# A Theoretical Framework of the Scaled Gaussian Stochastic Process in Prediction and Calibration

Mengyang Gu[y], Fangzheng Xie[z], and Long Wang[x]

**Abstract.** Model calibration or data inversion is one of the fundamental tasks in uncertainty quantication. In this work, we study the theoretical properties of the scaled Gaussian stochastic process (S-GaSP) for modeling the discrepancy between reality and the imperfect mathematical model. We establish an explicit connection between the Gaussian stochastic process (GaSP) and S-GaSP through the orthogonal series representation. The predictive mean estimator in the S-GaSP calibration model converges to reality at the same rate as the GaSP with suitable choices of the regularization and scaling parameters. We also show that the calibrated mathematical model in the S-GaSP calibration converges to the one that minimizes the $L_2$ loss between reality and the mathematical model, whereas the GaSP model with other, widely used covariance functions does not have this property. Numerical examples conrm the excellent nite sample performance of our approaches.

**Key words.** model misspecication, Bayesian prior, scaled Gaussian stochastic process prior, convergence, interpretability, orthogonal series representation

**MSC codes.** 62A01, 62F15, 62M20, 62P30

**DOI.** 10.1137/21M1409949

**1. Introduction.** Mathematical models are developed by scientists and engineers based on their expert knowledge of reproducing physical reality. With the rapid development of computational techniques in recent years, many mathematical models have been implemented in computer codes; these are often referred to as computer models or simulators.

Some parameters of mathematical models are unknown or unobservable in experiments. For example, the Kilauea volcano recently had one of its biggest eruptions in 2018. The location and volume of the magma chamber, as well as the magma supply and storage rate of this volcano, however, are unobservable. Field data, such as satellite interferograms and GPS measurements of the ground deformation, were used to estimate these parameters [1, 2]. Using eld observations to estimate parameters in a mathematical model and to identify the possible discrepancy between the mathematical model and the reality is widely known as model calibration or data inversion.

For any p-dimensional observable input $x \in \mathcal{X}$, denote $y^F(x) \in \mathbb{R}$ as the eld observation and $f^M(x; \theta) \in \mathbb{R}$ as a mathematical model with q-dimensional unobservable calibration

[y]Department of Statistics and Applied Probability, University of California, Santa Barbara, Santa Barbara, CA 93106-3110 USA (mengyang@pstat.ucsb.edu).

[z]Department of Statistics, Indiana University, Bloomington, Bloomington, IN 47408 USA (fxie@iu.edu).

[x]Department of Applied Mathematical and Statistics, Johns Hopkins University, Baltimore, MD 21218 USA (long.wang@jhu.edu).

parameters $2$ . Furthermore, let $y^R(x) = E[y^F(x)]$ represent the reality. A routinely used framework for calibrating the imperfect mathematical model is [15, 4]

$$(1) \qquad y^F(x) = f^M(x; ) + (x) + ;$$

where  is the noise and () is a discrepancy function between the reality and the mathematical model. Since the mathematical model is often developed by experts, we assume that the mean and trend of the observations are already included in the mathematical model.

The discrepancy function  was modeled as a Gaussian stochastic process (GaSP), in [15] and the framework has been widely studied in recent years [8, 14, 18]. This approach is referred to as GaSP calibration. Dene the prediction error by the $L_2$ loss

$$(2) \qquad \text{prediction error}(; ) = \int_{x \in X} [y^R(x) \quad \hat{y}^R_{;}(x)]^2 dx;$$

where $\hat{y}^R_{;}(x)$ is the prediction (e.g., predictive mean) of the reality, and the subscript $(; )$ implies the dependence on both the mathematical model and the discrepancy function. Since both the mathematical model and the discrepancy function are jointly estimated in GaSP calibration, the prediction error was found to be smaller when compared with the prediction error from the mathematical model or nonparametric regression alone [15].

It was shown in follow-up studies, however, that the calibrated mathematical model in GaSP calibration can be far away from the reality, which results in an identiability problem between the calibration parameters and the discrepancy function [3, 19, 26]. This is because the variability of the observations in GaSP calibration with some frequently used kernels, such as the Matern kernel, can be explained mostly by the estimated discrepancy function instead of the calibrated mathematical model.

Note that the true value of the parameter  in (1) is not well dened due to the inclusion of an unknown function . However, one can evaluate the distance between the reality and the calibrated computer model via some frequently used loss functions. A few recent studies measure the goodness of calibration in terms of the $L_2$ loss between the calibrated mathematical model and reality [25, 26, 29]. Thus, the calibration error may be dened by

$$(3) \qquad \text{calibration error}() = \int_{x \in X} [y^R(x) \quad f^M(x; )]^2 dx:$$

These studies seek to nd the $L_2$ minimizer of  that minimizes the calibration error, i.e., $_{L_2}$ := argmin$_2$ calibrationerror(). To estimate the $L_2$ minimizer, a few two-step ap-proaches were developed.　　　In [25], for instance, the reality is rst estimated through a nonparametric regression model without the assistance of the mathematical model. The calibration parameters are then estimated by minimizing the $L_2$ loss between the mathematical model and the estimator of the reality. For some complex applications, however, it is crucial to jointly estimate the reality and the calibration parameters, as the mathematical model developed based on the experts' knowledge can be very helpful for predicting the reality.

A recent model for the discrepancy function, called the scaled Gaussian stochastic process (S-GaSP) [12], was constructed such that both the prediction error and the calibration error are small. We call this approach the S-GaSP calibration. Note that the calibration error can be written as $\int_{x \in X} (x)^2 dx$. In S-GaSP, the prior distribution of the calibration error has

more probability mass near small values, reecting one's preference for smaller values of the $L_2$ norm of the discrepancy function. In contrast to the two-step approaches that seek $L_2$ minimizers, both GaSP and S-GaSP dene a sampling model of the discrepancy, such that the uncertainty of the parameters can be obtained exactly. Though the computational strategy of the S-GaSP calibration was studied in [12], theoretical properties of this process are yet to be examined.

In this work, we study the theoretical properties of S-GaSP. The contribution of this article contains three parts. First, we establish the explicit connection between GaSP and S-GaSP through the orthogonal representation of the process, which describes the reproducing kernel Hilbert space associated with an S-GaSP. Second, we show that the maximum likelihood estimator of the calibration parameter and the predictive mean in GaSP can be obtained by minimizing a squared error loss with a penalty term depending on the $L_2$ norm and the reproducing kernel Hilbert space norm. Third, we show that the predictive mean from S-GaSP converges to the reality at the same rate as the one from GaSP with suitable choices of the regularization and scaling parameters. Last but not least, using the same regularization and scaling parameters, the calibration parameters in S-GaSP converge to $L_2$, whereas GaSP calibration with other, widely used kernels does not enjoy this property. To the best of our knowledge, this work represents the rst eort in the literature showing that the joint estimation of the calibration parameters and the model discrepancy by S-GaSP (essentially a GaSP with a transformed kernel) allows these two convergence properties to hold at the same time, bridging the gap between the joint estimation procedure and other two-step estimation methods.

Although the two convergence properties discussed in this work can be achieved using the aforementioned two-step approaches [25, 29], nite sample studies suggest that jointly estimating the calibration parameters and the discrepancy as in S-GaSP calibration leads to high predictive accuracy. Additionally, since the sampling model is fully specied, uncertainties of the parameters in S-GaSP calibration can be naturally assessed through posterior distribu-tions via a Bayesian approach. Numerical results in section 5 suggest the maximum likelihood estimator (MLE) and Bayesian estimation are often similar, yet Bayesian estimation is more robust than numerical optimization.

The rest of this paper is organized as follows. In section 2, we introduce S-GaSP along with the orthogonal series representation and the joint estimation in calibration. Two convergence properties are discussed in section 3. In section 4, we introduce the discretized S-GaSP along with the parameter estimation under the Frequentist framework and Bayesian framework. A comparison between S-GaSP calibration and other alternatives is discussed in section 5 with numerical evidence provided in section 6. We conclude this work in section 7. The proofs of the theoretical results are given in the supplementary material (supplement.pdf [local/web 442KB]). We implement model calibration approaches in the RobustCalibration R package available on the Comprehensive R Archive Network (CRAN) [10], which allows for both MLE and Bayesian estimation by posterior sampling. The par-allel partial Gaussian stochastic process emulator for computer models with scalar-valued or vectorized output [11] can be called as a surrogate model from the package when the com-puter model is expensive. The source code of the examples in this article is available at https://github.com/UncertaintyQuantication/SGaSP-Theory.

**2. The scaled Gaussian stochastic process.** Denote $\zeta(\cdot) \sim \mathrm{GaSP}(0, \sigma^2 K(\cdot,\cdot))$ with variance $\sigma^2$ and correlation function $K(\cdot,\cdot)$ such that, for any inputs $\{x_i\}_{i=1}^n$, the marginal distribution $(\zeta(x_1), \ldots, \zeta(x_n))^\mathsf{T}$ follows a multivariate normal distribution with covariance $\mathrm{Cov}(\zeta(x_i), \zeta(x_j)) = \sigma^2 K(x_i, x_j)$. In order to have the mathematical model explain more variability, the S-GaSP prior introduced in [12] places more probability mass on a smaller random $L_2$ distance between the mathematical model and reality, as this measure is often used to quantify how well a mathematical model fits the reality. The S-GaSP calibration model is defined as the following hierarchical model :

$$
\begin{aligned}
& y^F(x) = f^M(x; \theta) + \zeta_z(x) + \epsilon; \\
& \zeta_z(x) = \zeta(x) \Big| \int_X \zeta^2(\zeta)d\zeta = Z \quad; \\
& \zeta(\cdot) \sim \mathrm{GaSP}(0, \sigma^2 K(\cdot,\cdot)); \\
& Z \sim p_Z(\cdot); \quad \epsilon \sim N(0, \sigma_0^2);
\end{aligned}
\tag{4}
$$

where, conditional on all parameters, the default choice of $p_Z(\cdot)$ is defined as

$$
p_Z(z) = \frac{g_Z(z)p(Z = z)}{\int_0^1 g_Z(t)p(Z = t)dt};
\tag{5}
$$

with $g_Z(z)$ being a nonincreasing scaling function and $p(Z = z)$ being the density of $Z$ at $z$ induced by a GaSP with mean 0 and covariance $\sigma^2 K(\cdot,\cdot)$.

We call $\zeta_z(\cdot)$ in (4) the scaled Gaussian stochastic process (S-GaSP). Given $Z = z$, S-GaSP becomes a GaSP constrained at the space $\int_{x \in X} \zeta^2(x)dx = z$. Note that $Z$ represents the $L_2$ distance between the reality and mathematical model. Given that $g(\cdot)$ is a nondecreasing function, the measure for $Z$ induced by S-GaSP has a larger prior probability mass near 0 than the one by GaSP, reflecting one's belief that the $L_2$ loss between the mathematical model and reality should be small.

It is easy to see that when $g_Z(\cdot)$ is a constant function, S-GaSP reduces to GaSP without any constraint. Conditioning on all parameters, we assume

$$
g_Z(z) = \frac{\lambda_z}{2\sigma^2}\exp\left\{-\frac{\lambda_z z}{2\sigma^2}\right\};
\tag{6}
$$

with a scaling parameter $\lambda_z$. We select $p_Z(\cdot)$ in (5) and $g_Z(\cdot)$ in (6) for computational reasons, as any marginal distribution of $\zeta_z$ still follows a multivariate normal distribution [12, Lemma 2.3]. Other scaling functions may also be used, but we do not pursue them in this study.

**2.1. Orthogonal series representation and marginal distribution.** Based on Karhunen–Loève theorem, GaSP with a stationary kernel admits the representation for any $x \in X$;

$$
\zeta(x) = \sum_{k=1}^{\infty} \sqrt{\lambda_k}\, Z_k \phi_k(x);
\tag{7}
$$

where $Z_k \overset{i.i.d}{\sim} N(0, 1)$, and $\lambda_k$ and $\phi_k(\cdot)$ are the $k$th eigenvalue and eigenfunction of the kernel $K(\cdot,\cdot)$, respectively. The S-GaSP can also be represented as an orthogonal series as shown below.

**Lemma 2.1 (Karhunen–Loeve expansion for the S-GaSP).** Assume $p_Z(\cdot)$ and $g_Z(\cdot)$ are defined in (5) and (6), respectively. For any $x \in \mathcal{X}$, the S-GaSP defined in (4) can be written as

$$\zeta_z(x) = \sum_{k=1}^{\infty} \sqrt{\frac{\lambda_k}{1 + \gamma_{z}\lambda_k}} Z_k \phi_k(x),$$

where $Z_k \overset{i.i.d}{\sim} N(0, 1)$, and $\lambda_k$ and $\phi_k(\cdot)$ are the $k$th eigenvalue and eigenfunction of the kernel $K(\cdot, \cdot)$, respectively.

The covariance function of S-GaSP can also be decomposed as an infinite orthogonal series, which is an immediate consequence of the fact that S-GaSP is indeed a GaSP with a transformed kernel (see Lemma 2.3 in [12] and Lemma 2.1).

**Corollary 2.2.** Assume $p_Z(\cdot)$ and $g_Z(\cdot)$ are defined in (5) and (6), respectively. The S-GaSP defined in (4) follows a multivariate normal distribution

$$[\zeta_z(x_1), \ldots, \zeta_z(x_n) \mid \sigma^2 R_z] \sim MN(0, \sigma^2 R_z),$$

where the $(i, j)$ entry of $R_z$ is

$$(8) \qquad K_z(x_i, x_j) = \sum_{k=1}^{\infty} \frac{\lambda_k}{1 + \gamma_z \lambda_k} \phi_k(x_i)\phi_k(x_j).$$

Corollary 2.2 implies that the $i$th eigenvalue of the kernel function $K_z(\cdot, \cdot)$ in S-GaSP is $\lambda_{z,k} := \lambda_k/(1 + \gamma_z\lambda_k)$, and the $k$th eigenfunction $\phi_k(\cdot)$ is the same as the one in GaSP. The form (8) does not give an explicit expression for the kernel in S-GaSP. Instead of truncating the series, one may discretize the integral $\int_{x \in \mathcal{X}} \zeta^2(x)dx$, which leads to an explicit expression of the covariance matrix, discussed in section 4.

To see the difference between GaSP prior and S-GaSP prior, we generate $n = 200$ equally spaced inputs $x_i \in [0, 1]$ for $i = 1, \ldots, 200$. The covariance matrix of these inputs is assumed to follow a unit variance Matern kernel with roughness parameter $5/2$, and the $(i, j)$ term is

$$(9) \qquad K(d_{i,j}) = \left(1 + \frac{\sqrt{5}d_{i,j}}{\gamma} + \frac{5d_{i,j}^2}{3\gamma^2}\right) \exp\left(-\frac{\sqrt{5}d_{i,j}}{\gamma}\right),$$

where $d_{i,j} = |x_i - x_j|$ for $i = 1, \ldots, n$ and $j = 1, \ldots, n$. The empirical eigenvalues are graphed as red symbols in the left panel of Figure 1, and the empirical CDF of the $L_2$ loss is shown in the right panel of Figure 1. First, when the correlation is large (i.e., large range parameter $\gamma$), the eigenvalues are typically more widely separated (left panel), and the prior probability of large $L_2$ loss $\int_\mathcal{R} \zeta^2(x)dx$ is large, even if the variance is assumed to be the same (right panel). Large correlation in GaSP is sometimes needed for accurate predictions; however, it simultaneously means the prior probability mass of the $L_2$ loss (between the reality and computer model) is large, leading to an identifiability problem. In comparison, the induced eigenvalues by the S-GaSP are graphed as blue symbols in the left panel of Figure 1, where the large eigenvalues are truncated yet the small eigenvalues remain almost the same.
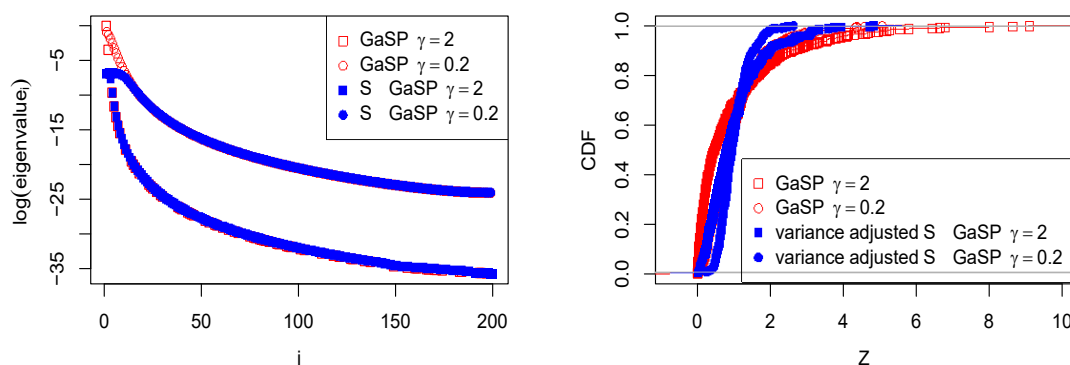
**Figure 1.** Eigenvalues and the cumulative distribution function (CDF) of the $L_2$ loss between GaSP and S-GaSP. The covariance is computed based on $n = 200$ inputs equally spaced from $[0, 1]$. In the left panel, the red symbols give the logarithm of approximated eigenvalues $\tilde{\lambda}_i / n$ by the Matern kernel in (9) with two range parameters, where $\tilde{\lambda}_i$ denotes the empirical eigenvalue of covariance matrix $\sigma^2 R$, with the $(i, j)$th term being $R_{i,j} = K(d_{i,j})$ from the Matern kernel in (9) and $\sigma^2 = 1$. The scaled eigenvalues $\tilde{\lambda}_{z,i} = \frac{\tilde{\lambda}_i/n}{n+z}$ by S-GaSP, with $z = 10^3$ are graphed by the blue symbols, for $i = 1, \ldots, n$. In the right panel, the empirical CDF of approximated $L_2$ loss $Z = \int_X \sigma^2(x) dx$ of GaSP and S-GaSP is shown. Each curve is computed based on 500 simulations, each containing $n = 200$ observations. The variance parameter $\sigma^2$ in S-GaSP is adjusted such that the summation of the eigenvalues is the same as the GaSP model with the Matern kernel.

Since we leave the variance as a free parameter estimated by the data, the spread of eigenvalues in S-GaSP is less extreme than that in GaSP. In the right panel of Figure 1, we plot the empirical distribution of the $L_2$ loss by S-GaSP, and the variance is adjusted such that the summation of the eigenvalues between GaSP and S-GaSP is the same. The S-GaSP has less prior mass on the large $L_2$ loss, reducing the identiability problem caused by large correlation.

The following Corollary 2.3 provides a decomposition of the $L_2$ loss $Z$ in S-GaSP, which follows from Lemma 2.1 in [12] and Corollary 2.2.

**Corollary 2.3.** Assume the same conditions in Lemma 2.1 hold. The distribution of $Z = \int_{x \in X} \sigma^2(x) dx$ induced by S-GaSP follows

$$Z \overset{d}{\sim} \sum_{k=1}^{\infty} \frac{\lambda_k}{1 + \frac{\lambda_k}{z}} \chi^2_k(1),$$

where $\{f_k(1)\}_{k=1}^{\infty}$ are independent chi-squared random variables with one degree of freedom.

Denote $\mathcal{H}$ and $\mathcal{H}_z$ as the reproducing kernel Hilbert space attached to GaSP with kernel $K(\cdot, \cdot)$ and S-GaSP with kernel $K_z(\cdot, \cdot)$, respectively. Let the native norm associated with $K(\cdot, \cdot)$ and $K_z(\cdot, \cdot)$ be $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{H}_z}$, respectively. We conclude this subsection by the explicit connection between the inner product of GaSP and that of S-GaSP.

**Lemma 2.4.** Assume $p_z(\cdot)$ and $g_z(\cdot)$ are dened as in (5) and (6), respectively. Let $h(\cdot) = \sum_{i=1}^{\infty} h_{ii}(\cdot)$ and $g(\cdot) = \sum_{i=1}^{\infty} g_{ii}(\cdot)$ be the elements in $\mathcal{H}$. It holds that

$$\langle h, g \rangle_{\mathcal{H}_z} = \langle h, g \rangle_{\mathcal{H}} + z \langle h, g \rangle_{L_2(X)}.$$

**2.2. Joint estimation in the S-GaSP calibration.** With the specication of $p_Z()$ in (5) and $g_Z()$ in (6), after marginalizing out $z = [z(x_1); \ldots; z(x_n)]$, the marginal distribution of the eld observations in (4) follows a multivariate normal distribution

(10) $$[y^F \mid ; ^2; _{0z}] \sim MN(f^M; ^2((n)^{-1}R_{0z} + I_n));$$

with the regularization parameter $_0 = ^2=(n^2)$ and the $(i; j)$ entry of $R_z$ as dened in (8).

Denote $L_z()$ as the likelihood for in (10). The joint estimator of $(; z())$ can be written as a penalized kernel ridge regression (KRR) estimator [28], where both the reproducing kernel Hilbert space (RKHS) norm and the $L_2$ norm of the discrepancy function are penalized simultaneously.

