

Modularized Calibration with laGP

Grant Hutchings

February 2022

1 Introduction

Consider the study of a physical process for which experiments are parameterized by a set of controllable inputs $\mathbf{x} = (x_1, \dots, x_{p_x})$. Suppose we have n noisy observations of this process at a number of input parameter settings $\mathbf{X}^F = [\mathbf{x}_1^F, \dots, \mathbf{x}_n^F]^T$. We denote a field observation at input setting \mathbf{x}^F as $y^F(\mathbf{x}^F)$ and the set of field observations $\mathbf{Y}^F(\mathbf{X}^F) = [y_1^F(\mathbf{x}_1^F), \dots, y_n^F(\mathbf{x}_n^F)]^T$. Field observations are often expensive and therefore limited, so scientists frequently turn to computer models to study the process. Suppose we have a computer model for the process which depends not only on \mathbf{x} , but also on a set of calibration parameters $\mathbf{t} = (t_1, \dots, t_{p_t})$. We will denote the response from our computer model as $y^M(\mathbf{x}, \mathbf{t})$ and the set of m computer model runs as $\mathbf{Y}^M(\mathbf{X}^M) = [y_1^F(\mathbf{x}_1^M), \dots, y_n^F(\mathbf{x}_n^M)]^T$. The goal of calibration is to determine the \mathbf{t}^* which is optimal for the computer model. Optimal can mean different things depending on if the computer model is an unbiased estimate of the field data or not. We will discuss this more in the following section.

2 Unbiased Calibration

Given a certain unknown \mathbf{t}^* , we expect our deterministic computer model to be an unbiased estimate of the field data. In other words, the field data is a noisy version of the simulator output.

$$y^F(\mathbf{x}^F) = y^M(\mathbf{x}^F, \mathbf{t}^*) + \epsilon; \epsilon \sim N(0, \sigma_{y^F}^2) \quad (1)$$

where $\sigma_{y^F}^2$ is the observation error variance. Calibration is the process of determining this \mathbf{t}^* . Computer models are often not sufficiently fast to use for calibration, as we would need to run the model at many different (\mathbf{x}, \mathbf{t}) pairs to determine \mathbf{t}^* . Instead we build a statistical model, or emulator, of the simulator using a Gaussian Process. This laGP emulator can be used for fast calibration.

2.1 Scalar Response

Let us first consider the case where the response from the field data and computer code are scalar as shown in Equation 1. Replacing our computer model with an laGP emulator $\eta(\mathbf{x}, \mathbf{t})$ we have

$$\begin{aligned} y^F(\mathbf{x}^F) &= \eta(\mathbf{x}^F, \mathbf{t}^*) + \epsilon; \epsilon \sim N(0, \sigma_{y^F}^2) \\ &\sim \mathcal{N}(\mu_\eta, \sigma_\eta^2 + \sigma_{y^F}^2) \end{aligned} \quad (2)$$

Before we fit $\eta(\mathbf{x}, \mathbf{t})$ to the simulator data \mathbf{Y}^M , there is some clever pre-computing we can do to facilitate fast calibration later. We follow the procedure of <https://arxiv.org/pdf/2005.00386.pdf>

where inputs (\mathbf{x}, \mathbf{t}) are scaled, or "stretched and compressed", by estimates of the global GP length-scale parameters. This means that future laGP fits using scaled inputs do not need to estimate length-scales, as they will be very near to 1. We will later show how this greatly reduces calibration time. We begin by fitting a full GP to a small subset of the available suite of computer model runs. This subset can be chosen randomly, or using some kind of space-filling algorithm. The size of this subset can be chosen based on available computational resources. Given a set of estimated length-scales $[\hat{\theta}_x, \hat{\theta}_t]$ we compute $[\mathbf{X}_{sc}^M, \mathbf{T}_{sc}^M] = [\mathbf{X}^M, \mathbf{T}^M] / \sqrt{[\hat{\theta}_x, \hat{\theta}_t]}$ and similarly for \mathbf{X}_{sc}^F .

With our stretched and compressed input space, we are ready to calibrate. Unbiased calibration is straightforward in that the optimal \mathbf{t}^* is one that maximizes the likelihood associated with the normal distribution in with Equation 2. A calibration algorithm is given below.

