

Stable likelihood computation for machine learning of linear differential operators with Gaussian processes

Chatrabgoun, O., Esmailbeigi, M., Cheraghi, M. & Daneshkhah, A

Author post-print (accepted) deposited by Coventry University's Repository

Original citation & hyperlink: Chatrabgoun, O, Esmailbeigi, M, Cheraghi, M & Daneshkhah, A 2022, 'Stable likelihood computation for machine learning of linear differential operators with Gaussian processes', International Journal for Uncertainty Quantification, vol. 12, no. 3, pp. 75-99. <https://doi.org/10.1615/Int.J.UncertaintyQuantification.2022038966>

DOI 10.1615/Int.J.UncertaintyQuantification.2022038966

ISSN 2152-5080

ESSN 2152-5099

Publisher: © BEGELL HOUSE Inc. 2023

Copyright © and Moral Rights are retained by the author(s) and/ or other copyright owners. A copy can be downloaded for personal non-commercial research or study, without prior permission or charge. This item cannot be reproduced or quoted extensively from without first obtaining permission in writing from the copyright holder(s). The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the copyright holders.

This document is the author's post-print version, incorporating any revisions agreed during the peer-review process. Some differences between the published version and this version may remain and you are advised to consult the published version if you wish to cite from it.

1 Stable Likelihood Computation for Machine Learning of Linear
2 Differential Operators with Gaussian Processes

3 **O. Chatrabgoun^{a,b}, M. Esmailbeigi^{a,*}, M. Cheraghi^a, A. Daneshkhah^c**

^a Faculty of Mathematical Sciences and Statistics, Malayer University,
Malayer 65719-95863, Iran

^b School of Computing, Electronics and Mathematics, Coventry University,
Coventry, CV1 5FB, UK

^c Centre for Computational Science and Mathematical Modelling, Coventry University,
Coventry, CV1 5FB, UK

4 **Abstract**

5 In many applied sciences, the main aim is to learn the parameters in the operational
6 equations which best fit the observed data. A framework for solving such problems is to em-
7 ploy Gaussian process (GP) emulators which are well-known as non-parametric Bayesian
8 machine learning techniques. GPs are among a class of methods known as kernel machines
9 which can be used to approximate rather complex problems by tuning their hyperparam-
10 eters. The maximum likelihood estimation (MLE) has widely been used to estimate the
11 parameters of the operators and kernels. However, the MLE-based and Bayesian infer-
12 ence in the standard form are usually involved in setting up a covariance matrix which
13 is generally ill-conditioned. As a result, constructing and inverting the covariance matrix
14 using the standard form will become unstable to learn the parameters in the operational
15 equations. In this paper, we propose a novel approach to tackle these computational com-
16 plexities and also resolve the ill-conditioning problem by forming the covariance matrix
17 using alternative bases via the Hilbert-Schmidt SVD (HS-SVD) approach. Applying this
18 approach yields a novel matrix factorization of the block-structured covariance matrix
19 which can be implemented stably by isolating the main source of the ill-conditioning. In
20 contrast to standard matrix decompositions which start with a matrix and produce the

*Corresponding author. E-mail address: m.esmailbeigi@malayeru.ac.ir

21 resulting factors, the HS-SVD is constructed from the Hilbert-Schmidt eigenvalues and
 22 eigenvectors without the need to ever form the potentially ill-conditioned matrix. We also
 23 provide stable MLE and Bayesian inference to adaptively estimate hyperparameters, and
 24 the corresponding operators can then be efficiently predicted at some new points using the
 25 proposed HS-SVD bases. The efficiency and stability of the proposed HS-SVD method
 26 will be compared with the existing methods by several illustrations of the parametric lin-
 27 ear equations, such as ordinary and partial differential equations, integro-differential and
 28 fractional order operators.

29 **Keyword:** Gaussian processes; Hilbert–Schmidt’s theory; Stable computation; Prob-
 30 abilistic machine learning; Uncertainty quantification.

31

32 1 Introduction

33 One of the major fields in applied sciences is to model different phenomena in terms of flexible
 34 operational equations [1, 2]. In other words, the researchers usually attempt to find a coherent
 35 form of flexible operational equations corresponding to the observed data to the effect that
 36 they best describe and govern them [3, 4]. Therefore, the necessity of existence of parametric
 37 operational equations is created, that is, the equations that have parameters and they increase
 38 the flexibility to cover the observed data. These parametric equations come from real-world
 39 mathematical modelling, and their parameters should be determined in terms of the observed
 40 data. Since nonlinear operators can be approximated by linear ones in many cases [5], an
 41 important category of these parametric operational equations is linear operations equation.
 42 Therefore, the major aim is to find their parameter which is known as a linear inverse problem
 43 (IP) [6].

44 To illustrate the key ingredients and focal point of this study, let us start by considering
 45 linear operational equations of the form

$$\mathcal{L}_{\mathbf{x}}^{\zeta} u(\mathbf{x}) = f(\mathbf{x}), \quad (1)$$

46 which models the relationship between $u(\mathbf{x})$ and $f(\mathbf{x})$ functions. Here, $f(\mathbf{x})$ is considered as a
 47 black-box forcing term, $\mathcal{L}_{\mathbf{x}}^{\zeta}$ is a linear operator equipped with parameter ζ and $u(\mathbf{x})$ denotes
 48 the latent solution. Take, for instance, the classical problem of heat conduction in a medium
 49 with unknown conductivity properties, albeit with an unknown thermal diffusivity coefficient
 50 ζ . The main objective for solving the linear parametric problems (1) is to optimally and

51 stably find the parameters ζ from the observed data using an implicit method. The linear
52 operational equations given by (1), have enormous practical benefits which are discussed in
53 [7, 8].

54 Raissi et al. [6] provide an innovative method for resolving such problems by employing
55 and adopting GPs [9, 10] within a Bayesian framework. Likewise, the Bayesian procedure
56 adopted here uses the GP as a flexible prior distribution over functions thereby providing ana-
57 lytical tractability. Once combined with the observed data, it will supply a fully probabilistic
58 approach to approximate the functions. Moreover, GPs are among a class of methods known
59 as kernel machines (as discussed in [5, 11, 12]) and have close similarities with regularisation
60 approaches [13, 14].

61 GP as a kernel-based non-parametric method relies on an appropriate selection of kernel
62 k . It is common to select a parametrized family of kernels. These kernels are parametrized by
63 one or more hyperparameters θ , which then need to be estimated on the basis of the observed
64 data. Indeed, the selection of the kernel k and the estimation of its hyperparameters have a
65 profound impact on the performance of the GP through the covariance matrix \mathbf{K} which is
66 constructed based on the selected kernel.

67 By applying the GP as a kernel-based method which includes the hyperparameters θ ,
68 the Eq. (1) can be considered as a surrogate model where the proposed GP serves as the
69 probabilistic approximation of the underlying problem. Given the observed data, the aim is
70 then to learn the hyperparameters θ and parameters ζ using a stable method. The learned
71 GP using the mentioned method can then be used to probabilistically solve the governing
72 linear operational equation given in (1) and predict the behaviour of this system in the future.
73 However, the goal pursued in many other studies addressing the similar inverse problem is
74 that the model parameters and hyperparameters are learned directly from the observed data
75 using various optimization techniques including Power Function method, Cross Validation
76 method, Trial and Error method, and the Contour-Pade algorithm [15, 16, 18]. Raissi et al. [6]
77 applied commonly done estimating the model parameters and hyperparameters directly from
78 the data by minimizing the negative log likelihood of the probabilistic model in GP regression
79 (See [9]), and then they used the usual optimization to fit a GP model to the parametric
80 operational problem given in Eq. (1). The main advantage of the method proposed in [6]
81 in comparison with the approaches mentioned above is that the optimal model parameters
82 and hyperparameters are all learned directly from the data by maximizing the joint marginal

83 log-likelihood of the probabilistic model. Therefore, minimizing the negative log likelihood
84 function estimates the (hyper)parameters that will most likely be used to produce the data
85 required for the linear operational problem. In constructing the likelihood function, the inverse
86 of the covariance matrix appeared in the quadratic term is used to learn from the training
87 data, while the log-determinant of the covariance matrix penalizes the model complexity.
88 If the covariance matrix is ill-conditioned, computing both the inverse and log-determinant
89 of the covariance matrix with standard methods (e.g., Cholesky factorization) will probably
90 be inaccurate. This would restrict us to use the maximum likelihood estimation (MLE) to
91 evaluate the validity of parameters and hyperparameters estimation and their posterior mean
92 accuracy.

93 This paper deals with stabilizing the likelihood computation of Gaussian process regression
94 for data of linear operational equations and subsequent optimization of (hyper-) parameters.
95 The main innovation is based on the novel matrix factorization of the block-structured co-
96 variance matrix which is generally unstable and ill-conditioned and needs to be inverted using
97 Hilbert-Schmidt SVD. These block-structured partitioning of the covariance matrix is imple-
98 mented stably without any major computational burdens by isolating the main source of the
99 ill-conditioning.

100 Without such isolating the covariance matrix is ill-conditioned such that the likelihood
101 computation and subsequent optimization of (hyper-) parameters is inaccurate. Therefore,
102 the main aim of the present paper is to stabilize ill-conditioning behaviour of the likelihood
103 computation through factorization of the covariance matrix (main source of instability) using
104 the Hilbert-Schmidt (HS) SVD as an alternative base. Whereas the standard and basic meth-
105 ods (e.g. Cholesky [6], RBF-QR [19] and weighted SVD [20]) for tackling this issue is still
106 suffering from the same presented instabilities in the likelihood function. Unlike the previous
107 methods, to stabilize the likelihood computation every positive definite kernel is represented in
108 terms of the (positive) eigenvalues and (normalized) eigenfunctions of an associated compact
109 integral operator without needing to decompose the covariance matrix in the same unstable
110 way [21]. Finally, using the properties such as orthogonality of eigenfunctions, the rapid de-
111 cay of eigenvalues for highly smooth kernels, and isolation of swiftly decaying eigenvalues (as
112 the main source of ill-conditioning in the covariance matrix) likelihood computation is imple-
113 mented in stabilized and stable manner. In fact, the principal key for computing likelihood
114 function is that it is not necessary to directly deal with the kernel in its closed form [16, 22]

115 because the composition of the covariance matrix can be obtained by the associate eigenvalues
 116 and eigenfunctions based on the Mercer’s expansion. Furthermore, due to the nature of the
 117 linear operators such as fractional-order operators and the inherent ill-conditioning of kernel-
 118 based approximation methods, we have included the fractional derivatives of the associate
 119 eigenfunctions of the kernel instead of the kernel directly [23, 24, 26]. Therefore, the source
 120 of the ill-conditioning in the posterior prediction process is analytically removed using the
 121 HS-SVD method proposed in this paper.

122 At the end, the efficiency and stability of the proposed HS-SVD method have been com-
 123 pared with the existing methods by several illustrations of the parametric linear equations,
 124 such as ordinary and partial differential equations, integro-differential and fractional order
 125 operators.

126 This paper is organized as follows. In Section 2 the challenges of the GPs as a data-driven
 127 algorithm for learning general parametric linear equations are presented. In Section 3, we
 128 introduce the new stable bases using the HS-SVD method and discuss how this method can
 129 tackle the computational challenges of the likelihood function. We illustrate several numerical
 130 examples in Section 4. Finally, some conclusions are presented in Section 5.

131 **2 Computational challenges in learning linear operators using** 132 **GPs**

133 First, we introduce the machine learning of the linear differential operators using GPs and
 134 its challenges in the computation of likelihood function and posterior mean. Without loss
 135 of generality, we assume the observed data was generated by a zero-mean GP, i.e., $\mu = 0$,
 136 although a nonzero mean can also be considered. In the following, the proposed data-driven
 137 algorithm presented by Raissi et al. [6] for learning general parametric linear equations of
 138 the form (1) corresponding to the differential operators is presented. The algorithm starts by
 139 assuming that $u(\mathbf{x})$ is GP with mean 0 and covariance function $k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$, i.e.,

$$u(\mathbf{x}) \sim \mathcal{GP}(0, k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})), \quad (2)$$

140 where $\boldsymbol{\theta}$ denotes the hyperparameters of the kernel k_{uu} . The key observation to make is
 141 that any linear transformation of a GP such as differentiation and integration is still a GP.

142 Consequently,

$$\mathcal{L}_{\mathbf{x}}^{\zeta} u(\mathbf{x}) = f(\mathbf{x}) \sim \mathcal{GP}(0, k_{ff}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta)), \quad (3)$$

143 with the following relationship between the kernels k_{uu} and k_{ff} ,

$$k_{ff}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}}^{\zeta} \mathcal{L}_{\mathbf{x}'}^{\zeta} k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}). \quad (4)$$

144 Furthermore, the covariance between $u(\mathbf{x})$ and $f(\mathbf{x}')$, and also the one between $f(\mathbf{x})$ and $u(\mathbf{x}')$
 145 are determined by $k_{uf}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}'}^{\zeta} k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$, and $k_{fu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}}^{\zeta} k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$,
 146 respectively. It goes without saying that the main purpose is to estimate the parameters ζ of
 147 the operator $\mathcal{L}_{\mathbf{x}}^{\zeta}$ and the hyperparameters of the kernels k_{ff} , k_{uf} , and k_{fu} .

148 There are some situations where it is reasonable to assume that the observations are noise-
 149 free, for example in computer simulations. Many scientific phenomena are investigated by com-
 150 plex computer models or codes. A feature of many computer experiments is that the output
 151 is deterministic i.e., rerunning the code with the same inputs gives identical observations[25].
 152 Due to the fact that our data in the present paper are taken from computer simulations, it
 153 is assumed $\mathbf{y} = \begin{bmatrix} \mathbf{y}_u \\ \mathbf{y}_f \end{bmatrix}$, such that $\mathbf{y}_u = u(\mathbf{X}_u)$, $\mathbf{y}_f = f(\mathbf{X}_f)$. Based on the aforementioned
 154 properties of the MLE, as pointed out by Raissi et al. in [6], “the hyperparameters $\boldsymbol{\theta}$ and
 155 more importantly the parameters ζ of the linear operator $\mathcal{L}_{\mathbf{x}}^{\zeta}$ can be trained by minimizing
 156 the negative log marginal likelihood (NLML)

$$\text{NLML}(\zeta, \boldsymbol{\theta}) = -\log p(\mathbf{y}|\zeta, \boldsymbol{\theta}), \quad (5)$$

157 where $p(\mathbf{y}|\zeta, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$, however (5) can be rewritten as

$$\text{NLML}(\zeta, \boldsymbol{\theta}) = \frac{1}{2} \log(|\mathbf{K}|) + \frac{1}{2} \mathbf{y}^{\top} \mathbf{K}^{-1} \mathbf{y} + \frac{N}{2} \log 2\pi, \quad (6)$$

158 and \mathbf{K} is given by

$$\mathbf{K} = \begin{bmatrix} k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta}) & k_{uf}(\mathbf{X}_u, \mathbf{X}_f; \boldsymbol{\theta}, \zeta) \\ k_{fu}(\mathbf{X}_f, \mathbf{X}_u; \boldsymbol{\theta}, \zeta) & k_{ff}(\mathbf{X}_f, \mathbf{X}_f; \boldsymbol{\theta}, \zeta) \end{bmatrix}.” \quad (7)$$

159 As pointed out by Raissi et al. in [6], “the marginal likelihood does not simply favor the models
 160 that fit the training data best. In fact, it induces an automatic trade-off between data-fit and
 161 model complexity. The likelihood function includes the \mathbf{K}^{-1} term in a quadratic framework
 162 which targets to fit the training data, while the log-determinant term $\log |\mathbf{K}|$ penalizes the

163 model complexity”. The most computationally intensive part of the training is associated with
 164 inverting dense covariance matrix \mathbf{K} , and computing determinant \mathbf{K} . This scales cubically
 165 with the number of observed data. This scaling is a well-known limitation of the GP, and
 166 even if \mathbf{K} is invertible it may still be numerically ill-conditioned [27]. These challenges make
 167 the results unreliable and reduce the validity of the method. Actually, when \mathbf{K} becomes ill-
 168 conditioned, it can lead to an ill-conditioned problem and computing with standard methods
 169 (e.g., Cholesky factorization) is probably inaccurate, leaving us unable to use the MLE to
 170 judge the validity of the method.