**Lemma 2.5.** *The maximum likelihood estimator* $^\wedge_{z;n} := \arg\max_z L_z()$ *and posterior mean* $_{z;n}() :\triangleq E[z() \mid y ; _{z;n}; ; z]$ *are the same as the following penalized KRR estimator:*

(11) $$^\wedge_{;n}; _{z;n}() = \arg\min_{()2H; 2} \frac{1}{n}\sum_{i=1}^{n} (y^F(x_i) - f^M(x_i; ) - (x_i))^2 + kk^2_{H_z}$$

*with* $kk^2_{H} = kk^2_{H} + _z kk^2_{L_2(X)}$.

The connection between KRR and the posterior mean in a GaSP regression model is well known. Here we establish an analogous connection in the S-GaSP calibration. In Lemma 2.5, both the $L_2$ norm and the native norm of the discrepancy function are penalized in S-GaSP calibration. When the discrepancy function is modeled as a GaSP, however, the $L_2$ norm of the discrepancy function is not a direct penalty (see the supplementary material supplement.pdf [local/web 442KB]). This property of the S-GaSP calibration is the key to guaranteeing that, under some regularity conditions, the estimated calibration parameters obtained from the joint estimation procedure (11) converges to the $L_2$ minimizer. A more detailed discussion is provided in section 3.

**3. Convergence properties of the S-GaSP calibration.** We discuss two convergence properties of the S-GaSP calibration in this section. First, the predictive mean estimator of the reality converges to the truth at the optimal rate with suitable choices of the regularization and scaling parameters. Second, the estimated calibration parameters by the S-GaSP calibration converge to $_{L_2}$ when sample size increases. These two properties are obtained by jointly estimating the discrepancy function and calibration parameters in (11).

**3.1. Convergence to the reality.** Let us rst consider the nonparametric regression,

(12) $$y(x_i) = f_0(x_i) + _i; \quad _i \overset{i.i.d}{\sim} N(0; _0); \quad ^2 \quad i = 1; \ldots; n;$$

where $f_0()$ denotes underlying truth. Here we place a zero-mean S-GaSP prior on the unknown truth with the default choices of $p_Z()$ and $g_Z()$ as in (5) and (6), respectively. This is a special case where the mathematical model is zero, i.e., $f^M(x; ) = 0$, and we will soon extend it to the general case when the mathematical model is not zero. For illustrative purposes, we follow [25] to assume that $x_1; \ldots; x_n$ are independently sampled from Unif$([0; 1]^p)$.

Assume the underlying truth $f_0() := E_y[y()]$ resides in the $p$-dimensional Sobolev space,

(13) $$W_2^m(X) = \left\{ f() = \sum_{k=1}^{\infty} f_k k() \in L_2(X) : \sum_{k=1}^{\infty} k^{2m=p} f_k^2 < 1 \right\};$$

with smoothness $m > p/2$; and $\{\phi_k(\cdot)\}_{k=1}^{\infty}$ being a sequence of the orthonormal basis of $L_2(\mathcal{X})$. For any integer vector $k = (k_1, \ldots, k_p)^T$ and a function $f(x_1, \ldots, x_p) : \mathcal{X} \to \mathbb{R}$, denote by $D^k$ the mixed partial derivative operator $D^k f(\cdot) := \partial^{|k|} f(\cdot) = \partial^{k_1} x_1 \cdots \partial^{k_p} x_p$ with $|k| = \sum_{i=1}^{p} k_i$. For any function in $W_2^m(\mathcal{X})$, we have $\|D^k f(\cdot)\|_{L_\infty(\mathcal{X})} < \infty$ for any $|k| < m$.

Recall that $\lambda = \lambda_0 = \tilde{\lambda}(n)$ in (10). By Lemma 2.5, the posterior mean estimator of $f(\cdot)$ with a S-GaSP prior is equivalent to the KRR estimator

$$(14) \qquad \hat{f}_{\lambda,\lambda_z;n} = \underset{f \in \mathcal{H}}{\arg\min} \left[ \frac{1}{n} \sum_{i=1}^{n} (y(x_i) - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 + \lambda \lambda_z \|f\|_{L_2(\mathcal{X})}^2 \right] : i=1$$

Recall that $\{\rho_k\}_{k=1}^{\infty}$ and $\{\phi_k\}_{k=1}^{\infty}$ are eigenvalues and eigenfunctions of the kernel $K(\cdot,\cdot)$ associated with $\mathcal{H}$, respectively. For all $k$, we assume the eigenvalues satisfy

$$(15) \qquad\qquad c k^{-2m/p} \leq \rho_k \leq C k^{-2m/p}$$

for some constants $c$ and $C > 0$. For all $k \in \mathbb{N}^+$ and $x \in \mathcal{X}$, we assume the eigenfunctions are bounded uniformly,

$$(16) \qquad\qquad \sup_{x \in \mathcal{X}} \sup_k |\phi_k(x)| \leq C_\phi ; k \geq 1$$

where $C_\phi > 0$ is a constant depending on the kernel $K(\cdot,\cdot)$.

We are now ready to state the convergence rate of the S-GaSP for the nonparametric regression model in (12).

**Theorem 3.1.** Assume the eigenvalues and eigenfunctions of $K(\cdot,\cdot)$ satisfy (15) and (16), respectively. Further assume $f_0 \in W_2^m(\mathcal{X})$; and denote $\gamma := (2m - p)^2/\{2m(2m + p)\}$. Consider the nonparametric regression model (12). For sufficiently large $n$, and any $\delta > 2$ and $C \in (0,1)$, with probability at least $1 - \exp\{-(\delta - 2)/3\} - \exp(-n^C)$, we have

$$\|\hat{f}_{\lambda,\lambda_z;n} - f_0\|_{L_2(\mathcal{X})}^2 \leq 2\delta^2 \left[ \|f_0\|_{L_2(\mathcal{X})} + \|f_0\|_{\mathcal{H}} + C_K \lambda_0^{\frac{p}{2m+p}} \right] \left( \frac{\|f_{\lambda,\lambda_z;n}\|_{\mathcal{H}}}{n} \right)$$

and

$$\|\hat{f}_{\lambda,\lambda_z;n} - f_0\|_{\mathcal{H}}^2 \leq 2\delta^2 \left[ \|f_0\|_{L_2(\mathcal{X})} + \|f_0\|_{\mathcal{H}} + C_K \lambda_0^{\frac{p}{\gamma}} \right]$$

by choosing $\lambda = n^{-2m/(2m+p)}$ and $\lambda_z = \lambda^{-1/2}$, where the constant $C_K$ depends on $K(\cdot,\cdot)$.

Our proof for Theorem 3.1 stems from [31], where the convergence rate was proved in a nonparametric regression model under the supremum norm. The proof of the convergence rate in Theorem 3.1 has two main differences in comparison to the proof of convergence of the nonparametric regression in [31]. First, $\lambda_z$ can go to infinity in Theorem 3.1, and consequently $\|\cdot\|_{\mathcal{H}_z}$ is unbounded, whereas $\|\cdot\|_{\mathcal{H}}$ was assumed bounded in proving the convergence of a nonparametric regression approach in [31]. Second, we generalize the proof to multivariate inputs. These two conditions make the proof more challenging. To solve this problem, we substantially modify the tools (e.g., see Lemma SM3.2 in the supplementary material) to prove Theorem 3.1.

The conditions in Theorem 3.1 can be relaxed in various ways. First, from the proof of Theorem 3.1, it is easy to see that if $\lambda = O(n^{-2m/(2m+p)})$ and $\lambda_z \sim O(\lambda^{1/2})$, the estimator still converges to the truth in the $L_2$ distance with the same rate $O(n^{-m/(2m+p)})$. Second, the design can be generalized to other space filling designs. Additionally, although the stationarity of the process is often assumed for computational purposes, it is not required in Theorem 3.1. Note that the regularity parameter and kernel parameters are held fixed in Theorem 3.1. We discuss the estimation of $\lambda$ and the parameters in the kernel function in section 4.1.

We are ready to discuss the convergence of estimating the reality in the calibration. The estimator for the reality in the S-GaSP calibration model is defined as

$$(17) \qquad \hat{y}^R_{\lambda,\lambda_z,n}(x) := f^M(x; \hat{\theta}_{\lambda,\lambda_z,n}) + \hat{\delta}_{\lambda,\lambda_z,n}(x)$$

for any $x \in \mathcal{X}$, where $(\hat{\theta}_{\lambda,\lambda_z,n}, \hat{\delta}_{\lambda,\lambda_z,n})$ is the estimator of the penalized KRR obtained by minimizing the loss in (11). The following Corollary 3.2 gives the convergence rate of the S-GaSP calibration model in predicting the reality. Similar to the extension for Theorem 3.1, the conditions in Corollary 3.2 can be relaxed by letting $\lambda = O(n^{-2m/(2m+p)})$ and $\lambda_z \sim O(\lambda^{1/2})$ to obtain the same convergence rate.

**Corollary 3.2.** Assume $y^R(\cdot) - f^M(\cdot; \theta) \in W_2^m(\mathcal{X})$ for any $\theta \in \Theta$ and $\sup_\theta \|y^R(\cdot) - f^M(\cdot; \theta)\|_{\mathcal{H}} < 1$. Let the eigenvalues and eigenfunctions of $K(\cdot, \cdot)$ satisfy (15) and (16), respectively. For sufficiently large $n$, any $\gamma > 2$, and $C \in (0, 1)$, with probability at least $1 - \exp\{-(\gamma - 2)/3\} - \exp(-n^C)$, we have

$$\|\hat{y}^R_{\lambda,\lambda_z,n}(\cdot) - y^R(\cdot)\|_{L_2(\mathcal{X})} \leq 2\sqrt{\gamma}\,\lambda^{\frac{p}{2}} \sup_\theta \|y^R(\cdot) - f^M(\cdot; \theta)\|_{L_2(\mathcal{X})}$$

$$+ \sqrt{\gamma}\,\lambda \sup_\theta \|y^R(\cdot) - f^M(\cdot; \theta)\|_{\mathcal{H}} + C_{K,0}\sqrt{\gamma}\,\lambda^{\frac{p}{2}} n^{-2m + \frac{m}{p} - \frac{m}{2}}$$

by choosing $\lambda = n^{-2m/(2m+p)}$ and $\lambda_z = \lambda^{1/2}$, where $\alpha = (2m - p)^2/(2m(2m + p))$ and $C_K$ is a constant only depending on the kernel.

In the example below, we illustrate the convergence using the function studied in [31], where $y^R(\cdot)$ lies in the Sobolev space $W_2^m(\mathcal{X})$ with $m = 3$ and $\mathcal{X} = [0, 1]$ when $J \to \infty$.

**Example 1.** Let the reality be $y^R(x) = 2\sum_{j=1}^{J} j^{-6}\cos(5(j - 0.5)x)\sin(5j)$, and consider $y^F(x) = y^R(x) + \epsilon$, where $\epsilon \sim N(0, 0.05^2)$ independently. Let $f^M(x; \theta) = \theta$. The goal is to predict $y^R(x)$ at $x \in [0, 1]$ and estimate $\theta$.

As a motivating example, we let $K(\cdot, \cdot)$ follow the Matern kernel in (9) with the range parameter $\gamma = 1$, as the reproducing kernel Hilbert space attached to the GaSP with this kernel is equal to Sobolev space $W_2^3(\mathcal{X})$. We test 50 configurations with the number of observations $n \in [\exp(5), \exp(10)]$, and the design points $\{x_i\}_{i=1}^n$ are equally spaced in $[0, 1]$. In each configuration, $N = 100$ simulation replicates are implemented, and we let $J = 100$ in each simulation. We then compute the average root of the mean squared error as follows:

$$(18) \qquad \mathrm{AvgRMSE}_{f^M + \delta} = \frac{1}{N}\sum_{i=1}^{N}\sqrt{\frac{1}{n}\sum_{j=1}^{n}(\hat{y}^R(x_j)_i - y^R(x_j))^2};$$
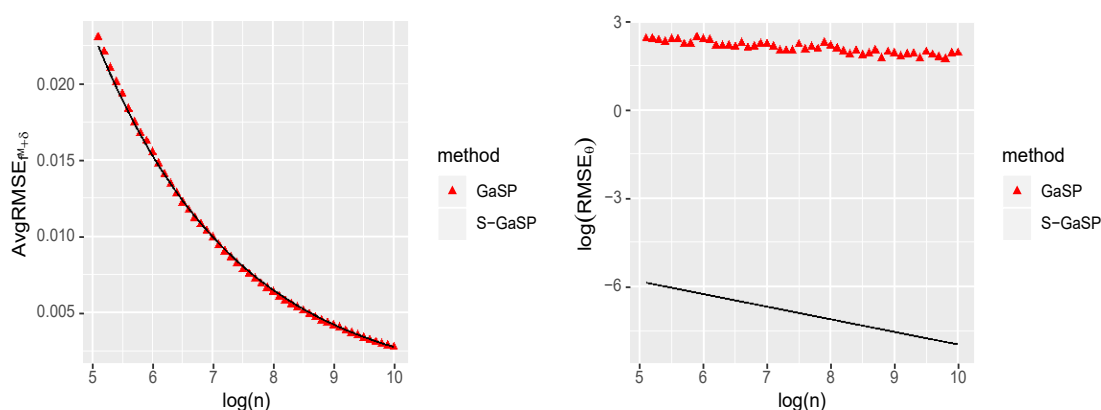
**Figure 2.** Prediction and calibration by the GaSP and discretized S-GaSP calibration models for Example 1. In the left panel, the $AvgRMSE_{fM+}$ of the GaSP calibration and that of the discretized S-GaSP calibration are graphed as the red triangles and blue dots, respectively; the black curve is $n^{-m=(2m+p)}=5$, representing the upper bound by Corollary 3.2 (up to a constant). In the right panel, the natural logarithm of the RMSE of the GaSP calibration and that of the discretized S-GaSP calibration are graphed as the red triangles and blue circles, respectively; the black line is $\log(n^{-m=(2m+p)}=40)$, the upper bound from Theorem 3.3 (up to a constant). $\sigma = n^{-2m=(2m+p)}10^{-4}$ and $z = \sigma^{1=2}$ are assumed.

where $\hat{y}_i^R(x_j)$ is an estimator of the reality at $x_j$ for $j = 1, \ldots, n$. The subscript $f^M+\delta$ indicates that both the calibrated mathematical model and the discrepancy can be used for prediction.

For both GaSP and S-GaSP calibration, the joint estimator, i.e., the posterior mean of the reality and the MLE of the calibration parameters discussed in Lemma 2.5, is implemented for each experiment at each configuration. In the left panel of Figure 2, the posterior mean estimator of the reality in both GaSP and S-GaSP calibration converges to the reality at the same rate, which matches the theoretical upper bound from Corollary 3.2. Here, for computational purposes we graph the results of discretized S-GaSP calibration, which replaces the integral $\int_{x\in\mathcal{X}} \delta^2(x)dx$ in the S-GaSP model in (4) by $(1=n)\sum_{i=1}^n \delta^2(x_i)$ in Figure 2. The discretized S-GaSP calibration is discussed in section 4.

To evaluate whether the calibrated mathematical model (here only a mean parameter) fits the data, we use the root of the mean squared error between the estimator of the calibration parameters and the $L_2$ minimizer, i.e., $RMSE = \sqrt{\frac{1}{N}\sum_{i=1}^N (\hat{\theta}_i - \theta_{L_2})^2}$, where $\hat{\theta}_i$ is an estimator of $\theta$ in the $i$th experiment.

Although GaSP and S-GaSP perform equally well in predictions for Example 1, and the estimator of the calibration parameter in the discretized S-GaSP calibration converges to the $L_2$ minimizer, that in the GaSP calibration does not converge to $\theta_{L_2}$, shown in the right panel of Figure 2. This problem is caused by the difference between the RKHS norm and the $L_2$ norm. As illustrated in Lemma 2.5, both the RKHS norm and $L_2$ norm of the discrepancy function are penalized in the S-GaSP calibration model, whereas the GaSP calibration model does not penalize the $L_2$ norm of the discrepancy function. In section 3.2, we further show that under some regularity conditions, the calibrated parameters in the S-GaSP calibration do

converge to the $L_2$ minimizer with the same choices of regularization parameter and scaling parameter as in Corollary 3.2.

**3.2. Convergence of calibration parameters.** We first list some regularity conditions for the convergence of calibration parameters. These conditions are also assumed in other studies for computer model calibration [25].

A1: $\theta_{L_2}$ is the unique $L_2$ minimizer and an interior point of $\Theta$.

A2: The Hessian matrix $\int_R \frac{\partial^2 (y^R(x) - f^M(x;\theta))^2}{\partial\theta\partial\theta^T} dx$ is invertible in a neighborhood of $\theta_{L_2}$.

A3: For all $j = 1, \ldots, q$, it holds that $\sup_\theta \left| \frac{\partial f^M(\cdot;\theta)}{\partial\theta_j} \right| < 1$.

A4: The function class $\{ y^R(\cdot) - f^M(\cdot;\theta) : \theta \in \Theta \}$ is Donsker. A5: $\sup_\theta \| y^R(\cdot) - f^M(\cdot;\theta) \|_H < 1$.

A6: The eigenvalues and eigenfunctions of $K(\cdot;\cdot)$ satisfy (15) and (16), respectively.

Assumptions A1–A3 are regularity conditions of $\theta_{L_2}$ and the mathematical model $y^M(\cdot;\theta)$ around $\theta_{L_2}$. Assumption A1 states that the $L_2$ loss minimizer is uniquely defined. If this condition is violated, reparameterization of calibration parameters may be required to reduce the intrinsic dimension of the parameter space. Assumptions A2–A3 ensure the Hessian matrix of the squared residuals is non-singular and the derivatives of computer models are bounded in the RKHS norm. The Donsker property in Assumption A4 is a standard requirement and ensures that the covering number of the function space for the discrepancy does not increase too fast [16]. Assumptions A5–A6 ensure the residuals between the reality and computer model are well-behaved such that the KRR estimator $\hat{z}_\theta$ converges to $y^R(\cdot) - f^M(\cdot;\theta)$ uniformly for each $\theta$ in terms of the $L_2$ loss.

Below, Theorem 3.3 guarantees the convergence of calibration parameters. As the calibration parameters and discrepancy function are estimated jointly in our approach, we extend the tools for proving the convergence of the two-step calibration approach [25] to prove Theorem 3.3. The detailed proof is given in the supplementary material (supplement.pdf [local/web 442KB]).

**Theorem 3.3.** Under assumptions A1–A6, the estimated parameters in (11) satisfy

$$\hat{\theta}_{z,n} = \theta_{L_2} + O_p(n^{-\frac{m}{2m+p}}),$$

provided that $\lambda = O(n^{-2m/(2m+p)})$ and $\lambda_z = O(\lambda^{1/2})$.

Note that $\lambda = O(n^{-2m/(2m+p)})$ and $\lambda_z = O(\lambda^{1/2})$ also guarantee that the posterior mean estimator in the S-GaSP calibration converges to the reality at the rate $O_p(n^{-m/(2m+p)})$ in terms of the $L_2$ loss. We briefly compare Theorem 3.3 with some existing results. The convergence rate of the calibration parameter is slower than $O(n^{-1/2})$ obtained in the two-step approach in both [25] and its Bayesian counterpart [30]. In [25] the authors established the $\sqrt{n}$-consistency of the $L_2$-calibration estimator, whereas a semi-parametric Bernstein-von Mises (BvM) theorem was developed in [30], providing a Bayesian analogy of the $\sqrt{n}$-consistency and the semi-parametric efficiency of the $L_2$ calibration method. In these approaches, the model of interest is a semiparametric model in the sense that the calibration parameter can be viewed as a functional of the unknown infinite-dimensional regression function $y^F$ by considering the functional below

$$(y^F) = \text{argmin} \| y^F(x) \quad f^M(x; ) \|^2_{L_2()} :$$

In contrast, the S-GaSP model does not incorporate the above explicit constraint directly. This is dierent from the work of [30], in which the authors formulated the prior specication by rst assigning a Gaussian process prior model on $y^F$ and then established the semiparametric BvM theorem of the posterior distribution of the induced posterior on $(y^F)$. In the current framework, such an exact constraint is not directly imposed on the relation between $y^F$ and . Due to the doubly-penalized least squares approach, given a realization of the discrepancy function (and hence, that of $y^F$), is not directly determined by $(y^F)$ but follows a distribution concentrated around $(y^F)$ a posteriori instead. Therefore, the semi-parametric BvM theorem no longer applies here. In addition, this extra layer of randomness results in a slower convergence rate than the parametric rate that can be obtained in a suciently regular semiparametric model [6].

Though the $O(n^{1=2})$ rate may be obtained by choosing $= O(n^{2m=(2m+p)})$ and $_z = O(^{1=2})$, we should be aware that, however, the $L_2$ minimizer is not the true calibration parameter but the one that minimizes the $L_2$ distance between the calibrated mathematical model and reality. In practice, the residuals between reality and the calibrated mathematical model with the $L_2$ minimizer may behave like noises, making them hard to be estimated by a nonprametric regression model alone. In comparison, the joint estimation of the discrepancy function and calibration parameters was found to have a smaller predictive error in numerical examples.

