asymptotically consistent or asymptotically optimal”. This means that when  $N$  is sufficiently large,  $\hat{\eta}_{\text{Sg}}$  and  $\hat{\eta}_{\text{Sy}}$  perform as well as  $\hat{\eta}_{\text{MSEg}}$  and  $\hat{\eta}_{\text{MSEy}}$ , respectively. It is also worth noting that even with increasing number of data the EB estimator  $\hat{\eta}_{\text{EB}}$  has another preference than to minimize MSEg and MSEy.

**Remark 11.** In contrast with  $W_g(P, \Sigma, \theta_0)$  and  $W_y(P, \Sigma, \theta_0)$ , a unique property of  $W_B(P, \theta_0)$  is that it does not depend on the limit  $\Sigma$  of  $\Phi^T \Phi / N$ . This can to some extent explain why the EB estimator is more robust than the SUREg and SUREy especially when  $\Phi^T \Phi$  is ill-conditioned. Interested readers can find experimental evidence for this in Pilonetto and Chiuso (2015). However, in contrast with the SUREg and SUREy, the EB estimator is not asymptotically optimal in the MSEg/MSEy sense.

**Remark 12.** The different expressions of the limit functions  $W_g(P(\eta), \Sigma, \theta_0)$ ,  $W_y(P(\eta), \Sigma, \theta_0)$ , and  $W_B(P(\eta), \theta_0)$  imply that the optimal hyperparameters  $\eta_g^*$ ,  $\eta_y^*$ , and  $\eta_B^*$  may be different. To check this, we consider a special case:  $P = \eta K$ , where  $\eta > 0$  and  $K$  is

fixed and positive definite. In this case, (40), (41) and (42) become

$$\begin{aligned} \eta_g^* &= \arg \min_{\eta \geq 0} \frac{\sigma^4}{\eta^2} \theta_0^T K^{-1} \Sigma^{-2} K^{-1} \theta_0 - \frac{2\sigma^4}{\eta} \text{Tr}(\Sigma^{-1} K^{-1} \Sigma^{-1}) \\ &= \frac{\theta_0^T K^{-1} \Sigma^{-2} K^{-1} \theta_0}{\text{Tr}(\Sigma^{-1} K^{-1} \Sigma^{-1})} \end{aligned}$$

$$\begin{aligned} \eta_y^* &= \arg \min_{\eta \geq 0} \frac{\sigma^4}{\eta^2} \theta_0^T K^{-1} \Sigma^{-1} K^{-1} \theta_0 - \frac{2\sigma^4}{\eta} \text{Tr}(\Sigma^{-1} K^{-1}) \\ &= \frac{\theta_0^T K^{-1} \Sigma^{-1} K^{-1} \theta_0}{\text{Tr}(\Sigma^{-1} K^{-1})} \end{aligned}$$

$$\begin{aligned} \eta_B^* &= \arg \min_{\eta \geq 0} \theta_0^T K^{-1} \theta_0 / \eta + \log \eta^n + \log \det(K) \\ &= \theta_0^T K^{-1} \theta_0 / n \end{aligned}$$

which shows that  $\eta_g^*$ ,  $\eta_y^*$  and  $\eta_B^*$  can be different. Clearly, when  $K = I_n$  and  $\Sigma = dI_n$  with  $d > 0$ ,  $\eta_g^* = \eta_y^* = \eta_B^*$ . For this case, the optimal value  $\eta_B^*$  has been given in Theorem 7.3 of Pilonetto and Chiuso (2015).

**Corollary 2.** Assume that  $\Phi^T \Phi / N \rightarrow dI_n$  almost surely with  $d > 0$  and  $P(\eta)$  is any positive definite parameterization of the kernel matrix. Then we have

$$\eta_g^* = \eta_y^* = \arg \min_{\eta \in \Omega} \theta_0^T P(\eta)^{-2} \theta_0 - 2 \text{Tr}(P(\eta)^{-1})$$

$$\eta_B^* = \arg \min_{\eta \in \Omega} \theta_0^T P(\eta)^{-1} \theta_0 + \log \det(P(\eta))$$

and further  $\eta_g^*$  and  $\eta_B^*$  are roots of the following system of equations, respectively:

$$\text{Tr} \left( P(\eta)^{-2} (P(\eta) - \theta_0 \theta_0^T) P(\eta)^{-1} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0, \quad i = 1, \dots, p$$

$$\text{Tr} \left( P(\eta)^{-1} (P(\eta) - \theta_0 \theta_0^T) P(\eta)^{-1} \frac{\partial P(\eta)}{\partial \eta_i} \right) = 0, \quad i = 1, \dots, p.$$

In addition, for the diagonal kernel matrix (32), we have

$$\eta_g^* = \eta_y^* = \eta_B^* = [(g_1^0)^2, \dots, (g_n^0)^2]^T.$$

In Theorem 1, we have considered the convergence of those hyperparameter estimators. In fact, we can further derive their corresponding convergence rate. To this end, we let  $\xi_N = o_p(a_N)$  denote that the sequence  $\{\xi_N/a_N\}$  for nonzero sequence  $\{a_N\}$  converges in probability to zero, i.e.,  $\forall \epsilon > 0, P(|\xi_N/a_N| > \epsilon) \rightarrow 0$  as  $N \rightarrow \infty$ , while  $\xi_N = O_p(a_N)$  denote that  $\{\xi_N/a_N\}$  is bounded in probability, i.e.,  $\forall \epsilon > 0, \exists L > 0$  such that  $P(|\xi_N/a_N| > L) < \epsilon, \forall N$ . Then we have the following theorem.

**Theorem 2.** Assume that  $\|\Phi^T \Phi / N - \Sigma\| = O_p(\delta_N)$ , where  $\|\cdot\|$  denotes the Frobenius norm for a square matrix,  $\delta_N \rightarrow 0$  as  $N \rightarrow \infty$  and  $P(\eta)$  is any positive definite parameterization of the kernel matrix. Then we have

$$\|\hat{\eta}_{\text{MSEg}} - \eta_g^*\| = O_p(\varpi_N), \quad \|\hat{\eta}_{\text{Sg}} - \eta_g^*\| = O_p(\mu_N) \quad (44a)$$

$$\|\hat{\eta}_{\text{MSEy}} - \eta_y^*\| = O_p(\varpi_N), \quad \|\hat{\eta}_{\text{Sy}} - \eta_y^*\| = O_p(\mu_N) \quad (44b)$$

$$\|\hat{\eta}_{\text{EEB}} - \eta_B^*\| = O_p(1/N), \quad \|\hat{\eta}_{\text{EB}} - \eta_B^*\| = O_p(1/\sqrt{N}) \quad (44c)$$

where

$$\varpi_N = \max(O_p(\delta_N), O_p(1/N)) \quad (45a)$$

$$\mu_N = \max(O_p(\delta_N), O_p(1/\sqrt{N})). \quad (45b)$$

Theorem 2 shows that the convergence rate of  $\hat{\eta}_{\text{EEB}}$  and  $\hat{\eta}_{\text{EB}}$  to  $\eta_B^*$  depends only on the fact  $\Phi^T \Phi / N \rightarrow \Sigma$  as  $N \rightarrow \infty$  ( $\Phi^T \Phi = O_p(N)$ ) but not on the rate  $\|\Phi^T \Phi / N - \Sigma\| = O_p(\delta_N)$ . Moreover, we have

- the convergence rate of  $\hat{\eta}_{\text{EEB}}$  to  $\eta_{\text{B}}^*$  is faster than that of  $\hat{\eta}_{\text{MSEg}}$  to  $\eta_{\text{g}}^*$  and that of  $\hat{\eta}_{\text{MSEy}}$  to  $\eta_{\text{y}}^*$ .
- the convergence rate of  $\hat{\eta}_{\text{EB}}$  to  $\eta_{\text{B}}^*$  is faster than that of  $\hat{\eta}_{\text{Sg}}$  to  $\eta_{\text{g}}^*$  and that of  $\hat{\eta}_{\text{Sy}}$  to  $\eta_{\text{y}}^*$ .
- the convergence rate of  $\hat{\eta}_{\text{MSEg}}$ ,  $\hat{\eta}_{\text{MSEy}}$  and  $\hat{\eta}_{\text{EEB}}$  to  $\eta_{\text{g}}^*$ ,  $\eta_{\text{y}}^*$  and  $\eta_{\text{B}}^*$ , respectively, is faster than that of  $\hat{\eta}_{\text{Sg}}$ ,  $\hat{\eta}_{\text{Sy}}$  and  $\hat{\eta}_{\text{EB}}$  to  $\eta_{\text{g}}^*$ ,  $\eta_{\text{y}}^*$  and  $\eta_{\text{B}}^*$ , respectively.

**Theorem 2** has the following corollary.

**Corollary 3.** Assume that  $\|\Phi^T \Phi / N - \Sigma\| = O_p(\delta_N)$ , where  $\delta_N \rightarrow 0$  as  $N \rightarrow \infty$  and  $P(\eta)$  is any positive definite parameterization of the kernel matrix. Then

$$\|\hat{\eta}_{\text{MSEg}} - \hat{\eta}_{\text{Sg}}\| = O_p(\mu_N) \quad (46a)$$

$$\|\hat{\eta}_{\text{MSEy}} - \hat{\eta}_{\text{Sy}}\| = O_p(\mu_N) \quad (46b)$$

$$\|\hat{\eta}_{\text{EEB}} - \hat{\eta}_{\text{EB}}\| = O_p(1/\sqrt{N}) \quad (46c)$$

where  $\mu_N$  is defined in (45b).

This corollary shows that the convergence rate of  $\|\hat{\eta}_{\text{EEB}} - \hat{\eta}_{\text{EB}}\|$  to zero is faster than that of  $\|\hat{\eta}_{\text{MSEg}} - \hat{\eta}_{\text{Sg}}\|$  and  $\|\hat{\eta}_{\text{MSEy}} - \hat{\eta}_{\text{Sy}}\|$  to zero.

## 6. Numerical simulation

In this section, we illustrate the theoretical results with numerical simulation.