171 After training the model and parameter estimation in the previous step, we predict the
 172 values $u(\mathbf{x})$ and $\mathcal{L}_{\mathbf{x}}^{\zeta}u(\mathbf{x}) = f(\mathbf{x})$ at a new test point \mathbf{x} by writing the posterior distributions

$$\begin{aligned} p(u(\mathbf{x})|\mathbf{y}) &= N(\bar{u}(\mathbf{x}), v_u^2(\mathbf{x})), \\ p(f(\mathbf{x})|\mathbf{y}) &= N(\bar{f}(\mathbf{x}), v_f^2(\mathbf{x})), \end{aligned} \tag{8}$$

173 such that

$$\begin{aligned} \bar{u}(\mathbf{x}) &= \mathbf{k}_u^{\top}(\mathbf{x})\mathbf{K}^{-1}\mathbf{y}, & v_u^2(\mathbf{x}) &= k_{uu}(\mathbf{x}, \mathbf{x}) - \mathbf{k}_u^{\top}(\mathbf{x})\mathbf{K}^{-1}\mathbf{k}_u(\mathbf{x}), \\ \bar{f}(\mathbf{x}) &= \mathbf{k}_f^{\top}(\mathbf{x})\mathbf{K}^{-1}\mathbf{y}, & v_f^2(\mathbf{x}) &= k_{ff}(\mathbf{x}, \mathbf{x}) - \mathbf{k}_f^{\top}(\mathbf{x})\mathbf{K}^{-1}\mathbf{k}_f(\mathbf{x}), \end{aligned} \tag{9}$$

174 and

$$\begin{aligned} \mathbf{k}_u^{\top}(\mathbf{x}) &= \mathbf{k}^{\top}(\mathbf{x}) = \begin{bmatrix} k_{uu}^{\top}(\mathbf{x}, \mathbf{X}_u) & k_{uf}^{\top}(\mathbf{x}, \mathbf{X}_f) \end{bmatrix}, \\ \mathbf{k}_f^{\top}(\mathbf{x}) &= \begin{bmatrix} k_{fu}^{\top}(\mathbf{x}, \mathbf{X}_u) & k_{ff}^{\top}(\mathbf{x}, \mathbf{X}_f) \end{bmatrix}, \end{aligned} \tag{10}$$

175 where for notational convenience the dependence of the kernels on hyperparameters and pa-
 176 rameters is dropped. Apart from $\bar{u}(\mathbf{x})$ and $\bar{f}(\mathbf{x})$, the posterior variances $v_u^2(\mathbf{x})$ and $v_f^2(\mathbf{x})$ can
 177 be used as good indicators of how confident one could be about the estimated parameters ζ
 178 of the linear operator $\mathcal{L}_{\mathbf{x}}^{\zeta}$ and posterior predictions made based on these parameters. How-
 179 ever, both of these representations can lead to severe numerical instability. To overcome these
 180 challenges, we address them in Section 3. Since the GP inherits the properties of its kernel, a
 181 brief review of the kernels and their hyperparameters is required.

182 Since everything hinges upon our selection of the kernel k_{uu} -though this kernel is not
 183 usually known-it is common to consider a parametrized family of kernels. As mentioned be-
 184 fore, the selection of these hyperparameters has a significant effect on the performance of the

Table 1: Some well-known kernels k_{uu} with their hyperparameters $\theta = (\varepsilon, \beta, a, b)$.

Name	Definition
Squared Exponential (SE)	$\exp\left(-\frac{1}{2}\sum_{i=1}^d\varepsilon_i(x_i-x'_i)^2\right)$
Multiquadrics (MQ)	$\sqrt{1+(\varepsilon r)^2}$
Periodic spline	$\sum_{n=1}^{\infty}\frac{2}{(4n^2\pi^2)^\beta}2\cos(2n\pi(x-x'))$
Generalized periodic spline	$\sum_{n=1}^{\infty}\frac{2}{(4n^2\pi^2+\varepsilon^2)^\beta}2\cos(2n\pi(x-x')), \quad x, x' \in [0, 1], \beta \in \mathbb{N}$
Chebyshev	$\left(2a(1-b)\frac{b(1-b^2)-2b(x^2+x'^2)+(1+3b^2)xx'}{(1-b^2)^2+4b(b(x^2+x'^2)-(1+b^2)xx')}\right)$ $+(1-a), \quad a \in (0, 1], b \in (0, 1)$

185 GP. Some common kernel families are included in Table 1 [9]. In Multiquadrics (MQ) kernel
186 $r = \|\mathbf{x} - \mathbf{x}'\|$ such that $\|\cdot\|$ is a norm on \mathbb{R}^d and usually the Euclidean norm [15]. The
187 hyperparameter b in Chebyshev kernel acts like a shape parameter, where $b \rightarrow 1$ yields more
188 peakier kernels with increased locality (and thus reduced interactions between kernels) and
189 $b \rightarrow 0$ yields a flatter kernel with increasingly global behavior which is less concentrated. Also,
190 the hyperparameter a is not that significant as long as $a \in (0, 1)$ as it just shifts and scales the
191 kernel vertically. However, setting $a = 1$ eliminates the vertical shift and therefore makes it
192 markedly more difficult to fit data with a nonzero mean. The existence of the hyperparameter
193 ε as shape parameter in squared exponential and multiquadrics kernel (or other hyperparam-
194 eters such as β in periodic spline) allows for flexibility to select a kernel supported by the data
195 without having to explore the endless selection of all positive definite kernels. Unfortunately,
196 this flexibility is often accompanied by the danger of severe ill-conditioning for small ε be-
197 cause of the increasing linear dependence of the vectors in the covariance matrix. Therefore,
198 computing the standard form of the covariance matrix is not a good idea [28, 29, 30]. The
199 ill-conditioning can be corrected using alternative bases and matrix factorization.

200 3 Overcoming computational challenges using HS-SVD bases

201 3.1 Stable computation in the posterior prediction

202 We now develop the HS-SVD method to create an alternative basis, to eliminate the ill-
203 conditioning of the covariance matrix \mathbf{K} . In the standard approach, we work with a basis

204 generated by the kernel "shifts" which are known as standard bases. The new basis derived
 205 from the eigenfunction expansion produces a linear system that is devoid of the ill-conditioning
 206 of the standard form of \mathbf{K} . The HS-SVD method has two advantages in practical application:
 207 (i) it is not necessary to form the covariance matrix, and (ii) it is not necessary to know the
 208 kernel and its derivatives in closed form [16, 22]. Therefore, we apply it to the covariance
 209 matrix \mathbf{K} . We define the truncated Hilbert–Schmidt expansion with M terms (or truncated
 210 Mercer series expansion) of our kernel k_{uu} in (2) as

$$k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \sum_{n=1}^M \lambda_n \varphi_n(\mathbf{x}) \varphi_n(\mathbf{x}') = \begin{bmatrix} \varphi_1(\mathbf{x}) & \dots & \varphi_M(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_M \end{bmatrix} \begin{bmatrix} \varphi_1(\mathbf{x}') \\ \vdots \\ \varphi_M(\mathbf{x}') \end{bmatrix}, \quad (11)$$

211 where the truncated Mercer series provides the best M -term approximation of the kernel in
 212 the mean-square error [16].

It should be noted that, (λ_n, φ_n) are orthonormal eigenpairs of a Hilbert–Schmidt integral operator $T_{k_{uu}} : L_2(\Omega, \rho) \rightarrow L_2(\Omega, \rho)$ defined as

$$(T_{k_{uu}} f)(\mathbf{x}) = \int_{\Omega} k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) f(\mathbf{x}') \rho(\mathbf{x}') d\mathbf{x}',$$

213 where $\Omega \subseteq \mathbb{R}^d$, ρ is a weight function and $\|k_{uu}\|_{L_2(\Omega \times \Omega, \rho \times \rho)} < \infty$ and also, M is chosen as the
 214 smallest value that satisfies $\lambda_M < \epsilon_{\text{mach}} \lambda_{n_u + n_f}$ and we will always assume that $M > n_u + n_f$
 215 such that n_u and n_f are the number of training points chosen from $u(\mathbf{x})$ and $f(\mathbf{x})$ functions,
 216 respectively. Here ϵ_{mach} is machine precision (assumed to be 10^{-16}). Therefore, the quadratic
 217 form (11) can be replaced by

$$k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \boldsymbol{\phi}(\mathbf{x})^\top \Lambda \boldsymbol{\phi}(\mathbf{x}'), \quad (12)$$

where

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \varphi_1(\mathbf{x}) \\ \vdots \\ \varphi_M(\mathbf{x}) \end{bmatrix}, \quad \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_M \end{bmatrix}.$$

218 The eigen-decomposition (12) provides an accurate approximation of the kernel k_{uu} without
 219 ever forming it. According to (12), the Hilbert–Schmidt eigen-decomposition of the vector
 220 $k_{uu}(\mathbf{x}, \mathbf{X}_u; \boldsymbol{\theta})^\top = [k_{uu}(\mathbf{x}, \mathbf{x}_1; \boldsymbol{\theta}) \dots k_{uu}(\mathbf{x}, \mathbf{x}_{n_u}; \boldsymbol{\theta})]$ and also the matrix $k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta})$ in (7)
 221 are as

$$k_{uu}(\mathbf{x}, \mathbf{X}_u; \boldsymbol{\theta})^\top = \boldsymbol{\phi}(\mathbf{x})^\top \Lambda \Phi^\top \Rightarrow k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta}) = \Phi \Lambda \Phi^\top, \quad (13)$$

where

$$\Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \vdots \\ \phi(\mathbf{x}_{n_u})^\top \end{bmatrix} = \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \dots & \varphi_M(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \varphi_1(\mathbf{x}_{n_u}) & \dots & \varphi_M(\mathbf{x}_{n_u}) \end{bmatrix}.$$

222 However, it is not recommended to directly use the decomposition (13) either because all of
 223 the ill-conditioning associated with matrix $k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta})$ is still present in the matrix Λ . We
 224 now use mostly standard numerical linear algebra to isolate some of the ill-conditioning and
 225 develop the HS-SVD. The key step in removing the ill-conditioning is to write the component
 226 matrices that appear in the eigen-decomposition of the matrix $k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta})$ in (13) using
 227 blocks $\Phi = (\Phi_1 \ \Phi_2)$ and $\Lambda = \begin{bmatrix} \Lambda_1 & \\ & \Lambda_2 \end{bmatrix}$ where

$$\Phi_1 = \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \dots & \varphi_{n_u}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \varphi_1(\mathbf{x}_{n_u}) & \dots & \varphi_{n_u}(\mathbf{x}_{n_u}) \end{bmatrix}, \Phi_2 = \begin{bmatrix} \varphi_{n_u+1}(\mathbf{x}_1) & \dots & \varphi_M(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \varphi_{n_u+1}(\mathbf{x}_{n_u}) & \dots & \varphi_M(\mathbf{x}_{n_u}) \end{bmatrix},$$

and $\Phi_1, \Lambda_1 \in \mathbb{R}^{n_u \times n_u}$, $\Phi_2 \in \mathbb{R}^{n_u \times (M-n_u)}$, $\Lambda_2 \in \mathbb{R}^{(M-n_u) \times (M-n_u)}$. Now, we can write the
 eigen-decomposition of the vector $k_{uu}(\mathbf{x}, \mathbf{X}_u; \boldsymbol{\theta})^\top$ in (13) as

$$k_{uu}(\mathbf{x}, \mathbf{X}_u; \boldsymbol{\theta})^\top = \underbrace{\phi(\mathbf{x})^\top \begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_2^\top \Phi_1^{-\top} \Lambda_1^{-1} \end{bmatrix}}_{\psi(\mathbf{x})^\top} \Lambda_1 \Phi_1^\top = \psi(\mathbf{x})^\top \Lambda_1 \Phi_1^\top,$$

228 and also, we can proceed to construct block matrix $k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta})$ as following

$$k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \boldsymbol{\theta}) = \Phi \Lambda \Phi^\top = \Phi \underbrace{\begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_2^\top \Phi_1^{-\top} \Lambda_1^{-1} \end{bmatrix}}_{\Psi} \Lambda_1 \Phi_1^\top = \Psi \Lambda_1 \Phi_1^\top. \quad (14)$$

229 The rate of decay of the Hilbert-Schmidt eigenvalues determines the smoothness of the
 230 kernel k_{uu} : the faster the eigenvalues decay, the smoother the kernel (and vice versa). More
 231 specifically, if the eigenvalues decay at an algebraic rate of $\mathcal{O}^{-\beta+1+\tau}$ with $\beta \in \mathbb{N}_0$ and arbi-
 232 trarily small $\tau > 0$, then the kernel will be a finite smooth kernel in C^β , and if the eigenvalues
 233 decay geometrically, then the kernel will be infinitely smooth, even analytic. Moreover, the
 234 smoothness of the kernel determines the rate of convergence of the kernel-based approximation

235 method. In this case the general rule of thumb is: the smoother the kernel, the faster the rate
 236 of convergence of the approximation method (for more information, see [16]).