The calibrated parameters of the GaSP calibration, on the other hand, typically do not converge to the $L_2$ minimizer. Let $\frac{@f^M(;)}{@_j} := \frac{@f^M(;)}{@_j}\Big|_{=\hat{}}$. A key dierence between the GaSP and the S-GaSP calibrations is stated in the following Corollary 3.4, which is an immediate consequence of the proof of Theorem 3.3.

**Corollary 3.4.** Under assumptions A1{A6, the estimator for the calibration parameters in the S-GaSP calibration in (11) satises

$$\frac{1}{_z}_{z;n}();\ \frac{@f^M(;_{;_z;n})}{@_j} +_{;_z;n}();\ \frac{@f^M(;_{;_z;n})}{@_j}\Big|_{L_2(X)} = 0: \quad \hat{}_{@_j}$$

Further assuming the mathematical model is dierentiable at $_{;n}$, where $_{;n}$ is the estimator of via replacement of the norm $\|\cdot\|_{H_z}$ by $\|\cdot\|_H$ in (11), we see that the estimator of the calibration parameters in the GaSP calibration satises

$$_{;n}();\ \frac{@f^M(;_{;n})}{@_j}\Big|_H = 0$$

for any $_j$, $j = 1; :::; q$.

To ensure the convergence of an estimator $\hat{}$ to the $L_2$ minimizer, one typical requirement is that $h_{L_2}();\ \frac{@f^M(;)}{@_j}|_{\hat{}}|_{L_2(X)} = o_p(1)$. It is easy to see that the S-GaSP satises this condition with $1=_z = o(1)$. However, because of the dierence between the RKHS norm and the $L_2$ norm, the estimated parameters $_{;n}$ in the GaSP calibration model can be far away from

the $L_2$ minimizer. As a result, the calibrated mathematical model may not fit the data in the GaSP calibration model, as found in previous studies [25, 29]. In Example 1, the estimated parameters in the discretized S-GaSP calibration converge to the $L_2$ minimizer when the sample size increases, whereas the estimated parameters in GaSP calibration with an unscaled kernel function do not converge, as shown in the right panel of Figure 2.

**4. Discretized scaled Gaussian stochastic process.** We address the computational issue in the S-GaSP calibration in this section. Instead of truncating the kernel function in (8) by the first several terms, we select $N_C$ distinct points to discretize the input space $[0, 1]^p$ and replace $\int_{\mathcal{X}} z(x)^2 dx$ by $(1/N_C) \sum_{i=1}^{N_C} z(x_i^C)^2$ in the S-GaSP model in (4).

Here we let the discretization points be the observed inputs, i.e., $x_i^C = x_i$ for $i = 1, \ldots, N_C$ and $N_C = n$. The discretized S-GaSP is to replace $z$ in (4) by

$$(19) \qquad z_d(x) = \left\{ z(x) \,\Big|\, \frac{1}{n}\sum_{i=1}^{n} z(x_i)^2 = Z_d \right\};$$

with density $p_{Z_d}(\cdot)$ as defined in (5). After marginalizing out $Z_d$, it follows from Lemma 2.4 in [12] that $z_d(\cdot)$ is a zero-mean GaSP with the covariance function

$$(20) \qquad \sigma^2 K_{z_d}(x_a, x_b) = \sigma^2 (K(x_a, x_b) - r^T(x_a)\tilde{R}^{-1}r(x_b))$$

for any $x_a, x_b \in \mathcal{X}$, where $\tilde{R} := R + nI_n/\lambda_z$. Denote the $i$th largest eigenvalues of $R$ and $\tilde{R}$ by $\tilde\rho$ and $\tilde\rho_{d,i,z}$. It is easy to see $\tilde\rho_{d,i,z} n = \tilde\rho_i / \{n(1 + \lambda_z/\eta)\}$, which coincides with the shrinkage of eigenvalues of a nondiscretized S-GaSP in (8), by using $\tilde\rho_i/n$ as an empirical approximation to $\rho_i$.

Recall $\eta = \sigma^2/(n\sigma_0^2)$. We have the following predictive distribution of the discretized S-GaSP calibration model.

**Theorem 4.1.** Assume $z_d(\cdot)$ in (19) with $p_{Z_d}(\cdot)$ and $g_{Z_d}(\cdot)$ defined as in (5) and (6), respectively. The predictive distribution of the field data at any $x \in \mathcal{X}$ by the discretized S-GaSP calibration model in (19) is a multivariate normal distribution

$$y^F(x) \mid y^F, \theta, \sigma^2, \lambda, \lambda_z \sim \mathcal{MN}(\hat\mu_{z_d}(x), \sigma^2((n\lambda)^{-1}K_0(x, x) + 1/\lambda_z));$$

where $\hat\mu_{z_d}(x) = f^M(x, \theta) + \dfrac{r^T(x)}{1+\lambda_z}\left(R + \dfrac{n}{1+\lambda_z}I_n\right)^{-1}(y^F - f^M)$; and $K_{z_d}(x, x) = K(x, x) - r^T(x)\left(I_n + \lambda R + \dfrac{\lambda n}{1+\lambda_z}I_n\right)^{-1}\dfrac{1}{(1+\lambda_z)}\tilde{R}_z^{-1}r(x)$, with $r(x) = (K(x, x_1), \ldots, K(x, x_n))^T$; and $\tilde{R}_z = R + \dfrac{n}{\lambda_z}I_n$ with the $(i, j)$ entry of $R$ being $K(x_i, x_j)$.

Theorem 4.1 indicates that the predictive mean of the reality in the discretized S-GaSP calibration model shrinks the posterior mean towards the mean function. When $\lambda_z = 0$, the shrinkage is zero, and the discretized S-GaSP becomes the GaSP.

Interestingly, when the observations contain no noise, the predictive mean and variance of the field data from the GaSP calibration model and the discretized S-GaSP calibration model are exactly the same, as stated in the following Lemma 4.2.

**Lemma 4.2.** Assume the conditions in Theorem 4.1 hold. If $\lambda_0 = 0$, the predictive distribution of the field data at any $x \in \mathcal{X}$ by the discretized S-GaSP model in (19) is a multivariate normal distribution with the predictive mean and variance as follows:

$$E[y^F(x) \mid y^F; \hat{\theta}; \hat{\sigma}; \hat{\gamma}_z] = f^M(x; \hat{\theta}) + r^T(x) R^{-1}(y^F - f^M);$$

$$V[y^F(x) \mid y^F; \hat{\theta}; \hat{\sigma}; \hat{\gamma}_z] = \hat{\sigma}^2 \left( K(x; x) - r^T(x) R^{-1} r(x) \right);$$

where $r(x) = (K(x; x_1); \ldots; K(x; x_n))^T$; and the $(i; j)$ entry of $R$ is $K(x_i; x_j)$.

The discretized S-GaSP with the discretization points on the observed input has the same computational complexity as GaSP, as computing the inverse and determinant of the covariance matrix in the likelihood function takes $O(n^3)$ operations, where $n$ is the number of field observations. The computational complexity between different calibration approaches will be compared in section 5.2.

**4.1. Parameter estimation.** We discuss the computational issue and the parameter estimation in this section. All the approaches are implemented in the RobustCalibration R package available on CRAN [10]. First, some of the mathematical models are numerical solutions of partial differential equations implemented as computer code, which is computationally expensive to run. In these cases, one often uses a statistical emulator to approximate the computer model based on a set of model runs [21, 23]. The GaSP emulator from the RobustGaSP R package is used to emulate the computer model for both scalar-value and vector-valued output when it is expensive to run in the RobustCalibration R package. One unique feature of the RobustCalibration R package is that the parallel partial Gaussian processes [11] from the RobustGaSP R package can be called to emulate computer models with massive numbers of coordinates.

We next discuss the estimation of the regularization parameters and the kernel parameters, which were held fixed in some studies. In practice, estimating these parameters rather than fixing them can improve the predictive performance. For any $x_a = (x_{a1}; \ldots; x_{ap})^T$ and $x_b = (x_{b1}; \ldots; x_{bp})^T$, the kernel is often assumed to have a product form [15],

$$(21) \qquad K(x_a; x_b) = \prod_{i=1}^{p} K_i(d_i);$$

where $d_i = |x_{ai} - x_{bi}|$ for $i = 1; \ldots; p$, and $K_i(\cdot)$ is a one-dimensional kernel function. One widely used kernel function is the Matern kernel [13]. When the roughness parameter is a half-integer, the Matern kernel has an explicit expression. For example, the Matern kernel with roughness parameter $5/2$ is given in (9). A good feature of the Matern kernel is that the sample path of a GaSP is $\lceil \nu_i \rceil - 1$ times differentiable, where $\nu_i$ is the roughness parameter.

Denote by $\gamma = (\gamma_1; \ldots; \gamma_p)^T$ the unknown range parameters in the covariance. The parameters in the discretized S-GaSP calibration model are the calibration parameters $\theta$, the variance parameter of the noise $\sigma^2$, the regularization parameter $\lambda = \sigma^2/(n^2 \sigma_0^2)$ the scaling parameter $\gamma_z$ and range parameters $\gamma$. The MLE of the variance of noise follows $\hat{\sigma}^2_{0;MLE} = S^2_{\gamma_z d}$ where $S^2_{\gamma_z d} = (y^F - f^M)^T \tilde{R}^{-1}_{\gamma_z d}(y^F - f^M)$ with $\tilde{R}^{-1}_{\gamma_z d} = \left( R^{-1} + \gamma_z I_n = n \right)^{-1} + n I_n$; and $\tilde{R}^{-1}_{\gamma_z d} = \gamma_z = (ng) + (R + n I_n = g)^{-1} = g^2$ with $g = \gamma_z + 1$. Marginalizing out $\gamma_z(\cdot)$ and

plugging $\hat{\sigma}^2_{0;MLE}$ into the likelihood of the discretized S-GaSP model in (19), one has the prole likelihood

$$\text{(22)} \qquad L_{Z_d}(\cdot;\cdot;\cdot;_z) \propto |R_{Z_d}|^{-1=2}(S^2)^{-\frac{n=2}{Z_d}}:$$

One may use an MLE approach to numerically maximize the prole likelihood in (22) to estimate the parameters. Note that $_z$ reects one's tolerance of how well a mathematical model should predict the reality without the discrepancy function, and thus this parameter may be chosen based on the experts' knowledge. Because of the conditions discussed in Theorem 3.3, $_z$ may be xed to be proportional to $^{-1=2}$ or be related to the sample size. For all numerical examples, the scaling parameter of the S-GaSP is $_z = (c_z|||)^{-1=2}$, where by default $c_z = 1$, $= ^2 = (^2 n)$; and $= (\tilde{}_1;:::;\tilde{}_{p})^T$, with $\tilde{}_i$ being the normalized range parameter (normalized by the length of each coordinate of the input variable). This choice is also implemented as a default choice in the RobustCalibration package, and users can specify this parameter as well. The inclusion of the range parameters is due to the confounding issue between the range parameter and the variance parameter of the process, whereas the ratio of these parameters can typically be estimated consistently from the data [32].

In the RobustCalibration package, we implement both the MLE and the Bayesian method for estimating the parameters and making predictions. For the MLE, the low-storage quasi-Newton method was used for optimization [17]. For the Bayesian method, we assume that the prior distribution is $(\cdot;\cdot;^2) \propto ()(_0^{\cdot})^{=2}$; with $= ^2=^2$ being the nugget parameter. Here $(\cdot;)$ is chosen as the joint robust prior for the kernel parameters [9], and $()$ may be specied by the user to reect the experts' knowledge. We assume a uniform distribution of $()$ in the numerical examples for demonstration purposes. In contrast to the MLE, the uncertainty in parameter estimation can be obtained by the posterior samples.

## 5. Comparison between dierent calibration approaches.

We compare a few calibration approaches in this section. One of the most popular frameworks for computer model calibration is the GaSP calibration approach [15], which models the discrepancy in (1) as a GaSP. The mathematical model and discrepancy function are jointly estimated under the Bayesian framework. The S-GaSP approach is an extended version of GaSP calibration and places more prior probability mass of the $L_2$ norm of the discrepancy near zero. Consequently, the calibrated mathematical models in the S-GaSP calibration t the data better than the ones in the GaSP calibration.

The S-GaSP calibration approach was inspired by a few pioneering approaches seeking to nd the $L_2$ minimizer of the calibration parameters [19, 25, 29]. The orthogonal Gaussian process proposed in [19] constrains the derivatives of the random $L_2$ norm of the discrepancy to be zeros, equivalently giving more prior probability mass of calibration parameters at the stationary points of the calibration parameters in terms of the $L_2$ loss. The S-GaSP model explores another transformation that avoids putting more prior probability mass at the local maxima of the $L_2$ loss of the discrepancy, and it has a closed-form likelihood function. The $L_2$ calibration was proposed in [25], where the reality is rst estimated by a KRR nonpara-metric regression, and the calibration parameters are then estimated by minimizing the $L_2$ loss between the estimated reality and mathematical model. The least squares (LS) calibra-tion is proposed in [29], where the calibration parameters are estimated by rst minimizing

the squared error between the mathematical model and observations, and a nonparamet-ric approach is then applied to estimate the residuals between the reality and mathematical model.

**5.1. Numerical comparison of predictive accuracy by the calibrated computer model and discrepancy.** We use the following example to illustrate that jointly estimating the discrepancy function and calibration parameters can be helpful for predicting the reality.

**Example 2.** Let $y^F(x) = y^R(x) + $ , where $\sim N(0; 0.05^2)$ independently and $y^R(x) = g_1(x)+g_2(x)$, with $g_1(x) = \sum_{j=1}^{\infty} j^{-1} \cos(5(j-0.5)x) \sin(5j)$ and $g_2(x) = \sum_{j=1}^{\infty} j^{-6} \cos(5(j-0.5)x) \sin(5j)$. Let the mathematical model be $f^M(x) = g_1$. The goal is to predict $y^R(x)$ and estimate . For any $x$, the first 100 terms in $g_1(x)$ and $g_2(x)$ were used in the computation.

We compare the GaSP, S-GaSP, $L_2$; and LS calibration approaches for Example 2. For both the GaSP and S-GaSP approaches, the calibration parameter, variance, and kernel parameters are estimated by the MLE. For the two-step approaches, a GaSP model is used as the nonparametric regression model in the first step of $L_2$ calibration and in the second step of the LS calibration. We use the Matern kernel function in (9) as the kernel function $K(\cdot ; \cdot)$ for all methods.

In the left panel of Figure 3, we found that the predictive errors by the GaSP and S-GaSP calibrations are considerably smaller than those of the LS and $L_2$ calibration approaches. This is because the reality contains $g_1(x)$, which is difficult to predict by a nonparametric regression approach alone. The mathematical model specified herein, however, can explain this term with calibration parameter close to 1. The estimated calibrated parameter in both GaSP and S-GaSP is indeed close to 1 in all experiments, which leads to better predictive performance.

The $L_2$ minimizer of the calibration parameters in this example is around 1.775. Note that the estimated calibration parameter S-GaSP calibration may not converge to the $L_2$ minimizer when sample size increases. This is because $||y^R(\cdot) - f^M(\cdot; )||_H$ is unbounded when $= 1$, which violates assumption A5. The calibrated computer model with calibration
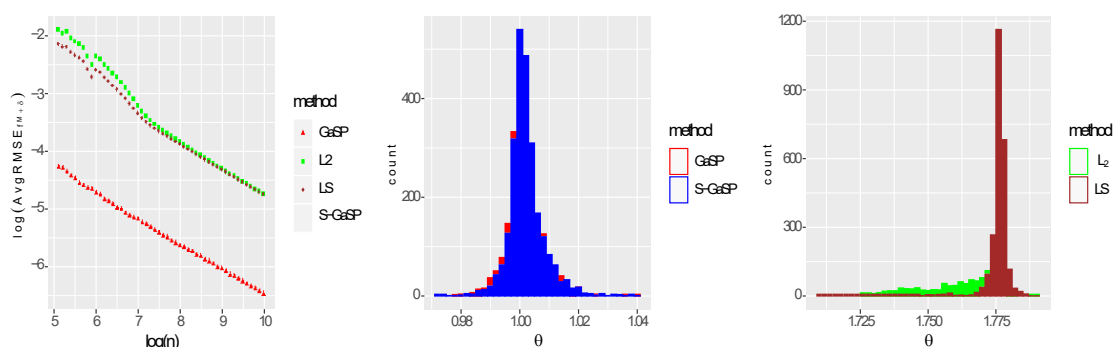


**Figure 3.** *Comparison of different approaches in Example 2. In the left panel, the logarithm of the $AvgRMSE_{f^M+\delta}$ for four calibration approaches is graphed at the logarithm of different sample sizes. The histogram of the estimated calibration parameter of each experiment for different approaches is given in the middle panel and the right panel.*

parameter being around 1 improves the predictive accuracy and interpretation, as the residual discrepancy term is a smooth term that is easy to predict.

The GaSP and S-GaSP calibration approaches are not always more accurate in predicting the reality than the two-step approaches. Indeed, for Example 5, to be discussed in section 5, the prediction in the $L_2$ calibration approach is the best among all methods, as the observations can be easily predicted through a nonparametric regression without the mathematical model. When the reality is complicated, the joint estimation strategy implemented in both GaSP and S-GaSP calibration approaches seems to have smaller predictive errors, which will be illustrated by a few more numerical examples.

As the sampling model of the observations is well dened in both GaSP calibration and S-GaSP calibration, a Bayesian approach can be implemented, and the uncertainty of the parameters can be naturally obtained by their posterior distributions. For the $L_2$ calibration, the asymptotic distribution of the estimator of the calibration parameter may be used to approximately quantify the uncertainty in parameter estimation [25]. A bootstrap approach is developed to assess the uncertainty of parameters for the LS calibration approach [29].

**5.2. Comparing computational complexity from dierent calibration approaches.** The computational cost of computer model calibration may come from two parts of this method: evaluating the computer model, and computing the likelihood or loss function in calibration. First, for a slow computer model, the GaSP emulator discussed in section 4.1 can be used to approximate the output from the computer model. Second, the computational complexity of all the methods ($L_2$, LS, GaSP, and S-GaSP calibrations) discussed in section 5.1 is $O(n^3)$ for computing the likelihood function or loss function in general, where n is the number of eld observations.

The MLE and the two-step approaches are faster than the Bayesian estimation by posterior samples, as they requires fewer iterations. However, there are two drawbacks to MLE and the two-step approaches. First, numerical optimization of the parameters could be unstable and could depend on the initial guesses when the numbers of calibration parameters and kernel parameters are large.           Second, the uncertainty of the parameters can be obtained naturally by a Bayesian method but not straightforwardly by the MLE and the two-step approaches.

Both the MLE and posterior sampling approaches are implemented in the Robust-Calibration package [10]. In the following, we will show results of both posterior sampling and MLE by the GaSP, S-GaSP, and no-discrepancy calibration approaches. The estimation by $L_2$ calibration and LS calibration will also be included for comparison.

**6. Numerical study.** In this section, we compare the numerical performance of several methods based on the prediction error and the calibration error in (2) and (3), respectively, both evaluated with respect to the $L_2$ loss. The prediction error is our primary consideration, because the out-of-sample prediction for the reality examines how well we can reproduce the reality. The calibration describes how well the calibrated mathematical model ts reality in terms of the $L_2$ loss [25]. The parameters in a mathematical model often have scientic interpretation, whereas a nonlinear discrepancy function might not be interpretable. Thus the discrepancy function is not used for prediction in the second criterion.

For all examples, we compute $\text{AvgRMSE}_{f^M+\delta}$ in (18) and $\text{AvgRMSE}_{f^M}$ through replacing prediction by the calibrated computer model output in (18). For assessing the uncertainty in predictions, we also compute the average length of the predictive interval and the proportion of the observations covered by the 95% predictive interval denfed as follows:

$$L_{CI}(95\%) = \frac{1}{Nn} \sum_{j=1}^{N} \sum_{i=1}^{n} \text{length} \{CI_{ij}(95\%)\};$$

$$P_{CI}(95\%) = \frac{1}{Nn} \sum_{j=1}^{N} \sum_{i=1}^{n} 1\{y_j^R(x_i) \in CI_{ij}(95\%)\};$$

where $CI_{ij}(95\%)$ is the 95% predictive interval; $N$ and $n$ are the total number of experiments and of held-out test data in each experiment, respectively. An efficient method should have small $\text{AvgRMSE}_{f^M+\delta}$ and $\text{AvgRMSE}_{f^M}$, short predictive interval, and $P_{CI}(95\%)$ close to the nominal 95%. For the real examples, we replace the test reality output by the held-out observations for out-of-sample predictions.

For GaSP, S-GaSP, and the no-discrepancy calibration (where the discrepancy is zero), the results by both posterior sampling and MLE are obtained by the RobustCalibration R package available from CRAN [10]. Five initial values were used to numerically optimize the calibration parameters and the kernel parameters in the MLE approach. Furthermore, we also implemented the $L_2$ calibration and LS calibration. The GaSP regression using the RobustGaSP R package is used to estimate the reality in the fIrst step of the $L_2$ calibration and estimate the residual (between the calibrated mathematical model and the reality) in the second step of the LS calibration. The kernel function $K(\cdot; \cdot)$ is assumed to be the Matérn covariance function (9) in all methods for demonstration purposes.