### 6.1. Test data-bank

The method in Chen et al. (2012) and Pillonetto and Chiuso (2015) is used to generate 1000 30th order test systems. Then for each test system, we consider four different test inputs:

- The first two test inputs are implemented by the MATLAB command `idinput` choosing the bandlimited white Gaussian noise with normalized bands  $[0, 0.6]$  and  $[0, 1]$ , respectively, and denoted by IT1 and IT2, respectively.
- The third and fourth test inputs are the white Gaussian noise of unit variance filtered by a second order rational transfer function  $1/(1 - aq^{-1})^2$  with  $a$  chosen to be 0.95 and 0.05, respectively, and denoted by IT3 and IT4, respectively.

To generate the data set, we simulate each system with one of the four test inputs to get the output, which is then corrupted by an additive white Gaussian noise. The signal-to-noise ratio (SNR), i.e., the ratio between the variance of the noise-free output and the noise, is uniformly distributed over  $[1, 10]$ , and is kept the same for the four test inputs.

Finally, in order to test the finite sample and asymptotic behaviour of the hyperparameter estimators, we consider data sets with different data lengths  $N = 500$  and  $8000$ , respectively.

### 6.2. Simulation setup

The performance of the RLS estimator (6b) is evaluated by the measure of fit (Ljung, 2012) defined as follows:

$$\text{Fit} = 100 \times \left( 1 - \frac{\|\hat{\theta} - \theta_0\|}{\|\theta_0 - \bar{\theta}_0\|} \right), \quad \bar{\theta}_0 = \frac{1}{n} \sum_{k=1}^n g_k^0$$

where  $n$  is set to 200. This fit is actually to evaluate the RLS estimator in the MSEg sense.

Here the unknown inputs  $u_{-1}, \dots, u_{-n+1}$  are not used (nonwindowed). The TC kernel (14c) is considered and its hyperparameter  $\eta = [c, \alpha]^T$  is estimated by using the estimators SUREg (18a),

**Table 1**

Average fits for 1000 test systems and data sets.

	MSEg	Sg	MSEy	Sy	EEB	EB
IT1						
$N = 500$	80.34	−2.4E9	78.07	53.83	77.98	77.26
$N = 8000$	90.63	−8.6E8	88.08	78.39	88.39	88.36
IT2						
$N = 500$	87.11	84.46	87.02	86.03	86.60	86.16
$N = 8000$	96.67	96.60	96.67	96.60	96.47	96.44
IT3						
$N = 500$	46.95	−2220	41.61	−146.4	39.47	39.03
$N = 8000$	57.67	−176.8	53.63	38.86	51.05	50.86
IT4						
$N = 500$	86.78	83.89	86.69	85.66	86.24	85.84
$N = 8000$	96.57	96.49	96.56	96.49	96.38	96.35

SUREy (18b), and EB (15a), respectively. For reference, we also consider their corresponding Oracle counterparts, i.e., the estimators MSEg (19a), MSEy (19b), and EEB (19c), respectively. The notations Sg, Sy, EB, MSEg, MSEy, and EEB are used to denote the corresponding simulation results, respectively.

### 6.3. Simulation results

The average fits are given in Table 1. The boxplots of the 1000 fits for IT1 and IT2 are displayed in Figs. 1–2, respectively. The boxplots for IT3 and IT4 are skipped because of their similarity with IT1 and IT2.

### 6.4. Findings

Firstly, for all tested cases and in terms of average accuracy and robustness, the Oracle estimators MSEg and MSEy (not implementable in practice) are better than Sg and Sy, respectively, while EB is just a little bit worse than but very close to its Oracle estimator EEB.

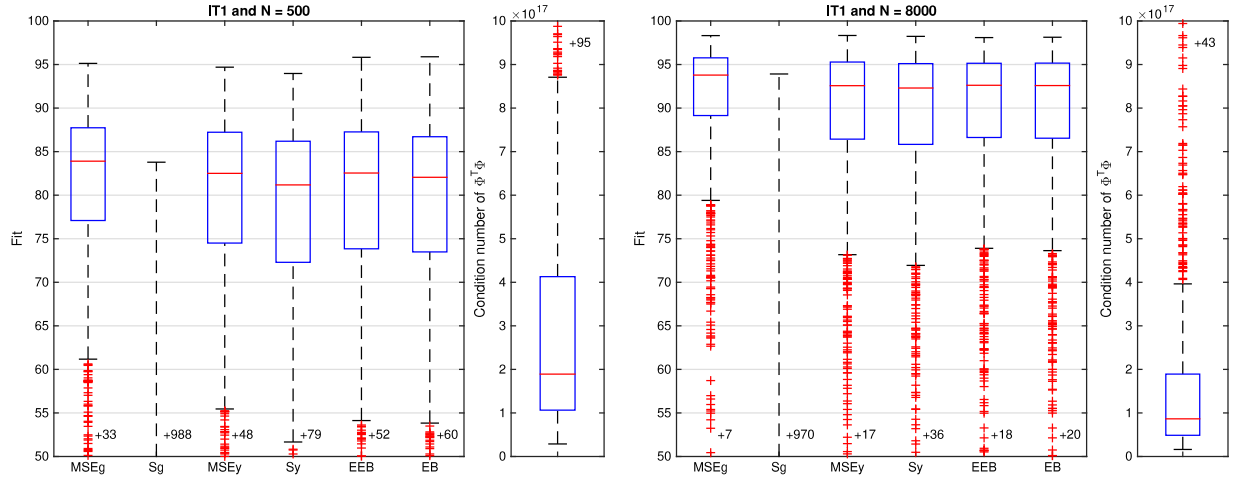
Secondly, we consider the cases with input IT1, where  $\Phi^T \Phi$  is very ill-conditioned for both  $N = 500$  and  $N = 8000$ . In this case and in terms of average accuracy and robustness, Sg performs badly because it depends on  $(\Phi^T \Phi)^{-1}$ . Moreover, Sy is better than Sg, but worse than EB.

Thirdly, we consider the case with input IT2 and  $N = 500$ , where  $\Phi^T \Phi$  is much better conditioned than the cases with input IT1. In this case and in terms of average accuracy and robustness, Sg behaves much better in contrast with the cases with input IT1. Moreover, EB and Sy are quite close though EB is a little bit better, and they are all better than Sg.

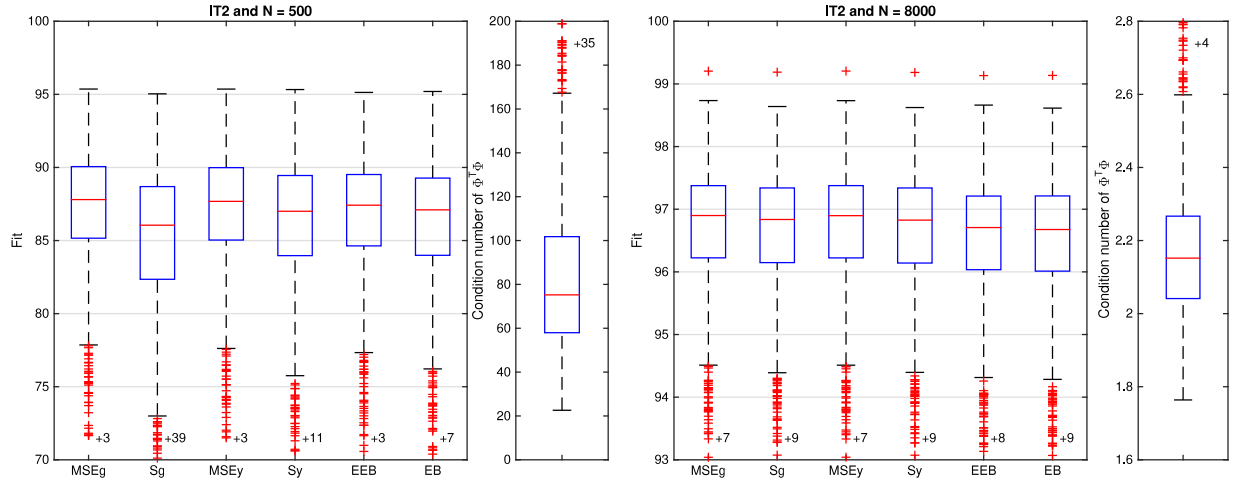
Lastly, we consider the case with input IT2 and  $N = 8000$ , where  $\Phi^T \Phi$  is very well-conditioned and in terms of average accuracy and robustness, Sg behaves much better in contrast with all the other cases, and performs as well as Sy and better than EB. Moreover, Sg and Sy are very close to the corresponding Oracle estimators MSEg and MSEy. These observations coincide with the results found in Theorem 1 and Corollary 2. Namely, Sg and Sy are asymptotically optimal but EB is not in the MSEg/MSEy senses and moreover, Sg and Sy give the same optimal hyperparameter estimate as their Oracle counterparts MSEg and MSEy, because the limit  $\Sigma = I_n$  of  $\Phi^T \Phi / N$  as  $N \rightarrow \infty$ . It can also be seen from Figs. 1 and 2 that the boxplots of EEB and EB are closer than that of MSEg and Sg and that of MSEy and Sy. This observation coincides with the result found in Corollary 3, that is, the convergence rate of  $\|\hat{\eta}_{\text{EEB}} - \hat{\eta}_{\text{EB}}\|$  to zero is faster than that of  $\|\hat{\eta}_{\text{MSEg}} - \hat{\eta}_{\text{Sg}}\|$  and  $\|\hat{\eta}_{\text{MSEy}} - \hat{\eta}_{\text{Sy}}\|$  to zero.

Based on the theoretical and simulation results, we have the following suggestions for choosing the hyperparameter estimators:





**Fig. 1.** Boxplot of the 1000 fits for the bandlimited white Gaussian noise input with the normalized band  $[0, 0.6]$  and boxplot of the condition numbers of the matrix  $\Phi^T \Phi$ : data lengths  $N = 500$  (left) and  $N = 8000$  (right).



**Fig. 2.** Boxplot of the 1000 fits for the bandlimited white Gaussian noise input with the normalized band  $[0, 1]$  and boxplot of the condition numbers of the matrix  $\Phi^T \Phi$ : data lengths  $N = 500$  (left) and  $N = 8000$  (right).