The Mercer's theorem guarantees the uniform convergence of the series (11) provided that $T_{k_{uu}}$ is a positive operator. Also, the Mercer series provides the best approximation for the selected kernel in terms of various metrics, in particular, the mean-square error [16]. Generally, our kernel k_{uu} is positive definite and thus $T_{k_{uu}}$ is a positive operator. It would be plausible to assume the linear operators \mathcal{L}^ζ are continuous (bounded) and due to the uniform convergence of the series (11), we can proceed to construct an approximation for $k_{uf}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}'}^\zeta k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$, $k_{fu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}}^\zeta k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ and $k_{ff}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}, \zeta) = \mathcal{L}_{\mathbf{x}}^\zeta \mathcal{L}_{\mathbf{x}'}^\zeta k_{uu}(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta})$ using Mercer series and construct block matrices of \mathbf{K} as following

$$\begin{aligned}
 k_{uf}(\mathbf{X}_u, \mathbf{X}_f; \boldsymbol{\theta}, \zeta) &= \Phi \underbrace{\begin{bmatrix} I_{n_f} \\ \Lambda_2 \Phi_{2, \mathcal{L}^{\mathbf{x}'}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^{-\top} \Lambda_1^{-1} \end{bmatrix}}_{\Psi_{\mathcal{L}^{\mathbf{x}'}}} \Lambda_1 \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^\top = \Psi_{\mathcal{L}^{\mathbf{x}'}} \Lambda_1 \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^\top, \\
 k_{ff}(\mathbf{X}_f, \mathbf{X}_f; \boldsymbol{\theta}, \zeta) &= \Phi_{\mathcal{L}^{\mathbf{x}}} \underbrace{\begin{bmatrix} I_{n_f} \\ \Lambda_2 \Phi_{2, \mathcal{L}^{\mathbf{x}'}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^{-\top} \Lambda_1^{-1} \end{bmatrix}}_{\Psi_{\mathcal{L}^{\mathbf{x}} \mathcal{L}^{\mathbf{x}'}}} \Lambda_1 \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^\top = \Psi_{\mathcal{L}^{\mathbf{x}} \mathcal{L}^{\mathbf{x}'}} \Lambda_1 \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^\top, \\
 k_{fu}(\mathbf{X}_f, \mathbf{X}_u; \boldsymbol{\theta}, \zeta) &= \Phi_{\mathcal{L}^{\mathbf{x}}} \underbrace{\begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_{2, \mathcal{L}^{\mathbf{x}'}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}'}}^{-\top} \Lambda_1^{-1} \end{bmatrix}}_{\Psi_{\mathcal{L}^{\mathbf{x}}}} \Lambda_1 \Phi_1^\top = \Psi_{\mathcal{L}^{\mathbf{x}}} \Lambda_1 \Phi_1^\top, \tag{15}
 \end{aligned}$$

237 where

$$\Phi_{\mathcal{L}^{\mathbf{x}}} = \begin{bmatrix} \mathcal{L}_1^{\mathbf{x}} \phi(\mathbf{x})^\top \\ \vdots \\ \mathcal{L}_{n_f}^{\mathbf{x}} \phi(\mathbf{x})^\top \end{bmatrix} = \begin{pmatrix} \Phi_{1, \mathcal{L}^{\mathbf{x}}} & \Phi_{2, \mathcal{L}^{\mathbf{x}}} \end{pmatrix}.$$

239 Indeed, to achieve more stability, the QR decomposition is used as

$$\Phi_{\mathcal{L}^{\mathbf{x}}} = \begin{pmatrix} \Phi_{1, \mathcal{L}^{\mathbf{x}}} & \Phi_{2, \mathcal{L}^{\mathbf{x}}} \end{pmatrix} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_{1, \mathcal{L}^{\mathbf{x}}} & \mathbf{R}_{2, \mathcal{L}^{\mathbf{x}}} \end{pmatrix} \Rightarrow \Phi_{2, \mathcal{L}^{\mathbf{x}}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}}}^{-\top} = \mathbf{R}_{2, \mathcal{L}^{\mathbf{x}}}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{R}_{1, \mathcal{L}^{\mathbf{x}}}^{-\top} = \mathbf{R}_{2, \mathcal{L}^{\mathbf{x}}}^\top \mathbf{R}_{1, \mathcal{L}^{\mathbf{x}}}^{-\top}. \tag{16}$$

240 For example, using the relation above, it can safely be concluded that the correction matrix

241 $[\Lambda_2 \Phi_{2, \mathcal{L}^{\mathbf{x}}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}}}^{-\top} \Lambda_1^{-1}]$ is as

$$[\Lambda_2 \Phi_{2, \mathcal{L}^{\mathbf{x}}}^\top \Phi_{1, \mathcal{L}^{\mathbf{x}}}^{-\top} \Lambda_1^{-1}] = [\Lambda_2 \mathbf{R}_{2, \mathcal{L}^{\mathbf{x}}}^\top \mathbf{R}_{1, \mathcal{L}^{\mathbf{x}}}^{-\top} \Lambda_1^{-1}].$$

242 Now, using the relations (14–15), the covariance matrix \mathbf{K} can be decomposed as follows

$$\begin{aligned}
\mathbf{K} &= \begin{bmatrix} \Psi \Lambda_1 \Phi_1^\top & \Psi_{\mathcal{L}^{x'}} \Lambda_1 \Phi_{1,\mathcal{L}^{x'}}^\top \\ \Psi_{\mathcal{L}^x} \Lambda_1 \Phi_1^\top & \Psi_{\mathcal{L}^x, \mathcal{L}^{x'}} \Lambda_1 \Phi_{1,\mathcal{L}^{x'}}^\top \end{bmatrix}, \\
&= \underbrace{\begin{bmatrix} \Psi & \Psi_{\mathcal{L}^{x'}} \\ \Psi_{\mathcal{L}^x} & \Psi_{\mathcal{L}^x, \mathcal{L}^{x'}} \end{bmatrix}}_{\Psi} \underbrace{\begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_1 \end{bmatrix}}_{\Lambda_1} \underbrace{\begin{bmatrix} \Phi_1^\top & 0 \\ 0 & \Phi_{1,\mathcal{L}^{x'}}^\top \end{bmatrix}}_{\Phi_1^\top} = \Psi \Lambda_1 \Phi_1^\top. \tag{17}
\end{aligned}$$

244 To demonstrate the usefulness of the HS-SVD, we write the posterior mean (9) as

$$\begin{aligned}
\bar{u}(\mathbf{x}) &= \mathbf{k}_u^\top(\mathbf{x}) \mathbf{K}^{-1} \mathbf{y} = \psi_u^\top(\mathbf{x}) \mathbf{b}, \\
\bar{f}(\mathbf{x}) &= \mathbf{k}_f(\mathbf{x})^\top \mathbf{K}^{-1} \mathbf{y} = \psi_f^\top(\mathbf{x}) \mathbf{b}, \tag{18}
\end{aligned}$$

245 where $\mathbf{b} = \Psi^{-1} \mathbf{y}$, $\psi_u^\top(\mathbf{x}) = [\psi(\mathbf{x})^\top \ \psi_{\mathcal{L}^{x'}}(\mathbf{x})^\top]$ and $\psi_f^\top(\mathbf{x}) = [\psi_{\mathcal{L}^x}(\mathbf{x})^\top \ \psi_{\mathcal{L}^x, \mathcal{L}^{x'}}(\mathbf{x})^\top]$ and also,

$$\begin{aligned}
\psi_{\mathcal{L}^x}(\mathbf{x})^\top &= \phi_{\mathcal{L}^x}(\mathbf{x})^\top \begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_2^\top \Phi_1^{-\top} \Lambda_1^{-1} \end{bmatrix}, \phi_{\mathcal{L}^x}(\mathbf{x})^\top = \mathcal{L}^x \phi(\mathbf{x})^\top, \\
\psi_{\mathcal{L}^{x'}}(\mathbf{x})^\top &= \phi(\mathbf{x})^\top \begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_{2,\mathcal{L}^{x'}}^\top \Phi_{1,\mathcal{L}^{x'}}^{-\top} \Lambda_1^{-1} \end{bmatrix}, \\
\psi_{\mathcal{L}^x, \mathcal{L}^{x'}}(\mathbf{x})^\top &= \phi_{\mathcal{L}^x}(\mathbf{x})^\top \begin{bmatrix} I_{n_u} \\ \Lambda_2 \Phi_{2,\mathcal{L}^{x'}}^\top \Phi_{1,\mathcal{L}^{x'}}^{-\top} \Lambda_1^{-1} \end{bmatrix}.
\end{aligned}$$

246 Now, the ill-conditioning due to the dangerous Λ_1^{-1} term, introduced by applying \mathbf{K}^{-1} , is
247 removed analytically through the Λ_1 term present in $\mathbf{k}_u^\top(\mathbf{x}) = \psi_u^\top(\mathbf{x}) \Lambda_1 \Phi_1^\top$ or $\mathbf{k}_f^\top(\mathbf{x}) =$
248 $\psi_f^\top(\mathbf{x}) \Lambda_1 \Phi_1^\top$. Also, it should be noted that computation of the posterior variances is sub-
249 ject to the same ill-conditioning as any expression involving \mathbf{K}^{-1} ; this ill-conditioning can be
250 similarly resolved with

$$\begin{aligned}
v_u^2(\mathbf{x}_0) &= k_{uu}(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}_u^\top(\mathbf{x}_0) \mathbf{K}^{-1} \mathbf{k}_u(\mathbf{x}_0) = k(\mathbf{x}_0, \mathbf{x}_0) - \psi_u^\top(\mathbf{x}_0) \Psi^{-1} \mathbf{k}_u(\mathbf{x}_0), \\
v_f^2(\mathbf{x}_0) &= k_{ff}(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}_f^\top(\mathbf{x}_0) \mathbf{K}^{-1} \mathbf{k}_f(\mathbf{x}_0) = \mathcal{L}_{\mathbf{x}_0}^\zeta \mathcal{L}_{\mathbf{x}_0}^\zeta k(\mathbf{x}, \mathbf{x}) - \psi_f^\top(\mathbf{x}_0) \Psi^{-1} \mathbf{k}_f(\mathbf{x}_0). \tag{19}
\end{aligned}$$

251 As a result, the main portion of the ill-conditioning can be resolved in the posterior variances
252 $v_u^2(\mathbf{x})$ and $v_f^2(\mathbf{x})$ using (19).

253 **3.2 Stable likelihood function computation**

254 For the likelihood function defined in (6), when \mathbf{K} becomes ill-conditioned (e.g., for small
 255 ε , since the use of “flatter” kernels will lead to more and more similar entries in the system
 256 matrix \mathbf{K} , and therefore to potential numerical problems due to the ill-conditioning of \mathbf{K}
 257 [16]), computing $\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}$ and $\det \mathbf{K}$ by standard methods (such as Cholesky factorization)
 258 may be inaccurate, leaving us unable to use the MLE to judge the validity of the small ε
 259 for our estimation, despite the fact that (18) would allow us to make posterior predictions
 260 accurately. Using the HS-SVD decomposition of the covariance matrix \mathbf{K} in (17), we can
 261 follow a similar strategy as in Subsection 3.1 to the stable computation of log likelihood
 262 function (6). Computing $\log \det \mathbf{K}$ is relatively straightforward using $\mathbf{K} = \mathbf{\Psi} \mathbf{\Lambda}_1 \mathbf{\Phi}_1^\top$, as

$$\log |\mathbf{K}| = \log |\mathbf{\Psi}| + \log |\mathbf{\Lambda}_1| + \log |\mathbf{\Phi}_1^\top|. \quad (20)$$

263 It should be noted that $\mathbf{\Lambda}_1$ is diagonal, and therefore the very small eigenvalues can be handled
 264 by taking their logarithms. A similar strategy will allow us to compute $\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}$. Applying
 265 $\mathbf{\Psi} \mathbf{b} = \mathbf{y}$ and the Hilbert–Schmidt SVD (17) to $\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}$ gives

$$\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} = (\mathbf{\Psi} \mathbf{b})^\top (\mathbf{\Psi} \mathbf{\Lambda}_1 \mathbf{\Phi}_1^\top)^{-1} \mathbf{\Psi} \mathbf{b} = \mathbf{b}^\top \mathbf{\Psi}^\top \mathbf{\Phi}_1^{-\top} \mathbf{\Lambda}_1^{-1} \mathbf{b}. \quad (21)$$

Now, we are in a situation to find the Hilbert–Schmidt decomposition of the negative log marginal likelihood in (6) using (20) and (21) as follows

$$\begin{aligned} \text{NLML}_{\text{HS}}(\zeta, \theta) &= \frac{1}{2} \mathbf{b}^\top \mathbf{\Psi}^\top \mathbf{\Phi}_1^{-\top} \mathbf{\Lambda}_1^{-1} \mathbf{b} \\ &\quad + \frac{1}{2} (\log |\mathbf{\Psi}| + \log |\mathbf{\Lambda}_1| + \log |\mathbf{\Phi}_1^\top|) + \frac{N}{2} \log 2\pi. \end{aligned} \quad (22)$$

The (hyper)parameters θ and ζ can be trained by employing a Quasi-Newton optimizer L-BFGS to minimize the negative log marginal likelihood [6]. To set the hyperparameters by minimizing the negative log marginal likelihood, we seek the partial derivatives of the marginal likelihood with respect to the (hyper)parameters. The partial derivatives of the marginal likelihood with respect to the (hyper)parameters can be calculated by the relation

$$\frac{\partial}{\partial \theta_j} \text{NLML}(\zeta, \theta) = -\frac{\partial}{\partial \theta_j} \log p(\mathbf{y} | \zeta, \theta) = -\frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \mathbf{K}^{-1}) \frac{\partial \mathbf{K}}{\partial \theta_j} \right),$$

where $\boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{y}$ [9]. Equivalent relation for the partial derivatives of Hilbert–Schmidt decomposition of the negative log marginal likelihood in (22) is as

$$\frac{\partial}{\partial \theta_j} \text{NLML}_{\text{HS}}(\zeta, \theta) = -\frac{1}{2} \text{tr} \left(\boldsymbol{\beta} (\mathbf{b} \mathbf{b}^\top \boldsymbol{\beta}^\top - \mathbf{\Psi}^{-1}) \frac{\partial (\mathbf{\Psi} \mathbf{\Lambda}_1 \mathbf{\Phi}_1^\top)}{\partial \theta_j} \right),$$

266 where $\beta = (\Lambda_1 \Phi_1^\top)^{-1}$.

267 In the following, we explain how the HS-SVD alternative bases can reduce the instabilities
 268 in (hyper)parameter estimation and posterior prediction using the GPs. The first step in
 269 creating a stable basis is to study the structure of the system in (13) and asking the question
 270 Why is this system ill-conditioned?

271 It would seem that because the eigenfunctions φ_n are orthogonal the matrix Φ in relation
 272 (13) should be relatively well behaved. Therefore, the ill-conditioning appears to instead
 273 originate in the diagonal matrix Λ which contains block matrices Λ_1 and Λ_2 , whose values are
 274 the eigenvalues of \mathbf{K} and its 2-norm condition number is as

$$\text{cond}(\Lambda) = \frac{\lambda_1}{\lambda_M}. \quad (23)$$

275 This condition number is not directly relevant, since the entire Λ matrix is never inverted, but
 276 it serves to give an idea of the delicate nature of the \mathbf{K} matrix.