1. For a given $\mathbf{t}_{sc} = \mathbf{t} / \sqrt{\hat{\theta}_t}$, fit an laGP model to the simulator data with inputs $[\mathbf{X}_{sc}^M, \mathbf{T}_{sc}^M]$ and make predictions at $[\mathbf{X}_{sc}^F, \mathbf{t}_{sc}]$. As our inputs are scaled, length-scales are fixed at 1, and a small nugget is added for numerical stability.

$$\begin{aligned} \eta(\mathbf{X}^F, \mathbf{t}) &= \text{laGP}(X = [\mathbf{X}_{sc}^M, \mathbf{T}_{sc}^M], \\ Z &= \mathbf{Y}^M, \\ XX &= [\mathbf{X}_{sc}^F, \mathbf{t}_{sc}], \\ d &= 1, \\ g &= 1e - 7) \end{aligned} \quad (3)$$

2. Compute the residuals for the predictions: $\mathbf{r} = \mathbf{Y}^F(\mathbf{X}_{sc}^F) - \boldsymbol{\mu}_\eta(\mathbf{X}_{sc}^F, \mathbf{t}_{sc}) \in \mathbb{R}^n$.
3. Compute the log-likelihood of these residuals subject to Equation 2 with $\hat{\sigma}_{y^F}^2 = \frac{1}{n} \mathbf{r}^T \mathbf{r}$.

$$ll \propto -\frac{1}{2} \left(\log(|\boldsymbol{\Sigma}_\eta + \hat{\sigma}_{y^F}^2 \mathbf{I}_n|) + \mathbf{r}^T (\boldsymbol{\Sigma}_\eta + \hat{\sigma}_{y^F}^2 \mathbf{I}_n)^{-1} \mathbf{r} \right) \quad (4)$$

where $\boldsymbol{\Sigma}_\eta \in \mathbb{R}^{n \times n}$ is a diagonal matrix of prediction variances returned by laGP.

Steps 1-3 can be run over a grid of possible \mathbf{t} values, or can be implemented as an objective function in a more sophisticated optimizer in order to find the \mathbf{t} which maximizes the likelihood.

2.2 Functional Response

We follow a similar approach to Higdon '08 to extend this procedure observations and computer simulations with functional response. Suppose the simulation response is of dimension p and observed data is of dimension q . Throughout the following section we will index experiments with the index $i = 1, \dots, n$, simulations with the index $j = 1, \dots, m$, and orthogonal bases, or principal components with $k = 1, \dots, n_{pc}$. We now have

$$\mathbf{y}^F(\mathbf{x}^F) = \mathbf{y}^M(\mathbf{x}^F) + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}_n, \sigma_{y^F}^2 \mathbf{I}_q) \quad (5)$$

where $\mathbf{y} \in \mathbb{R}^q$. In the above equation, and throughout the documents, $\mathbf{y}^M(\mathbf{x}^F)$ implies not only that we are evaluating the simulator at the observed input \mathbf{x}^F , but also on the observed support.

If we have m computer model runs, stack the response vectors in columns to produce the matrix $\mathbf{Y}^M = [\mathbf{y}_1^M | \dots | \mathbf{y}_m^M] \in \mathbb{R}^{p \times m}$ and perform an SVD decomposition $\mathbf{Y}^M = \mathbf{U} \mathbf{D} \mathbf{W}^T = \mathbf{B}^M \mathbf{W}^T$ where

$\mathbf{B}^M = \mathbf{U}\mathbf{D}$. In practice we approximate the response using only the first n_{pc} basis vectors and associated weights.

$$\mathbf{Y}^M \approx \mathbf{B}^M[:, 1 : n_{pc}] \mathbf{W}^T [1 : n_{pc}, :] \quad n_{pc} \leq \min(p, q) \quad (6)$$

To cast field observations into this basis, we take $\mathbf{B}^F = \mathbf{B}^M$ if the field observations and simulations are on the same support. If not, we compute basis vectors \mathbf{b}^F by interpolating the columns of \mathbf{B}^M at the observed support. We can then compute the matrix of weights, \mathbf{Z}^T as

$$\begin{aligned} \mathbf{Y}^F &= \mathbf{B}^F \mathbf{Z}^T \\ \implies (\mathbf{B}^F)^T \boldsymbol{\Sigma}_y \mathbf{Y}^F &= (\mathbf{B}^F)^T \boldsymbol{\Sigma}_y \mathbf{B}^F \mathbf{Z}^T \\ \implies \mathbf{Z}^T &= ((\mathbf{B}^F)^T \boldsymbol{\Sigma}_y \mathbf{Y}^F)^{-1} (\mathbf{B}^F)^T \boldsymbol{\Sigma}_y \mathbf{B}^F \end{aligned}$$

where $\boldsymbol{\Sigma}_y = \hat{\gamma}_y \mathbf{I}_n$ and $\hat{\gamma}_y$ is the estimated observation precision (Go back to Hidgon to make sure I understand why the precision is included).