**6.1. Simulated example. Example 3.** Let $y^F(x) = y^R(x) + \epsilon$, where $x = (x_1, x_2) \in [0,1]^2$, $y^R(x) = \sin(0.2x_1)x_2 + \sin(2x_1)x_2 + 1$; and $\epsilon_i \overset{iid}{\sim} N(0, 0.1^2)$. The mathematical model is $f^M(x; \theta) = \sin(\theta_1 x_1)x_2 + \theta_2$ with $\theta_1 \in [0, 10]$ and $\theta_2 \in R$.

We fIrst consider a simulated study in Example 3 and test two confIgurations, where the sample sizes of the observations are taken to be $n = 25$ and $n = 50$, respectively. For each confIguration, we test $N = 100$ experiments with $n = 2500$ held-out reality points, where the inputs are equally spaced in each interval. The input variable in each experiment is generated from the maximin Latin hypercube design [23]. Here we call the deSolve R package [24] to numerically solve the dIfferential equation at each sampled or iterative parameter value. Emulating the vectorized output based on a GaSP emulator built in the RobustCalibration package is faster than directly calling the numerical solver for this example. It is shown in [10] that the calibration results based on the numerical solver and the emulator are quite similar for this example.

The predictive errors by dIfferent approaches are given in Table 1. First, we found that the $\text{AvgRMSE}_{f^M+\delta}$ of all methods that include a model of discrepancy is better than $\text{AvgRMSE}_{f^M}$, as the mathematical model is imperfect compared to the reality. Second, the GaSP and S-GaSP calibration approaches perform better than the two-step LS and $L_2$ calibration methods in terms of $\text{AvgRMSE}_{f^M+\delta}$. Although around 95% of the held-out test data

**Table 1**

Predictive performance by dierent methods for Example 3. We show performance of the GaSP, S-GaSP, and no-discrepancy (N-D) calibration approaches based on posterior sampling (PS) and MLE in estimating parameters. Also included are the $L_2$ and LS calibration approaches. The smallest error is highlighted in bold.

| n = 25 | GaSP | GaSP | S-GaSP | S-GaSP | N-D | N-D | $L_2$ | LS |
|---|---|---|---|---|---|---|---|---|
| | PS | MLE | PS | MLE | PS | MLE | | |
| $AvgRMSE_{fM+}$ | .0556 | .0686 | .0558 | .0656 | / | / | .102 | .0702 |
| $AvgRMSE_{fM}$ | .143 | .139 | .131 | .134 | .130 | .131 | .131 | .131 |
| $L_{CI}(95\%)$ | .206 | .207 | .209 | .207 | 0.186 | / | .320 | .213 |
| $P_{CI}(95\%)$ | .931 | .851 | .932 | .840 | 0.415 | / | .871 | .847 |
| n = 50 | GaSP | GaSP | S-GaSP | S-GaSP | N-D | N-D | $L_2$ | LS |
| | PS | MLE | PS | MLE | PS | MLE | | |
| $AvgRMSE_{fM+}$ | .0404 | .0439 | .0402 | .0418 | / | / | .0655 | .0437 |
| $AvgRMSE_{fM}$ | .144 | .140 | .130 | .133 | .128 | .128 | .128 | .128 |
| $L_{CI}(95\%)$ | .153 | .151 | .151 | .137 | 0.124 | / | .245 | .154 |
| $P_{CI}(95\%)$ | .938 | .912 | .935 | .900 | 0.269 | / | .949 | .918 |

are covered by the 95% predictive interval in all methods, the average lengths of the predic-tive intervals by the GaSP and S-GaSP calibrations are shorter than those of the two-step approaches.

Second, the predictive error of the MLE approach by GaSP and S-GaSP is a little worse than that of the Bayesian approach by posterior sampling when the sample size is small, as the numerical optimization is less stable. Furthermore, as the uncertainty of calibration and the kernel parameters were not taken into account in the MLE, the proportion of the held-out data covered in the 95% interval ($P_{CI}(95\%)$) is typically much smaller than 95%.

For the $L_2$ calibration approach, the mathematical model is not used in the rst step, and the parameters in the mathematical model are estimated in the second step to minimize the $L_2$ loss. Here, the high-requency term ($\sin(2x_1)x_2$) makes the prediction in the rst step by nonparametric regression less accurate than the joint estimation of the calibration parameters and the discrepancy function. For the LS approach, the residuals between the tted mathematical model and the reality can be more dicult to predict by a GaSP model than a joint model of the mathematical model and discrepancy, as the residuals contain more changes of monotonicity locally than the reality.

The estimated calibration parameters of Example 3 are graphed in Figure 4. In the left panel, the estimate of $_1$ using the LS and $L_2$ methods is close to the $L_2$ minimizer (graphed as the solid line), whereas the estimate of $_1$ using the GaSP and S-GaSP is, in fact, closer to 2. This is because the model complexity is naturally built into the calibration: an estimated $_1$ that is close to 2 makes the prediction better, since the high-frequency term can be approximately explained by the calibrated mathematical model. The estimate of $_1$ being close to 2 is important to have better predictive performance, as shown in Table 1. The estimate of $_2$ is close to the $L_2$ minimizer by the S-GaSP, $L_2$; and LS calibration approaches, whereas GaSP calibration leads to an estimated $_2$ that has a large variability. Consequently, the GaSP calibration approach produces the larger calibration error ($AvgRMSE_{fM}$) compared to other calibration methods, as illustrated in Table 1.
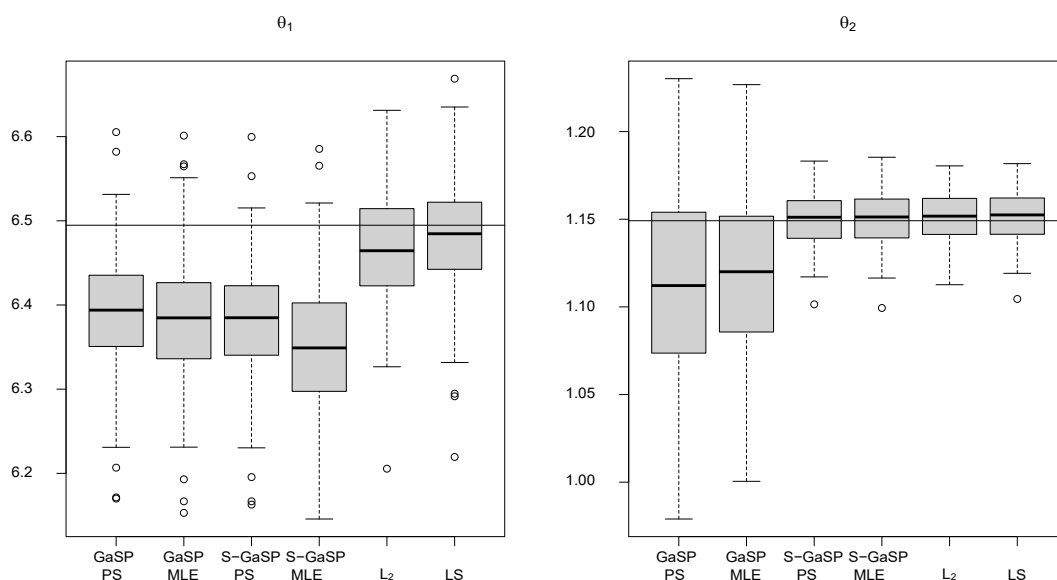
**Figure 4.** Boxplots of the posterior mean of calibration parameters from the GaSP with posterior sampling (PS), MLE, and from the S-GaSP, $L_2$, and LS calibration approaches for each experiment in Example 3. The no-discrepancy calibration with both PS and MLE are almost the same as $L_2$ and LS calibration, so they are not shown here. The solid lines are the $L_2$ minimizer, which is around 6:48 and 1:15 for $\theta_1$ and $\theta_2$, respectively.

Example 3 indicates that the calibrated mathematical model that minimizes the $L_2$ loss between the reality and mathematical model might not always be the optimal choice for predictions. In this example, when $\theta_1$ is estimated to be close to 2, the predictive accuracy can be improved signicantly. The other parameter, $\theta_2$, is a mean parameter. The two types of errors are both small when $\theta_2$ is estimated to be close to the $L_2$ minimize and when $\theta_1$ is estimated to be close to 2, which is achieved in the S-GaSP calibration.

**6.2. Chemical system interaction. Example 4.** Consider the system interaction between two chemical substances $y_1$ and $y_2$:

$$\dot{y}_1(t) = -10^{\theta_1 - 3} y_1(t);$$
$$\dot{y}_2(t) = 10^{\theta_1 - 3} y_1(t) - 10^{\theta_2 - 3} y_2(t);$$

where 2 repeated observations of the second chemical substance are measured at each of the 6 time points $t = \{10, 20, 40, 80, 160, 320\}$ in [7]. The goal is to estimate $\theta_1$ and $\theta_2$ and to predict the values of the chemical substances across time.

We consider $(\theta_1, \theta_2) \in [0.5, 1.5]^2$. As the number of observations is limited, we rst perform a leave-one-out comparison by holding out two repeated observations for prediction at each time point. We replace the reality in each criterion by the held-out observations to test each approach. The predictive performance of the leave-one-out comparison for Example 4 is given in Table 2.

First, the predictive errors ($AvgRMSE_{fM+}$) of GaSP or S-GaSP with posterior sampling are smaller than those of the MLE approach, as the Bayesian method is more stable, in

**Table 2**
Predictive performance by dierent methods for Example 4. The smallest error is highlighted in bold.

|  | GaSP, PS | GaSP, MLE | S-GaSP, PS | S-GaSP, MLE |
|---|---|---|---|---|
| $AvgRMSE_{fM_+}$ | 5.19 | 6.47 | 5.39 | 6.50 |
| $AvgRMSE_{fM}$ | 6.10 | 7.18 | 5.70 | 6.44 |
| $L_{CI}(95\%)$ | 26.5 | 18.1 | 29.2 | 18.9 |
| $P_{CI}(95\%)$ | 91:7% | 83:3% | 91:7% | 83:3% |
|  | N-D, PS | N-D, MLE | $L_2$ | LS |
| $AvgRMSE_{fM_+}$ | / | / | 10.2 | 6.39 |
| $AvgRMSE_{fM}$ | 5.84 | 6.22 | 6.95 | 6.22 |
| $L_{CI}(95\%)$ | 23.7 | 15.5 | 52.5 | 23.7 |
| $P_{CI}(95\%)$ | 91:7% | 75% | 91:7% | 100% |

particular in a small sample scenario. The $L_2$ approach is less accurate in prediction in this example (shown in the upper panels in Figure 5), as predicting the reality by nonparametric regression can be inaccurate due to the small number of observations. On the other hand, the GaSP calibration with either posterior sampling or MLE has a large calibration error, based on predictions by the calibrated mathematical model alone. The $L_2$ calibration and no-discrepancy calibration have small $AvgRMSE_{fM}$. The S-GaSP calibration has a relatively small predictive error under both predictive criteria.

The upper panels in Figure 5 give posterior samples of the calibration parameters and MLE by GaSP, S-GaSP, and no-discrepancy calibration for Example 4. The posterior samples by S-GaSP and no-discrepancy calibration are very similar, and the MLEs by these two methods are close. This is reasonable, as the calibrated computer model ts the observations reasonably well, shown in the lower left panel. The posterior samples and MLE by GaSP are slightly dierent from those by S-GaSP or no-discrepancy calibration. Consequently, the calibrated computer model by GaSP underestimates the output, shown in the lower left panel in Figure 5.

**6.3. Ion channel experiments.** In this example, we consider calibrating the mathematical model for the sodium ion channels using real observations from the whole cell voltage clamp experiments [20].

**Example 5.** The data sets consist of 19 observations of normalized current needed to maintain the membrane potential at 35mV over time [19]. Denote by the input variable x the natural logarithm of time. The mathematical model has the following expression:

$$f^M(x; ) = e^T_4 \exp[\exp(x)A()]e_4;$$

where $e_i$ is a column vector with 1 at the ith element and 0 for the rest of the components, the rst exp is the matrix exponential, $= (_1; _2; _3)_T$, and

$$A() = \begin{pmatrix} 0 & 0 & 0 & 1 \\ _{23}_1 & 0 & 0 & B \\ B_2 & _{12} & 0_1 C & C \\ @_0 & _2 & _{2}_1 & _1 A \\ 0 & 0 & _2 & _1 \end{pmatrix}$$

The ranges of the calibration parameters considered here are $_i \in [0; 10]$ for $i = 1; 2; 3$.
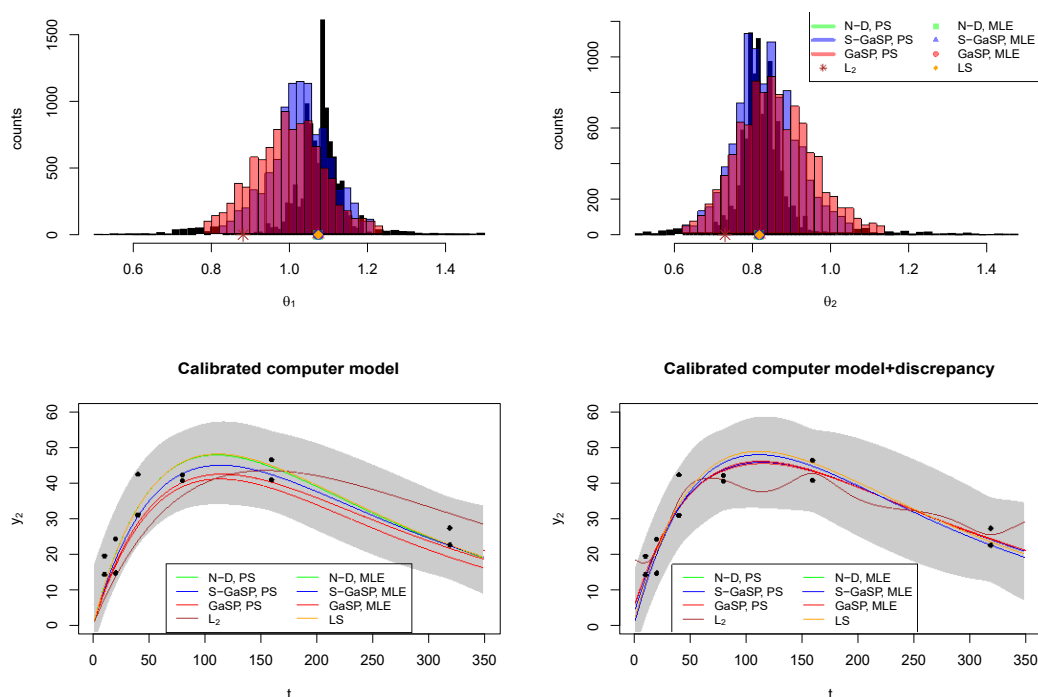
**Figure 5.** *In the upper panels, histograms are posterior samples of two calibration parameters from the GaSP, S-GaSP, and no-discrepancy calibration approaches. Point estimates of MLE, $L_2$, and LS are plotted on the x-coordinate. In the lower left panel, calibrated computer models by dierent approaches were graphed, and the shared area is the 95% posterior interval from GaSP calibration. In the lower right panel, the predictive distribution of the summation of calibrated computer models and discrepancy is plotted, and the shaded area is the 95% posterior predictive interval from S-GaSP calibration. The black circles are observations.*

We rst consider leave-one-out comparisons for the ion channel experiment. The predictive errors by dierent approaches are given in Table 3. Here, since we have a moderately large number of observations for the one-dimensional input variable, GaSP regression by the RobustGaSP R package has high accuracy in predicting the reality even without the computer model. Including the mathematical model does not improves the prediction accuracy in this example. Therefore, the predictive error by the $L_2$ calibration is the smallest, as the GaSP regression without the mathematical model is used for predicting the reality in the rst step. Between posterior sampling and the MLE approach, the posterior sampling typically has a smaller predictive error, and the predictive interval covers more held-out test points, which are consistent with previous examples.

The posterior samples and estimation of the reality by dierent calibration approaches are graphed in Figure 6. For better visualization, we reduce the posterior samples by 10 and only graph one-tenth of the total posterior samples. In this example, the posterior samples of the S-GaSP are closer to the ones with the no-discrepancy calibration. For this reason, the calibrated computer model by the S-GaSP calibration (graphed as the dashed blue curve in the lower right panel) ts the observations well, whereas the calibrated computer

Table 3

Predictive performance by dierent methods for Example 5.

| | GaSP, PS | GaSP, MLE | S-GaSP, PS | S-GaSP, MLE |
|---|---|---|---|---|
| $AvgRMSE_{fM+}$ | $1.77 \times 10^{3}$ | $1.97 \times 10^{3}$ | $1.52 \times 10^{3}$ | $1.36 \times 10^{3}$ |
| $AvgRMSE_{fM}$ | $1.13 \times 10^{2}$ | $1.75 \times 10^{2}$ | $5.62 \times 10^{3}$ | $5.76 \times 10^{3}$ |
| $L_{CI}(95\%)$ | $5.64 \times 10^{3}$ | $3.56 \times 10^{3}$ | $9.58 \times 10^{3}$ | $7.10 \times 10^{3}$ |
| $P_{CI}(95\%)$ | $94.7\%$ | $84.2\%$ | $100\%$ | $73.6\%$ |

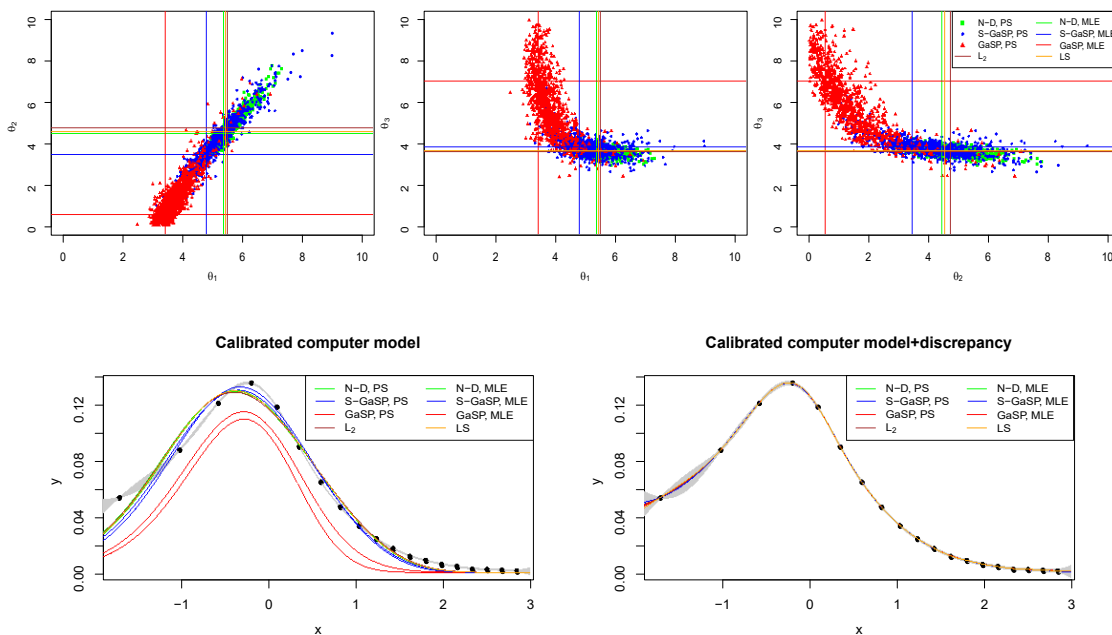| | N-D, PS | N-D, MLE | $L_2$ | LS |
|---|---|---|---|---|
| $AvgRMSE_{fM+}$ | / | / | $8.26 \times 10^{4}$ | $2.17 \times 10^{3}$ |
| $AvgRMSE_{fM}$ | $6.45 \times 10^{3}$ | $6.51 \times 10^{3}$ | $4.83 \times 10^{3}$ | $6.49 \times 10^{3}$ |
| $L_{CI}(95\%)$ | $2.58 \times 10^{2}$ | $2.24 \times 10^{2}$ | $9.45 \times 10^{3}$ | $5.80 \times 10^{3}$ |
| $P_{CI}(95\%)$ | $84.2\%$ | $84.2\%$ | $84.2\%$ | $84.2\%$ |



Figure 6. The posterior samples from no-discrepancy, S-GaSP, and GaSP calibrations are given in the upper panels. The crossings of each colored line are MLEs and $L_2$ and LS calibration. The t by the calibrated computer model by dierent approaches is plotted in the lower left panel, and the predictive mean by combining the calibrated computer model and discrepancy is plotted in the lower right panel. The black dots are observa-tions. The 95% predictive intervals by GaSP and S-GaSP with posterior sampling (PS) are plotted in the lower left and lower right panels, respectively.

model by the GaSP calibration (graphed as the dashed red curve in the lower right panel) under-estimates the values of the observations. This identiability issue of the GaSP calibration in this example was also reported in [19]. In contrast, the calibrated mathematical model ts the observations reasonably well in the S-GaSP calibration. Finally, when using

the calibrated computer model and discrepancy function, all methods t the observations well as shown in the lower right panel in Figure 6.

