# Bayes and Biased Estimators Without Hyper-parameter Estimation: Comparable Performance to the Empirical-Bayes-Based Regularized Estimator

Yue Ju, Bo Wahlberg, Håkan Hjalmarsson

Abstract—Regularized system identification has become a significant complement to more classical system identification. Kernel-based regularized estimators often perform better than the maximum likelihood estimator in terms of minimizing mean squared error (MSE), but often require hyper-parameter estimation. This paper focuses on ridge regression and the regularized estimator by employing the empirical Bayes hyperparameter estimator. We utilize the excess MSE expressions to develop a family of generalized Bayes estimators and a family of closed-form biased estimators. They have the same excess MSE as the empirical-Bayes-based regularized estimator, but eliminate the need for hyper-parameter estimation.

#### I. PRELIMINARIES AND PROBLEM STATEMENT

We consider the following linear regression model,

$$Y = \Phi \theta + E$$
,

where  $\boldsymbol{Y} \in \mathbb{R}^N$  is the measurement output vector with Nbeing the sample size,  $\boldsymbol{\Phi} \in \mathbb{R}^{N \times n}$  is a lower triangular matrix consisting of inputs, and  $\boldsymbol{\theta} \in \mathbb{R}^n$  is the model parameter vector to be estimated. The measurement noise vector  $\boldsymbol{E} \in \mathbb{R}^N$  is assumed to follow  $\mathcal{N}(\boldsymbol{0}, \sigma^2 \mathbf{I}_N)$ . Given an estimator  $\boldsymbol{\hat{\theta}} \in \mathbb{R}^n$  of the unknown parameter  $\boldsymbol{\theta}$ , we evaluate its average performance by its mean squared error (MSE):  $\text{MSE}(\boldsymbol{\hat{\theta}}) = \mathbb{E}(\|\boldsymbol{\hat{\theta}} - \boldsymbol{\theta}_0\|_2^2)$ , where  $\boldsymbol{\theta}_0 \in \mathbb{R}^n$  is the "true" value of  $\boldsymbol{\theta}$  and the expectation  $\mathbb{E}$  is with respect to the measurement noise  $\boldsymbol{E}$ . The smaller its MSE, the better its performance.

# A. Maximum likelihood and regularized estimators

One classical estimator of  $\boldsymbol{\theta}$  is the ML estimator given by  $\hat{\boldsymbol{\theta}}^{\mathrm{ML}} = (\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^{\top} \boldsymbol{Y}$ . It is well-known that  $\hat{\boldsymbol{\theta}}^{\mathrm{ML}}$  is unbiased but may have a large variance, which will result in large  $\mathrm{MSE}(\hat{\boldsymbol{\theta}}^{\mathrm{ML}}) = \sigma^2 \mathrm{Tr}[(\boldsymbol{\Phi}^{\top} \boldsymbol{\Phi})^{-1}]$ . To achieve a better bias-variance trade-off, we consider the following regularized ridge regression estimator,

$$\hat{\boldsymbol{\theta}}^{\mathsf{R}}(\hat{\boldsymbol{\eta}}) = [\boldsymbol{\Phi}^{\top}\boldsymbol{\Phi} + (\sigma^2/\hat{\boldsymbol{\eta}})\mathbf{I}_n]^{-1}\boldsymbol{\Phi}^{\top}\boldsymbol{Y}, \qquad (1)$$

where  $\hat{\eta}$  is a hyper-parameter estimator. A commonly used hyper-parameter estimator is the EB one [4], given by

$$\hat{\eta}_{\text{EB}} = \operatorname*{arg\,min}_{\eta > 0} \mathscr{F}_{\text{EB}}(\eta),$$
$$\mathscr{F}_{\text{EB}}(\eta) = \mathbf{Y}^{\top} \mathbf{Q}(\eta)^{-1} \mathbf{Y} + \log \det(\mathbf{Q}(\eta))$$

with  $\mathbf{Q}(\eta) = \eta \mathbf{\Phi} \mathbf{\Phi}^{\top} + \sigma^2 \mathbf{I}_N$ . Correspondingly,  $\hat{\boldsymbol{\theta}}^{\mathsf{R}}(\hat{\eta}_{\mathrm{EB}})$  in the form of (1) will be referred to as the EB-based regularized estimator in this paper.

# B. Problem statement

Since  $MSE(\hat{\theta}^{R}(\hat{\eta}_{EB}))$  is analytically intractable in finitesample scenarios, we apply a high-order asymptotic quantity: the excess MSE (XMSE) [2]. It can be used to quantify the difference between  $MSE(\hat{\theta}^{R}(\hat{\eta}_{EB}))$  and  $MSE(\hat{\theta}^{ML})$  for large sample sizes. Based on [2, Theorem 2], we can derive the XMSE of  $\hat{\theta}^{R}(\hat{\eta}_{EB})$ .