In fact, the rate of decay of the Hilbert-Schmidt eigenvalues determines the smoothness of the kernel k_{uu} . Therefore, in the course of using more smooth kernels, the problem of instability will be raised more seriously according to relation the (23). This connection between the ill-conditioning of the system and the smoothness of the kernel has been studied in [16]. It would appear then that the presence of the Λ matrix is the main source of ill-conditioning. The key step in removing the ill-conditioning is to write the component matrices that appear in the eigen-decomposition of \mathbf{K} as (17). Therefore, instead of solving the standard systems (6) and (9) which have the potential of yielding inaccurate and unreliable results, we now solve the transformed systems (18), (19) and (22), which are more numerically stable. In addition, we always make sure to simultaneously apply the matrices Λ_1 and Λ_2 required in the computation of the corrector matrices $[\Lambda_2 \Phi_2^\top \Phi_1^{-\top} \Lambda_1^{-1}]$ and $[\Lambda_2 \Phi_{2,\mathcal{L}^x}^\top \Phi_{1,\mathcal{L}^x}^{-\top} \Lambda_1^{-1}]$. For example, in MATLAB software, $[\Lambda_2 \Phi_2^\top \Phi_1^{-\top} \Lambda_1^{-1}]$ is computed as follows

$$\text{bsxfun}(@rdivide, \Lambda_2, \Lambda_1) .* (\Phi_2 / \Phi_1').$$

277 This minimizes the chances of producing an overflow or underflow error caused, respectively,
 278 by dividing the small eigenvalues at the end of Λ_1 or multiplying by the even tinier ones in
 279 Λ_2 . On the other hand, using the QR decomposition (16) is sometimes preferable to directly
 280 invert $\Phi_{1,\mathcal{L}^x}^\top$, because it may be a more stable computation than LU factorization, which may
 281 be preferable depending on the scale of the various eigenfunctions [31]. Therefore, the use

282 of alternative HS-SVD bases (with or without the optional QR step) allows us to isolate the
 283 ill-conditioning primarily in the Λ_1 factor and resolve the corrector matrices safely and recom-
 284 mended for the stable computation of optimal parameters and operators posterior predictions,
 285 especially in the kernel flat-limit.

286 4 Numerical experiments

287 The main purpose of the numerical results is to point out the efficiency, validity and stability
 288 of the method presented in this paper to estimate (hyper)parameters and operators posterior
 289 prediction in a numerically stable way as presented here. In the following, we present a series
 290 of numerical experiments that demonstrate the effectiveness of our approach. We have imple-
 291 mented the process of (hyper)parameters estimation by employing an L-BFGS optimization
 292 method and posterior prediction of various differential and integral linear operators in new
 293 points as well. Using various figures, the numerical stability of calculation of MLE in the HS-
 294 SVD method in comparison with the direct method is shown. Also, to show the accuracy and
 295 numerical stability of posterior mean of the HS-SVD method in comparison with the standard
 296 method, the maximum absolute error (absolute error) graphs of the posterior means and the
 297 condition number of the covariance matrices for both methods are presented. The absolute
 298 error between the exact function $u(\mathbf{x})$ and the predicted mean $\bar{u}(\mathbf{x})$ is described as below:

$$\max_{1 \leq i \leq N} |u(\mathbf{x}_i) - \bar{u}(\mathbf{x}_i)|.$$

299 It should also be noted that in the following examples, we have ignored the constant term
 300 $\frac{N}{2} \log 2\pi$ in the calculation of the likelihood function (6). In fact, in calculating of NLML,
 301 we have used the relations $Dmle = NLML(\boldsymbol{\zeta}, \boldsymbol{\theta}) - \frac{N}{2} \log 2\pi$ for standard (direct) computation
 302 (labeled in the figures with Direct Likelihood) and $HSmle = NLML_{HS}(\boldsymbol{\zeta}, \boldsymbol{\theta}) - \frac{N}{2} \log 2\pi$ for HS-
 303 SVD computation (labeled in the figures with HS-SVD Likelihood) as the likelihood criterion.
 304 Note that in all figures and tables to compare "Direct Likelihood" and "HS-SVD likelihood"
 305 the NLML (Negative Log Marginal Likelihood) quantity is used such that a higher direct
 306 likelihood means a lower NLML.

307 **Example 1** Consider the one dimensional fractional equation

$$\mathcal{L}_x^\alpha u(x) = {}_{-\infty}^{RL} D_x^\alpha u(x) - u(x) = f(x),$$

308 where $\alpha \in \mathbb{R}$ and ${}_{-\infty}^{RL}D_x^\alpha$ is defined in the Riemann-Liouville sense [32]. As pointed out
309 by Raissi et al. in [6], "Fractional operators often arise in modeling anomalous diffusion
310 processes. Their non-local behavior poses serious computational challenges as it involves
311 costly convolution operations for resolving the underlying non-Markovian dynamics". It should
312 be noted that ${}_{-\infty}^{RL}D_x^\alpha k_{uu}$, ${}_{-\infty}^{RL}D_y^\alpha [{}_{-\infty}^{RL}D_x^\alpha k_{uu}]$ and ${}_{-\infty}^{RL}D_x^\alpha \varphi(x)$ were obtained using generalized
313 Gauss-Laguerre quadrature method, involving a weight function of the form $x^{\alpha_{gGL}} e^{-x}$ for
314 $\alpha_{gGL} > -1$ as

$$\int_0^\infty f(x) dx = \int_0^\infty x^{\alpha_{gGL}} e^{-x} [e^x x^{-\alpha_{gGL}} f(x)] dx \simeq \sum_i^n w_i e^{x_i} x_i^{-\alpha_{gGL}} f(x_i). \quad (24)$$

315 We use the Golub-Welsch algorithm to find the nodes, but we compute the weights by evalu-
316 ating the generalized Gauss-Laguerre polynomial at these nodes for higher relative accuracy.
317 In practice, it is essential for α_{gGL} to match the fractional part of power of the monomial
318 in the integrand f , as the remainder yields a smooth function (for more information, see
319 [23]). However, the presented stable machine learning method using the GP overcomes these
320 computational challenges, and we can seamlessly handle all such linear cases without any
321 modifications. In this example, we have used a generalized periodic spline kernel on a set of
322 one-dimensional data in the interval $[0, 1]$. This kernel has eigenvalues and eigenfunctions

$$\lambda_n = \begin{cases} ((2j\pi)^2 + \varepsilon^2)^{-\beta} & n = 2j - 1, \\ ((2j\pi)^2 + \varepsilon^2)^{-\beta} & n = 2j, \end{cases}$$

$$\varphi_n(x) = \begin{cases} \sqrt{2} \sin(2j\pi x) & n = 2j - 1, \\ \sqrt{2} \cos(2j\pi x) & n = 2j, \end{cases}$$

324 for $j = 1, 2, \dots$. It should also be emphasized that we have used the roots of the squared
325 exponential kernel eigenfunctions as training points labelled as "Roots". We create data
326 values $\{x_u, y_u\}, \{x_f, y_f\}$ by sampling the test function

$$u(x) = \frac{1}{2} e^{-2\pi i x} \left(\frac{(2\pi + i)e^{4i\pi x}}{-1 + (2i\pi)\sqrt{2}} + \frac{2\pi - i}{-1 + (-2i\pi)\sqrt{2}} \right)$$

327 and $f(x) = 2\pi \cos(2\pi x) - \sin(2\pi x)$, for $\alpha = \sqrt{2}$ which the training points chosen in the interval
328 $[0, 1]$. It should be noted that to illustrate the instability and computational challenges of the
329 direct approach and the efficiency of the HS-SVD approach, various number of Roots data

330 points are chosen in the interval $[0, 1]$. In Table 2, the absolute error, condition number,
 331 the optimal values of α and likelihood criterion are presented for different values of (n_u, n_f)
 332 of Roots data points. It should be made clear here that in Table 2, the optimal values of
 333 (hyper)parameters are obtained and then the absolute error for the posterior mean of the
 334 fractional operators in new points using the optimal (hyper)parameters are reported. The
 335 likelihood criterion is evaluated for α spaced uniformly in $[.1, 5]$, ε spaced logarithmically
 336 in $[.1, 10]$ and $\beta = 3$ using both the direct approach (labeled Direct Likelihood) based on
 337 Cholesky decomposition and HS-SVD method (labeled HS-SVD Likelihood) in Fig. 1 and
 338 and specifically 2 which provides more accurate and more stable results. A similar pattern
 339 is observed for other values of β . In addition, as the number of data points increases, the
 340 instability of the direct method increases. This issue is not brought to the forefront here
 341 because it seems to be redundant. In fact, we have demonstrated the training process to learn
 342 the optimal α and ε parameters simultaneously using both direct and HS-SVD techniques, in
 343 Figs. 1 and 2. we have also made posterior predictions at $N_{eval}=100$ evenly spaced points
 344 in the domain to calculate the error value logarithmically for different values of α and ε with
 345 $\beta = 3$ to show the validity of the HS-SVD method. Due to the error figures, we have found out
 346 that the HS-SVD method correctly determines a region for an "optimal" (hyper) parameter
 347 estimate α and ε , while we have noticed instability in the direct approach. It is clear that by
 348 increasing the number of training points, the direct approach in parameters estimation and
 349 operators posterior prediction loses accuracy and suffers a complete breakdown because \mathbf{K}^{-1}
 350 is too ill-conditioned. Absolute error between the true fractional order α and the estimated
 351 one (top) and also between the exact function $u(x)$ and the predicted mean $\bar{u}(x)$ (middle) in
 352 the logarithmic scale as a function of the total number of training points for $u(x)$ and $f(x)$
 353 denoted by n_u and n_f and condition number of covariance matrix \mathbf{K} and matrix Ψ (bottom)
 354 using both approaches demonstrated in Fig. 3. As Figs. 1,2 and 3 show, in contradiction
 355 with the notion that increasing the number of data points leads to an increase in the accuracy
 356 of calculations, we see that this does not happen in direct method because of ill-conditioning
 357 covariance matrix \mathbf{K}^{-1} . While in HS-SVD method, by increasing the number of data points,
 358 in addition to maintaining stability, the accuracy of the results also increases. It is apparent
 359 that the stably computed likelihood parametrization criterion, HSmle, identifies a region for
 360 an "optimal" (hyper)parameters estimate that matches the region of smallest error. According
 361 to the Figs. 1, 2 and 3, we can see the HS-SVD algorithm learns the parameter α and ε to have

362 "optimal" values, while the direct maximum likelihood estimator does not precisely locate it.
 363 That's because the MLE direct computation loses accuracy and suffers a complete breakdown
 364 because \mathbf{K}^{-1} is too ill-conditioned. Furthermore, that indicates the satisfactory performance
 365 of the method presented in this paper for different trainings.

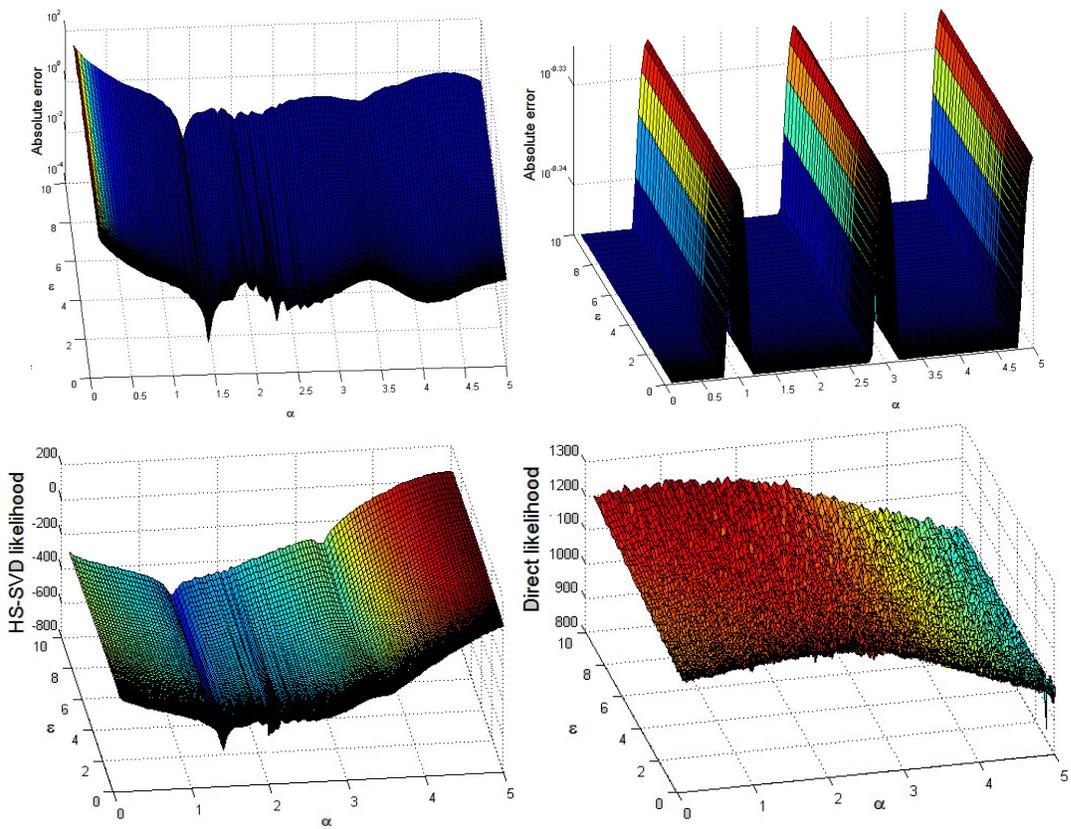


Figure 1: Comparison of the Negative Log Marginal likelihood (NLML) criterion and the error of the posterior mean computed by the direct method (right) and HS-SVD method (left) for $\beta = 3$ and different values ϵ and α . The top row shows the error of the posterior mean based on $n_u = n_f = 10$ Roots data points using generalized periodic spline kernel displayed. The bottom row displays the corresponding likelihood estimates for Example 1.