- (i) When the regression matrix is well-conditioned and the data is sufficiently long, the two SUREs should be used since they are asymptotically optimal;
- (ii) When the regression matrix is ill-conditioned or the data is short, the EB estimator should be used.

## 7. Conclusions

Kernel matrix design and hyperparameter estimation are two core issues for the kernel based regularization methods. In contrast with the former issue, there are few results reported for the latter issue. In this paper, we focused on the latter issue and studied the properties of several hyperparameter estimators including the empirical Bayes (EB) estimator, two Stein's unbiased risk estimators (SURE) and their corresponding Oracle counterparts, with an emphasis on the asymptotic properties of these hyperparameter estimators. Our major results are the following:

- The first order optimality conditions of these hyperparameter estimators are put in similar forms that better expose their relation and lead to several insights on these hyperparameter estimators.

- As the number of data goes to infinity, the two SUREs converge to the best hyperparameter minimizing the corresponding mean square error, respectively, while the more widely used EB estimator converges to another best hyperparameter minimizing the expectation of the EB estimation criterion. This indicates that the two SUREs are asymptotically optimal in the corresponding MSE sense but the EB estimator is not.
- The convergence rate of two SUREs is slower than that of the EB estimator, and moreover, unlike the two SUREs, the EB estimator is independent of the convergence rate of  $\Phi^T \Phi / N$  to its limit.

The results enhance our understanding about these hyperparameter estimators and are one step forward towards the goal of building a theory of the hyperparameter estimation for the kernel-based regularization methods.

## Acknowledgments

This work is supported in part by the National Natural Science Foundation of China under contract Nos. 61773329 and 61603379,

the Thousand Youth Talents Plan funded by the central government of China, the Shenzhen Research Projects Ji-20170189 and Ji-20160207 funded by the Shenzhen Science and Technology Innovation Council, the Presidential Fund PF. 01.000249 and the start-up grant 2014.0003.23 funded by the Chinese University of Hong Kong, Shenzhen, a research grant for junior researchers under contract No. 2014-5894 funded by Swedish Research Council, the National Key Basic Research Program of China (973 Program) under contract No. 2014CB845301, and the Presidential Fund of the Academy of Mathematics and Systems Science, CAS under contract No. 2015-hwxyqncr-mbq.

## Appendix A

Appendix A contains the proof of the results in the paper, for which the technical lemmas are placed in Appendix B. The proofs of Propositions 1, 5, 7, 8 and Corollaries 1, 2, and 3 are straightforward and thus omitted.

### A.1. Proof of Proposition 2

Under the setting  $P^{-1} = \beta A / \sigma^2$ , the MSEg (11a) of the RLS estimator (6b) is a function of  $\beta$  for a given  $A$ :

$$\text{MSEg}(\beta) = \text{Bias}(\beta) + \text{Var}(\beta) \text{ where} \quad (\text{A.1})$$

$$\begin{aligned} \text{Bias}(\beta) &= \beta^2 \theta_0^T A (\Phi^T \Phi + \beta A)^{-1} (\Phi^T \Phi + \beta A)^{-1} A \theta_0 \\ \text{Var}(\beta) &= \sigma^2 \text{Tr}((\Phi^T \Phi + \beta A)^{-1} \Phi^T \Phi (\Phi^T \Phi + \beta A)^{-1}). \end{aligned}$$

Note that  $\text{MSEg}(0) = \sigma^2 \text{Tr}((\Phi^T \Phi)^{-1})$  corresponds to the MSEg of the LS estimator (5b). The derivatives of  $\text{Bias}(\beta)$  and  $\text{Var}(\beta)$  with respect to  $\beta$  are as follows:

$$\begin{aligned} \frac{d\text{Bias}(\beta)}{d\beta} &= 2\beta \theta_0^T A (\Phi^T \Phi + \beta A)^{-1} (\Phi^T \Phi + \beta A)^{-1} A \theta_0 \\ &\quad - 2\beta^2 \theta_0^T A (\Phi^T \Phi + \beta A)^{-1} A (\Phi^T \Phi + \beta A)^{-1} \\ &\quad \times (\Phi^T \Phi + \beta A)^{-1} A \theta_0 \end{aligned} \quad (\text{A.2})$$

$$\begin{aligned} \frac{d\text{Var}(\beta)}{d\beta} &= -2\sigma^2 \text{Tr}((\Phi^T \Phi + \beta A)^{-1} A (\Phi^T \Phi + \beta A)^{-1} \\ &\quad \times \Phi^T \Phi (\Phi^T \Phi + \beta A)^{-1}) \end{aligned} \quad (\text{A.3})$$

where the formula  $\frac{dC^{-1}(\beta)}{d\beta} = -C^{-1}(\beta) \frac{dC(\beta)}{d\beta} C^{-1}(\beta)$  for an invertible matrix  $C(\beta)$  is used. Then we have

$$\begin{aligned} \left. \frac{d\text{Bias}(\beta)}{d\beta} \right|_{\beta \rightarrow 0^+} &= 0 \\ \left. \frac{d\text{Var}(\beta)}{d\beta} \right|_{\beta \rightarrow 0^+} &= -2\sigma^2 \text{Tr}((\Phi^T \Phi)^{-1} A (\Phi^T \Phi)^{-1}) < 0. \end{aligned}$$

Therefore, we have  $\left. \frac{d\text{MSEg}(\beta)}{d\beta} \right|_{\beta \rightarrow 0^+} < 0$ . This means that  $\text{MSEg}(\beta) < \text{MSEg}(0)$  in some small right neighbourhood of the origin  $\beta = 0$ .

Under the assumption that  $A$  is positive definite, denote

$$M(\beta) \triangleq E(\hat{\theta}^R - \theta_0)(\hat{\theta}^R - \theta_0)^T.$$

We first prove  $M(0) - M(\beta) > 0$  for  $0 < \beta < 2\sigma^2/(\theta_0^T A \theta_0)$ . A straightforward calculation gives

$$\begin{aligned} M(0) - M(\beta) &= \sigma^2 (\Phi^T \Phi)^{-1} - \sigma^2 (\Phi^T \Phi + \beta A)^{-1} \Phi^T \Phi (\Phi^T \Phi + \beta A)^{-1} \\ &\quad - \beta^2 (\Phi^T \Phi + \beta A)^{-1} A \theta_0 \theta_0^T A (\Phi^T \Phi + \beta A)^{-1} \\ &= \beta (\Phi^T \Phi + \beta A)^{-1} (\sigma^2 [2A + \beta A (\Phi^T \Phi)^{-1} A] - \beta A \theta_0 \theta_0^T A) \\ &\quad \times (\Phi^T \Phi + \beta A)^{-1}. \end{aligned}$$

As a result, to prove  $M(0) - M(\beta) > 0$ , it suffices to show

$$\sigma^2 [2A + \beta A (\Phi^T \Phi)^{-1} A] - \beta A \theta_0 \theta_0^T A > 0 \quad (\text{A.4})$$

which is true if  $2\sigma^2 I_n - \beta A^{1/2} \theta_0 \theta_0^T A^{1/2} > 0$  due to

$$\begin{aligned} \sigma^2 [2A + \beta A (\Phi^T \Phi)^{-1} A] - \beta A \theta_0 \theta_0^T A \\ &> 2\sigma^2 A - \beta A \theta_0 \theta_0^T A \\ &= A^{1/2} (2\sigma^2 I_n - \beta A^{1/2} \theta_0 \theta_0^T A^{1/2}) A^{1/2} > 0 \end{aligned}$$

where  $A^{1/2}$  is the symmetric and positive definite square root of  $A$  when  $A$  is positive definite. In addition, the eigenvalues of  $A^{1/2} \theta_0 \theta_0^T A^{1/2}$  are  $\theta_0^T A \theta_0$  and zero (with multiplicity  $n - 1$ ). This shows  $2\sigma^2 I_n - \beta A^{1/2} \theta_0 \theta_0^T A^{1/2} > 0$  for  $0 < \beta < 2\sigma^2/(\theta_0^T A \theta_0)$ .

Note that  $\text{MSEg}(\beta) = \text{Tr}(M(\beta))$ . One has proved that  $M(0) - M(\beta)$  is positive definite if  $0 < \beta < 2\sigma^2/(\theta_0^T A \theta_0)$ , so we have  $\text{MSEg}(0) - \text{MSEg}(\beta) = \text{Tr}(M(0) - M(\beta)) > 0$ .

The proof for the MSEy (11b) is similar to that for the MSEg (11a) by using the connection (10).

**Remark A1.** When  $\beta \rightarrow \infty$ , from the MSEg (A.1) we have

- (1)  $\text{Bias}(\beta) \rightarrow \theta_0^T \theta_0$  and  $\frac{d\text{Bias}(\beta)}{d\beta} \rightarrow 0$ ,
- (2)  $\text{Var}(\beta) \rightarrow 0$  and  $\frac{d\text{Var}(\beta)}{d\beta} \rightarrow 0$ ,
- (3)  $\text{MSEg}(\beta) \rightarrow \theta_0^T \theta_0$  and  $\frac{d\text{MSEg}(\beta)}{d\beta} \rightarrow 0$ .

### A.2. Proof of Proposition 3

To prove (24a), let us set

$$\mathcal{F}_{S_{g1}}(P) = \sigma^4 Y^T Q^{-T} \Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y$$

$$\mathcal{F}_{S_{g2}}(P) = \sigma^2 \text{Tr}(2R^{-1} - (\Phi^T \Phi)^{-1}).$$

By (B.1) and (B.4), the derivative of  $\mathcal{F}_{S_{g1}}(P)$  is

$$\begin{aligned} \frac{\partial \mathcal{F}_{S_{g1}}(P)}{\partial P} &= \sigma^4 \sum_{i,j} (2\Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y Y^T)_{ij} \frac{\partial (Q^{-1})_{ij}}{\partial P} \\ &= -2\sigma^4 \sum_{i,j} (\Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y Y^T)_{ij} \Phi^T Q^{-T} J_{ij} Q^{-T} \Phi \\ &= -2\sigma^4 \Phi^T Q^{-T} \Phi (\Phi^T \Phi)^{-2} \Phi^T Q^{-1} Y Y^T Q^{-T} \Phi \end{aligned} \quad (\text{A.5})$$

and using (B.16) implies the derivative of  $\mathcal{F}_{S_{g2}}(P)$

$$\begin{aligned} \frac{\partial \mathcal{F}_{S_{g2}}(P)}{\partial P} &= 2\sigma^2 \sum_{i=1}^n \frac{\partial (R^{-1})_{ii}}{\partial P} \\ &= 2\sigma^4 P^{-T} R^{-T} R^{-T} P^{-T} = 2\sigma^4 H^{-T} \bar{H}^{-T}. \end{aligned} \quad (\text{A.6})$$

Combining (A.5) with (A.6) derives (24a).