7. Concluding remarks. We have introduced a theoretical framework of a scaled Gaussian stochastic process (S-GaSP) for calibration and prediction. We showed that under certain routinely used assumptions, the predictive mean of the S-GaSP calibration model converges to the reality as fast as the GaSP calibration with some suitable choices of the regularization parameter and scaling parameter. The estimated calibration parameters in the S-GaSP calibration converge to the $L_2$ minimizer with the same choices of the regularization parameter and scaling parameter, whereas those in the GaSP calibration typically do not converge to the $L_2$ minimizer. The results rely on the orthogonal series representation of the processes studied in this work. The GaSP, S-GaSP, and no-discrepancy calibration approaches were implemented in the RobustCalibration R package, where the parameters can be estimated by both posterior sampling and MLE. Furthermore, for computationally expensive computer models, a GaSP emulator for scalar-valued and vectorized output can be called in RobustCalibration for accelerating the computation.

The numerical studies indicate that jointly estimating the calibration parameters and the discrepancy function in GaSP and S-GaSP calibration can improve the predictive accuracy if the reality is too complicated to be predicted precisely by a nonparametric regression model alone. The mathematical model, which typically contains some information about the trend and the shape of the reality function, can be helpful in predictions. We also empirically found that the calibrated mathematical model that minimizes the $L_2$ distance between the reality and the mathematical model may not always be the best for reducing the predictive error. This is because the residuals may behave like noises, which could be dicult to estimate accurately by nonparametric models. The S-GaSP calibration gives predictions that are at least as accurate as the GaSP calibration, and the calibrated mathematical model by the S-GaSP calibration ts the real observations more closely than the ones by the GaSP calibration in almost all examples.

We outline a few extensions of the S-GaSP model below. First, from the theoretical perspective, we did not study the convergence of the discretized S-GaSP, whereas the numerical studies indicate that the convergence rate from the discretized S-GaSP is the same as from the S-GaSP. Second, we illustrated that the S-GaSP calibration can be implemented in both Frequentist and Bayesian ways. The contraction rate of the S-GaSP under the Bayesian framework, however, was not studied. The studies of the contraction rate of the GaSP regression may be extended to achieve this goal [5], and the adaptive approach with respect to the smoothness level of latent function in [27] may be extended in model calibration. Third, we x the scaling parameter in the S-GaSP calibration as a function of the sample size, whereas historical information may be used to develop a reasonable prior for this parameter [22]. Furthermore, the S-GaSP calibration framework may be extended to include a model of correlated noises, as the eld observations may contain time series and images [2].

## REFERENCES

[1] K. R. Anderson, I. A. Johanson, M. R. Patrick, M. Gu, P. Segall, M. P. Poland, E. K. Montgomery-Brown, and A. Miklius, Magma reservoir failure and the onset of caldera collapse at Klauea volcano in 2018, Science, 366 (2019), eaaz182.

[2] K. R. Anderson and M. P. Poland, Abundant carbon in the mantle beneath Hawai'i, Nat. Geosci., 10 (2017), pp. 704{708.

[3] P. D. Arendt, D. W. Apley, and W. Chen, Quantication of model uncertainty: Calibration, model discrepancy, and identiability, J. Mech. Des., 134 (2012), 100908.

[4] M. J. Bayarri, J. O. Berger, R. Paulo, J. Sacks, J. A. Cafeo, J. Cavendish, C.-H. Lin, and J. Tu, A framework for validation of computer models, Technometrics, 49 (2007), pp. 138{154.

[5] A. Bhattacharya, D. Pati, and D. Dunson, Anisotropic function estimation using multi-bandwidth Gaussian processes, Ann. Stat., 42 (2014), pp. 352{381.

[6] P. J. Bickel, C. A. Klaassen, P. J. Bickel, Y. Ritov, J. Klaassen, J. A. Wellner, and Y. Ritov, Ecient and Adaptive Estimation for Semiparametric Models, Vol. 4, Johns Hopkins University Press, Baltimore, 1993.

[7] G. Box and G. Coutie, Application of digital computers in the exploration of functional relationships, Proc. IEE-Part B Radio Electron. Eng., 103 (1956), pp. 100{107.

[8] M. Goldstein and J. Rougier, Probabilistic formulations for transferring inferences from mathematical models to physical systems, SIAM J. Sci. Comput., 26 (2004), pp. 467{487, https://doi.org/10.1137/S106482750342670X.

[9] M. Gu, Jointly robust prior for Gaussian stochastic process in emulation, calibration and variable selection, Bayesian Anal., 14 (2019), pp. 857{885.

[10] M. Gu, Robust Calibration: Robust Calibration of Computer Models in R, preprint, https://arxiv.org/abs/2201.01476, 2022.

[11] M. Gu and J. O. Berger, Parallel partial Gaussian process emulation for computer models with massive output, Ann Appl. Stat., 10 (2016), pp. 1317{1347.

[12] M. Gu and L. Wang, Scaled Gaussian stochastic process for computer model calibration and prediction, SIAM/ASA J. Uncertain. Quantif., 6 (2018), pp. 1555{1583, https://doi.org/10.1137/17M1159890.

[13] M. S. Handcock and M. L. Stein, A Bayesian analysis of Kriging, Technometrics, 35 (1993), pp. 403{410.

[14] D. Higdon, J. Gattiker, B. Williams, and M. Rightley, Computer model calibration using high-dimensional output, J. Am. Stat. Assoc., 103 (2008), pp. 570{583.

[15] M. C. Kennedy and A. O'Hagan, Bayesian calibration of computer models, J. R. Stat. Soc. Ser. B, 63 (2001), pp. 425{464.

[16] M. R. Kosorok, Introduction to Empirical Processes and Semiparametric Inference, Springer, 2008.

[17] D. C. Liu and J. Nocedal, On the limited memory BFGS method for large scale optimization, Math. Program., 45 (1989), pp. 503{528.

[18] R. Paulo, G. Garca-Donato, and J. Palomo, Calibration of computer models with multivariate output, Comput. Stat. Data Anal., 56 (2012), pp. 3959{3974.

[19] M. Plumlee, Bayesian calibration of inexact computer models, J. Am. Stat. Assoc., 112 (2017), pp. 1274{1285.

[20] M. Plumlee, V. R. Joseph and H. Yang, Calibrating functional parameters in the ion channel models of cardiac cells, J. Am. Stat. Assoc., 111 (2016), pp. 500{509.

[21] J. Sacks, W. J. Welch, T. J. Mitchell, and P. H. Wynn, Design and analysis of computer experiments, Statist. Sci., 4 (1989), pp. 409{435.

[22] J. M. Salter, D. B. Williamson, J. Scinocca, and V. Kharin, Uncertainty quantication for com-puter models with spatial output using calibration-optimal bases, J. Am. Stat. Assoc., (2019), pp. 1{24.

[23] T. J. Santner, B. J. Williams, and W. I. Notz, The Design and Analysis of Computer Experiments, Springer-Verlag, 2003.

[24] K. Soetaert, T. Petzoldt, and R. W. Setzer, Package deSolve: Solving initial value dierential equations in R, J. Stat. Softw., 33 (2010), pp. 1{25.

[25] R. Tuo and C. J. Wu, Ecient calibration for imperfect computer models, Ann. Stat., 43 (2015), pp. 2331{2352.

[26] R. Tuo and C. F. J. Wu, A theoretical framework for calibration in computer models: Parametrization, estimation and convergence properties, SIAM/ASA J. Uncertain. Quantif., 4 (2016), pp. 767{795, https://doi.org/10.1137/151005841.

[27] A. W. Van der Vaart and J. H. Van Zanten, Adaptive Bayesian estimation using a Gaussian random eld with inverse gamma bandwidth, Ann. Stat., 37 (2009), pp. 2655{2675.

[28] G. Wahba, Spline Models for Observational Data, CBMS-NSF Reg. Conf. Ser. Appl. Math. 59, SIAM, 1990.

[29] R. K. Wong, C. B. Storlie, and T. Lee, A frequentist approach to computer model calibration, J. R. Stat. Soc. Ser. B , 79 (2017), pp. 635{648.

[30] F. Xie and Y. Xu, Bayesian projected calibration of computer models, J. Am. Stat. Assoc., 116 (2021), pp. 1965{1982.

[31] Y. Yang, A. Bhattacharya, and D. Pati, Frequentist coverage and sup-norm convergence rate in Gaussian process regression, preprint, https://arxiv.org/abs/1708.04753, 2017.

[32] H. Zhang, Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics, J. Am. Stat. Assoc., 99 (2004), pp. 250{261.

# Supplementary Materials for A theoretical framework of the scaled Gaussian stochastic process

## Mengyang Gu, Fangzheng Xie and Long Wang

All the formulas in this supplementary materials are cross-referenced in the main body of the article. We rst give a brief introduction of the Gaussian stochastic process model and reproducing kernel Hilbert space in Section S2. The proof for Section 2 is given in Section S3. The proof for Theorem 3.1 and two auxiliary lemmas are provided in Section S4. Section S5 encloses the proof for Theorem 3.3 and provides additional results regarding the convergence of S-GaSP calibration when kernel parameters are estimated.

**S1. Background: Gaussian stochastic process.** All the formulas in this supplementary materials are cross-referenced in the main body of the article. We rst give a brief introduction of the Gaussian stochastic process model and reproducing kernel Hilbert space in Section S2. The proof for Section 2 is given in Section S3. The proof for Theorem 3.1 and two auxiliary lemmas are provided in Section S4. Section S5 encloses the proof for Theorem 3.3 and provides additional results regarding the convergence of S-GaSP calibration when kernel parameters are estimated.

**S2. Background: Gaussian stochastic process.** Assume the mean and trend of the reality are properly modeled in the mathematical model. Consider to model the unknown discrepancy function in the calibration model (1) via a real-valued zero-mean Gaussian stochastic process ($\cdot$) on a p-dimensional input domain $X$,

$$\text{(S1)} \qquad \qquad \delta(\cdot) \sim \text{GaSP}(0, \sigma^2 K(\cdot, \cdot));$$

where $\sigma^2$ is a variance parameter and $K(x_a, x_b)$ is the correlation for any $x_a, x_b \in X$, parameterized by a kernel function. For simplicity, we assume $X = [0, 1]^p$ in this work.

For any $\{x_1, \ldots, x_n\}$, the outputs $(\delta(x_1), \ldots, \delta(x_n))^T$ follow a multivariate normal distribution

$$\text{(S2)} \qquad \qquad [\delta(x_1), \ldots, \delta(x_n) \mid \sigma^2, R] \sim \text{MN}(0, \sigma^2 R);$$

where the $(i, j)$ entry of $R$ is $K(x_i, x_j)$. Some frequently used kernel functions include the power exponential kernel and the Matern kernel. We defer the issue of estimating the param-

eters in the kernel function in Section [4] and assume $K(\cdot, \cdot)$ is known for now.

The reproducing kernel Hilbert space (RKHS), denoted as $\mathcal{H}$, attached to the Gaussian stochastic process $\text{GaSP}(0, \sigma^2 K(\cdot, \cdot))$, is the completion of the space of all functions

$$x \mapsto \sum_{i=1}^{k} w_i K(x_i, x), \quad w_1, \ldots, w_k \in \mathbb{R}; \ x_1, \ldots, x_k; x \in \mathcal{X}; k \in \mathbb{N},$$

with the inner product

$$\left\langle \sum_{i=1}^{k} w_i K(x_i, \cdot), \sum_{j=1}^{m} w_j K(x_j, \cdot) \right\rangle_{\mathcal{H}} = \sum_{i=1}^{k}\sum_{j=1}^{m} w_i w_j K(x_i, x_j).$$

For any function $f(\cdot) \in \mathcal{H}$, denote $\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle_{\mathcal{H}}}$ the RKHS norm or the native norm. Because the evaluation maps in RKHS are bounded linear, it follows from the Riesz representation theorem that for each $x \in \mathcal{X}$ and $f(\cdot) \in \mathcal{H}$, one has $f(x) = \langle f(\cdot), K(\cdot, x) \rangle_{\mathcal{H}}$.

Denote $L_2(\mathcal{X})$ the space of square-integrable functions $f : \mathcal{X} \to \mathbb{R}$ with $\int_{x \in \mathcal{X}} f^2(x)dx < \infty$. We denote $\langle f, g \rangle_{L_2(\mathcal{X})} := \int_{x \in \mathcal{X}} f(x)g(x)dx$ the usual inner product in $L_2(\mathcal{X})$. By Mercer's theorem, there exists an orthonormal sequence of continuous eigenfunctions $\{\phi_k\}_{k=1}^{\infty}$ with a sequence of non-increasing and non-negative eigenvalues $\{\rho_k\}_{k=1}^{\infty}$ such that

$$(S3) \qquad\qquad K(x_a, x_b) = \sum_{k=1}^{\infty} \rho_k \phi_k(x_a)\phi_k(x_b),$$

for any $x_a, x_b \in \mathcal{X}$.

The RKHS $\mathcal{H}$ contains all functions $f(\cdot) = \sum_{k=1}^{\infty} f_k \phi_k(\cdot) \in L_2(\mathcal{X})$ with $f_k = \langle f, \phi_k \rangle_{L_2(\mathcal{X})}$ and $\sum_{k=1}^{\infty} f_k^2 / \rho_k < \infty$. For any $g(\cdot) = \sum_{k=1}^{\infty} g_k \phi_k(\cdot) \in \mathcal{H}$ and $f(\cdot)$, the inner product can be represented as $\langle f, g \rangle_{\mathcal{H}} = \sum_{k=1}^{\infty} f_k g_k / \rho_k$. For more properties of the RKHS, we refer to Chapter 1 of [11] and Chapter 11 of [2].

## S2.1. The equivalence between the maximum likelihood estimator and the kernel ridge regression estimator in calibration.

Assume one has a set of observations $\mathbf{y}^F := \left(y^F(x_1), \ldots, y^F(x_n)\right)^T$ and mathematical model outputs $\mathbf{f}^M_\theta := (f^M(x_1, \theta), \ldots, f^M(x_n, \theta))^T$, where $\theta = (\theta_1, \ldots, \theta_q)^T \in \mathbb{R}^q$ is a q-dimensional vector of the calibration parameters. Denote the regularization parameter $\lambda := \sigma^2/(n\sigma_0^2)$. For the calibration model (1) with $\delta$ modeled as a GaSP in (S1), the marginal distribution of $\mathbf{y}^F$ follows a multivariate normal

after marginalizing out $\delta$,

(S4) $$[\mathbf{y}^F \mid \mu; \sigma^2; \phi]_\delta \sim MN(\mathbf{f}^M; \sigma^2((n\lambda)^{-1}\mathbf{R} + \mathbf{I}_n)):$$

Let $L(\phi)$ be the likelihood for $\phi$ in (S4) given the other parameters in the model. For any given $\lambda$, the maximum likelihood estimator (MLE) of $\phi$ is denoted as

(S5) $$\hat{\phi}_{\lambda;n} := \underset{\phi}{\arg\max}\, L(\phi):$$

Conditioning on the observations, $\hat{\phi}_{\lambda;n}$ and $\lambda$, the predictive mean of the discrepancy function at any $x \in \mathcal{X}$ has the following expression

(S6) $$\hat{\delta}_{\lambda;n}(x) := E[\delta(x) \mid \mathbf{y}^F; \hat{\phi}_{\lambda;n}; \lambda] = \mathbf{r}^T(x)(\mathbf{R} + n\lambda \mathbf{I}_n)^{-1}\mathbf{y}^F - \mathbf{f}^M_{\hat{\phi}_{\lambda;n}}$$

with $\mathbf{r}(x) = (K(x_1; x); \cdots; K(x_n; x))^T$ and $\mathbf{I}_n$ being the n-dimensional identity matrix.

It is well-known that the predictive mean in (S6) can be written as the estimator for the kernel ridge regression (KRR). In the following lemma, we show that $(\hat{\phi}_{\lambda;n}; \hat{\delta}_{\lambda;n}(\cdot))$ is equivalent to the KRR estimator.

Lemma S1. The maximum likelihood estimator $\hat{\phi}_{\lambda;n}$ defined in (S5) and predictive mean estimator $\hat{\delta}_{\lambda;n}(\cdot)$ defined in (S6) can be expressed as the estimator of the kernel ridge regression as follows

$$(\hat{\phi}_{\lambda;n}; \hat{\delta}_{\lambda;n}(\cdot)) = \underset{\phi, \delta}{\arg\min}\, \ell_{\lambda;n}(\phi; \delta); \quad 2; \delta(\cdot) \in \mathcal{H}$$

where

(S7) $$\ell_{\lambda;n}(\phi; \delta) = \frac{1}{n}\sum_{i=1}^{n}\{y^F(x_i) - f^M(x_i; \phi) - \delta(x_i)\}^2 + \lambda\|\delta\|^2_{\mathcal{H}}:$$

Proof of Lemma S1. By the representer lemma [6, 11], for any $\phi \in \Theta$ and $x \in \mathcal{X}$, one has

(S8) $$\hat{\delta}_{\lambda;n;\phi}(x) = \sum_{i=1}^{n} w_i(\phi) K(x_i; x):$$

Denote $\mathbf{w} = (w_1(\phi); \cdots; w_n(\phi))^T$. Since $\langle K(x_i; \cdot); K(x_j; \cdot)\rangle_{\mathcal{H}} = K(x_i; x_j)$, (S7) becomes to 3

nd and $w$ that minimize

(S9)
$$\frac{1}{n}(y^F - f^M - Rw)^T(y^F - f^M - Rw) + w^T Rw.$$

For any , solving the minimization for (S9) with regard to $w$ gives

(S10)
$$\hat{w} = (R + nI_n)^{-1}(y^F - f^M).$$

Then plugging $\hat{w}$ into (S9), based on the Woodbury matrix identity, one has

$$\frac{1}{n}(y^F - f^M - R\hat{w})^T(y^F - f^M - R\hat{w}) + \hat{w}^T R\hat{w}$$

$$= \frac{1}{n}(y^F - f^M)^T[(I_n - R(R + nI_n)^{-1})^T(I_n - R(R + nI_n)^{-1})]$$

$$\quad + (y^F - f^M)^T(R + nI_n)^{-1}R(R + nI_n)^{-1}(y^F - f^M)$$

(S11)
$$= (y^F - f^M)^T(R + nI_n)^{-1}(y^F - f^M);$$

which shows that the minimizer of on right-hand side of (S11) is the same as the MLE of in (S5). Finally, plugging the estimator $\hat{}_{;n}$ into (S9), the result follows from the KRR estimator of () in (S8) with the weights in (S10). ■

Although modeling the discrepancy function by the GaSP typically improves the prediction accuracy of the reality, the penalty term of (S7) only contains $\|\cdot\|_H$ to control the complexity of the discrepancy. As the RKHS norm is not equivalent to the $L_2$ norm, the calibrated computer model could deviate a lot from the best performed mathematical model in terms of the $L_2$ loss [8]. In Section 2, we introduce the scaled Gaussian stochastic process that predicts the reality as accurately as the GaSP with the aid of the mathematical model, but has more prior mass on the small $L_2$ distance between the reality and mathematical model. As a consequence, the KRR estimator of the new model penalizes both $\|\cdot\|_H$ and $\|\cdot\|_{L_2(X)}$ simultaneously.

### S3. Proof for Section 2.

Proof of Lemma 2.1. By Karhunen-Loeve expansion, we have

$$(x) = \sum_{i=1}^{\infty} \sqrt{\rho_i} Z_i \phi_i(x)$$

with $Z_i \overset{i.i.d}{\sim} N(0,1)$. Denote $W_k = \sum_{i=k+1}^{\infty} \rho_i Z_i^2$ for any $k \in N^+$. From the definition of 4

$Z = \int_{x \in X} \psi^2(x)dx$ and $\int_{x \in X} \phi_i^2(x)dx = 1$ for any $i \in \mathbb{N}^+$, it is straightforward to see that

(S12)
$$Z = \sigma^2(\lambda_1 Z_1^2 + \ldots + \lambda_k Z_k + W_k^2):$$

In the following expressions, we are conditioning on all parameters and they are dropped for simplicity. From the construction of

$$\zeta_z(x) = \mu(x) + \int_{x \in X}^{Z} \psi^2(x)dx = Z$$

and $Z \sim p_Z(\cdot)$, the joint density of $(Z_1; \ldots; Z_k; W_k)$ in the S-GaSP can be expressed as

$p_z(Z_1 = z_1; \ldots; Z_k = z_k; W_k = w_k)$

$= \int_{Z^0}^{Z_1} p(Z_1 = z_1; \ldots; Z_k = z_k; W_k = w_k \mid Z = z) p_Z(Z = z)dz$

$\int_0^{Z_1} \frac{p(Z = z; Z_1 = z_1; \ldots; Z_k = z_k; W_k = w_k)}{p(Z = z)} g_Z(Z = z)p(Z = z)dz$

$\propto p(Z_1 = z_1; \ldots; Z_k = z_k)p(W_k = w_k)$

$\int_0^{Z_1} p(Z = z \mid Z_1 = z_1; \ldots; Z_k = z_k; W_k = w_k) \exp\left(\frac{z_z}{2\sigma^2}\right) dz$

$\propto p(Z_1 = z_1; \ldots; Z_k = z_k)p(W_k = w_k)$

$\int_0^{Z_1} \mathbb{1}\left[z = \sigma^2(\lambda_1 z^2 + \ldots + \lambda_k z^2 + w_k)\right]_k \exp\left(z^{\frac{z}{2\sigma^2}}\right) dz \, !\#$

$\propto \exp\left(\frac{1}{2}\sum_{i=1}^{X^k} z_i^2\right) p(W_k = w_k) \exp\left(\frac{z}{2}\sum_{i=1}^{X^k} \lambda_i z_i^2 + w_k\right)$

$= \prod_{i=1}^{Y^k} \exp\left(\frac{1}{2}(1 + \sigma_{zi})z_i^2\right) p(W_k = w_k)\exp(\sigma_z w_k = 2);$

where $\mathbb{1}$ in the fourth step is a Dirac delta function.