Lemma 1: Assume that  $\lim_{N\to\infty} \mathbf{\Phi}^{\top} \mathbf{\Phi}/N = \mathbf{I}_n$ . We then have  $\operatorname{XMSE}(\hat{\boldsymbol{\theta}}^{\mathsf{R}}(\hat{\eta}_{\operatorname{EB}})) = (-n^2 + 4n)(\sigma^2)^2/\|\boldsymbol{\theta}_0\|_2^2$ .

- In this work, we consider the following two problems.
- 1) Is it possible to design a generalized Bayes estimator<sup>1</sup> that has the same XMSE as  $\hat{\theta}^{R}(\hat{\eta}_{EB})$ ? It is free of hyper-parameters, thereby eliminating the computational cost associated with estimating such parameters.
- 2) Although a generalized Bayes estimator does not need any hyper-parameter, it often needs to be computed using sampling methods. The question thus arises whether it is possible to design a biased estimator in closed form that has the same XMSE as  $\hat{\theta}^{R}(\hat{\eta}_{EB})$ .

# II. BAYES AND BIASED ESTIMATORS

We first design generalized Bayes estimators that have the same XMSE as  $\hat{\theta}^{R}(\hat{\eta}_{EB})$ .

**Theorem 1:** If the weighting function of  $\hat{\theta}^{\mathrm{Bayes, EB}}$  is

$$\pi(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^{2-n} (C_1 \|\boldsymbol{\theta}\|_2 + C_2 \|\boldsymbol{\theta}\|_2^{-1})^2, \qquad (3)$$

where  $C_1, C_2 \in \mathbb{R}$  are arbitrary constants, we then have  $\text{XMSE}(\hat{\theta}^{\text{Bayes},\text{EB}}) = \text{XMSE}(\hat{\theta}^{\text{R}}(\hat{\eta}_{\text{EB}})).$ 

We then design the following biased estimators that have the same XMSE as  $\hat{\theta}^{R}(\hat{\eta}_{EB})$  and  $\hat{\theta}^{Bayes,EB}$ . **Theorem 2:** If  $\hat{\theta}^{Biased,EB} = \hat{\theta}^{ML} + (1/N)\boldsymbol{b}_{N}(\hat{\theta}^{ML})$  and

**Theorem 2:** If  $\hat{\boldsymbol{\theta}}^{\text{Biased},\text{EB}} = \hat{\boldsymbol{\theta}}^{\text{ML}} + (1/N)\boldsymbol{b}_N(\hat{\boldsymbol{\theta}}^{\text{ML}})$  and  $\boldsymbol{b}_N(\hat{\boldsymbol{\theta}}^{\text{ML}}) = \sigma^2 N(\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$ 

$$\times \left[2 - n + \frac{2(C_1 \|\hat{\boldsymbol{\theta}}^{\mathrm{ML}}\|_2 - C_2 \|\hat{\boldsymbol{\theta}}^{\mathrm{ML}}\|_2^{-1})}{C_1 \|\hat{\boldsymbol{\theta}}^{\mathrm{ML}}\|_2 + C_2 \|\hat{\boldsymbol{\theta}}^{\mathrm{ML}}\|_2^{-1}}\right] \frac{\hat{\boldsymbol{\theta}}^{\mathrm{ML}}}{\|\hat{\boldsymbol{\theta}}^{\mathrm{ML}}\|_2^2}, (4)$$

where  $C_1, C_2 \in \mathbb{R}$  are arbitrary constants, then we have  $\text{XMSE}(\hat{\theta}^{\text{Based},\text{EB}}) = \text{XMSE}(\hat{\theta}^{\text{Bayes},\text{EB}}) = \text{XMSE}(\hat{\theta}^{\text{R}}(\hat{\eta}_{\text{EB}})).$ 

## **III. NUMERICAL SIMULATION**

We generate 100 collections of test systems and inputoutput data. For each collection, 1) we generate  $\tilde{\theta}_0$  as a realization of  $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  and scale  $\theta_0 = m_{\theta} \tilde{\theta}_0$  such that

<sup>&</sup>lt;sup>1</sup>For a generalized Bayes estimator, its nonnegative weighting function  $\pi(\theta)$  can be improper, i.e.,  $\int \pi(\theta) d\theta = +\infty$ .

 $\|\boldsymbol{\theta}_0\|_2 = 1$ ; 2) generate  $\{\tilde{u}(t)\}_{t=1}^N$  as independent realizations of  $\mathcal{N}(0,1)$  and set  $\sigma^2 = 1$ ; 3) scale  $u(t) = m_u \tilde{u}(t)$  such that the sample SNR, which is the ratio between the sample variance of  $\boldsymbol{\Phi}\boldsymbol{\theta}_0$  and the measurement noise variance  $\sigma^2$ , is 5; 4) corrupt the noise-free output  $\boldsymbol{\Phi}\boldsymbol{\theta}_0$  with  $N_{\text{MC}} = 200$  additive independent noise realizations of  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ , to obtain  $N_{\text{MC}}$  measurement output sequences  $\{y(t)\}_{t=1}^N$ .