The proof of (24b) and (24c) is similar by using Lemma B1 and so it is omitted.

### A.3. Proof of Proposition 4

Recall that  $R = \Phi^T \Phi + \sigma^2 P^{-1}$  and  $V$  is the noise vector. It follows from (6b) that

$$\begin{aligned} \hat{\theta}^R - \theta_0 &= R^{-1} \Phi^T Y - \theta_0 \\ &= -\sigma^2 R^{-1} P^{-1} \theta_0 + R^{-1} \Phi^T V \\ &= -\sigma^2 H^{-1} \theta_0 + R^{-1} \Phi^T V \end{aligned}$$

which derives

$$\begin{aligned} \text{MSEg}(P) &= \sigma^4 \theta_0^T H^{-T} H^{-1} \theta_0 + \sigma^2 \text{Tr}(R^{-1} \Phi^T \Phi R^{-T}) \\ &= \text{MSEg1}(P) + \text{MSEg2}(P). \end{aligned}$$

For the term  $\text{MSEg1}(P)$ , using the formulae (B.1) and (B.4) gives

$$\begin{aligned} \frac{\partial \text{MSEg1}(P)}{\partial P} &= \sigma^4 \sum_{i,j} (2H^{-1}\theta_0\theta_0^T)_{ij} \frac{\partial (H^{-1})_{ij}}{\partial P} \\ &= \sigma^4 \sum_{i,j} (2H^{-1}\theta_0\theta_0^T)_{ij} (-H^{-T}J_{ij}H^{-T}\Phi^T\Phi) \\ &= -2\sigma^4 H^{-T}H^{-1}\theta_0\theta_0^T H^{-T}\Phi^T\Phi \\ &= -2\sigma^4 H^{-T}H^{-1}\theta_0\theta_0^T \Phi^T Q^{-T}\Phi. \end{aligned} \quad (\text{A.7})$$

By using the formulae (B.6) and (B.16), one derives

$$\begin{aligned} \frac{\text{MSEg2}(P)}{\partial P} &= \sigma^2 \sum_{i,j} (2R^{-1}\Phi^T\Phi)_{ij} \frac{\partial (R^{-1})_{ij}}{\partial P} \\ &= \sigma^2 \sum_{i,j} (2R^{-1}\Phi^T\Phi)_{ij} (\sigma^2 P^{-T}R^{-T}J_{ij}R^{-T}P^{-T}) \\ &= 2\sigma^4 P^{-T}R^{-T}R^{-1}\Phi^T\Phi R^{-T}P^{-T} \\ &= 2\sigma^4 H^{-T}H^{-1}P\Phi^T Q^{-T}\Phi. \end{aligned} \quad (\text{A.8})$$

Combining (A.7) with (A.8) implies the conclusion (26a).

The proof of (26b) and (26c) is similar by using Lemma B1 and so it is omitted.

#### A.4. Proof of Proposition 6

Let us first consider the EB estimator. Under the assumptions that  $P(\eta) = \eta I_n$  and  $\Phi^T\Phi = N I_n$ , we have  $S^{-1} = 1/(\eta + \sigma^2/N)I_n$ . Further, by using (28f) and (23), we obtain

$$\frac{d\mathcal{F}_{\text{EB}}(P(\eta))}{d\eta} = \frac{nN^2}{(\eta N + \sigma^2)^2} \left( \eta + \frac{\sigma^2}{N} - \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} \right).$$

By using the Lagrange multiplier, the resulting Lagrangian is

$$\mathcal{L}_{\text{EB}}(\eta, \lambda) = \mathcal{F}_{\text{EB}}(P(\eta)) - \lambda \eta$$

where  $\lambda$  is the Lagrange multiplier. Thus the corresponding Karush–Kuhn–Tucker (KKT) conditions (Boyd & Vandenberghe, 2004) are

$$\begin{aligned} \frac{nN^2}{(\eta N + \sigma^2)^2} \left( \eta + \frac{\sigma^2}{N} - \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} \right) - \lambda &= 0 \\ \lambda \eta &= 0, \quad \lambda \geq 0, \quad \eta \geq 0 \end{aligned}$$

and its solution is

$$\begin{cases} \hat{\eta} = \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} - \frac{\sigma^2}{N}, \hat{\lambda} = 0 & \text{if } \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} \geq \frac{\sigma^2}{N} \\ \hat{\eta} = 0, \hat{\lambda} = \frac{nN^2}{\sigma^4} \left( \frac{\sigma^2}{N} - \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} \right) & \text{if } \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} < \frac{\sigma^2}{N}. \end{cases}$$

Therefore, we obtain the hyperparameter  $\eta$  estimated by the EB estimator is

$$\max\left(0, \frac{(\hat{\theta}^{\text{LS}})^T \hat{\theta}^{\text{LS}}}{n} - \frac{\sigma^2}{N}\right). \quad (\text{A.9})$$

Under the same setting  $\Phi^T\Phi = N I_n$ , the KKT conditions corresponding to the SUREg and SUREy estimators derive that the arguments minimizing  $\mathcal{F}_{\text{Sg}}(P(\eta))$  and  $\mathcal{F}_{\text{Sy}}(P(\eta))$  under the constraint  $\eta \geq 0$  are also (A.9), respectively.

Likewise, by using (28a), (28c), and (28e) we have all the hyperparameters minimizing the  $\text{MSEg}(P(\eta))$ ,  $\text{MSEy}(P(\eta))$ , and  $\text{EEB}(P(\eta))$  satisfy the equation

$$\text{Tr}(\eta I_n - \theta_0 \theta_0^T) = 0$$

and its solution is  $\theta_0^T \theta_0 / n$ .

#### A.5. Proof of Proposition 9

Under the assumptions that  $\Phi^T\Phi/N \rightarrow \Sigma > 0$  and the white noise  $v(t)$ , we have  $(\Phi^T\Phi)^{-1} = O_p(1/N) \rightarrow 0$ ,  $S^{-1} \rightarrow P^{-1}$ ,  $NR^{-1} \rightarrow \Sigma^{-1}$ ,  $R^{-1}\Phi^T\Phi \rightarrow I_n$ , and  $\hat{\theta}^{\text{LS}} \rightarrow \theta_0$  almost surely as  $N \rightarrow \infty$ .

Let us first prove (36a). Using (27), we rewrite  $\text{MSEg}(P)$  in (11a) as follows:

$$\begin{aligned} \text{MSEg}(P) &= \sigma^4 \theta_0^T S^{-1} (\Phi^T\Phi)^{-2} S^{-1} \theta_0 \\ &\quad + \sigma^2 \text{Tr}(R^{-1}\Phi^T\Phi R^{-1}). \end{aligned}$$

The limit

$$\begin{aligned} &N^2(R^{-1}\Phi^T\Phi R^{-1} - (\Phi^T\Phi)^{-1}) \\ &= -\sigma^2 N^2 R^{-1} (2P^{-1} + \sigma^2 P^{-1} (\Phi^T\Phi)^{-1} P^{-1}) R^{-1} \\ &\rightarrow -2\sigma^2 \Sigma^{-1} P^{-1} \Sigma^{-1} \end{aligned} \quad (\text{A.10})$$

yields that

$$\begin{aligned} &N^2(\text{MSEg}(P) - \sigma^2 \text{Tr}((\Phi^T\Phi)^{-1})) \\ &= \sigma^4 \theta_0^T S^{-1} (N^2 (\Phi^T\Phi)^{-2}) S^{-1} \theta_0 \\ &\quad + \sigma^2 N^2 \text{Tr}(R^{-1}\Phi^T\Phi R^{-1} - (\Phi^T\Phi)^{-1}) \\ &\rightarrow \sigma^4 \theta_0^T P^{-1} \Sigma^{-2} P^{-1} \theta_0 - 2\sigma^4 \text{Tr}(\Sigma^{-1} P^{-1} \Sigma^{-1}) \\ &= W_g(P, \Sigma, \theta_0). \end{aligned} \quad (\text{A.11})$$

To prove (36b), note that the first term of  $\mathcal{F}_{\text{Sg}}(P)$  can be rewritten as  $\sigma^4 (\hat{\theta}^{\text{LS}})^T S^{-1} (\Phi^T\Phi)^{-2} S^{-1} \hat{\theta}^{\text{LS}}$ . Thus one derives

$$\begin{aligned} &N^2(\mathcal{F}_{\text{Sg}}(P) - \sigma^2 \text{Tr}((\Phi^T\Phi)^{-1})) \\ &= \sigma^4 (\hat{\theta}^{\text{LS}})^T S^{-1} N^2 (\Phi^T\Phi)^{-2} S^{-1} \hat{\theta}^{\text{LS}} \\ &\quad + 2\sigma^2 N^2 \text{Tr}(R^{-1} - (\Phi^T\Phi)^{-1}) \\ &\rightarrow W_g(P, \Sigma, \theta_0) \end{aligned} \quad (\text{A.12})$$

where we use the limit

$$\begin{aligned} N^2(R^{-1} - (\Phi^T\Phi)^{-1}) &= -\sigma^2 NR^{-1} P^{-1} N (\Phi^T\Phi)^{-1} \\ &\rightarrow -\sigma^2 \Sigma^{-1} P^{-1} \Sigma^{-1}. \end{aligned}$$