After integrating out $W_k$, it is clear that $Z_i$'s are independently distributed as $N(0; 1=(1 + \sigma_{zi}))$ under the measure induced by the S-GaSP. Since $k$ is arbitrary, we have

$$\zeta_z(x) = \mu(x) + \sum_{i=1}^{X^d} \sqrt{\frac{\lambda_i}{1 + \sigma_{zi}}} Z_{\sqrt{i}}(x)$$

with $Z_i \overset{i.i.d}{\sim} N(0; 1)$, from which the proof is complete.  ■

Proof of Lemma 2.4.  First note that for any $x_a; x_b \in X$, we have $K(x_a; x_b) = \sum_{i=1}^{P_1} \lambda_i \phi_i(x_a)\phi_i(x_b)$ 5  ▌

and $K_z(x_a; x_b) = \frac{1}{P}\sum_{i=1}^{P}\frac{1}{\lambda_{z;i}}\phi_i(x_a)\phi_i(x_b)$ with $\lambda_{z;i} = \lambda_i=(1+z_i)$. For $h(\cdot) = \sum_{i=1}^{P} h_i\phi_i(\cdot) \in \mathcal{H}$ and $g(\cdot) = \frac{1}{P}\sum_{i=1}^{P} g_i\phi_i(\cdot) \in \mathcal{H}$, one has

$$\langle h; g\rangle_{\mathcal{H}_z} = \sum_{i=1}^{P}\frac{1}{\lambda_{z;i}}h_i g_i = \sum_{i=1}^{P}\frac{1}{\lambda_i}h_i g_i + z\sum_{i=1}^{P} h_i g_i = \langle h; g\rangle_{\mathcal{H}} + z\langle h; g\rangle_{L_2}: \quad\blacksquare$$

**Proof of Lemma 2.5.** We show below that for any $\lambda \geq 2$ and any $x \in X$, one has

(S13)
$$\hat{\mu}_{\lambda;z;n}(x) = \sum_{i=1}^{n} w_{z;i}(\lambda)K_z(x_i; x):$$

For any $\mu(\cdot) \in \mathcal{H}$, decomposing it into the linear combination of the basis $\{K_z(x_i; \cdot)\}_{i=1}^{n}$ and the orthogonal complement $v(\cdot)$ gives

$$\mu(\cdot) = \sum_{i=1}^{n} \tilde{w}_{z;i}(\lambda)_z K_z(x_i; \cdot) + v(\cdot); $$

where $\langle v(\cdot); {}_z K_z(\cdot; x_i)\rangle_{\mathcal{H}_z} = 0$ for $i = 1; \ldots; n$.

To evaluate $\mu(\cdot)$ at $x_j$ for any $j = 1; \ldots; n$, we have

$$\mu(x_j) = \left\langle \sum_{i=1}^{n} \tilde{w}_{z;i}(\lambda)_z K_z(x_i; \cdot) + v(\cdot); K_z(x_j; \cdot)\right\rangle_{\mathcal{H}_z}$$
$$= \sum_{i=1}^{n} \tilde{w}_{z;i}(\lambda)_z K_z(x_i; x_j);$$

which is independent from $v(\cdot)$. Hence the rst term on right-hand side of (11) is also independent from $v(\cdot)$. For the second term on right-hand side of (11), since $v(\cdot)$ is orthogonal to $\{K_z(x_i; \cdot)\}_{i=1}^{n}$, plugging in the decomposition of $\mu(\cdot)$, we have

$$\|\mu\|_{\mathcal{H}_z}^2 = \left(\left\|\sum_{i=1}^{n} \tilde{w}_{z;i}(\lambda)_z K_z(x_i; \cdot)\right\|_{\mathcal{H}_z} + \|v\|_{\mathcal{H}_z}^2\right)$$
$$\geq \left\|\sum_{i=1}^{n} \tilde{w}_{z;i}(\lambda)_z K_z(x_i; \cdot)\right\|_{\mathcal{H}_z}^2:$$

Thus choosing $v(\cdot) = 0$ does not change the rst term on right-hand side of (11), but also minimizes the second term on right-hand side of (11). Letting $w_{z;i}(\lambda) = \tilde{w}_{z;i}(\lambda)_z$, we have proved (S13). The rest of the proof can be derived similarly as the proof for Lemma S1, so it is omitted here. $\quad\blacksquare$

**S4. Proof for Section 3.1.** We prove Theorem 3.1 in this Section. Two auxiliary lemmas used for the proof of Theorem 3.1 are given after the proof.

**Proof for Theorem 3.1.** Define a new inner product on $H$ as

(S14)
$$\langle f, g\rangle = (1 + \mu^p)\langle f, g\rangle_{L_2(X)} + \langle f, g\rangle_H$$

Let $f = \sum_{k=1}^{\infty} f_k \phi_k$ and $g = \sum_{k=1}^{\infty} g_k \phi_k$ be elements in $H$. Then

$$= (1 + \mu^p) \sum_{k=1}^{\infty} f_k g_k + \sum_{k=1}^{\infty} \frac{f_k g_k}{\lambda_k} = \sum_{k=1}^{\infty} \left(1 + \mu^p + \frac{1}{\lambda_k}\right) f_k g_k.$$

By letting $\tilde\lambda_k$ by $\tilde\lambda_k^{-1} = 1 + \mu^p + \lambda_k^{-1}$, we can define a new reproducing kernel

(S15)
$$K(x; x^0) = \sum_{k=1}^{\infty} \tilde\lambda_k \phi_k(x)\phi_k(x^0)$$

Since $c^{-1}k^{-2m/p} \le \lambda_k \le C^{-1}k^{-2m/p}$ and $|\phi_i(\cdot)| < C$ for some positive constants $c$, $C$ and $\tilde C$, bounding the sums by integrals, we have

$$\sup_{x;x^0} K(x; x^0) \le \tilde C^2 \sum_{k=1}^{\infty} \frac{1}{1 + ck^{2m/p}} \le \tilde C^2 \sum_{k=1}^{\infty} \int_{k-1}^{k} \frac{1}{1 + cx^{2m/p}}dx \frac{1}{2^{p/2m}}$$

$$= \tilde C^2 c^{-p/2m} \int_0^{\infty} \frac{(c)^{p/2m}}{1 + f(c)^{p/2m}x^{2m/p}}dx =$$

$$\tilde C^2 c^{-p/2m}\mu^{p/2m} \int_0^{\infty} \frac{1}{1 + x^{2m/p}}dx.$$

Thus

(S16)
$$\sup_{x;x^0} K(x; x^0) \le C_K^2 \mu^{p/(2m)};$$

for some constant $C_K$ depending on $K$. Define the following linear operators $F : H \to H$ and $P : H \to H$ via

$$(Fg)(x) = \int_X g(x^0)K(x; x^0)dx^0; \quad \text{and} \quad (Pg)(x) = g(x) - (Fg)(x);$$

7

Clearly, we have

$$\langle f, Fg\rangle = \sum_{k=1} \langle f, \phi_k\rangle_{L_2(X)}\langle g, \phi_k\rangle_{L_2(X)} = \langle f, g\rangle_{L_2(X)};$$

(S17) $\qquad \langle f, Pg\rangle = \langle f, g\rangle \qquad \langle f, Fg\rangle = \rho\langle f, g\rangle_{L_2(X)} + \langle f, g\rangle_H:$

Denote the loss function

$$\ell_n(f) = \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2 + \rho\|f\|_{L_2(X)}^2 + \|f\|_H^2;$$

and the estimator $\hat{f}_n := \operatorname{argmin}_{f\in H}\ell_n(f)$. Let $D\ell_n(f) : H \to H$ be the Frechet derivative of $\ell_n$ evaluated at $f$. Clearly, for any $g \in H$,

$$D\ell_n(f)g = \frac{2}{n}\sum_{i=1}^{n}\left(f(x_i) - y_i\right)\langle hK(x_i;\cdot); g\rangle_i + 2\langle hPf; g\rangle *$$

(S18) $$= \left\langle \frac{2}{n}\sum_{i=1}^{n}(f(x_i) - y_i)K(x_i;\cdot) + 2(Pf)(\cdot); g(\cdot)\right\rangle:$$

It follows that $D\ell_n(\hat{f}_n)g = 0$ for all $g \in H$, and hence, $S_n(\hat{f}_n)(\cdot) = 0$, where

$$S_n(f)(\cdot) = \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))K(x_i;\cdot) - (Pf)(\cdot):$$

Define $S(f)(\cdot) = E_{y;x}(S_n(f)(\cdot))$. Then

$$S(f)(\cdot) = \int_{x\in X}(f_0(x) - f(x))K(x;\cdot)dx - (Pf)(\cdot)$$

$$= (F(f_0 - f))(\cdot) - (Pf)(\cdot) = (Ff_0)(\cdot) - f(\cdot);$$

and therefore, $S(Ff_0)(\cdot) = 0$. Let $f = \hat{f}_n - Ff_0$. By the definitions of $S_n$ and $S$, we

8

have

$$\langle S_n(\hat{f}_n) - S(\hat{f}_n), \phi \rangle - \langle S_n(Ff_0) - S(Ff_0)g, \phi \rangle$$
$$= \langle S_n(\hat{f}_n) - S_n(Ff_0), \phi \rangle - \langle S(\hat{f}_n) - S(Ff_0), \phi \rangle$$
$$= \frac{1}{n} \sum_{i=1}^{X^n} \langle Ff_0(x_i) - \hat{f}_n(x_i) \rangle K(x_i; \cdot) + \langle P(Ff_0)(\cdot) - \hat{f}_n(\cdot), \phi \rangle$$
$$+ \langle \hat{f}_n(\cdot) - (Ff_0)(\cdot), \phi \rangle$$
$$= \frac{1}{n} \sum_{i=1}^{X^n} f(x_i) K(x_i; \cdot) - \langle (Pf)(\cdot) + (f)(\cdot) \rangle$$
$$= \frac{1}{n} \sum_{i=1}^{X^n} f(x_i) K(x_i; \cdot) + E_x f f(x) K(x; \cdot)g : i=1$$

On the other hand, $S_n(\hat{f}_n)(\cdot) = S(Ff_0)(\cdot) = 0$ and $S(\hat{f}_n)(\cdot) = \hat{}(Ff_0)(\cdot) - \hat{f}_n(\cdot) = \hat{} f(\cdot)$. Therefore,

$$\langle S_n(\hat{f}_n) - S(\hat{f}_n), \phi \rangle - \langle S_n(Ff_0) - S(Ff_0)g, \phi \rangle = \langle f(\cdot), S_n(Ff_0)(\cdot) \rangle:$$

Dene the event

$$A_n(t) = \left( \left\| \frac{1}{n} \sum_{i=1}^{n} g(x_i)K(x_i; \cdot) - E_x f g(x)K(x; \cdot)g \right\| < t\|g\|_K \text{ for all } g \in H \right) :$$

Applying Lemma S2 on $g(\cdot)=\|g\|_K$,

$$P_x f A_n(t)g \geq 1 - 2\exp \frac{p(6m - p)=(4m^2) nt^2}{K}$$

for some constant $K > 0$. The deviation threshold t will be specied later, and from now we consider data points $(x_i; y_i)_{i=1}^n$ over the event $A_n(t)$.

Over the event $A_n(t)$, we have

$$\langle S_n(\hat{f}_n) - S(\hat{f}_n), \phi \rangle - \langle S_n(Ff_0) - S(Ff_0)g, \phi \rangle$$
$$= \frac{1}{n} \sum^{X^n} f(x_i)K(x_i; \cdot) + E_x f f(x)K(x; \cdot)g \, i=1$$
$$= \langle f(\cdot), S_n(Ff_0)(\cdot) \rangle;$$

implying that

$$\|f - S_n(F f_0)\| = \left\| \frac{1}{n}\sum_{i=1}^{n} f(x_i)K(x_i; \cdot) - E_x f f(x)K(x; \cdot)g \right\|$$

$$\le t\|f\|.$$

Now we proceed to bound $\|S_n(F f_0)\|$. Write

$$\|S_n(F f_0)\| = \left\| \frac{1}{n}\sum_{i=1}^{n} f y_i - F f_0(x_i)g K(x_i; \cdot) - (P F f_0)(\cdot) \right\|$$

$$\le \left\| \frac{1}{n}\sum_{i=1}^{n} f f_0(x_i) - F f_0(x_i)g K(x_i; \cdot) - f F(f_0 - F f_0)g(\cdot) \right\|$$

$$+ \left\| \frac{1}{n}\sum_{i=1}^{n} {}_i K(x_i; \cdot) \right\|$$

$$= \left\| \frac{1}{n}\sum_{i=1}^{n} f f_0(x_i) - F f_0(x_i)g K(x_i; \cdot) - E_x[f f_0(x) - F f_0(x)g K(x; \cdot)] \right\|$$

$$+ \left\| \frac{1}{n}\sum_{i=1}^{n} {}_i K(x_i; \cdot) \right\|$$

$$\le t\|f_0 - F f_0\| + \left\| \frac{1}{n}\sum_{i=1}^{n} {}_i K(x_i; \cdot) \right\|;$$

where the last inequality is due to the construction of the event $A_n(t)$. To bound the second term of the preceding display, we let $= [K(x_i; x_j)]_{nn}$ and $= [{}_1; \ldots; {}_n]^T$. By the Hanson-Wright inequality [7], for all $x > 0$, we have

$$P_x \quad {}^T \quad {}^2 \quad {}_0 \left[ tr() + 2\sqrt{tr() x + 2\|k\|_F x^2} \right] \quad e^{-x^2}:$$

Since by the Cauchy-Schwarz inequality,

$$\text{tr}(\Sigma) = \sum_{i=1}^{n} \| K(x_i; \cdot) \|_{\mathcal{H}}^2 = \sum_{i=1}^{n} K(x_i; x_i) \le C^2 n^{p=(2m)}; $$

$$\text{tr}(\Sigma^2) \le \sum_{i=1}^{n}\sum_{j=1}^{n} \| K(x_i; \cdot) \|_{L_2(X)} \| K(x_j; \cdot) \|_{L_2(X)} \le C^4 n^2 \overset{p=m}{\phantom{.}}; $$

$$\| \Sigma \|_F = \sqrt{\text{tr}(\Sigma^2)} \le C^2 n^{p=(2m)}; $$

it follows that

(S19)    $$ \text{tr}(\Sigma) + 2\sqrt{\text{tr}(\Sigma)x} + 2\| \Sigma \|_F x \le 2 C^2 n^{p=(2m)}(1 + 2x + 2x^2): $$

Set the event $B_n$ to be

$$ B_n = \left( \frac{1}{n} \sum_{i=1}^{n} e_i K(x_i; \cdot) < \delta_0 C_K n^{1=2} \frac{p=(4m)}{\phantom{x}} 1=2 \right); $$

where $\delta = 2 + 3x^2$. Since $1 + 2x + 2x^2 \le 2 + 3x^2 = \delta$, by taking $x = {}^p(\frac{\phantom{-}}{\phantom{-}} - 2)=3$, we have $P(B_n) \ge 1 - \exp(-(\delta - 2)=3)$ for any $\delta > 2$. Putting all pieces obtained above together, we have

$$ \| f \| \le \| f - \hat{S}_n(F f_0) \| + \| \hat{S}_n(F f_0) \| $$
$$ \le t \| f \| + t \| f_0 \| - \| F f_0 \| + \delta_0 C_K n^{1=2} \frac{p=(4m)}{\phantom{x}} 1=2 \quad \text{(S20)} $$
$$ = t \| f \| + t \| P f_0 \| + \delta_0 C_K n^{1=2} \frac{p=(4m)}{\phantom{x}} 1=2; $$

over the event $A_n(t) \setminus B_n$. Now take $\delta = n^{2m=(2m+p)}$. Choose any $C \in (0; 1)$ and let $t = \sqrt{\frac{\phantom{x}}{\kappa}} = (n^{(1-C)} \log(2))$. Then, for suciently large $n$,

$$ P_x\{A_n(t)\} \ge 1 - 2\exp\left(\frac{nt^2}{\kappa} - 1\right) \ge \exp\left(-n^C\right); $$

where $\kappa = (2m - p)^2 = (2m(2m + p))$, and therefore,

$$ \| f \| \le \| P f_0 \| + 2\delta_0 C_K n^{m=(2m+p)1=2}; $$

with probability at least

$$P f A_n(t) \setminus B_n) g = 1 \quad P f A_n^c(t) [ B_n^c g$$
$$1 \quad P f A_n^c(t) g \quad P(B_n^c) = 1 \quad \text{expf} ( 2)=3g \quad \text{expf} n^C g$$

for suciently large n. Observe that

$$kPf_0k \overset{2}{=} \sum_{k=1}^{X^1} (1 \quad k)hf_0; ki_{L_2(X)}k() \quad \overset{2}{=} \sum_{k=1}^{X^1} \frac{(1 \quad k)^2}{k} hf_0; ki_{L_2(X)}^2$$

$$= \sum_{k=1}^{X^1} \frac{(\overset{p}{} + = )k^2}{1 + \quad p+=_k} hf_0; ki_{L_2(X)}^2 \quad 2 + 2(=k)^2 \frac{hf_0; ki_{L_2(X)} 1}{1 + = k} \quad 1 \quad 2$$

$$2 \sum_{k=1} hf_0; ki_{L_2(X)_2} + 2 \sum_{k=1} \frac{X}{k} 1 hf_0; ki_{L_2(X)}$$

$$= 2kf_0k_{L_2(X)}^2 + 2kf_0k_H^2$$

$$2n^{2m=(2m+p)} kf_0k_{L_2(X)} + kf_0k_H^2:$$

Hence, we proceed to compute

$$k\hat{f}_n \quad f_0k_{L_2(X)} \quad kf_n^\wedge \quad f_0k$$
$$kf_n^\wedge \quad Ff_0k + kFf_0 \quad f_0k$$
$$= kfk + kPf_0k$$
$$2 \overset{p}{} 2(kf_0k_{L_2(X)} + kf_0k_H) + 2_0C_K^{1=2} n^{m=(2m+p)}$$

with probability at least 1 expf ( 2)=3g exp n^C for suciently large n. The bound for $kf_{;z;n^\wedge} \quad f_0k_H$ follows immediately by the denition of jj jj, completing the proof. ∎

The following the Lemma S1 is Theorem 3.6 in [5], which is needed for the proof of Lemma S2.

Lemma S1. Let $(X_j)_{j=0}^1$ be a sequence of random elements in a Hilbert space H with norm $k k_H$. Suppose that $(X_j)_{j=0}^1$ forms a martingale in the sense that $E(X_j j X_0; :::; X_{j 1}) = X_j$ a.s., and that the dierence sequence $(D_j)_{j=1}^1 = (X_j \quad X_{j 1})_{j=1}^1$ satises $kD_jk^2 \quad b_j^2$ a.s. and $P_{j=1}^1 b_j^2 \quad b^2$. Then for any t 0,

$$P \quad \sup_{j 1} kX_jk_H \quad t \quad 2\exp \quad \frac{t^2}{2b^2}:$$

12

The following maximum inequality for functional empirical processes in the Sobolev space $W_2^m(X;1)$, which generalizes Lemma 5.1 in [12] to multivariate functions, is of fundamental importance to the proof of Theorem 3.1.