For each collection of test system and input-output data, we will perform  $N_{\rm MC} = 200$  Monte Carlo (MC) simulations. The regularized estimator  $\hat{\theta}^{\rm R}(\hat{\eta}_{\rm EB})$  will be implemented using [1, Algorithm 2]. The Bayes estimator  $\hat{\theta}^{\rm Bayes,EB}$  will be approximated using the sampling method. The average performance of an estimator  $\hat{\theta}$  will be measured by the sample mean of  $\|\hat{\theta} - \theta_0\|_2$  over 200 MC simulations, referred to as the sample MSE( $\hat{\theta}$ ). As its relative version, the average FIT( $\hat{\theta}$ ) [3] is given by the sample mean of FIT( $\hat{\theta}$ ) =  $100 \times (1 - \|\hat{\theta} - \theta_0\|_2 / \|\theta_0 - \bar{\theta}_0\|_2)$  with  $\bar{\theta}_0 = \frac{1}{n} \sum_{k=1}^{n} [\theta_0]_k$ over 200 MC simulations. The better  $\hat{\theta}$  performs, the smaller its sample MSE, while the larger its average FIT.

We consider the following two settings: 1) n = 1 N = 5, and  $M_s = 200, 2$ ) n = 5, N = 15, and  $M_s = 500$ . From Fig. 1-2, we can observe that for n = 1, N = 5,  $\hat{\theta}^{\text{ML}}$  outperforms  $\hat{\theta}^{\text{R}}(\hat{\eta}_{\text{EB}})$ , while for n = 5, N = 15, the performance of  $\hat{\theta}^{\text{ML}}$ is worse, which confirms the discussions after Lemma 1; for at least one combination of  $C_1$  and  $C_2$ ,  $\hat{\theta}^{\text{Bayes,EB}}$  with (3) and  $\hat{\theta}^{\text{Biased,EB}}$  with (4) perform similarly to  $\hat{\theta}^{\text{R}}(\hat{\eta}_{\text{EB}})$ ; among different combinations of  $C_1$  and  $C_2$ , for n = 1,  $\hat{\theta}^{\text{Bayes,EB}}$  and  $\hat{\theta}^{\text{Biased,EB}}$  with  $C_1 = 1, C_2 = 0$  perform the best; while for n = 5,  $\hat{\theta}^{\text{Bayes,EB}}$  and  $\hat{\theta}^{\text{Biased,EB}}$  with  $C_1 = 0, C_2 = 1$  perform the best.



Fig. 1: Sample means of the sample MSE and the average FIT for n = 1 and N = 5.

For larger n and N, the influence of different ratios of  $C_1$  and  $C_2$  on the performance of  $\hat{\theta}^{\text{Bayes},\text{EB}}$  and  $\hat{\theta}^{\text{Biased},\text{EB}}$  becomes weaker. In Fig. 3 and Table<sup>2</sup> I, we consider n = 80, N = 360, and  $M_s = 5 \times 10^3$ . We can observe that the



Fig. 2: Sample means of the sample MSE and the average FIT for n = 5 and N = 15.

performance of  $\hat{\theta}^{R}(\hat{\eta}_{EB})$  is quite close to that of  $\hat{\theta}^{Bayes,EB}$ and  $\hat{\theta}^{Biased,EB}$ , while its computing time is over twice that of  $\hat{\theta}^{Bayes,EB}$  and over 500 times that of  $\hat{\theta}^{Biased,EB}$ .



Fig. 3: Sample means of the sample MSE and the average FIT for n = 80 and N = 360.

TABLE I: Sample means of the sample MSE, the average FIT, and the total computing time for n = 80 and N = 360.

	$\hat{oldsymbol{ heta}}^{ extsf{R}}(\hat{\eta}_{ extsf{EB}})$	$\hat{oldsymbol{ heta}}^{\mathrm{Bayes,EB}}$	$\hat{ heta}^{\mathrm{Biased},\mathrm{EB}}$
sample MSE	$5.37 \times 10^{-2}$	$5.43 \times 10^{-2}$	$5.37 \times 10^{-2}$
average FIT	76.81	76.69	76.81
computing time (s)	$1.81 \times 10^3$	$6.68  imes 10^2$	3.51

### REFERENCES

- T. Chen and L. Ljung. Implementation of algorithms for tuning parameters in regularized least squares problems in system identification. *Automatica*, 49:2213–2220, 2013.
- [2] Y. Ju, B. Wahlberg, and H. Hjalmarsson. Excess mean squared error of empirical Bayes estimators. arXiv e-prints: 2503.11863, 2025.
- [3] L. Ljung. System Identification Toolbox for Use with MATLAB. The Math Works, 1995.
- [4] G. Pillonetto, T. Chen, A. Chiuso, G. De Nicolao, and L. Ljung. *Regularized System Identification: Learning Dynamic Models from Data*. Springer Nature, 2022.

<sup>&</sup>lt;sup>2</sup>For  $\hat{\theta}^{\mathrm{Bayes, EB}}$  and  $\hat{\theta}^{\mathrm{Biased, EB}}$ , we first calculate the total computing time, and the sample means of the average FIT and the sample MSE over 100 collections of test systems and data. Then, we calculate the sample means of these three statistics over different combinations of  $C_1$  and  $C_2$ .