Similarly, we can rewrite  $\text{MSEy}(P)$  as

$$\begin{aligned} \text{MSEy}(P) &= \sigma^4 \theta_0^T S^{-1} (\Phi^T\Phi)^{-1} S^{-1} \theta_0 + N \sigma^2 \\ &\quad + \text{Tr}(R^{-1}\Phi^T\Phi R^{-1}\Phi^T\Phi) \end{aligned} \quad (\text{A.13})$$

and hence the assertion (36c) is proved by

$$\begin{aligned} &N(\text{MSEy}(P) - (n+N)\sigma^2) \\ &= \sigma^4 \theta_0^T S^{-1} N (\Phi^T\Phi)^{-1} S^{-1} \theta_0 \\ &\quad + \sigma^2 N \text{Tr}(R^{-1}\Phi^T\Phi R^{-1}\Phi^T\Phi - I_n) \\ &\rightarrow W_y(P, \Sigma, \theta_0) \end{aligned} \quad (\text{A.14})$$

where we use the formulae

$$\begin{aligned} &N(R^{-1}\Phi^T\Phi R^{-1}\Phi^T\Phi - I_n) \\ &= -\sigma^2 NR^{-1} [2P^{-1} + \sigma^2 P^{-1} (\Phi^T\Phi)^{-1} P^{-1}] R^{-1} \Phi^T\Phi \\ &\rightarrow -2\sigma^2 \Sigma^{-1} P^{-1}. \end{aligned}$$

To prove (36d), we need some identities. A straightforward calculation shows that

$$Q(I_N - \Phi(\Phi^T\Phi)^{-1}\Phi^T)Q = \sigma^4(I_N - \Phi(\Phi^T\Phi)^{-1}\Phi^T).$$

This means that

$$\sigma^4 Q^{-1}(I_N - \Phi(\Phi^T\Phi)^{-1}\Phi^T)Q^{-1} = I_N - \Phi(\Phi^T\Phi)^{-1}\Phi^T$$

and hence we derive

$$\begin{aligned} & \sigma^4 Y^T Q^{-2} Y + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y - Y^T Y \\ &= \sigma^4 Y^T Q^{-1} \Phi (\Phi^T \Phi)^{-1} \Phi^T Q^{-1} Y. \end{aligned}$$

It follows from (B.11) and (B.14) that

$$\begin{aligned} & N[\mathcal{F}_{\text{Sy}}(P) + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y - Y^T Y - 2n\sigma^2] \\ &= N[\sigma^4 Y^T Q^{-1} \Phi (\Phi^T \Phi)^{-1} \Phi^T Q^{-1} Y + 2\sigma^2 \text{Tr}(R^{-1} \Phi^T \Phi - I_n)] \\ &= N[\sigma^4 (\hat{\theta}^{\text{LS}})^T S^{-1} (\Phi^T \Phi)^{-1} S^{-1} \hat{\theta}^{\text{LS}} + 2\sigma^2 \text{Tr}(R^{-1} \Phi^T \Phi - I_n)] \\ &\rightarrow W_y(P, \Sigma, \theta_0) \end{aligned} \quad (\text{A.15})$$

where we use the limit

$$N(R^{-1} \Phi^T \Phi - I_n) = -\sigma^2 N R^{-1} P^{-1} \rightarrow -\sigma^2 \Sigma^{-1} P^{-1}.$$

Similarly, we need two identities to prove (36e). Using the Sylvester's determinant identity  $\det(I_n + AB) = \det(I_n + BA)$  derives

$$\det(Q) = \sigma^{2(N-n)} \det(\Phi^T \Phi) \det(P + \sigma^2 (\Phi^T \Phi)^{-1})$$

which implies

$$\begin{aligned} & \log \det(Q) - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \\ &= \log \det(S) \rightarrow \log \det(P). \end{aligned} \quad (\text{A.16})$$

Starting with the identity  $I_N = \sigma^2 Q^{-1} + \Phi P \Phi^T Q^{-1}$  gives

$$\begin{aligned} & \sigma^2 \text{Tr}(Q^{-1}) = N - \text{Tr}(\Phi P \Phi^T Q^{-1}) \\ &= N - \text{Tr}(R^{-1} \Phi^T \Phi) \rightarrow N - n. \end{aligned}$$

Therefore, the limit (36e) is proved by

$$\begin{aligned} & \text{EEB}(P) - (N - n) - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \\ &= \theta_0^T S^{-1} \theta_0 + (\sigma^2 \text{Tr}(Q^{-1}) - (N - n)) \\ &+ \log \det(Q) - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \\ &\rightarrow \theta_0^T P^{-1} \theta_0 + \log \det(P) = W_B(P, \theta_0). \end{aligned} \quad (\text{A.17})$$

At last, we finish the proof by checking (36f). The identity

$$Q(I_N - \Phi (\Phi^T \Phi)^{-1} \Phi^T) / \sigma^2 = I_N - \Phi (\Phi^T \Phi)^{-1} \Phi^T$$

implies that

$$\begin{aligned} & Y^T Q^{-1} Y + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y / \sigma^2 - Y^T Y / \sigma^2 \\ &= Y^T Q^{-1} \Phi (\Phi^T \Phi)^{-1} \Phi^T Y. \end{aligned} \quad (\text{A.18})$$

It follows from (A.16), (A.18), and (B.11) that

$$\begin{aligned} & \mathcal{F}_{\text{EB}}(P) + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y / \sigma^2 - Y^T Y / \sigma^2 \\ &- (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \\ &= Y^T Q^{-1} Y + Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T Y / \sigma^2 - Y^T Y / \sigma^2 \\ &+ \log \det(Q) - (N - n) \log \sigma^2 - \log \det(\Phi^T \Phi) \\ &= Y^T Q^{-1} \Phi (\Phi^T \Phi)^{-1} \Phi^T Y + \log \det(S) \\ &\rightarrow W_B(P, \theta_0). \end{aligned} \quad (\text{A.19})$$

#### A.6. Proof of Theorem 1

Firstly, we prove  $\hat{\eta}_{\text{MSEg}} \rightarrow \eta_g^*$  as  $N \rightarrow \infty$ . Define

$$\overline{\text{MSEg}}(P) \triangleq N^2 (\text{MSEg}(P) - \sigma^2 \text{Tr}((\Phi^T \Phi)^{-1})). \quad (\text{A.20})$$

Clearly, we have  $\hat{\eta}_{\text{MSEg}}$  also minimizes  $\overline{\text{MSEg}}(P(\eta))$ , i.e.,

$$\hat{\eta}_{\text{MSEg}} = \underset{\eta \in \Omega}{\text{argmin}} \overline{\text{MSEg}}(P(\eta)).$$

Under Assumption 2, there exists a compact set

$$\overline{\Omega} \subset \Omega \quad (\text{A.21})$$

containing  $\eta_g^*$  such that  $0 < d_1 \leq \|P(\eta)\| \leq d_2 < \infty$  for all  $\eta \in \overline{\Omega}$ . Then by Lemma B3 in Appendix B, to prove  $\hat{\eta}_{\text{MSEg}} \rightarrow \eta_g^*$  as  $N \rightarrow \infty$ , it suffices to show that  $\overline{\text{MSEg}}(P(\eta))$  converges to  $W_g(P(\eta), \Sigma, \theta_0)$  almost surely and uniformly in  $\overline{\Omega}$ , as  $N \rightarrow \infty$ .

It follows from (A.11) and (A.10) that

$$\begin{aligned} & \overline{\text{MSEg}}(P(\eta)) - W_g(P, \Sigma, \theta_0) \\ &= \sigma^4 Z_1 + 2\sigma^4 \text{Tr}(Z_2) - \sigma^6 \text{Tr}(Z_3) \end{aligned} \quad (\text{A.22})$$

$$\begin{aligned} Z_1 &= \theta_0^T S^{-1} (N^2 (\Phi^T \Phi)^{-2}) S^{-1} \theta_0 - \theta_0^T P^{-1} \Sigma^{-2} P^{-1} \theta_0 \\ Z_2 &= \Sigma^{-1} P^{-1} \Sigma^{-1} - N^2 R^{-1} P^{-1} R^{-1} \\ Z_3 &= -N^2 R^{-1} P^{-1} (\Phi^T \Phi)^{-1} P^{-1} R^{-1}. \end{aligned} \quad (\text{A.23})$$

For the term  $Z_1$ , we have

$$\begin{aligned} Z_1 &= \theta_0^T (S^{-1} - P^{-1}) (N^2 (\Phi^T \Phi)^{-2}) S^{-1} \theta_0 \\ &+ \theta_0^T P^{-1} (N^2 (\Phi^T \Phi)^{-2} - \Sigma^{-2}) S^{-1} \theta_0 \\ &+ \theta_0^T P^{-1} \Sigma^{-2} (S^{-1} - P^{-1}) \theta_0 \end{aligned} \quad (\text{A.24})$$

where

$$S^{-1} - P^{-1} = -\sigma^2 S^{-1} (\Phi^T \Phi)^{-1} P^{-1}. \quad (\text{A.25})$$

Note that  $\Phi^T \Phi / N \rightarrow \Sigma$  implies that  $\|N(\Phi^T \Phi)^{-1}\| = O_p(1)$ . Then further noting that  $d_1 \leq \|P(\eta)\| \leq d_2$  and  $\|S(\eta)^{-1}\| < \|(P(\eta))^{-1}\| \leq 1/d_1$  for  $\eta \in \overline{\Omega}$ , we have  $Z_1$  converges to zero almost surely and uniformly in  $\overline{\Omega}$ .

For the term  $Z_2$ , we have

$$\begin{aligned} & \Sigma^{-1} P^{-1} \Sigma^{-1} - N^2 R^{-1} P^{-1} R^{-1} \\ &= (\Sigma^{-1} - N R^{-1}) P^{-1} \Sigma^{-1} + N R^{-1} P^{-1} (\Sigma^{-1} - N R^{-1}). \end{aligned}$$

The assertions  $N R^{-1} \rightarrow \Sigma^{-1}$  and  $\|N R^{-1} - \Sigma^{-1}\| = O_p(1)$  yield that  $Z_2$  converges to zero almost surely and uniformly in  $\overline{\Omega}$ . Finally, by noting  $(\Phi^T \Phi)^{-1} \rightarrow 0$  as  $N \rightarrow \infty$  it is easy to see that  $Z_3$  also converges to zero almost surely and uniformly. Making use of these facts shows that  $\overline{\text{MSEg}}(P(\eta))$  converges to  $W_g(P(\eta), \Sigma, \theta_0)$  almost surely and uniformly in  $\overline{\Omega}$  and hence, by Lemma B3,  $\hat{\eta}_{\text{MSEg}} \rightarrow \eta_g^*$  as  $N \rightarrow \infty$  almost surely.