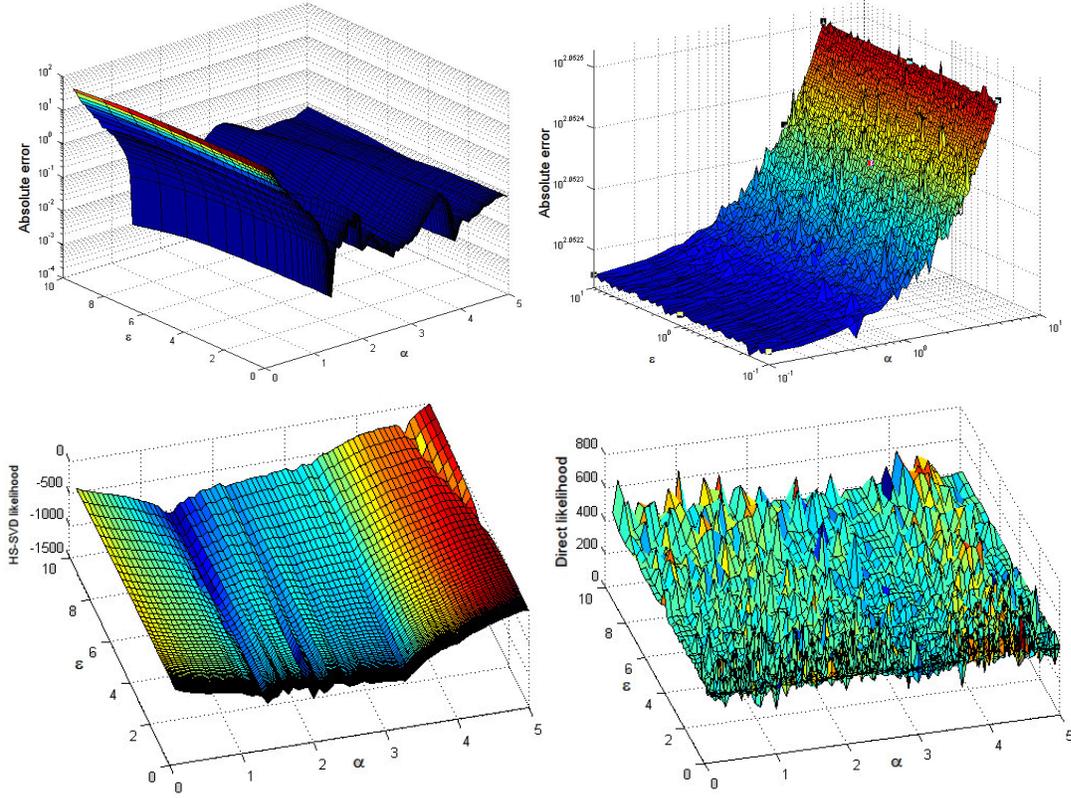


Figure 2: Comparison of the Negative Log Marginal likelihood (NLML) criterion and the error of the posterior mean computed by the direct method (right) and HS-SVD method (left) for $\beta = 3, n_u = n_f = 20$ and different values ε and α for Example 1.

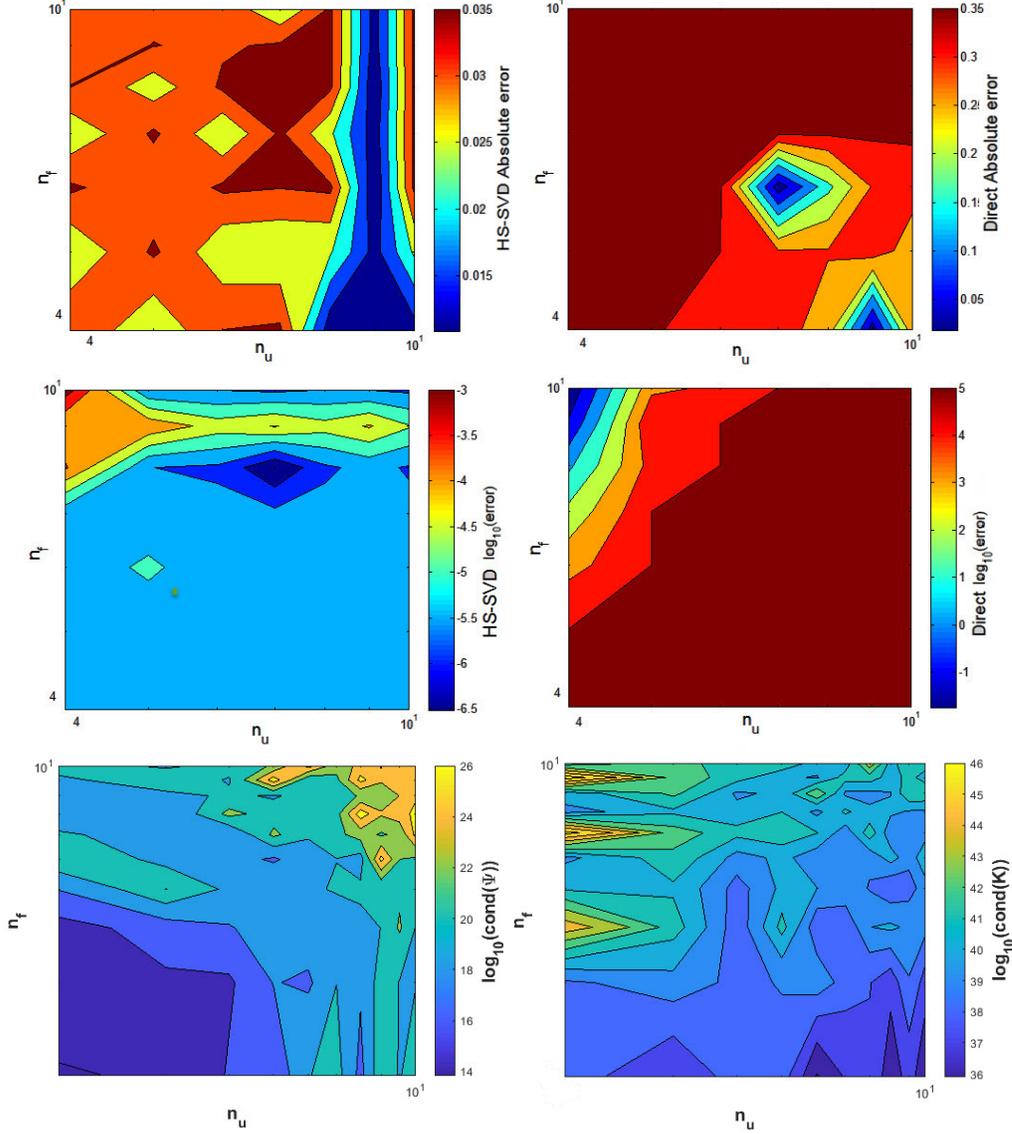


Figure 3: Absolute error between the true fractional order α and the estimated one (top), the exact function $u(x)$ and the predicted mean $\bar{u}(x)$ in the logarithmic scale as a function of the total number of training points for $u(x)$ and $f(x)$ (middle), denoted by n_u and n_f is shown with both methods. Condition number of covariance matrix \mathbf{K} defined in (17) and matrix Ψ (bottom) with $\beta = 3$ in the logarithmic scale is demonstrated using Roots data points for Example 1.

Table 2: Parameter estimation α and absolute error (error) of fractional operators posterior prediction for optimal (hyper)parameters obtained with direct and HS-SVD methods using different number of Roots data points for Example 1.

(n_u, n_f)	HS-SVD method				Direct method			
	α	HSmle	cond(Ψ)	error	α	Dmle	cond(\mathbf{K})	error
(5, 5)	4.8651	-166.3783	5.5679e+16	0.4518	5.000	568.6083	2.7907e+41	4.7185e-01
(10, 10)	2.0353	-577.3515	2.3678e+18	0.0938	5.0000	939.1081	1.4683e+44	4.7185e-01
(20, 20)	1.4043	-1.0783e+03	1.4259e+22	0.0057	3.8715	195.8455	7.2055e+45	1.1283e+02

366 **Example 2** Consider the one dimensional fractional equation

$$\mathcal{L}_x^\alpha u(x) = {}_0^C D_x^{75} u(x) + \alpha u(x) = f(x),$$

367 where $\alpha \in \mathbb{R}$ and ${}_0^C D_x^\alpha$ are defined in the Caputo sense [32]. In this example, we have used
368 periodic spline kernel on a set of one-dimensional data. Also, the fractional derivatives of the
369 kernel and ${}_0^C D_x^{75} k_{uu}$, ${}_0^C D_y^{75} [{}_0^C D_x^{75} k_{uu}]$ and ${}_0^C D_x^{75} \varphi(x)$ are approximated based on power series
370 expansion with Maple software by the command “fracdiff” as “fracdiff($k_{uu}(x, y)$, α , method =
371 series, method-options = [about= a]”. The optional parameters for this method to be specified
372 in method-options, are about= a and order = o. The value of a specifies the point on how
373 to expand the series and o specifies the accuracy or order of the series. For convenience, the
374 order o = 20 and the starting point of the interval a = 0 are considered.

The periodic spline kernel has eigenvalues and eigenfunctions

$$\lambda_n = \begin{cases} (2j\pi)^{(-2\beta)} & n = 2j - 1, \\ (2j\pi)^{(-2\beta)} & n = 2j, \end{cases}$$

$$\varphi_n(x) = \begin{cases} \sqrt{2} \sin(2j\pi x) & n = 2j - 1, \\ \sqrt{2} \cos(2j\pi x) & n = 2j, \end{cases}$$

375 for $\beta \in \mathbb{N}$, $j = 1, 2, \dots$. We simulate data values $\{x_u, y_u\}, \{x_f, y_f\}$ by sampling the test
376 function $u(x) = \sin(2\pi x)$ and

$$f(x) = 6.9320 \sqrt[4]{x} {}_1F_2(1.0; 0.625, 1.125; -9.86960 x^2) + 2 \sin(2\pi x)$$

377 for $\alpha = 2$ and $n_u = n_f = 10$ Chebyshev data points are chosen in the interval $[0, 1]$ where
 378 ${}_1F_2$ is the generalized hypergeometric function. Likelihood criterion is obtained for the values
 379 of $\beta \in \{1, 2, \dots, 10\}$ and α spaced logarithmically in $[10^{-4}, 10]$ using both the direct method
 380 (labelled **Direct Likelihood** in Fig. 4) based on Cholesky decomposition, and the more elaborate
 381 for relation (22) which provide more stable results (labelled **HS-SVD Likelihood**). The data is
 382 then used to make posterior predictions at `Neval` =100 evenly spaced points in the domain,
 383 and the absolute error compared to $u(x)$ displayed in Fig. 5. As Fig. 4 shows, the posterior
 384 mean (top row), as well as the likelihood criterion (bottom row), can be stably and reliably
 385 computed with the help of the HS-SVD (left column)-as compared to the direct approach,
 386 displayed in the right column, and computed using the standard methods such as the Cholesky
 387 decomposition. According to Fig. 4, we can observe that the HS-SVD algorithm correctly
 388 learns the parameter α and β to have "optimal" values, while the direct MLE is near the
 389 optimal error and it does not exactly locate it. Also, According to Fig. 5, we realize that the
 390 HS-SVD method determines a more precise region for "optimal" (hyper) parameters α and β .
 391 In addition, in Table 3 the absolute error and likelihood criterion for different values of (n_u, n_f)
 392 of Chebyshev data points are presented. In Table 3, the optimal value of the parameter α is
 393 presented, and the absolute error for the posterior mean of the fractional operators in new
 394 data points using the optimal (hyper) parameters are reported.

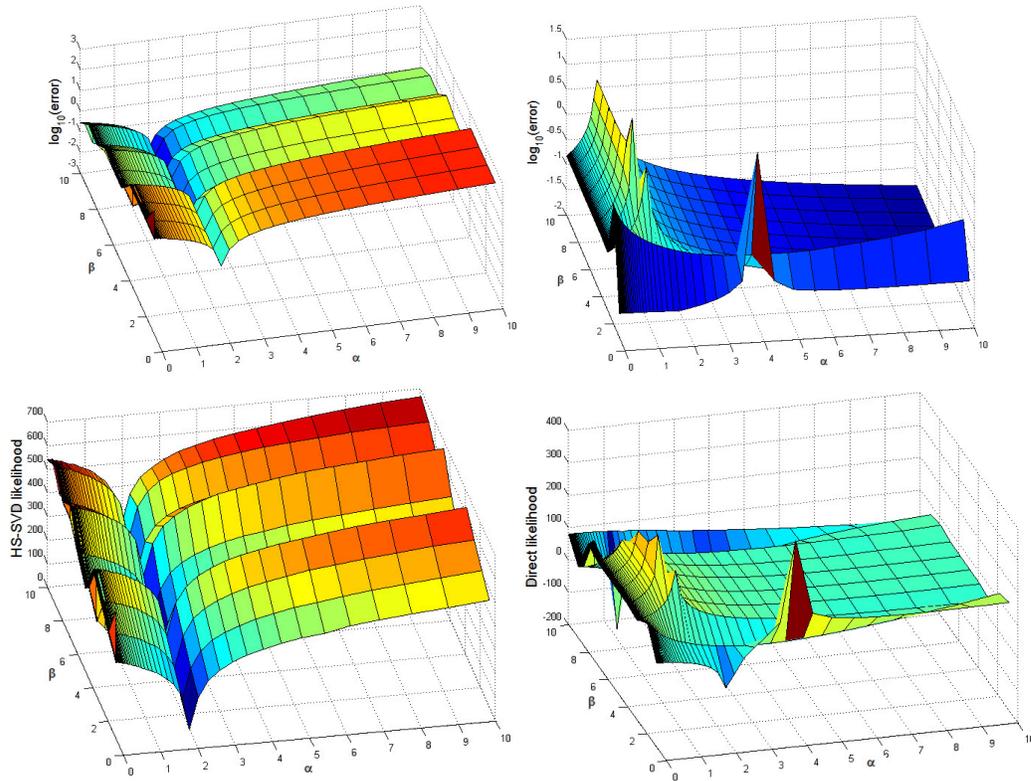


Figure 4: Comparison of Negative Log Marginal likelihood (NLML) criterion computed with both methods. The top row shows the error of the posterior mean based on Chebyshev data points using periodic spline kernel. The bottom row displays the corresponding likelihood estimates.

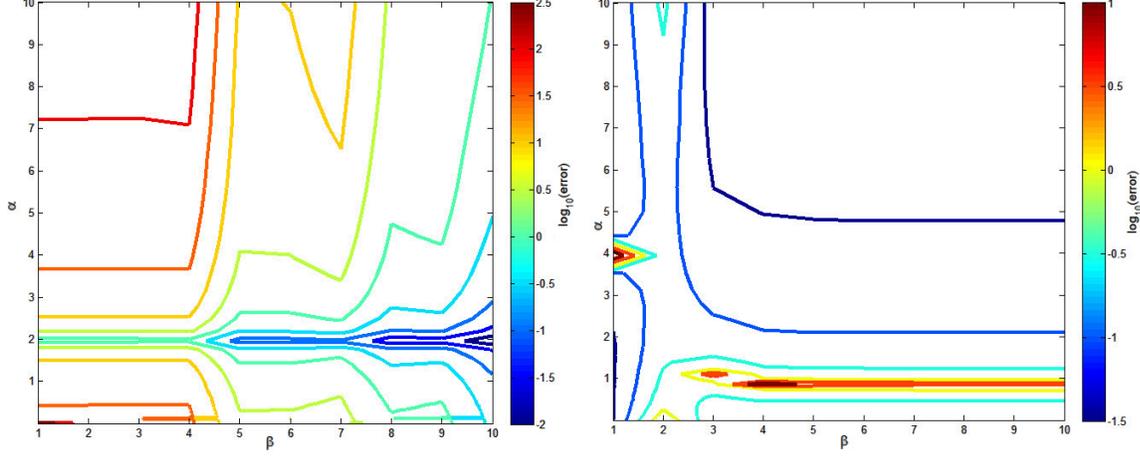


Figure 5: Absolute error for Example 2 using periodic spline kernel computed via the HS-SVD approach (left) and direct approach (right).

Table 3: Parameters estimation α and absolute error of fractional operators posterior prediction for optimal (hyper) parameters obtained with direct and HS-SVD methods by Chebyshev data points for Example 2.