Lemma S2. Denote $W_2^m(X;1) = \{f \in W_2^m(X) : \|f\| \le 1\}$. Suppose $x_1, \ldots, x_n$ are independently and uniformly drawn from $X$. Then there exists some constant $\kappa$ depending on the kernel $K$, such that for any $t > 0$,

$$
P_x\left\{ \sup_{g \in W_2^m(X;1)} \left| \frac{1}{n}\sum_{i=1}^{n} [g(x_i)K(x_i;\cdot)] - E_x[g(x)K(x;\cdot)] \right| \ge t \right\}
$$
$$
\le 2\exp\left( -\frac{d(6m-d)=(4m^2)nt^2}{\kappa} \right):
$$

Proof of Lemma S2. We follow the argument used in the proof of Lemma 6.1 in [13]. Denote

$$
\{Z_n(g)\}(\cdot) = \frac{1}{n}\sum_{i=1}^{n}[g(x_i)K(x_i;\cdot) - E_x\{g(x)K(x;\cdot)\}]:
$$

Fix $g, h \in H$, $n$, and $\cdot$, consider the following sequence of martingale $(X_j)_{j=0}^\infty$ in $H$:

$$
X_j = \begin{cases} 0; & \text{if } j = 0; \\ j\{Z_j(g) - Z_j(h)\}g; & \text{if } j = 1,\ldots,n \\ X_n; & \text{if } j \ge n+1: \end{cases}
$$

Clearly, for $j = 1,\ldots,n$,

$$
(X_j - X_{j-1})(\cdot) = \{g(x_j) - h(x_j)\}gK(x_j;\cdot) - E_x[\{g(x_j) - h(x_j)\}gK(x_j;\cdot)]
$$

and $X_j - X_{j-1} = 0$ for $j \ge n+1$. Observe that

$$
\|K(x_j;\cdot)\| = \sqrt{\langle K(x_j;\cdot), K(x_j;\cdot)\rangle} = \sqrt{K(x_j;x_j)} \le C_K^{d=(4m)}
$$

with probability one. Therefore, with probability one, we have

$$
\|X_j - X_{j-1}\|^2 \le 4C_K^{2\,d=(2m)}\|g - h\|_{L_1}^2
$$

for $j = 1,\ldots,n$, and hence, we invoke the bounded difference inequality for martingales in

13

Banach space (Lemma S1) to derive

$$P\left(\|Z_n(g) - Z_n(h)\| \geq t\right) = P\left(\|nfZ_n(g) - Z_n(h)g\| \geq nt\right)$$

$$\leq P\left(\sup_{j \geq 1} \|jfZ_j(g) - Z_j(h)g\| \geq nt\right)$$

$$\leq 2\exp\left(-\frac{nt^2}{8C_K \ 2^{d=(2m)}\|g - h\|_{L_1}^2}\right).$$

Applying Lemma 8.1 in [3], we obtain the following bound

$$(S21) \qquad \|\|Z_n(g) - Z_n(h)\|\|_{\psi_2} \leq \frac{24C_K^{\frac{p}{2}} \ d^{2}=(4m)}{\sqrt{p n}} \|g - h\|_{L_1};$$

where $\|\cdot\|_{\psi_2}$ is the Orlicz norm associated with $\psi_2(s) = \exp(s^2) - 1$.

Now let $\psi = \psi\log(3=2)g^{1=2}$ and set $\phi(x) = \psi_2(x)$. Clearly, $\phi(1) = 1=2$, and $\phi(x)\phi(y) \leq \phi(xy)$ for any $x; y \geq 1$. Applying Lemma 8.2 in [3], the Orlicz norm of the maximum of finitely many random variables can be bounded by the maximum of these Orlicz norms as follows:

$$\left\|\max_{1 \leq i \leq k}(\xi_i)\right\|_{\psi} = \left\|\max_i \phi^{-1}(k) \max_{1 \leq i \leq k} \|\xi_i\| \right\| = \phi_2^{-1}(k) \max_{1 \leq i \leq k} \|\xi_i\|_{\psi_2};$$

namely,

$$(S22) \qquad \left\|\max_{1 \leq i \leq k}\right\|_{\psi_2} \leq \psi_2^{-1}(k) \max_{1 \leq i \leq k} \|\xi_i\|_{\psi_2};$$

where $\{\xi_i\}_{i=1}^k$ are finitely many random variables.

Next we apply the "chaining" argument. Let $\varepsilon > 0$ be some constant to be determined later. Construct a sequence of function classes $(G_j)_{j=0}^1$ in $H(1)$ satisfying the following conditions:

(i) For any $G_j$ and any $h_j; g_j \in G_j$, $\|h_j - g_j\|_{L_1} \geq \varepsilon=2^j$, and $G_j$ is maximal in the sense that for any $g_j \notin G_j$, there exists some $h_j \in G_j$ such that $\|h_j - g_j\| < \varepsilon=2^j$.

(ii) For any $G_{j+1}$, and any $g_{j+1} \in G_{j+1}$, select a unique element $g_j \in G_j$ such that $\|g_{j+1} - g_j\|_{L_1} \leq \varepsilon=2^j$. Thus, there exists a finite sequence $(g_0; g_1; \ldots; g_{j+1})$ such that $\|g_i - g_{i+1}\|_{L_1} \leq \varepsilon=2^i$ for $i = 0; \ldots; j$, and $g_i \in G_i$.

Therefore, for any $g_{j+1}; h_{j+1} \in G_{j+1}$ with $\|g_{j+1} - h_{j+1}\|_{L_1} \leq \varepsilon$, there exists two sequences

14

$(g_i)_{i=0}^{j+1}$, $(h_i)_{i=0}^{j+1}$, such that $g_i, h_i \in G_i$, $\max\{\|g_i - g_{i+1}\|_{L_1}; \|h_i - h_{i+1}\|_{L_1}\} \le 2^i$, and that

$$\|g_0 - h_0\|_{L_1} \le \sum_{i=0}^{j}(\|g_i - g_{i+1}\|_{L_1} + \|h_i - h_{i+1}\|_{L_1}) + \|h_{j+1} - g_{j+1}\|_{L_1}$$

$$\le \sum_{i=0}^{j} 2 \cdot \frac{1}{2^i} + \epsilon \le 5\epsilon;$$

and hence, by (S21) one has

$$\text{(S23)} \qquad \|\|Z_n(g_0) - Z_n(h_0)\|\|_2 \le \frac{5^{\frac{p}{24}} C_K^2}{p \cdot n} \quad p=(4m)\epsilon:$$

We also notice that $G_j \subseteq H(1) \cap \{f \in H : \|f\|_H \le \frac{1}{2}\}$, and therefore, the cardinality of $G_j$ can be bounded by the metric entropy of $\{f \in H : \|f\|_H \le \frac{1}{2}\}$, which is known in the literature [1]:

$$\log |G_j| \le \log N_{[]}(\epsilon = 2^j; \{f \in H : \|f\|_H \le \frac{1}{2}\}; \|\cdot\|_{L_1}) \le c_0 \cdot \frac{1}{2^j}^{p=m} ;$$

where $c_0$ is some absolute constant.

Now suppose $g, h$ are arbitrary functions in $H(1)$ such that $\|g - h\|_{L_1} \le 2$. For any $j \ge 2$, there exists $g_j, h_j \in G_j$ such that

$$\max\{\|g_j - g\|_{L_1}; \|h_j - h\|_{L_1}\} \le 2^j;$$

15

and hence, $\|g_j - h_j\|_{L_1} \leq \varepsilon''$. Therefore, for any $j \geq 2$,

$$\sup_{g,h \in 2W_2^m(X;1);\|g-h\|_{L_1} \leq \varepsilon''} \|Z_n(g) - Z_n(h)\|_2$$

$$\leq \sup_{g,h \in 2W_2^m(X;1);\|g-h\|_{L_1} \leq \varepsilon''} \|Z_n(g) - Z_n(g_j)\|_2 + \|Z_n(g_j) - Z_n(h_j)\|_2$$

$$+ \|Z_n(h_j) - Z_n(h)\|_2$$

$$\leq \frac{2^p \overline{24}C_K^2}{p} \left( \frac{d=(4m)}{n} \right) \max \{\|g - g_j\|_{L_1}; \|h - h_j\|_{L_1}\} + \sup_{g,h \in 2W(X;1);\|g-h\| \leq \varepsilon''} \|Z_n(g_j) - Z_n(h_j)\|_2$$

$$\leq \frac{2^p \overline{24}C_K^2 m}{\overline{n}} \frac{d=(4m)\varepsilon''}{2^j} + 1 \max_{g_j,h_j \in G_j;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|Z_n(g_j) - Z_n(h_j)\|_2 :$$

We focus on the second term of the preceding display. Fix $j \geq 2$, for any $g_j, h_j \in G_j$, consider the finite sequences $(g_0; g_1; \ldots; g_j)$ and $(h_0; h_1; \ldots; h_j)$ such that $g_i; h_i \in G_i$ and $\|g_i - g_{i+1}\|_{L_1} \leq \varepsilon'' = 2^i$, $i = 1; \ldots; j - 1$. Invoking the inequality (S22), we have

$$\max_{g_j,h_j \in G_j;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|Z_n(g_j) - Z_n(h_j)\|_2$$

$$\leq \max_{g_j,h_j \in G_j;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|\{Z_n(g_j) - Z_n(h_j)\} - \{Z_n(g_0) - Z_n(h_0)\}\|_2$$

$$+ \max_{g_0,h_0 \in G_0;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|Z_n(g_0) - Z_n(h_0)\|_2$$

$$\leq \max_{g_j,h_j \in G_j;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|\{Z_n(g_j) - Z_n(h_j)\} - \{Z_n(g_0) - Z_n(h_0)\}\|_2$$

$$+ \frac{2^p}{\underline{}} \log \overline{(1 + jG_0 G_j)} \max_{(g_0;h_0) \in G_0 G_0;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|\|Z_n(g_0) - Z_n(h_0)\| \|_2 :$$

Clearly, the second term can be bounded by inequality (S23):

$$\frac{2^p}{\underline{}} \log \overline{(1 + jG_0 G_j)} \max_{(g_0;h_0) \in G_0 G_0;\|g_j - h_j\|_{L_1} \leq \varepsilon''} \|\|Z_n(g_0) - Z_n(h_0)\| \|_2$$

$$\leq \frac{10^p \overline{24}C_K^2}{p} \left( \frac{p=(4m)\varepsilon''}{n} \right) \frac{q}{\log \overline{1 + \exp} \left( 2c_0 \frac{p=2m)\varepsilon''}{} \right) p=m} ;$$

16

since

$$|G_0 \ G_0| = |G_0|^2 \ \exp(2 \log N_{[]}(\varepsilon; \|f\| \ 1g; \| \ \|_{L_1})) \ \text{(S24)}$$

$$\exp(2c_0 \ {}^{p=(2m)}\varepsilon \ {}^{p=m});$$

it suces to bound the rst term. Write

$$\max_{\substack{g_j; h_j 2 G_j; \|g_j \ h_j\|_{L_1} \varepsilon \\ j \ 1}} \|f Z_n \ (g_j) \ Z_n \ (h_j)g \ f Z_n \ (g_0) \ Z_n \ (h_0)g\|^2$$

$$2 \ \sum_{i=0} \max_{(g_i;g_{i+1}) 2 G_i G_{i+1}; \|g_i \ g_{i+1}\|_{L_1} \varepsilon=2^i} \|Z_n \ (g_{i+1}) \ Z_n \ (g_i)\|^2$$

$$2 \ \sum_{i=0}^{\lceil 1 2p} \log(1 + |G_i| |G_{i+1}|)$$

$$\max_{(g_i;g_{i+1}) 2 G_i G_{i+1}; \|g_i \ g_{i+1}\|_{L_1} \varepsilon=2^i} \|\|Z_n \ (g_{i+1}) \ Z_n \ (g_i)\| \ \|^2$$

$$\frac{4^p \ 24C_{K}^2}{\frac{p}{\sqrt{n}}} \ \sum_{i=0}^{p=(4m) \ \lceil 1} \sqrt{\log \ 1 + \exp \ 2c_0 \ {}^{p=(2m)}_{(}(\varepsilon=2^i) \ {}^{p=m} \ \varepsilon=2^i};$$

where inequalities (S21) and (S22) are applied. Bounding the sum by integral, we have

$$\sum_{i=0}^{\lceil 1} \sqrt{\log \ 1 + \exp \ 2c_0 \ {}^{p=}_{(}{}^{2m)}_{p=m} \ i}$$
$$(\varepsilon=2^i) \quad \varepsilon=2$$

$$\sum_{i=0}^{\lceil 1} \int_{\varepsilon=2^{i+1}}^{\varepsilon=2^i} \frac{q}{\log f 1 + \exp(2c_0 \ {}^{p=(2m)}x \ {}^{p=m})g} dx$$

$$\int_{\varepsilon}^{\varepsilon} \frac{q}{\log f 1 + \exp(2c_0 \ {}^{p=(2m)}x \ {}^{p=m})g} dx: 0$$

Putting all pieces above together, we obtain the following bound:

$$\sup_{g; h 2 W_{2n}(X;1); \|g \ h\|_{L_1} \varepsilon} \|Z_n(g) \ Z_n(h)\|^2$$

$$. \ \frac{{}^{p=(4m)}}{\frac{p}{\sqrt{n}} \ \frac{\varepsilon}{2^j}} + \int_0^{\varepsilon} \frac{q}{\log f 1 + \exp(2c_0 \ {}^{p=(2m)}x \ {}^{p=m})g} dx$$

$$+ \ \varepsilon \ \sqrt{\log f 1 + \exp(2c_0 \ {}^{p=}_{(}{}^{2m)}\varepsilon \ {}^{p=m})g}:$$

By taking $j \to 1$, we can let the first term in the squared bracket tend to 0, and hence,

$$
\sup_{g,h \in W_2^m(X;1); \|g-h\|_{L_1} \leq \varepsilon} \|Z_n(g) - Z_n(h)\|
$$

$$
\lesssim \frac{\varepsilon^{p-(4m)}}{n} \int_0^{\varepsilon^{-q}} \log(1 + \exp(2c_0 \varepsilon^{-(2m)} x^{-m})) \, dx
$$

$$
+ \left[ \log(1 + \exp(2c_0 \varepsilon^{-(2m)} \varepsilon^{-m})) \right]
$$

$$
\cdot \frac{\varepsilon^{p-(4m)}}{n} \int_0^{\varepsilon^{-q}} \log(1 + \exp(2c_0 \varepsilon^{-(2m)} x^{-m})) \, dx \to 0
$$

Now we take $h = 0$, which implies $Z_n(h) = 0$ by the construction of $Z_n$. Furthermore, by the property of reproducing kernel $K$ and the Cauchy-Schwarz inequality,

$$
\|g - h\|_{L_1} \leq \sup_{x \in X} |g(x)| = \sup_{x \in X} |\langle g(\cdot), K(x, \cdot) \rangle|
$$

$$
\leq \sup_{x \in X} \|g\| \sqrt{\langle K(x,\cdot), K(x,\cdot) \rangle} \leq C_K \varepsilon^{p-(4m)}.
$$

Taking $\varepsilon = C_K \varepsilon^{p-(4m)}$, we obtain

$$
\sup_{g \in W_2^m(X;1)} \|Z_n(g)\|
$$

$$
\leq \sup_{g \in W_2^m(X;1); \|g\|_{L_1} \leq \varepsilon} \|Z_n(g)\|
$$

$$
\leq \sup_{g,h \in W_2^m(X;1); \|g-h\|_{L_1} \leq \varepsilon} \|Z_n(g) - Z_n(h)\|
$$

$$
\lesssim n^{1/2} \varepsilon^{p-(4m)} \int_0^{\varepsilon^{-q}} \log(1 + \exp(2c_0 \varepsilon^{-(2m)} x^{-m})) \, dx \to 0
$$

$$
\lesssim n^{1/2} \varepsilon^{p(6m-p)=(8m^2)}.
$$

Hence, invoking Lemma 8.1 in [3], we finally obtain

$$
\mathbb{P}\left( \sup_{g \in W_2^m(X;1)} \|Z_n(g)\| > t \right) \leq 2 \exp\left( -\frac{nt^2}{2\varepsilon^{p(6m-p)=(4m^2)}} \right);
$$

for some absolute constant $\varepsilon_K$ depending on $K$ only, completing the proof. $\blacksquare$

**S5. Proof for Section 3.2.** Denote $\hat{\theta}_z := \hat{\theta}_{;z;n}$, $\ell_z(\theta) :\overset{\Delta}{=} \ell_{;z;n}(\theta)$ and $\ell_z(\theta;\cdot) := \ell_{;z;n}(\theta;\cdot)$ in (11).

We need the following Corollary S1 and Lemma S2 to prove theorem 3.3. Corollary S1 is a direct consequence of Theorem 3.1. We repeatedly use the fact that for any $f(\cdot) \in L_2(X)$, there exists a constant $C$ such that $\|f\|_{L_2(X)} \le C \|f\|_H$ in the following proof.

**Corollary S1.** Denote $\hat{\theta}_z = \text{argmin}_{\theta \in H} \ell_z(\theta;\cdot)$ for each $z \in \Theta$. Under the Assumptions A1 to A6, for suciently large $n$ and any $\gamma > 2$ and $C \in (0;1)$, with probability at least $1 - \exp\{-(\gamma-2)=3\} - \exp\{-n^C\}$, one has

$$\sup_{z \in \Theta} \|\hat{\theta}_{z;\theta}(\cdot) - (y^R(\cdot) - f^M(\cdot;\theta))\|_{L_2(X)}$$
$$2^{\frac{p}{2}} \cdot 2 \sup_{z \in \Theta} \|y^R(\cdot) - f^M(\cdot;\theta)\|_{L_2(X)}$$
$$+ \sup_{z \in \Theta} \|y^R(\cdot) - f^M(\cdot;\theta)\|_H + C_K \theta_0^{1=2} \cdot n^{-\frac{2m+\overline{p}-m}{2};2}$$

and

$$\sup_{z \in \Theta} \|\hat{\theta}_{z;\theta}(\cdot)\|_H \le (2^{\frac{p}{2}} \cdot 2 + 1) \sup_{z \in \Theta} \|(y^R(\cdot) - f^M(\cdot;\theta))\|_H$$
$$+ 2^{\frac{p}{2}} \cdot 2 \sup_{z \in \Theta} \|y^R(\cdot) - f^M(\cdot;\theta)\|_{L_2(X)} + 2^{\frac{p}{2}} \cdot C_K \theta_0^{1=2}$$

by choosing $\theta = n^{-2m=(2m+p)}$ and $\theta_z = \theta^{1=2}$, where $C_K$ is a constant depending on the kernel $K(\cdot;\cdot)$.

**Lemma S2.** Under assumptions A1 to A6,
(i) it holds that

$$\sup_{z \in \Theta} \left| \frac{1}{n} \sum_{i=1}^{n} (y^R(x_i) - f^M(x_i;\theta) - \hat{\theta}_{z;\theta}(x_i))^2 \right.$$
$$\left. - \int_{x \in X} (y^R(x) - f^M(x;\theta) - \hat{\theta}_{\lambda}(x))^2 dx \right| = o_p(n^{-1=2});$$

and

$$\sup_{z \in \Theta} \left| \frac{1}{n} \sum_{i=1}^{n} (y^R(x_i) - f^M(x_i;\theta) - \hat{\theta}_{z;\theta}(x_i))_i \right| = o_p(n^{-1=2});$$

19

(ii) for any $j = 1, \ldots, q$, one has

$$\frac{1}{n} \sum_{i=1}^{n} (y^R(x_i) - f^M(x_i; \hat{z}) - \hat{z}(x_i)) \frac{\partial f^M(x_i; \hat{z})}{\partial_j}$$

$$= \int_{x \in X} (y^R(x) - f^M(x; \hat{z}) - \hat{z}(x)) \frac{\partial f^M(x; \hat{z})}{\partial_j} dx + o_p(n^{-1/2}):$$

Proof. Denote

$$W_2^m(H; B) := \left\{ f() = \sum_{j=1}^{\infty} f_j() \in L_2(X) : \sum_{j=1}^{\infty} j^{2m=p} f_j^2 \le B^2 \right\};$$

and

$$s_i^2(;) := (y^R(x_i) - f^M(x_i;) - (x_i))^2;$$

$$u_i(;) := (y^R(x_i) - f^M(x_i;) - (x_i))_i;$$

$$r_i^2(;) := (y^R(x_i) - f^M(x_i;) - (x_i)) \frac{\partial f^M(x;)}{\partial_j}$$

for $(;) \in W_2(X; B^m)$ and some $B > 0$ that will be specied later. Dene the empirical processes

$$\hat{s}(;) := \frac{1}{n} \sum_{i=1}^{n} f s_i^2(;)_i - E_{x_i}[s_i^2(;)]g;$$

$$\hat{u}(;) := \frac{1}{n} \sum_{i=1}^{n} f u_i(;) - E_{x_i;i}^i[u_i(;)]g; \quad i=1$$

$$\hat{r}(;) := \frac{1}{n} \sum_{i=1}^{n} f r_i(;) - E_{x_i}[r_i(;)]g;$$

where

$$E_{x_i}[s_i^2(;)] = \int_{x \in X} (y^R(x) - f^M(x;) - (x))^2 dx;$$

$$E_{x_i;i}[u_i(;)] = \int_{x \in X} (y^R(x) - f^M(x;) - (x))_i dx = 0;$$

$$E_{x_i}[r_i(;)] = \int_{x \in X} (y^R(x) - f^M(x;) - (x)) \frac{\partial f^M(x;)}{\partial_j} dx:$$