Secondly, we prove that  $\hat{\eta}_{\text{Sg}} \rightarrow \eta_g^*$  as  $N \rightarrow \infty$  and the proof is similar to that of  $\hat{\eta}_{\text{MSEg}} \rightarrow \eta_g^*$  as  $N \rightarrow \infty$ . Define

$$\overline{\mathcal{F}}_{\text{Sg}}(P(\eta)) \triangleq N^2 (\mathcal{F}_{\text{Sg}}(P(\eta)) - \sigma^2 \text{Tr}((\Phi^T \Phi)^{-1})).$$

Then, we have

$$\hat{\eta}_{\text{Sg}} = \underset{\eta \in \Omega}{\text{argmin}} \overline{\mathcal{F}}_{\text{Sg}}(P(\eta)). \quad (\text{A.26})$$

It follows from (A.12) that

$$\begin{aligned} & \overline{\mathcal{F}}_{\text{Sg}}(P(\eta)) - W_g(P, \Sigma, \theta_0) = \sigma^4 Z'_1 + 2\sigma^4 \text{Tr}(Z'_2) \\ Z'_1 &= (\hat{\theta}^{\text{LS}})^T S^{-1} N^2 (\Phi^T \Phi)^{-2} S^{-1} \hat{\theta}^{\text{LS}} - \theta_0^T P^{-1} \Sigma^{-2} P^{-1} \theta_0 \\ Z'_2 &= \Sigma^{-1} P^{-1} \Sigma^{-1} - N R^{-1} P^{-1} N (\Phi^T \Phi)^{-1}. \end{aligned}$$

For the terms  $Z'_1$  and  $Z'_2$ , we have

$$\begin{aligned} Z'_1 &= (\hat{\theta}^{\text{LS}} - \theta_0)^T S^{-1} N^2 (\Phi^T \Phi)^{-2} S^{-1} \hat{\theta}^{\text{LS}} \\ &+ \theta_0^T (S^{-1} - P^{-1}) N^2 (\Phi^T \Phi)^{-2} S^{-1} \hat{\theta}^{\text{LS}} \\ &+ \theta_0^T P^{-1} (N^2 (\Phi^T \Phi)^{-2} - \Sigma^{-2}) S^{-1} \hat{\theta}^{\text{LS}} \\ &+ \theta_0^T P^{-1} \Sigma^{-2} (S^{-1} - P^{-1}) \hat{\theta}^{\text{LS}} \\ &+ \theta_0^T P^{-1} \Sigma^{-2} P^{-1} (\hat{\theta}^{\text{LS}} - \theta_0) \end{aligned} \quad (\text{A.27})$$

$$\begin{aligned} Z'_2 &= (\Sigma^{-1} - N R^{-1}) P^{-1} \Sigma^{-1} \\ &+ N R^{-1} P^{-1} (\Sigma^{-1} - N (\Phi^T \Phi)^{-1}). \end{aligned} \quad (\text{A.28})$$



Then, noting that  $\hat{\theta}^{LS} \rightarrow \theta_0$ ,  $S^{-1} \rightarrow P^{-1}$ ,  $N(\Phi^T \Phi)^{-1} \rightarrow \Sigma^{-1}$ ,  $NR^{-1} \rightarrow \Sigma^{-1}$  almost surely as  $N \rightarrow \infty$ , and  $\|NR^{-1}\| = O_p(1)$ ,  $\|\hat{\theta}^{LS}\| = O_p(1)$ , and  $d_1 \leq \|P(\eta)\| \leq d_2$ ,  $\|S(\eta)^{-1}\| < \|(P(\eta))^{-1}\| \leq 1/d_1$ , for  $\eta \in \bar{\Omega}$ , one can show that each term of (A.27) and (A.28), and thus both  $Z'_1$  and  $Z'_2$  converge to zero almost surely and uniformly in  $\bar{\Omega}$ . Therefore,  $\mathcal{F}_{sg}(P(\eta))$  converges to  $W_g(P, \Sigma, \theta_0)$  almost surely and uniformly in  $\bar{\Omega}$ . It then follows from Lemma B3 that  $\hat{\eta}_{sg} \rightarrow \eta_g^*$  almost surely as  $N \rightarrow \infty$ .

The proof of (43b) and (43c) can be done similarly and thus is omitted. The first order optimality conditions of  $\eta_g^*$ ,  $\eta_y^*$ , and  $\eta_b^*$  can be derived in a similar way as Proposition 4 and thus is omitted. This completes the proof.

#### A.7. Proof of Theorem 2

We first prove that  $\|\hat{\eta}_{MSEg} - \eta_g^*\| = O_p(\varpi_N)$ .

Noting (A.11), the  $i$ th elements of the gradient vectors of  $\overline{MSEg}(P(\eta))$  and  $W_g(P(\eta), \Sigma, \theta_0)$  with respect to  $\eta$  are, respectively, for  $1 \leq i \leq p$ ,

$$\begin{aligned} \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta_i} &= 2\sigma^4 N^2 \theta_0^T S^{-1} (\Phi^T \Phi)^{-2} \frac{\partial S^{-1}}{\partial \eta_i} \theta_0 \\ &\quad + 2\sigma^2 N^2 \text{Tr} \left( \frac{\partial R^{-1}}{\partial \eta_i} \Phi^T \Phi R^{-1} \right) \\ \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta_i} &= 2\sigma^4 \theta_0^T P^{-1} \Sigma^{-2} \frac{\partial P^{-1}}{\partial \eta_i} \theta_0 \\ &\quad - 2\sigma^4 \text{Tr} \left( \Sigma^{-1} \frac{\partial P^{-1}}{\partial \eta_i} \Sigma^{-1} \right). \end{aligned} \quad (\text{A.29})$$

Using the identity  $\frac{\partial R^{-1}}{\partial \eta_i} = -R^{-1} \frac{\partial R}{\partial \eta_i} R^{-1} = -\sigma^2 R^{-1} \frac{\partial P^{-1}}{\partial \eta_i} R^{-1}$ , we see their difference is

$$\frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta_i} - \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta_i} = 2\sigma^4 (\Upsilon_1 + \text{Tr}(\Upsilon_2))$$

$$\begin{aligned} \text{where } \Upsilon_1 &= \theta_0^T S^{-1} (N^2 (\Phi^T \Phi)^{-2}) \frac{\partial S^{-1}}{\partial \eta_i} \theta_0 \\ &\quad - \theta_0^T P^{-1} \Sigma^{-2} \frac{\partial P^{-1}}{\partial \eta_i} \theta_0, \\ \Upsilon_2 &= \Sigma^{-1} \frac{\partial P^{-1}}{\partial \eta_i} \Sigma^{-1} - NR^{-1} \frac{\partial P^{-1}}{\partial \eta_i} R^{-1} \Phi^T \Phi NR^{-1}. \end{aligned}$$

Noting  $\|N(\Phi^T \Phi)^{-1} - \Sigma^{-1}\| = O_p(\delta_N)$ ,  $\|S^{-1} - P^{-1}\| = O_p(1/N)$ ,  $\|\frac{\partial S^{-1}}{\partial \eta_i} - \frac{\partial P^{-1}}{\partial \eta_i}\| = O_p(1/N)$ ,  $\|R^{-1} \Phi^T \Phi - I_n\| = O_p(1/N)$ ,  $\|NR^{-1} - \Sigma^{-1}\| = O_p(\delta_N)$ , and  $d_1 \leq \|P(\eta)\| \leq d_2$  and  $\|S(\eta)^{-1}\| < \|(P(\eta))^{-1}\| \leq 1/d_1$  for  $\eta \in \bar{\Omega}$  yields

$$|\Upsilon_1| = O_p(\varpi_N), \quad |\text{Tr}(\Upsilon_2)| = O_p(\varpi_N) \quad (\text{A.30})$$

uniformly in  $\bar{\Omega}$ , where  $\bar{\Omega}$  is defined in (A.21). Therefore,

$$\left\| \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} - \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta} \right\| = O_p(\varpi_N)$$

uniformly for any  $\eta \in \bar{\Omega}$ . Since  $\hat{\eta}_{MSEg}$  and  $\eta_g^*$  minimize  $\overline{MSEg}(P)$  and  $W_g(P, \Sigma, \theta_0)$ , respectively, we have

$$\left. \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} \right|_{\eta=\hat{\eta}_{MSEg}} = 0 \text{ and } \left. \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta} \right|_{\eta=\eta_g^*} = 0.$$

It follows that

$$\left. \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} \right|_{\eta=\eta_g^*} = O_p(\varpi_N).$$

In addition, by using (A.29), the  $(i, j)$ -element of the Hessian matrix of  $W_g(P(\eta), \Sigma, \theta_0)$  is

$$\begin{aligned} &\frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta_i \partial \eta_j} \\ &= 2\sigma^4 \theta_0^T P^{-1} \Sigma^{-2} \frac{\partial^2 P^{-1}}{\partial \eta_i \partial \eta_j} \theta_0 + 2\sigma^4 \theta_0^T \frac{\partial P^{-1}}{\partial \eta_j} \Sigma^{-2} \frac{\partial P^{-1}}{\partial \eta_i} \theta_0 \\ &\quad - 2\sigma^4 \text{Tr} \left( \Sigma^{-1} \frac{\partial^2 P^{-1}}{\partial \eta_i \partial \eta_j} \Sigma^{-1} \right). \end{aligned} \quad (\text{A.31})$$

The Hessian matrix  $\frac{\partial^2 \overline{MSEg}(P(\eta))}{\partial \eta \partial \eta^T}$  of  $\overline{MSEg}(P(\eta))$  is omitted here for simplicity. Then, it can be shown that

$$\left\| \frac{\partial^2 \overline{MSEg}(P(\eta))}{\partial \eta \partial \eta^T} - \frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta \partial \eta^T} \right\| = o_p(1)$$

uniformly for any  $\eta \in \bar{\Omega}$ . Applying the Taylor expansion to  $\frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta}$  yields

$$\begin{aligned} 0 &= \left. \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} \right|_{\eta=\hat{\eta}_{MSEg}} = \left. \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} \right|_{\eta=\eta_g^*} \\ &\quad + \left. \frac{\partial^2 \overline{MSEg}(P(\eta))}{\partial \eta \partial \eta^T} \right|_{\eta=\bar{\eta}} (\hat{\eta}_{MSEg} - \eta_g^*) \end{aligned}$$

where  $\bar{\eta}$  lies between  $\hat{\eta}_{MSEg}$  and  $\eta_g^*$ .