	HS-SVD method				Direct method			
	α	HSmle	cond(Ψ)	error	α	Dmle	cond(\mathbf{K})	error
(5, 5)	1.9307	-46.6439	6.1427e+03	0.0152	0.7543	-26.5178	1.8353e+15	11.4753
(10, 10)	1.9307	-83.5802	2.6829e+05	0.0108	0.0910	-71.1816	7.8813e+16	0.5588
(20, 20)	1.9698	-125.819	1.9509e+07	0.0015	1.2068	-63.9671	4.0022e+17	0.3404

395 **Example 3** Consider the following differential equation,

$$\mathcal{L}_x^\alpha u(x) = \frac{d^2}{dx^2} u(x) + \frac{\alpha x}{x^2 + 1} \frac{d}{dx} u(x) + u(x) = f(x).$$

Note that the functions $u(x) = -2 \cos(4\pi x) + \sin(4\pi x)$ and

$$f(x) = 32\pi^2 \cos(4\pi x) - 16\pi^2 \sin(4\pi x) + \frac{\alpha x (8\pi \sin(4\pi x))}{x^2 + 1} \\ \frac{\alpha x (4\pi \cos(4\pi x))}{x^2 + 1} - 2 \cos(4\pi x) + \sin(4\pi x)$$

satisfy the equation. We create data values $\{x_u, y_u\}, \{x_f, y_f\}$ by sampling the test function $u(x)$ and $f(x)$ with $n_u = n_f = 40$ Chebyshev data points are chosen in the interval $[-1, 1]$ and $\alpha = 6$. In this example we have used Chebyshev kernel with eigenvalues and eigenfunctions

$$\lambda_n = \begin{cases} 1 - a & n = 0, \\ \frac{a(1-b)b^n}{b} & n = 1, 2, \dots \end{cases}$$

$$\varphi_n(x) = \sqrt{2 - \delta_{n0}} T_n(x),$$

396 where T_n are Chebyshev polynomials of degree n . Likelihood criterion is obtained for values
397 of $a = .5$ and b spaced uniformly in $[10^{-4}, .9]$ using both methods in Fig. 6. The data is
398 then used to make predictions at `Neval` =100 evenly spaced points in the domain, and the
399 absolute errors compared to $u(x)$ are displayed in Figs. 6 and 7. As Fig. 6 shows, the posterior
400 mean and the likelihood criterion can stably and reliably be computed with the help of the
401 HS-SVD -as compared to the direct approach, computed with the standard methods. It is
402 apparent that the stably computed likelihood parametrization criterion identifies a region for
403 an "optimal" (hyper) parameters estimate α and b that matches the region of the smallest
404 error in Figs. 6 and 7. Figure 8 indicates Maximum error between exact function and predicted
405 mean $\bar{u}(x)$, exact function and predicted mean $\bar{f}(x)$ in the logarithmic scale as a function of
406 the total number of training points $u(x)$ and $f(x)$, denoted by n_u and n_f , with both methods
407 and also, the condition number of covariance matrix \mathbf{K} and matrix Ψ based on Chebyshev
408 data points. It also reveals that by increasing the number of data points the accuracy of
409 calculations increases. Needless to say, this does not hold true in the direct method. As far
410 as the HS-SVD method is concerned, an increase in the number of data points leads to an
411 increase in the accuracy of the results and also the maintenance of the stability. Of course,
412 whenever \mathbf{K} is not severely ill-conditioned (usually this is true for kernels with a low level of
413 smoothness such as Matern kernels or compactly supported Wendland kernels) it is easier and
414 more convenient to work with its Cholesky factorization $\mathbf{K} = \mathbf{L}\mathbf{L}^\top$ as a fundamental tool in
415 matrix computations.

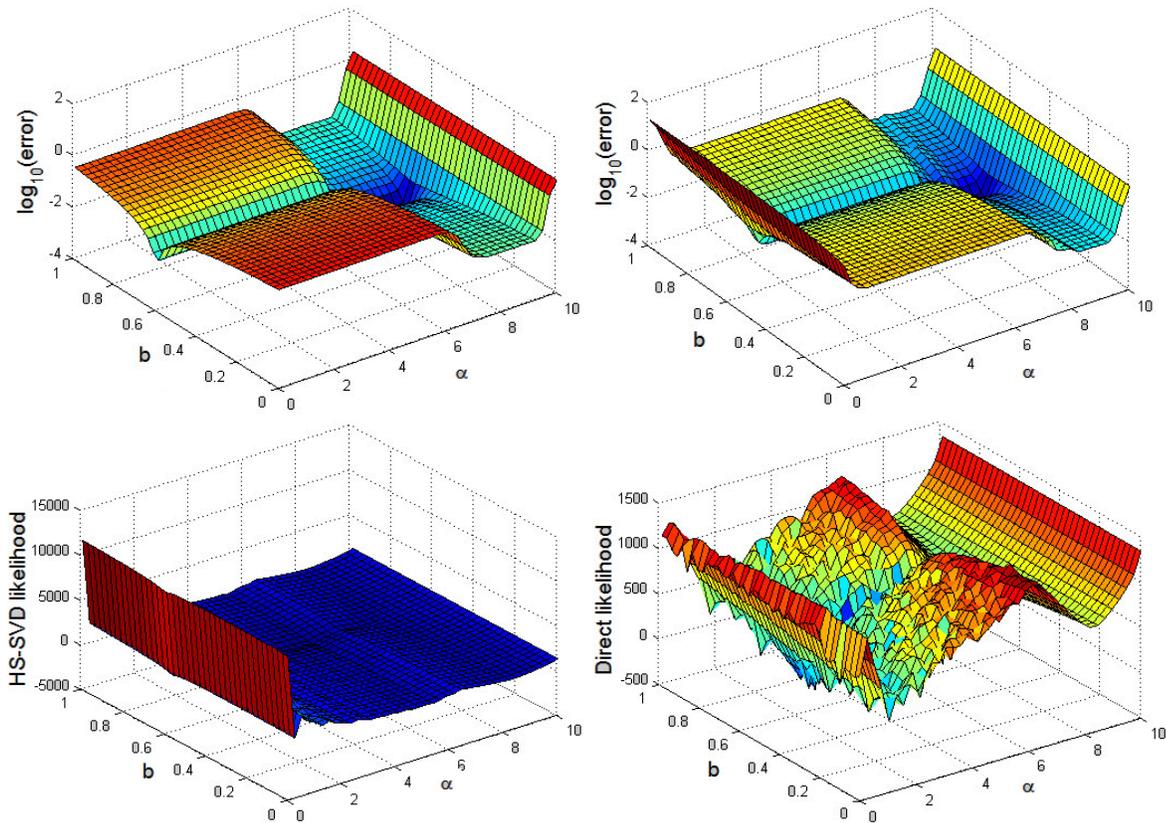


Figure 6: Comparison of Negative Log Marginal likelihood (NLML) criterion computed with both methods. The top row shows the error of the posterior mean based on Chebyshev data points using Chebyshev kernel. The bottom row displays the corresponding likelihood estimates.

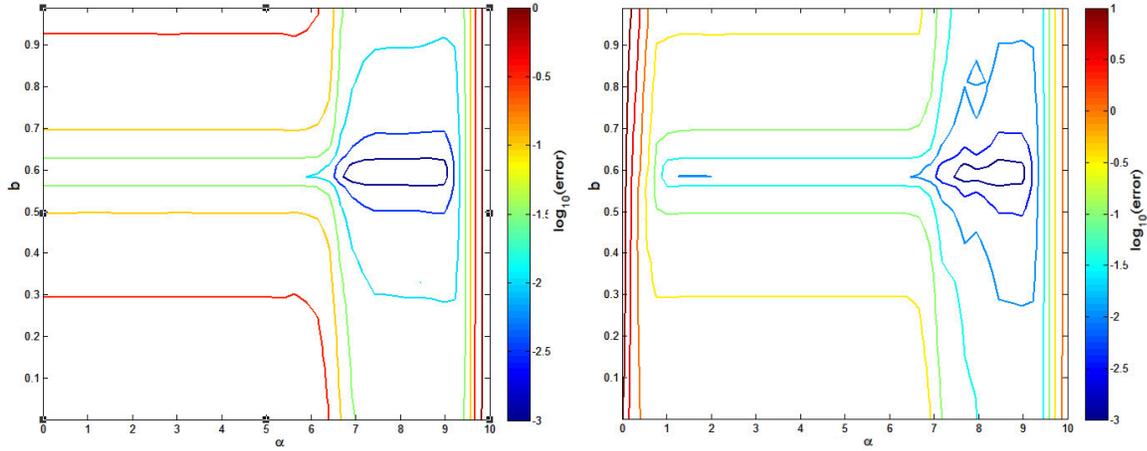


Figure 7: Absolute error for Example 3 using Chebyshev kernel computed via the HS-SVD approach (left) and direct approach (right).

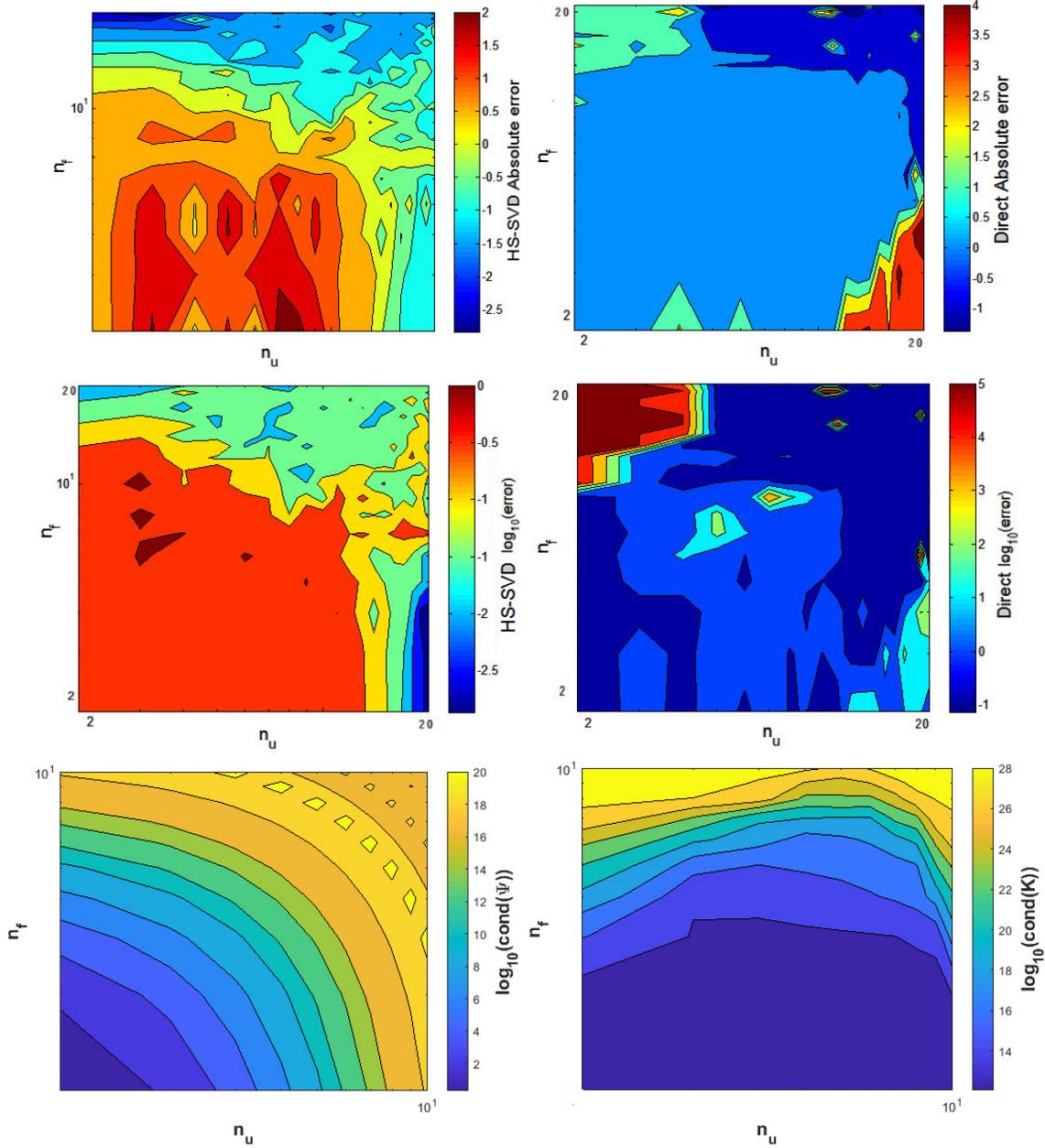


Figure 8: Maximum error between exact function $u(x)$ and predicted mean $\bar{u}(x)$ (top), exact function $f(x)$ and predicted mean $\bar{f}(x)$ (middle) in the logarithmic scale as a function of the total number of training points $u(x)$ and $f(x)$, denoted by n_u and n_f , are demonstrated with HS-SVD and direct methods. The condition number of covariance matrix \mathbf{K} and matrix Ψ (bottom) in the logarithmic scale are demonstrated using Chebyshev data points for Example 3.

416 **Example 4** Consider the following integro-differential equation,

$$\mathcal{L}_x^\alpha u(x) = \frac{d}{dx}u(x) + 2u(x) + \alpha \int_0^x u(t)dt = f(x).$$

417 The functions $u(x) = \sin(2\pi x)$ and

$$f(x) = 2\pi \cos(2\pi x) + 2\sin(2\pi x) - \alpha \frac{\cos(2\pi x) - 1}{2\pi}$$

418 satisfy the equation. We have used the data $\{x_u, y_u\}, \{x_f, y_f\}$ generated from $u(x)$ and $f(x)$
419 with $n_u = n_f = 20$ Chebyshev data points chosen in the interval $[0, 1]$ for $\alpha = 3$. We
420 have demonstrated the effectiveness of the HS-SVD method using the generalized periodic
421 spline kernel. Likelihood criterion is evaluated for $\beta = 3$ and α and ε uniformly spaced in
422 $[.01, 10]$. The data is then used to make predictions at `Neval=100` evenly spaced points in the
423 domain, and the absolute error is displayed in Fig. 9. Also, likelihood criterion is evaluated
424 for the values of ε logarithmically spaced in $[10^{-2}, 10^2]$ in Fig. 10. This data is then used
425 to make predictions for different values β at `Neval=100` evenly spaced points in the domain,
426 and the absolute errors are displayed in Fig. 10. It is apparent that the HS-SVD method
427 suffers no ill-conditioning. The maximum likelihood estimator is near the "optimal" error,
428 though it does not precisely locate it. It is clear that by increasing the values of β the MLE
429 direct computation loses accuracy and suffers a complete breakdown because \mathbf{K}^{-1} is too ill-
430 conditioned. Also, the absolute errors between the true parameter α and the estimated one,
431 between the exact functions $u(x)$ and $f(x)$ and the predicted means $\bar{u}(x)$ and $\bar{f}(x)$ and also
432 the condition number of covariance matrix \mathbf{K} and matrix Ψ using Chebyshev data points are
433 demonstrated in Fig. 11. In Table 4, using the optimal (hyper) parameters the absolute error,
434 the condition number, the optimal values of α and the likelihood criterion are presented for
435 different values of (n_u, n_f) .