20

By Assumptions A3 and A4, the function classes $f@f^M(;)=@_j : 2 g$ and $F = fy^R()$ $f^M(;)$; $2 g$ are Donsker. Note that, by denition, $W^m(X; B_2)$ is also Donsker. Since both $W_2(X; B)$ and $F$ are uniformly bounded, the function classes

$$f(y^R() \quad f^M(;) \quad ())^2 : 2 ; 2 W^m(X; B_2)g; \qquad \text{and}$$

$$(y^R() \quad f^M(;) \quad ())^{@f} \frac{;^M}{@()_j} : 2 ; 2 W^m(X; B)_2$$

are also Donsker classes. Furthermore, letting $f_;(; x) = (y^R(x) \quad f^M(x;) \quad (x))$, observe that for any $(_1;_1)$ and $(_2;_2)$, the distance

$$E_0 (f^1_{;_1} \quad f^2_{;_2})^2 {}^{1=2}$$

$$= {}_0 k f^M(;_1) \quad _1() \quad f^M(;_2) + _2()k_{L_2(X)}$$

$$_0 \quad kf^M(;_1) \quad f^M(;_2)k_{L_2(X)} + k_1() \quad _2()k_{L_2(X)}$$

can be bounded by the $L_2(X)$-distance of functions in $ff^M(;) : 2 g$ and $() 2 W^m(X; B)$. In addition, by Assumption A4 $ff^M(;) : 2 g$ and $W^m_2(X; B)$ are Donsker classes, it follows that the function class

$$ff_; 2 C(R \quad X) : 2 ; 2 W_2(X; B)g^m$$

is also Donsker, since its metric entropy can be upper bounded by those of $ff^M(;) : 2 g$ and $W^m_2(X; B)$. By Theorem 2.4 in [4], for any $t_1 > 0$ and any $B > 0$, there exists $t_2; t2^{\zeta}; t_{00}^{0} > 0$ such that

(S25) $\quad \limsup_{n ! 1} P \left( \sup_{kk_H B; 2; ky^R() \ f^M(;) \ ()k_{L_2(X)} t_2} j \hat{s}^2(;)j > t_1 \right) < t_1;$

(S26) $\quad \limsup_{n ! 1} P \left( \sup_{kk_H B; 2; ky^R() \ f^M(;) \ ()k_{L_2(X)} t_2} j_0(; )j > t_1 \right) < t_1;$

(S27) $\quad \limsup_{n ! 1} P \left( \sup_{kk_H B; 2; ky^R() \ f^M(;) \ ()k_{L_2(X)} t^{00}} j(;)j > t_1 \right) < t_1:$

Note that by Corollary S1, $\sup_2 k_{z;} \hat{K}_H^2$ is asymptotically tight, and therefore for any $" > 0$, there exists $B_0 > 0$ and some integer $N 2 N_+$, both depending on , such that $P(\sup_2 k_{z;} \hat{K}_H > B_0) "=3$ for all $n > N$. Now take $B = B_0$, $t_1 = =3$. Then we can choose $t_2$ to be a value that satises (S25), $t_{0 2}$ satisfying (S26), and $t_{00}$ satisfying (S27). By

21

Corollary S1 and Assumption A5, $\sup k_{z;}(\hat{\cdot})\ (y^R(\cdot)\ f^M(\cdot;\cdot))k_{L_2(X)} = O_P(n^{\ 2m=(2m+p)})$, and hence there exists $t_3 > 0$, depending on and $n$, such that for all $n > N$, it holds that

$$P\left\{\sup_2 ky^R(\cdot)\ f^M(\cdot;\cdot)\ \hat{z;}(\cdot)k_{L_2(X)}\ t_3\right\} < "=3:$$

Without loss of generality, we may require $t_3\ \min\{t_2; t2; \frac{t^0}{0\hat{\underline{z}}}g$ by taking suciently large $n$. Then for suciently large $n$, we obtain

$$P\left\{\sup_{2\ 0} j\hat{s}^2(\cdot;_{z;})\hat{j} > "\right\}$$

$$P @\ \sup_{k_{z;}k_H\ \hat{B_0};\ 2;\ ky^R(\cdot)\ f^M(\cdot;\cdot)\ _{z;}(\cdot)k_{L_2(X)}t_2}\ \hat{}\ j\hat{s}^2(\cdot;_{z;})\hat{j} > t_1\ A\ _{\sup}^{\quad 1}$$

$$+\ P\left\{\sup_2 ky^R(\cdot)\ f^M(\cdot;\cdot)\ \hat{z;}(\cdot)k_{L_2(X)} > t_2\right\}$$
$$+\ P\left\{\sup_{2\ \hat{}} k_{z;}k_H > B_0\right\}$$

$$<\ "=3 +\ "=3 +\ "=3 =\ ";$$

and similarly,

$$P\left\{\sup_2 j(\cdot;_{z;})\hat{j} > "\right\} <\ \quad\text{and}\quad P\left\{\sup_2 j\dot{(\cdot;_{z;})}\hat{j} > "\right\} < :$$

Therefore,

$$p\frac{1}{n}\sup_2 j\hat{s}(;_{z;})\hat{j}$$

$$=\ \sup_2\ \frac{1}{n}\sum_{j=1}^{X^n} (y^R(x_i)\ f^M(x_i;\cdot)\ \hat{z;}(x_i))^2$$
$$Z_{x2X} (y^R(x)\ f^M(x;\cdot)\ _{\hat{}}(x))^2 dx = o_p(n^{\ 1=2});$$

and

$$p\frac{1}{n}\sup_2 j\dot{(\cdot;_{z;})}\hat{j} =\ \sup_2\ \frac{1}{n}\sum_{i=1}^{X^n}(y^R(x_i)\ f^M(x_i;\cdot)\ \hat{z;}(x_i))_i = o_p(n^{\ 1=2});$$

completing the proof of (i). The proof of (ii) can be completed by observing that

$$\frac{1}{n}\sum_{i=1}^{n}(y^R(x_i) - f^M(x_i; \hat{z}) - \hat{z}(x_i))\frac{@f^M(x_i; \hat{z})}{@_j}$$

$$\int_{x\in X}(y^R(x) - f^M(x; \hat{z}) - \hat{z}(x))\frac{@f^M(x; \hat{z})}{@_j}dx$$

$$= p\frac{1}{n}j(\hat{z}; \wedge \wedge)j_{z,z} \cdot p\frac{1}{n_2}\sup_z j(; z;)j = o_p(\hat{n}^{1=2}):$$ ∎

Proof for Theorem 3.3. Without loss of generality, it suces to prove the case when $z = 1=2$. For the general case when $z = O(1=2)$, the proof follows similarly. We rst show $z \to^p \hat{L_2}$. By the denition of $z$, $L_2$, and the theory of M-estimators (see, Theorem 5.7 in [10]), it suces to show that $1=2(\ell_z(; z;) - \wedge - 2)\to^p ky^R() - f^M(;)k^2_{2(X)}$ uniformly for each $2$. Note that

$$\ell_z(\hat{z}(); )$$

$$= \frac{1}{n}\sum_{i=1}^{n}(y^R(x_i) - f^M(x_i;) - \hat{z}(x_i))^2 + \frac{1}{n}\sum_{i=1}^{n}ilon^2_i$$

$$+ \frac{2}{n}\sum_{i=1}^{n}(y^R(x_i) - f^M(x_i;) - \hat{z}(x_i))_i + kz;k^2 \to^p k_{\mathcal{H}}; k_{L_2(X)}^{\wedge}_{i=1}{}^2$$

$$:= A_n + B_n + C_n + D_n + E_n:$$

For $A_n$, by Lemma S2 (i) and Corollary (S1), one has

(S28)  $$\sup_2 \frac{1}{n}\sum_{i=1}^{n}(y^R(x_i) - f^M(x_i;) - \hat{z}(x_i))^2 = o_p(n^{1=2})$$

Since $E[B_n] = {}^2_0$ and $V[B_n] = O(n^{-1})$, Chebyshev's inequality implies $(1=n)\sum_{i=1}^{n}ilon^2 \to^p {}^2_0$ $+_0 O_p(n^{1=2})$ for $B_n$. For $C_n$, Lemma S2 (i) guarantees that

$$\sup_2 \frac{2}{n}\sum_{i=1}^{n}(y^R(x_i) - f^M(x_i;) - \hat{z}(x_i))_i = o_p(n^{1=2})$$

Since $= O(n^{2m=(2m+p)})$, by the asymptotic tightness of $\sup k_z; k_{\mathcal{H}}^{\wedge}$ (Corollary S1), one 23

has $\sup_{2} \|k_{z;} k_{H}\|\hat{} = o_{p}^{2}(n^{-1=2})$. By putting the above all pieces together, we obtain

(S29)    $$\sup_{2} \|^{-1=2}(\ell_{z}(z;(\cdot); ) \qquad )_{0}^{2} \quad k_{A;} \|k_{L_{2}(X)}^{2} = O_{p}((n)^{-1=2}):$$

For any , by the Cauchy-Schwarz inequality, one has

$$\|k_{z;}\hat{}k_{L_{2}(X)}^{2} \qquad ky^{R}() \quad f^{M}(\cdot;)\|k_{L_{2}(X)}^{2}$$

$$k(z;\hat{(\cdot)} \qquad (y^{R}() \quad f^{M}(\cdot;))\|k_{L_{2}(X)} \|k_{z;}(\hat{} + y^{R}() \qquad f^{M}(\cdot;)\|k_{L_{2}(X)}$$

Recall that

$$\sup_{2} \|k(\hat{z};(\cdot) \qquad (y^{R}() \quad f^{M}(\cdot;))\|k_{L_{2}(X)} = O_{p}(n^{-m=(2m+d)})$$

by Corollary S1 and Assumption A4. Using Assumptions A4 and the asymptotic tightness of $\sup \|k_{z;}\hat{}k_{H}$ (Corollary S1), one has

$$\|k_{z;}^{\wedge}(\cdot) + y^{R}() \qquad f^{M}(\cdot;)\|k_{L_{2}(X)}$$

$$\|k_{z;}(\hat{)}\|k_{L_{2}(X)} + \sup_{2} \|ky^{R}() \qquad f^{M}(\cdot;)\|k_{L_{2}(X)}$$

$$C\|k_{z;}(\hat{)}k_{H} + \sup_{2} \|ky^{R}() \qquad f^{M}(\cdot;)\|k_{H} = O_{p}(1):$$

Thus

$$\sup_{2} \|k_{z;}\hat{}k^{2}_{2(X)} \qquad \|ky^{R}() \quad f^{M}(\cdot;)\|k^{2}_{2(X)} = O_{p}(n^{-m=(2m+d)});$$

and hence,

$$\sup_{2} \|^{-1=2}(\ell_{z}(\cdot; z;) \quad \hat{} \qquad )_{0}^{2} \quad \|ky^{R}() \qquad f^{M}(\cdot;)\|k^{2}_{2(X)} = o_{p}(1);$$

from which we conclude $\hat{z}\to^{p} L_{2}$.

Next we derive the convergence rate of $\hat{z}$. Apply the Frechet derivative on $\ell_{z}$ with regard to $(\cdot)$ and the partial derivative on $\ell_{z}$ with regard to $_{j}$, $j = 1; :::; q$. For any $g(\cdot) \in H$, $\hat{z}$

and $\hat{\gamma}_z$ satisfy

$$0 = \frac{2}{n}\sum_{i=1}^{X^n}(y_i^F - f(x_i;\hat{z}) - \hat{\gamma}_z(x_i))g(x_i) + 2\langle h_z(\cdot), g(\cdot)\rangle_H$$
$$+ 2\rho\langle \underline{h}_z(\cdot), g(\cdot)\rangle_{L_2(X)};$$
(S30)

$$0 = \frac{2}{n}\sum_{i=1}^{X^n}(y_i^F - f(x_i;\hat{z}) - \hat{\gamma}_z(x_i))\frac{@f^M(x_i;\hat{z})}{@_j}:$$
(S31)

Choosing $g(\cdot) = \frac{@f^M(\cdot;\hat{z})}{@_j}$ and plugging (S31) into (S30), one has

(S32)
$$\rho\left[\langle\hat{\gamma}_z(\cdot), \frac{@f^M(\cdot;z)}{@_j}\rangle_z\right]_H^* + \left[\langle\hat{\gamma}_z(\cdot), \frac{@f^M(\cdot;z)}{@_j}\rangle_z\right]_{L_2(X)}^* = 0:$$

Substituting (S31) into (S32) and by Lemma S2 (ii), we have

$$0 = \frac{1}{n}\sum_{i=1}^{X^n}(y_i^F - f(x_i;z) - \hat{\gamma}_z(x_i))\frac{@f^M(x_i;z)}{@_j}$$

$$= \left[\int (y^R(x) - f^M(x;\hat{z}))\frac{@f^M(x;z)}{@_j}dx + \langle\hat{\gamma}_z(\cdot), \frac{@f^M(\cdot;z)}{@_j}\rangle_z\right]_{L_2(X)}^*$$

$$\cdot \frac{1}{n}\sum_{i=1}^{X^n}\frac{@f^M(x_i;z)}{@_j} + o_p(n^{-1=2})$$

$$= \left[\int \frac{@(y^R(x) - f^M(x;z))^2}{@_j}dx - \rho\langle\hat{\gamma}_z(\cdot), \frac{@f^M(\cdot;z)}{@_j}\rangle_H\right]^*$$

$$\cdot \frac{1}{n}\sum_{i=1}^{X^n}\frac{@f^M(x_i;z)}{@_j} + o_p(n^{-1=2}):$$

Applying Taylor expansion to the rst term on the right-hand side at $_{L_2}$, for any $j = 1;\ldots;q$, we obtain

$$\left(\int \frac{@^2(y^R(x) - f^M(x;\hat{z}))^2}{@_j@}dx\right)^T(\hat{z} - _{L_2})$$

$$= \int \frac{@^2(y^R(x) - f^M(x;_{L_2}))^2}{@_j@}dx + o_p(1)^T(\hat{z} - _{L_2})$$

(S33)
$$= \left[\rho\langle\hat{\gamma}_z(\cdot), \frac{@f^M(\cdot;z)}{@_j}\rangle_H\right]^* + \frac{1}{n}\sum_{i=1}^{X^F}\frac{@f^M(x_i;z)}{@_j} + o_p(n^{-1=2});$$

where $\tilde{z}$ lies within the q dimensional rectangle between $_{L_2}$ and $_z$. Observe that Corollary S1 and assumption A3 imply

$$z; \hat{\wedge} \frac{@f^M(;z)^{\wedge +}}{@_j}\Big|_H k_z k_H \frac{@\hat{f}^M(;L_2)}{@_j} + o_p(1) \Big|_H = O_p(1):$$

Now we consider the second term. Dene the empirical process

$$G_n() = \frac{1}{\cancel{p}_n} X_i^n \frac{@f^M(x_i;)}{@_j} {}_i \frac{@f^M(x_i;L_2)}{@_j} ;{}_{i=1}$$

and denote

$$f(;x) = \frac{@f^M(x,)}{@_j} \frac{@f^M(x;L_2)}{@_j} :$$

Since

$$E_{ilon;x}[f_1(;x) \quad f_2(;x)]^2$$

$$= E_{ilon;x} \left[ {}^2 \frac{@f^M(x;)_1}{@_j} \frac{@f^M(x;)_2}{@_j} \right]^{\#}$$

$$= \cancel{0}^2 \left[ \frac{@f^M(x;)_1}{@_j} \frac{@f^M(x;)_2}{@_j} \right]_{L_2(X)} ;$$

therefore the function class $ff(;x) 2 C(R \ X) : 2 g$ is Donsker by Assumption A3, and hence, $G_n()$ converges weakly to a tight Gaussian stochastic process, denoted by $G()$. W.l.o.g., we may take $G()$ a version that has uniformly continuous sample paths (see Chapter 6 in [9]). Since $G_n(L_2) = 0$ for all n, it follows that $G(L_2) = 0$. By the consistency of $_z$ and the continuous mapping theorem [10], $G_n(z) = G(L_2) + o_p(1) = \hat{} o_p(1)$. Therefore,

$$\frac{1}{n} X_i^n \frac{@f^M(x_i;z)^{\wedge}}{@_j} = \frac{1}{\cancel{p}_n} G_n(z) \cancel{4} \frac{1}{n} X_{i=1}^n \frac{@f^M(x_i;L_2)}{@_j} = O_p(n^{1=2}){}_{i=1}$$

To sum up,

$$\int^Z \frac{@^2(y^R(x) \quad f^M(x;L_2))^2}{@@^T} dx + o_p(1)(_z \quad ^{\wedge} {}_{L_2})$$

$$= O_p(n^{m=(2m+p)}) + O_p(n^{1=2}) + o_p(n^{1=2}) = O_p(n^{m=(2m+p)});$$

completing the proof. ∎

S6. Proof and additional results for Section 4. The identities in the Lemma S1 are used repeatedly in the proof of the Theorem 4.1 and Lemma 4.2.

26

Lemma S1. Denote $\sigma^2 R_{z_d}$ the covariance matrix of $(z_d(x_1), \ldots, z_d(x_n))^T$, where the $(i, j)$ entry being $\sigma^2 K_{z_d}(x_i, x_j)$ defined in (20). Denote $r_{z_d}(x) = (K_{z_d}(x, x_1), \ldots, K_{z_d}(x, x_n))^T$ for any $x \in X$. One has the following identities

(S34)
$$R_{z_d}^{-1} = R^{-1} + \frac{z}{n} I_n;$$

(S35)
$$r_{z_d}^T(x) = \frac{n}{z} r^T(x) R^{-1} = r^T(x) R^{-1} R_{z_d};$$

for any $x \in X$.

Proof. By the definitions of $R_{z_d}$ and the Woodbury Identity, one has

$$R_{z_d} = R - R \tilde{R}^{-1} R = R \left( I_n - \left( I_n + \frac{z}{n} R \right)^{-1} \right)$$

$$= R \left( \frac{z}{n} R + I_n \right)^{-1} = \left( R^{-1} + \frac{z}{n} I_n \right)^{-1};$$

from which (S34) follows.

Equation (S35) can be shown similarly by noting $r_{z_d}^T(x) = r^T(x) - r^T(x) \tilde{R}^{-1} R$ and the Woodbury Identity. ∎

Proof of Theorem 4.1. The predictive mean is as follows

$$\hat{z}(x) = E[y^F(x) \mid y^F; \hat{\theta}^2; \hat{\theta}]$$

$$= f^M(x; \hat{\beta}) + r_{z_d}(x)^T (R_{z_d} + n I_n)^{-1}(y^F - f^M)$$

$$= f^M(x; \hat{\beta}) + r(x)^T R^{-1} R_{z_d} (R_{z_d} + n I_n)^{-1}(y^F - f^M)$$

$$= f^M(x; \hat{\beta}) + r(x)^T R^{-1} \left( I_n + n \left( R^{-1} + \frac{z}{n} I_n \right) \right)^{-1}(y^F - f^M)$$

$$= f^M(x; \hat{\beta}) + \frac{r(x)^T}{(1 + z)} \left( R^{-1} + \frac{n}{1 + z} I_n \right)^{-1}(y^F - f^M);$$

where the last two equalities follow from (S35) and (S34), respectively.

The predictive variance can be obtained using (S35) and (S34) as follows

$$K_z(x, x) = K_{z_d}(x, x) - r_{z_d}^T(x)(R_{z_d} + n I_n)^{-1} r_{z_d}(x)$$

$$= K(x, x) - r^T(x) \tilde{R}^{-1} r(x)$$

$$- (1 + z)^{-1} r(x)^T \left( R^{-1} + \frac{n}{1 + z} I_n \right)^{-1} \frac{R^{-1}}{z} r(x) \quad \blacksquare$$

27

from which the result follows.

Proof of Lemma 4.2. When $\sigma_0^2 = 0$, the predictive mean is as follows

$$E[y^F(x) \mid y^F; \beta; \sigma^2; \sigma_0^2] = f^M(x;\beta) + r_{z_d}(x)^T R^{-1}(y^F_{z_d} - f^M)$$

$$= f^M(x;\beta) + r(x)^T R^{-1} R_{z_d} R_{z_d}^{-1}(y^F - f^M)$$

$$= f^M(x;\beta) + r(x)^T R^{-1}(y^F - f^M):$$

The predictive variance can be obtained similarly. ∎

## REFERENCES

[1] D. E. Edmunds and H. Triebel, Function spaces, entropy numbers, dierential operators, vol. 120, Cambridge University Press, 2008.

[2] S. Ghosal and A. Van der Vaart, Fundamentals of nonparametric Bayesian inference, vol. 44, Cambridge University Press, 2017.

[3] M. R. Kosorok, Introduction to empirical processes and semiparametric inference., Springer, 2008.

[4] E. Mammen and S. Van de Geer, Penalized quasi-likelihood estimation in partial linear models, The Annals of Statistics, (1997), pp. 1014{1035.

[5] I. Pinelis, Optimum bounds for the distributions of martingales in banach spaces, The Annals of Probability, (1994), pp. 1679{1706.

[6] C. E. Rasmussen, Gaussian processes for machine learning, MIT Press, 2006.

[7] M. Rudelson, R. Vershynin, et al., Hanson-wright inequality and sub-gaussian concentration, Electronic Communications in Probability, 18 (2013).

[8] R. Tuo and C. J. Wu, Ecient calibration for imperfect computer models, Annals of Statistics, 43 (2015), pp. 2331{2352.

[9] S. A. Van de Geer, Empirical Processes in M-estimation, vol. 6, Cambridge university press, 2000.

[10] A. W. Van der Vaart, Asymptotic statistics, vol. 3, Cambridge university press, 2000.

[11] G. Wahba, Spline models for observational data, vol. 59, SIAM, 1990.

[12] Y. Yang, A. Bhattacharya, and D. Pati, Frequentist coverage and sup-norm convergence rate in Gaussian process regression, arXiv preprint arXiv:1708.04753, (2017).

[13] Y. Yang, Z. Shang, and G. Cheng, Non-asymptotic theory for nonparametric testing, arXiv preprint arXiv:1702.01330, (2017).