Clearly,

$$\left. \frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta \partial \eta^T} \right|_{\eta=\eta_g^*} = O_p(1).$$

Then under Assumption 2, we have  $\left. \frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta \partial \eta^T} \right|_{\eta=\eta_g^*}$  is positive definite. For sufficiently large  $N$ ,  $\bar{\eta}$  would be close to  $\eta_g^*$ . In this case, we also have  $\left. \frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta \partial \eta^T} \right|_{\eta=\bar{\eta}}$  is positive definite. Then it follows that

$$\begin{aligned} &\hat{\eta}_{MSEg} - \eta_g^* \\ &= - \left( \left. \frac{\partial^2 \overline{MSEg}(P(\eta))}{\partial \eta \partial \eta^T} \right|_{\eta=\bar{\eta}} \right)^{-1} \left. \frac{\partial \overline{MSEg}(P(\eta))}{\partial \eta} \right|_{\eta=\eta_g^*} \\ &= O_p(1) O_p(\varpi_N) = O_p(\varpi_N). \end{aligned}$$

Now, we prove  $\|\hat{\eta}_{sg} - \eta_g^*\| = O_p(\mu_N)$  and the proof is similar to that of  $\|\hat{\eta}_{MSEg} - \eta_g^*\| = O_p(\varpi_N)$ . By (A.12), the  $i$ th element of gradient vector of  $\mathcal{F}_{sg}(P(\eta))$  is

$$\begin{aligned} \frac{\partial \mathcal{F}_{sg}(P(\eta))}{\partial \eta_i} &= 2\sigma^4 (\hat{\theta}^{LS})^T S^{-1} N^2 (\Phi^T \Phi)^{-2} \frac{\partial S^{-1}}{\partial \eta_i} \hat{\theta}^{LS} \\ &\quad + 2\sigma^2 N^2 \text{Tr} \left( \frac{\partial R^{-1}}{\partial \eta_i} \right). \end{aligned} \quad (\text{A.32})$$

Using the identity  $\frac{\partial R^{-1}}{\partial \eta_i} = -\sigma^2 R^{-1} \frac{\partial P^{-1}}{\partial \eta_i} R^{-1}$ , we see

$$\frac{\partial \mathcal{F}_{sg}(P(\eta))}{\partial \eta_i} - \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta_i} = 2\sigma^4 \Upsilon'_1 + 2\sigma^4 \text{Tr}(\Upsilon'_2)$$

$$\begin{aligned} \text{where } \Upsilon'_1 &= (\hat{\theta}^{LS})^T S^{-1} N^2 (\Phi^T \Phi)^{-2} \frac{\partial S^{-1}}{\partial \eta_i} \hat{\theta}^{LS} \\ &\quad - \theta_0^T P^{-1} \Sigma^{-2} \frac{\partial P^{-1}}{\partial \eta_i} \theta_0 \end{aligned} \quad (\text{A.33})$$

$$\Upsilon'_2 = \Sigma^{-1} \frac{\partial P^{-1}}{\partial \eta_i} \Sigma^{-1} - NR^{-1} \frac{\partial P^{-1}}{\partial \eta_i} NR^{-1}.$$

Since  $\Phi^T \Phi / N \rightarrow \Sigma$  and  $v(t)$  is a white noise, we have  $\|\hat{\theta}^{LS} - \theta_0\| = O_p(1/\sqrt{N})$ . Then noting that  $\|N(\Phi^T \Phi)^{-1} - \Sigma^{-1}\| = O_p(\delta_N)$ ,  $\|S^{-1} - P^{-1}\| = O_p(1/N)$ ,  $\|\frac{\partial S^{-1}}{\partial \eta_i} - \frac{\partial P^{-1}}{\partial \eta_i}\| = O_p(1/N)$ ,  $\|NR^{-1} - \Sigma^{-1}\| = O_p(\delta_N)$ , and  $\|NR^{-1}\| = O_p(1)$ ,  $\|\hat{\theta}^{LS}\| = O_p(1)$ , and  $d_1 \leq \|P(\eta)\| \leq d_2$  and  $\|(P(\eta))^{-1}\| < \|(P(\eta))^{-1}\| \leq 1/d_1$  for  $\eta \in \bar{\Omega}$ , yields

$$|\gamma'_1| = \max(O_p(1/\sqrt{N}), O_p(1/N), O_p(\delta_N)) = O_p(\mu_N)$$

$$|\text{Tr}(\gamma'_2)| = O_p(\delta_N)$$

uniformly in  $\bar{\Omega}$ . It follows that

$$\left\| \frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta} - \frac{\partial W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta} \right\| = O_p(\mu_N)$$

uniformly for any  $\eta \in \bar{\Omega}$ . This implies

$$\frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta} \Big|_{\eta=\eta_g^*} = O_p(\mu_N). \quad (\text{A.34})$$

Similarly, one can obtain the Hessian matrix of  $\bar{\mathcal{F}}_{\text{sg}}(P(\eta))$  and can show that

$$\left\| \frac{\partial^2 \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta \partial \eta^T} - \frac{\partial^2 W_g(P(\eta), \Sigma, \theta_0)}{\partial \eta \partial \eta^T} \right\| = o_p(1) \quad (\text{A.35})$$

uniformly for any  $\eta \in \bar{\Omega}$ . Applying the Taylor expansion of  $\frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta}$  shows

$$0 = \frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta} \Big|_{\eta=\hat{\eta}_{\text{sg}}} = \frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta} \Big|_{\eta=\eta_g^*} + \frac{\partial^2 \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta \partial \eta^T} \Big|_{\eta=\tilde{\eta}} (\hat{\eta}_{\text{MSEg}} - \eta_g^*)$$

where  $\tilde{\eta}$  lies between  $\hat{\eta}_{\text{sg}}$  and  $\eta_g^*$ . For sufficiently large  $N$ , we have

$$\begin{aligned} & \hat{\eta}_{\text{sg}} - \eta_g^* \\ &= - \left( \frac{\partial^2 \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta \partial \eta^T} \Big|_{\eta=\tilde{\eta}} \right)^{-1} \frac{\partial \bar{\mathcal{F}}_{\text{sg}}(P(\eta))}{\partial \eta} \Big|_{\eta=\eta_g^*} \\ &= O_p(1) O_p(\mu_N) = O_p(\mu_N). \end{aligned}$$

The proof of (44b) and (44c) can be done in a similar way and thus is omitted. This completes the proof.

## Appendix B

This appendix contains the technical lemmas used in the proof in Appendix A.

### B.1. Matrix differentials and related identities

This section introduces the differentiation of a function  $f(X)$  where  $X$  is a matrix. It is assumed that  $X$  has no special structure, i.e., that the elements of  $X$  are independent. For convenience and readability, the formulae used in the paper are stated in the following lemmas.

**Lemma B1** (Petersen & Pedersen, 2012). Assume that  $b$  is a column vector, and  $A$ ,  $B$  and  $X$  are matrices with compatible dimensions. Then we have

$$\frac{\partial b^T X^T A X b}{\partial X} = (A + A^T) X b b^T \quad (\text{B.1})$$

$$\frac{\partial b^T X^{-1} b}{\partial X} = -X^{-T} b b^T X^{-T} \quad (\text{B.2})$$

$$\frac{\partial \log|\det(X)|}{\partial X} = X^{-T} \quad (\text{B.3})$$

$$\frac{\partial (X^{-1})_{kl}}{\partial X_{ij}} = -(X^{-1})_{ki} (X^{-1})_{jl} \quad (\text{B.4})$$

$$\frac{\partial \text{Tr}(A X^{-1} B)}{\partial X} = -(X^{-1} B A X^{-1})^T \quad (\text{B.5})$$

$$\frac{\partial \text{Tr}(A X B X^T A^T)}{\partial X} = A^T A X (B + B^T) \quad (\text{B.6})$$

where  $(\cdot)_{ij}$  denotes the  $(i, j)$ th element of a matrix.

**Lemma B2.** We have the following identities:

$$\sum_{ij} (A)_{ij} J_{ij} = A \quad (\text{B.7})$$

$$Y - \Phi \hat{\theta}^R = \sigma^2 Q^{-1} Y \quad (\text{B.8})$$

$$\hat{\theta}^{LS} - \hat{\theta}^R = \sigma^2 (\Phi^T \Phi)^{-1} \Phi^T Q^{-1} Y \quad (\text{B.9})$$

$$A(I_N + BA)^{-1} = (I_N + AB)^{-1} A \quad (\text{B.10})$$

$$\Phi^T Q^{-1} \Phi = S^{-1} \Phi^T Q^{-1} Y = S^{-1} \hat{\theta}^{LS} \quad (\text{B.11})$$

$$\Phi^T Q^{-T} Q^{-T} \Phi = S^{-T} (\Phi^T \Phi)^{-1} S^{-T} \quad (\text{B.12})$$

$$\Phi^T Q^{-T} Q^{-1} Y = S^{-T} (\Phi^T \Phi)^{-1} S^{-1} \hat{\theta}^{LS} \quad (\text{B.13})$$

$$I_N - \sigma^2 Q^{-1} = \Phi P \Phi^T Q^{-1} = Q^{-1} \Phi P \Phi^T = \Phi R^{-1} \Phi^T \quad (\text{B.14})$$

$$\frac{\partial (Q^{-1})_{ij}}{\partial P} = -\Phi^T Q^{-T} J_{ij} Q^{-T} \Phi \quad (\text{B.15})$$

$$\frac{\partial (R^{-1})_{ij}}{\partial P} = \sigma^2 P^{-T} R^{-T} J_{ij} R^{-T} P^{-T} \quad (\text{B.16})$$

where  $J_{ij}$  is the matrix whose  $(i, j)$ -element is one and zero for all other elements.