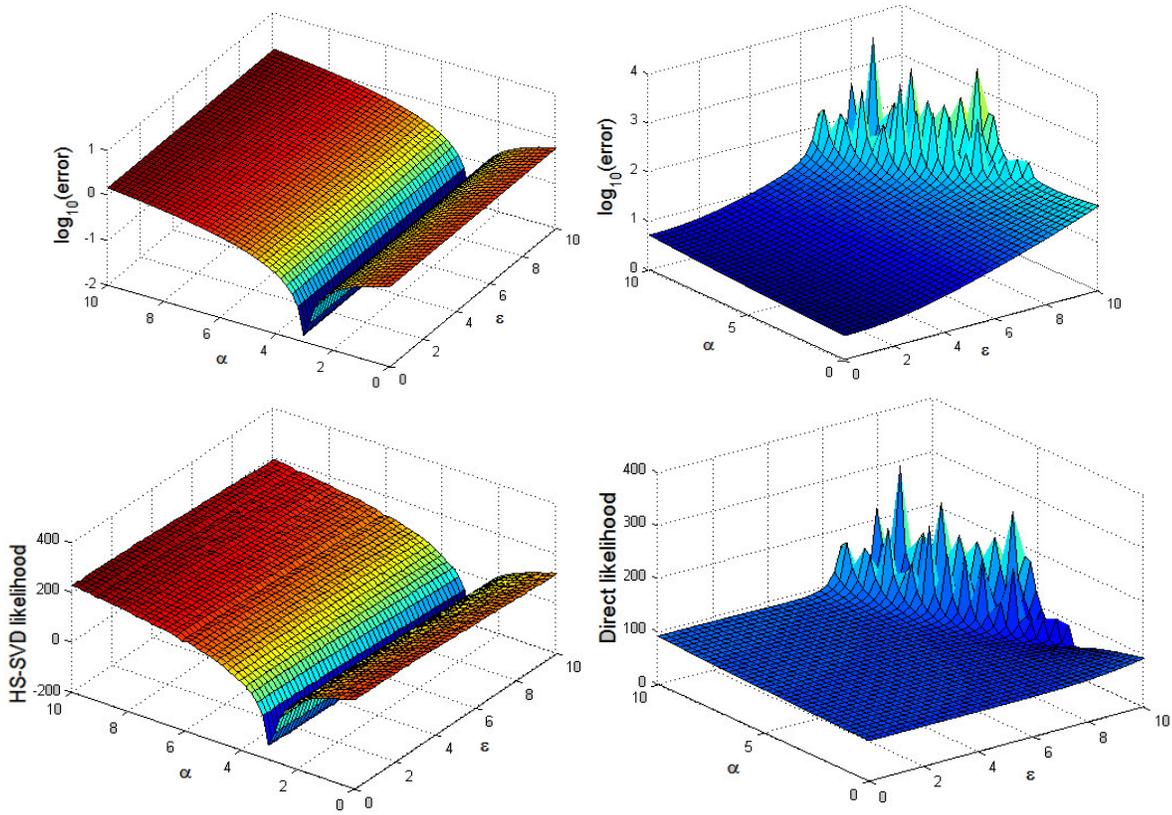


Figure 9: Comparison of Negative Log Marginal likelihood (NLML) criterion computed with both methods for Example 4. The top row shows the error of the posterior mean based on Chebyshev data points using Chebyshev kernel. The bottom row displays the corresponding likelihood estimates.

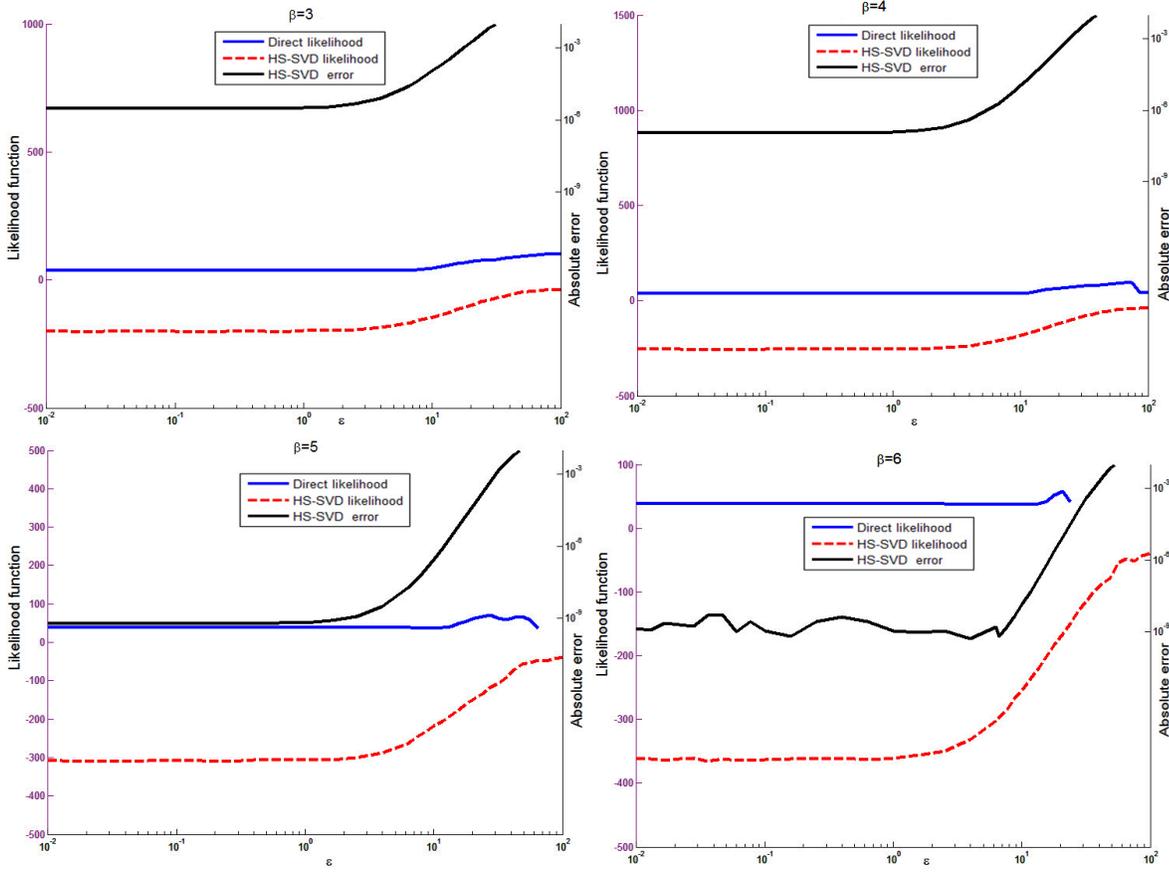


Figure 10: Comparison of the absolute error and MLE estimators of the optimal shape parameter ϵ for Example 4 using generalized periodic spline kernel computed via both approaches for different values of β .

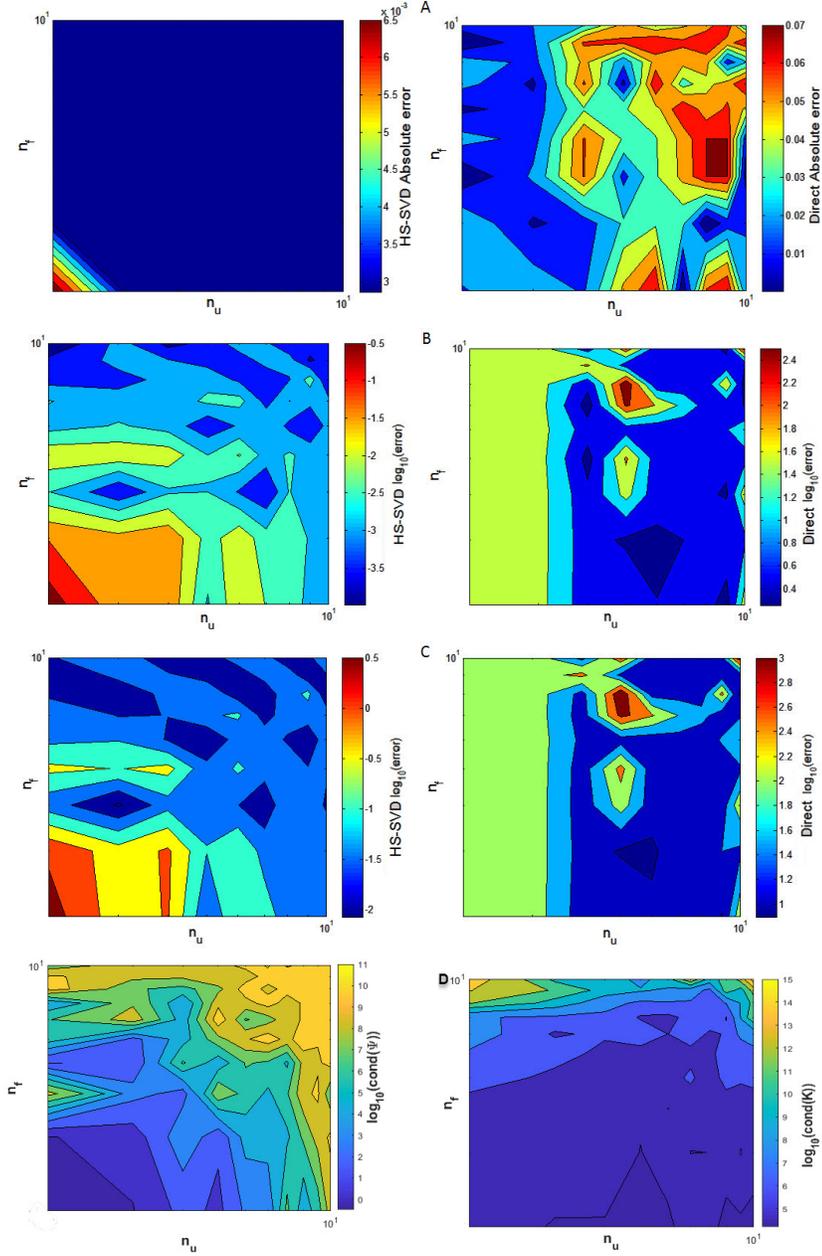


Figure 11: The absolute error between the true parameter α and the estimated one (A), between exact function $u(x)$ and predicted mean $\bar{u}(x)$ in the logarithmic scale (B), between exact function $f(x)$ and predicted mean $\bar{f}(x)$ in the logarithmic scale (C) are demonstrated with both methods. The condition number of covariance matrix \mathbf{K} and matrix Ψ in the logarithmic scale is demonstrated using Chebyshev data points for Example 4 (D).

Table 4: Parameter estimation α and the absolute error of operator posterior prediction for optimal (hyper)parameters obtained with both methods using Chebyshev data points for Example 4.

(n_u, n_f)	HS-SVD method				Direct method			
	α	HSmle	cond(Ψ)	error	α	Dmle	cond(\mathbf{K})	error
(5, 5)	2.98	-225.93	1.5863e+05	0.033	3.0665	-32.7392	5.6795e+09	120.0447
(10, 10)	3.008	-565.1892	2.0447e+11	0.0085	2.9365	-62.1680	2.4611e+13	590.8914
(20, 20)	3.0016	-309.2043	5.3467e+16	0.0013	3.8837	23.6461	1.5684e+20	35.3956

436 **Example 5 (Transport Equation)** Consider the following differential equation,

$$\mathcal{L}_{(x,t)}^\zeta u(x,t) = \frac{\partial u(x,t)}{\partial t} + \zeta \frac{\partial u(x,t)}{\partial x} = f(x,t).$$

437 The functions $u(x,t) = \exp(-x) \sin(2\pi t)$ and

$$f(x,t) = 2\pi \exp(-x) \cos(2\pi t) - \zeta \exp(-x) \sin(2\pi t)$$

satisfy the equation. We have used the data $\{\mathbf{x}_u = (x_u, t_u), \mathbf{y}_u\}, \{\mathbf{x}_f = (x_f, t_f), \mathbf{y}_f\}$ generated by $\mathbf{y}_u = u(x_u, t_u)$ and $\mathbf{y}_f = f(x_u, t_u)$ with $n_u = n_f = 32$ Halton data points chosen in the interval $[0, 1]^2$ for $\zeta = 1$. We have also demonstrated the effectiveness of the HS-SVD method using Squared Exponential kernel with eigenvalues and eigenfunctions

$$\lambda_n = \sqrt{\frac{\alpha^2}{\alpha^2 + \delta^2 + \varepsilon^2}} \left(\frac{\varepsilon^2}{\alpha^2 + \delta^2 + \varepsilon^2} \right)^{n-1} \quad n = 1, 2, \dots$$

$$\varphi_n(x) = \gamma_n e^{-\delta^2 x^2} H_{n-1}(\alpha \beta x),$$

where the H_n are Hermite polynomials of degree n , and

$$\beta = \left(1 + \left(\frac{2\varepsilon}{\alpha} \right)^2 \right)^{\frac{1}{4}}, \quad \gamma_n = \sqrt{\frac{\beta}{2^{n-1} \Gamma(n)}}, \quad \delta^2 = \frac{\alpha^2}{2} (\beta^2 - 1),$$

are constants such that they are defined in terms of the shape parameter ε and the parameter α in the weight function $\rho(x) = \frac{\alpha}{\sqrt{\pi}} e^{-\alpha^2 x^2}$ of the Hilbert-Schmidt integral operator. The multivariate case is easily obtained using the tensor product form of Squared Exponential kernel, i.e., for d-variate functions we have

$$\boldsymbol{\lambda}_n = \prod_{j=1}^2 \lambda_{n_j} = \prod_{j=1}^2 \sqrt{\frac{\alpha_j^2}{\alpha_j^2 + \delta_j^2 + \varepsilon_j^2}} \left(\frac{\varepsilon_j^2}{\alpha_j^2 + \delta_j^2 + \varepsilon_j^2} \right)^{n_j-1},$$

and

$$\varphi_n(\mathbf{x}) = \prod_{j=1}^2 \varphi_{n_j}(x_j) = \prod_{j=1}^2 \gamma_{n_j} e^{-\delta_j^2 x_j^2} H_{n_j-1}(\alpha_j \beta_j x_j),$$

438 where $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$. For more on tensor product kernels we point the reader to [16, 22].

439 Note that this formulation allows us to take different shape parameters ε_j and integral
 440 weights α_j for different space dimensions (i.e., k_{uu} may be an anisotropic kernel), or we can
 441 take them all equal, i.e., $\alpha_j = \alpha$ and $\varepsilon_j = \varepsilon$, $j = 1, 2$ (and then k_{uu} is isotropic or radial) [22].
 442 In the example, we restrict ourselves to using the same α_j and ε_j in all dimensions. Likelihood
 443 criterion is evaluated by a fixed value of $\alpha = 3$ and a grid of 625 different values of $[\zeta, \varepsilon]$ with
 444 each component uniformly spaced in $[.01, 10]$. The data is then used to make predictions at
 445 `Neval=81` evenly spaced points in the domain. The absolute error and likelihood criterion are
 446 displayed in Fig. 12. In Table 5, the absolute error, the condition number, the optimal values
 447 of ζ and the likelihood criterion are presented for different values of (n_u, n_f) .