**Proof.** The identities (B.7)–(B.14) can be verified by a straightforward calculation. Using (B.4) gives

$$\begin{aligned} \frac{\partial (Q^{-1})_{ij}}{\partial P_{st}} &= \sum_{a,b} \frac{\partial (Q^{-1})_{ij}}{\partial Q_{ab}} \frac{\partial Q_{ab}}{\partial P_{st}} \\ &= - \sum_{a,b} (Q^{-1})_{ia} (Q^{-1})_{bj} \Phi_{as} (\Phi^T)_{bt} \\ &= - \sum_{a,b} (\Phi^T)_{sa} (Q^{-T})_{ai} (Q^{-T})_{jb} \Phi_{bt} \\ &= - (\Phi^T Q^{-T})_{si} (Q^{-T} \Phi)_{jt} \end{aligned}$$

which implies (B.15). While (B.16) can be proved in a similar way.  $\square$

### B.2. Convergence result for extremum estimators

**Lemma B3** (Ljung, 1999, Theorem 8.2). Assume that

- (1)  $M(\eta)$  is a deterministic function that is continuous in  $\eta \in \Omega$  and minimized at the set

$$\mathcal{D} = \underset{\eta \in \Omega}{\text{argmin}} M(\eta) = \{\eta | \eta \in \Omega, M(\eta) = \min_{\eta' \in \Omega} M(\eta')\}$$

where  $\Omega$  is a compact subset of  $\mathbb{R}^p$ .

- (2) A sequence of functions  $\{M_N(\eta)\}$  converges to  $M(\eta)$  almost surely and uniformly in  $\Omega$  as  $N$  goes to  $\infty$ .

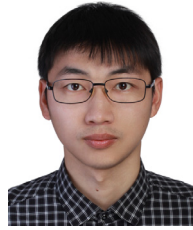
Then  $\hat{\eta}_N = \arg \min_{\eta \in \Omega} M_N(\eta)$  converges to  $\mathcal{D}$  almost surely, namely,

$$\inf_{\eta^* \in \mathcal{D}} \|\hat{\eta}_N - \eta^*\| \rightarrow 0, \text{ as } N \rightarrow \infty.$$

## References

- Aravkin, A., Burke, J. V., Chiuso, A., & Pillonetto, G. (2012a). On the estimation of hyperparameters for empirical bayes estimators: Maximum marginal likelihood vs minimum mse. In *Proceeding of the IFAC symposium on system identification* (pp. 125–130), Brussels, Belgium.

- Aravkin, A., Burke, J. V., Chiuso, A., & Pillonetto, G. (2012b). On the MSE properties of empirical Bayes methods for sparse estimation. In *Proceeding of the IFAC symposium on system identification* (pp. 965–970), Brussels, Belgium.
- Aravkin, A., Burke, J. V., Chiuso, A., & Pillonetto, G. (2014). Convex vs non-convex estimators for regression and sparse estimation: the mean squared error properties of ARD and GLasso. *Journal of Machine Learning Research (JMLR)*, 15, 217–252.
- Boyd, S., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge University Press.
- Carli, F. P., Chen, T., & Ljung, L. (2017). Maximum entropy kernels for system identification. *IEEE Transactions on Automatic Control*, 62, 1471–1477.
- Chen, T. Continuous-time DC kernel-a stable generalized first-order spline kernel. *IEEE Transactions on Automatic Control*, to appear in 2019. <https://doi.org/10.1109/TAC.2018.2825365>.
- Chen, T. (2018b). On kernel design for regularized LTI system identification. *Automatica*, 90, 109–122.
- Chen, T., Andersen, M. S., Ljung, L., Chiuso, A., & Pillonetto, G. (2014). System identification via sparse multiple kernel-based regularization using sequential convex optimization techniques. *IEEE Transactions on Automatic Control*, 59, 2933–2945.
- Chen, T., Ardeschiri, T., Carli, F. P., Chiuso, A., Ljung, L., & Pillonetto, G. (2016). Maximum entropy properties of discrete-time first-order stable spline kernel. *Automatica*, 66, 34–38.
- Chen, T., Ohlsson, H., & Ljung, L. (2012). On the estimation of transfer functions, regularizations and Gaussian processes—Revisited. *Automatica*, 48, 1525–1535.
- Chen, T., & Pillonetto, G. On the stability of reproducing kernel Hilbert spaces of discrete-time impulse responses. *Automatica*, accepted on April 18, 2018, to appear.
- Dinuzzo, F. (2015). Kernels for linear time invariant system identification. *SIAM Journal on Control and Optimization*, 53, 3299–3317.
- Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12, 55–67.
- Ljung, L. (1999). *System identification: Theory for the user*. Upper Saddle River, NJ: Prentice-Hall.
- Ljung, L. (2012). *System identification toolbox for use with MATLAB* (8th ed.). Natick, MA: The MathWorks, Inc..
- Ljung, L., Singh, R., & Chen, T. (2015). Regularization features in the system identification toolbox. In *Proceedings of the IFAC symposium on system identification* (pp. 745–750), Beijing, China.
- Marconato, A., Schoukens, M., & Schoukens, J. (2016). Filter-based regularisation for impulse response modelling. *IET Control Theory & Applications*, 11, 194–204.
- Petersen, K. B., & Pedersen, M. S. (2012). *The Matrix Cookbook*. <http://matrixcookbook.com>.
- Pillonetto, G., Chen, T., Chiuso, A., Nicolao, G. De., & Ljung, L. (2016). Regularized linear system identification using atomic, nuclear and kernel-based norms: The role of the stability constraint. *Automatica*, 69, 137–149.
- Pillonetto, G., & Chiuso, A. (2015). Tuning complexity in regularized kernel-based regression and linear system identification: The robustness of the marginal likelihood estimator. *Automatica*, 58, 106–117.
- Pillonetto, G., Chiuso, A., & De Nicolao, G. (2011). Prediction error identification of linear systems: A nonparametric gaussian regression approach. *Automatica*, 47, 291–305.
- Pillonetto, G., & De Nicolao, G. (2010). A new kernel-based approach for linear system identification. *Automatica*, 46, 81–93.
- Pillonetto, G., Dinuzzo, F., Chen, T., De Nicolao, G., & Ljung, L. (2014). Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50, 657–682.
- Theobald, C. M. (1974). Generalizations of mean square error applied to ridge regression. *Journal of the Royal Statistical Society. Series B. Statistical Methodology*, 36, 103–106.
- Zorzi, M. (2017). On the robustness of the Bayes and Wiener estimators under model uncertainty. *Automatica*, 83, 133–140.
- Zorzi, M., & Chiuso, A. (2017). The harmonic analysis of kernel functions. arXiv preprint arXiv:1703.05216.



**Biqiang Mu** received the Bachelor of Engineering degree in Material Formation and Control Engineering from Sichuan University in 2008 and the Ph.D. degree in Operations Research and Cybernetics from the Academy of Mathematics and Systems Science, Chinese Academy of Sciences in 2013. He was a postdoc at the Wayne State University from 2013 to 2014 and also at the Western Sydney University from 2015 to 2016. He is a postdoc in the Division of Automatic Control, Department of Electrical Engineering, Linköping University since 2016 and he is also an assistant professor at the Academy of Mathematics and Systems Science, Chinese Academy of Sciences. His research interests include system identification (data-driven modelling and analysis), machine learning, and their applications.



**Tianshi Chen** received his Bachelor's degree and Master's degree both from Harbin Institute of Technology in 2001 and 2005, respectively. He received his Ph.D. degree in Automation and Computer-Aided Engineering from The Chinese University of Hong Kong in December 2008. From April 2009 to December 2015, he was working in the Division of Automatic Control, Department of Electrical Engineering, Linköping University, Linköping, Sweden, first as a Postdoc (April 2009–March 2011) and then as an Assistant Professor (April 2011–December 2015). In May 2015, he received the Youth Talents Award of the Thousand Talents Plan of China, and in December 2015, he returned to China and joined the Chinese University of Hong Kong, Shenzhen, as an Associate Professor. He has also been holding a visiting associate professor position at Linköping University since 2016.

He has been mainly working in the area of system identification (data-driven modelling and analysis), statistical signal processing, machine learning, data science, nonlinear control, and their applications. He has participated in several projects in Sweden, Europe and China. He is an Associate Editor for *Automatica* (2017–present), *System & Control Letters* (2017–present), and *IEEE Control Systems Society Conference Editorial Board* (2016–present).

He has been mainly working in the area of system identification (data-driven modelling and analysis), statistical signal processing, machine learning, data science, nonlinear control, and their applications. He has participated in several projects in Sweden, Europe and China. He is an Associate Editor for *Automatica* (2017–present), *System & Control Letters* (2017–present), and *IEEE Control Systems Society Conference Editorial Board* (2016–present).



**Lennart Ljung** received his Ph.D. in Automatic Control from Lund Institute of Technology in 1974. Since 1976 he is Professor of the chair of Automatic Control in Linköping, Sweden. He has held visiting positions at Stanford and MIT and has written several books on System Identification and Estimation. He is an IEEE Fellow, an IFAC Fellow and an IFAC Advisor as well as a member of the Royal Swedish Academy of Sciences (KVA), a member of the Royal Swedish Academy of Engineering Sciences (IVA), an Honorary Member of the Hungarian Academy of Engineering and a Foreign Associate of the US National Academy of Engineering (NAE).

He has received honorary doctorates from the Baltic State Technical University in St Petersburg, from Uppsala University, Sweden, from the Technical University of Troyes, France, from the Catholic University of Leuven, Belgium and from Helsinki University of Technology, Finland. In 2002 he received the Quazza Medal from IFAC, in 2003 he received the Hendrik W. Bode Lecture Prize from the IEEE Control Systems Society, and he was the recipient of the IEEE Control Systems Award for 2007.