Table 5: Parameter estimation ζ and the absolute error of operator posterior prediction for optimal (hyper)parameters obtained with both methods using Halton data points for Example 5.

(n_u, n_f)	HS-SVD method				Direct method			
	ζ	HSmle	cond(Ψ)	error	ζ	Dmle	cond(\mathbf{K})	error
(8, 8)	0.7	-60.2578	2.3783e+04	0.0577	10	0.2091	6.9541e+12	1.0475
(16, 16)	0.8577	-130.6324	5.6463e+04	0.0205	3.0665	-32.7392	5.6795e+14	0.5317
(32, 32)	1.0771	-987.4220	5.3892e+10	.002524	7.4296	-495.7323	2.7838e+17	1.03254

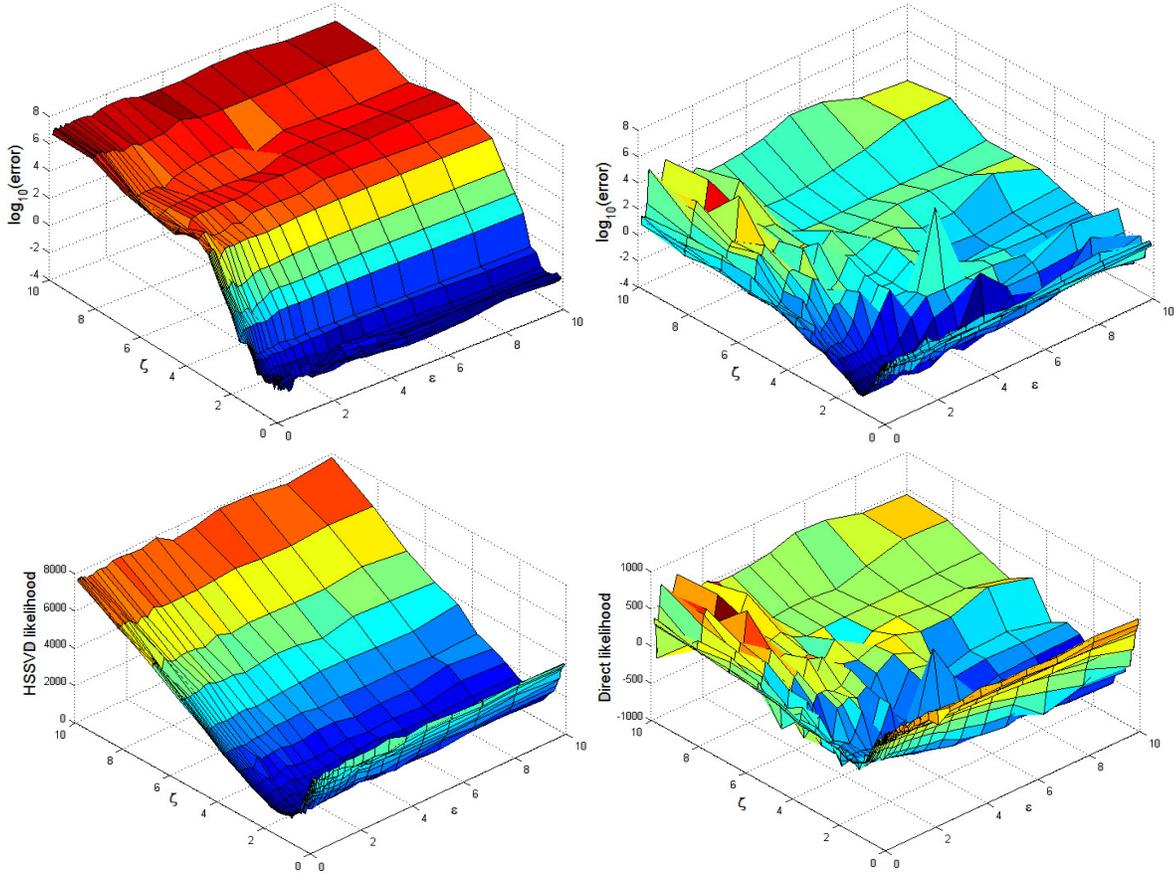


Figure 12: Comparison of Negative Log Marginal likelihood (NLML) criterion computed for $\alpha = 1$ with both methods for Example 5. The top row shows the error of the posterior mean based on Halton data points using Squared Exponential kernel. The bottom row displays the corresponding likelihood estimates.

448 5 Conclusion

449 In this paper, we made the unified framework to deal with parametric linear operational
 450 equation in (1) which was probabilistically approximated by employing the GPs, and was made
 451 computationally more stable and reliable by developing a novel computational strategy for
 452 more adaptive parameters and hyperparameters learning leading to more accurately predicting

453 operators at some unseen operational data points. The standard computational strategies
 454 suggested for solving the above linear inverse problems would usually become severely ill-
 455 conditioned to estimate the model parameters and hyperparameters, particularly when the
 456 number of data points increase, and the flat kernels (e.g., the squared exponential kernel
 457 with a small shape parameter, ε) are used. It is evident that by increasing the number of
 458 the observed data, the direct approach to estimating parameters and hyperparameters and
 459 predicting the operators would become less accurate, and it does not then correctly identify
 460 a region for "optimal" (hyper)parameters that match the region of the smallest error. As
 461 a result, it will suffer a complete breakdown, because \mathbf{K}^{-1} is too ill-conditioned. In this
 462 paper, we proposed an alternative computational approach using the HS-SVD at which the
 463 computation of the likelihood function becomes more stable, and the determination of the
 464 MLEs for optimal posterior predictions is now possible. It is thus apparent that the HS-SVD
 465 method correctly identifies a region for "optimal" (hyper)parameters that match the region of
 466 the smallest error. The proposed approach was validated by illustrating it in several benchmark
 467 problems and various kernels with different attributes. The numerical illustrations confirm the
 468 stability of the proposed method, particularly when the number of the observed data points
 469 increase, and the flat kernels are used when the standard computational strategies reviewed
 470 above would be unable to handle. The absolute error, condition number, the optimal values
 471 of (hyper)parameters and also likelihood criterion are presented for different number of data
 472 points for both approaches of direct and HS-SVD. We found out that the HS-SVD method
 473 correctly determines a region for an "optimal" (hyper) parameters estimate, while instability
 474 was witnessed in the direct approach. It was clear that by increasing the number of training
 475 points, the direct approach in parameters estimation and operators posterior prediction loses
 476 accuracy and suffers a complete breakdown because \mathbf{K}^{-1} is too ill-conditioned. Needless
 477 to say, the process of reducing the condition number, by the HS-SVD approach, also varies
 478 from one problem to another, depending on the type and nature of the model, the number
 479 of training points and even the type of the training points (e.g., Chebyshev points, Halton
 480 points, random points and and so forth). Therefore, the ill-conditioning improvement in the
 481 various examples may be different.

482 Future researchers are highly recommended to investigate the technique introduced in this
 483 paper for the problem of learning nonlinear operational equations as discussed in [5]. The main
 484 problem in nonlinear operational equations is that the stated condition "linear transformation

485 of a GP such as differentiation and integration remains GP” is no longer valid, and a linear
486 approximation of a nonlinear operator must be used. Furthermore, the results gained from
487 this study are limited to those special kernels for which a Hilbert-Schmidt SVD is available (
488 the collection of positive definite kernels and their known Mercer series are presented in [16]).
489 Therefore, understanding how a Mercer series with numerically computed eigenvalues and
490 eigenfunctions affects the quality of these computations can allow this strategy to be applied to
491 a wider range of kernels (which is currently limited by the availability of the Mercer series [17]).

492 The work in this paper is limited to low dimension in \mathbf{x} , as a tensor product basis is used
493 in terms of eigenfunctions. This removes an important advantage of kernel methods, which
494 are formally dimension-independent. Despite this limitation, we still have the advantage over
495 using regression in a spline basis that we can see the parameters of the operational equations
496 as kernel hyperparameters and estimate them using a machine learning strategy. At the same
497 time, in our future studies, we are trying to generalize the proposed method in this paper to
498 high dimensional problems without being dependent on the tensor product.

499 References

- 500 [1] Bender, E. A., *An Introduction to Mathematical Modeling*, New York: DOVER PUBLI-
501 CATIONS INC., Mineola, 1978.
- 502 [2] Vries, K. de, Nikishova, A., Czaja, B., Zavodszky, G. and Hoekstra, A. G., Inverse un-
503 certainty quantification of a cell model using a Gaussian process metamodel ,*Int. J.*
504 *Uncertainty Quantif.*, **10**(4), pp. 333-349, 2020.
- 505 [3] Narayan, A., Yan, L. and Zhou, T., Optimal design for kernel interpolation: Applications
506 to uncertainty quantification, *J. Comput. Phys.*, **430**, pp. 1-20, 2021.
- 507 [4] Qin, T., Chen, Z., Jakeman, J. D. and Xiu, D., Deep learning of parameterized equations
508 with applications to uncertainty quantification, *Int. J. Uncertainty Quantif.*, **11**(2), pp.
509 63-82, 2021.
- 510 [5] Raissi, M., Perdikaris, P. and Em Karniadakis, G., Numerical Gaussian Processes for
511 Time-Dependent and Nonlinear Partial Differential Equations, *SIAM J. Sci. Comput.*,
512 **40**(1), pp. 172-198, 2018.

- 513 [6] Raissi, M., Perdikaris, P. and Em Karniadakis, G., Machine learning of linear differential
514 equations using Gaussian processes, *J. Comput. Phys.*, **348**(1), pp. 683-693, 2017.
- 515 [7] Mueller, J. L. and Siltanen, S., *Linear and Nonlinear Inverse Problems with Practical*
516 *Applications*, Computational Science and Engineering, vol. 10. Philadelphia, PA, USA:
517 SIAM, 2012.
- 518 [8] Neto, F. D. M. and Neto, A. J. d. S., *An Introduction to Inverse Problems with Applica-*
519 *tions*, Springer-Verlag Berlin Heidelberg, 2013.
- 520 [9] Williams, C. K. and Rasmussen, C. E., *Gaussian processes for machine learning*, the MIT
521 Press, 2006.
- 522 [10] Murphy, K. P., *Machine learning: a probabilistic perspective*, MIT press, 2012.
- 523 [11] Vapnik, V., *The nature of statistical learning theory*, Springer Science & Business Media,
524 2013.
- 525 [12] Guo, M. and Hesthaven, J. S., Reduced order modeling for nonlinear structural analysis
526 using Gaussian process regression, *Comput. Meth. Appl. Mech. Eng.*, **341**(1), pp. 807-826,
527 2018.
- 528 [13] Tikhonov, A., Solution of incorrectly formulated problems and the regularization method,
529 *Soviet Math. Dokl.*, **5**, pp. 1035-1038, 1963.
- 530 [14] Poggio, T. and Girosi, F., Networks for approximation and learning, *Proceedings of the*
531 *IEEE*, **78**, pp. 1481-1497, 1990.
- 532 [15] Fasshauer, G. E., *Meshfree Approximation Methods with Matlab*, Interdisciplinary Math-
533 ematical Sciences, Vol. 6, World Scientific Publishing, Singapore, 2007.
- 534 [16] Fasshauer, G. E. and McCourt, M., *Kernel-based Approximation Methods using MAT-*
535 *LAB*, *Interdisciplinary Mathematical Sciences*, World Scientific Publishing, Singapore,
536 2015.
- 537 [17] McCourt, M., and Fasshauer, G. E., *Stable Likelihood Computation for Gaussian Ran-*
538 *dom Fields*, *Recent Applications of Harmonic Analysis to Function Spaces, Differential*
539 *Equations, and Data Science. Applied and Numerical Harmonic Analysis*, Birkher, Cham,
540 2015.

- 541 [18] Kohn, R., Ansley, C. F. and Tharm, D., The performance of cross-validation and maxi-
542 mum likelihood estimators of spline smoothing parameters, *J. Am. Stat. Assoc.*, **86**(416),
543 pp. 1042-1050, 1991.
- 544 [19] Fornberg, B., Larsson, E. and Flyer, N., Stable Computations with Gaussian Radial Basis
545 Functions, *SIAM J. Sci. Comput.*, **33**(2), pp. 869-892, 2011.
- 546 [20] Pazouki, M. and Schaback, R., Bases for kernel-based spaces, *J. Comput. Appl. Math.*,
547 **236**(4), pp. 575-588, 2011.
- 548 [21] Quang, M. H., Niyogi, P. and Yao, Y., *Mercer's Theorem, Feature Maps, and Smoothing*,
549 In International Conference on Computational Learning Theory (COLT), 2007.
- 550 [22] Fasshauer, G. E. and McCourt, M., Stable evaluation of Gaussian RBF interpolants,
551 *SIAM J. Sci. Comput.*, **34**(2), pp. 737-762, 2012.
- 552 [23] Gulian, M., Raissi, M., Perdikaris, P. and Em Karniadakis, G., Machine learning of space-
553 fractional differential equations, *SIAM J. Sci. Comput.*, **41**(4), pp. 2485-2509, 2019.
- 554 [24] Esmailbeigi, M., Chatrabgoun, O. and Cheraghi, M., The Role of Hilbert-Schmidt SVD
555 basis in Hermite-Birkhoff interpolation in fractional sense, *Comput. Appl. Math.*, **38**(82),
556 pp. 1-20, 2019, DOI: 10.1007/s40314-019-0849-x.
- 557 [25] Sacks, J., Welch, W. J., Mitchell, T. J. and Wynn, H. P., Design and Analysis of Computer
558 Experiments, *Statistical Science*, **4**(4), pp. 409-435, 1989.
- 559 [26] Esmailbeigi, M., Chatrabgoun, O. and Cheraghi, M., Fractional Hermite interpolation
560 using RBFs in high dimensions over irregular domains with application, *J. Comput. Phys.*,
561 **375**, pp. 1091-1120, 2018.
- 562 [27] Rasmussen, C. E. and Ghahramani, Z., *Occam's razor*, Adv. Neural. Inform. Process.
563 Syst., pp. 294-300, 2001.
- 564 [28] Aronszajn, N., Theory of reproducing kernels, *Transactions of the American mathematical*
565 *society*, **68**, pp. 337-404, 1950.
- 566 [29] Saitoh, S., *Theory of reproducing kernels and its application*, Longman, 1988.

- 567 [30] Berliet, A. and Thomas-Agnan, C., *Reproducing kernel Hilbert spaces in probability and*
568 *statistics*, Springer Science & Business Media, 2011.
- 569 [31] Datta, B. N., *Numerical Linear Algebra and Applications*, second edition, SIAM, Philadel-
570 phia, 2010.
- 571 [32] Podlubny, I., *Fractional differential equations: an introduction to fractional derivatives,*
572 *fractional differential equations, to methods of their solution and some of their applica-*
573 *tions*, Academic press, 1998.