With our data properly decomposed, our emulator consists of n_{pc} independent laGP models on the basis weight vectors, the rows of \mathbf{W}^T which can be combined through the basis vectors to reproduce the functional response. We will denote the emulator as $\eta(\mathbf{x}, \mathbf{t}) = \mathbf{B}\mathbf{w}(\mathbf{x}, \mathbf{t})$ where $\mathbf{w}(\mathbf{x}, \mathbf{t}) \in \mathbb{R}^{n_{pc}}$ contains a basis weight for each of the k components.

Calibration is quite similar to the scalar case where we use predictions from the emulator at the observed inputs. Given the unbiased simulator assumption, we have

$$\mathbf{y}_i^F(\mathbf{x}^F) = \mathbf{B}^F \mathbf{w}_i(\mathbf{x}^F, \mathbf{t}^*) + \boldsymbol{\varepsilon}_i$$

where

$$\begin{aligned} \mathbf{w}_i &\sim \mathcal{N}(\boldsymbol{\mu}_{wi}, \boldsymbol{\Sigma}_{wi}) \\ \boldsymbol{\varepsilon}_i &\sim \mathcal{N}(\mathbf{0}, \sigma_{y^F}^2 \mathbf{I}_q) \\ \implies \mathbf{B}^F \mathbf{z}_i &\sim \mathbf{B}^F \mathcal{N}(\boldsymbol{\mu}_{wi}, \boldsymbol{\Sigma}_{wi}) + \mathcal{N}(\mathbf{0}, \sigma_{y^F}^2 \mathbf{I}) \\ \implies \mathbf{B}^F \mathbf{z}_i &\sim \mathcal{N}(\mathbf{B}^F \boldsymbol{\mu}_{wi}, \mathbf{B}^F \boldsymbol{\Sigma}_{wi} (\mathbf{B}^F)^T + \sigma_{y^F}^2 \mathbf{I}) \\ \implies \mathbf{z}_i &\sim \mathcal{N}(\boldsymbol{\mu}_{wi}, \boldsymbol{\Sigma}_{wi} + \sigma_{y^F}^2 ((\mathbf{B}^F)^T \mathbf{B}^F)^{-1}) \end{aligned} \quad (7)$$

This is the likelihood we will use for calibration. Once again we begin by fitting n_{pc} independent full GP models to a subset of the computer model runs to estimate length-scale parameters for each component of the model. With a set of estimated length-scales for each of the n_{pc} models, inputs can be stretched and compressed (sc) independently for each model. Calibration proceeds similarly to the scalar response case.

1. For a given \mathbf{t} , fit n_{pc} independent laGP models to the computer model data with scaled inputs $[\mathbf{X}_{sc}^M, \mathbf{T}_{sc}^M]$ and predict at $[\mathbf{X}_{sc}^F, \mathbf{t}_{sc}]$. As our inputs are stretched and compressed, length-scales are fixed at 1, and a small nugget is used. laGP returns the mean and variance of the sampling distribution

$$\begin{aligned} \boldsymbol{\mu}_{wk}, \boldsymbol{\sigma}_{wk}^2 &= \text{laGP}(X = [\mathbf{X}_{sc}^M, \mathbf{T}_{sc}^M]_k, \\ &\quad Z = \mathbf{W}^T[k, :], \\ &\quad XX = [\mathbf{X}_{sc}^F, \mathbf{t}_{sc}]_k, \\ &\quad d = 1, \\ &\quad g = 1e - 7); \quad k = 1, \dots, n_{pc} \end{aligned} \quad (8)$$

Subscripts on input matrices are to indicate that inputs are scaled differently for each k . Store prediction means and variances in matrices $\mathbf{M}_w = [\boldsymbol{\mu}_{w1} | \dots | \boldsymbol{\mu}_{wn_{pc}}]^T \in \mathbb{R}^{n_{pc} \times n}$, $\mathbf{S}_w = [\boldsymbol{\sigma}_{w1}^2 | \dots | \boldsymbol{\sigma}_{wn_{pc}}^2]^T \in \mathbb{R}^{n_{pc} \times n}$

2. Compute the residual matrices:

$$\begin{aligned}\mathbf{R}_w &= \mathbf{Z}^T - \mathbf{M}_w \\ \mathbf{R}_y &= \mathbf{Y}^F - \mathbf{B}^F \mathbf{M}_w\end{aligned}$$

3. Estimate error variance: Store \mathbf{R}_y in a vector \mathbf{r}_y and estimate

$$\hat{\sigma}_y^2 = \frac{1}{n * p} \mathbf{r}_y^T \mathbf{r}_y$$

4. Compute the log likelihood for each experiment from the normal distribution in Equation 7.

$$ll \propto \sum_{i=1}^n -\frac{1}{2} (\log(|\boldsymbol{\Sigma}_i|) + \mathbf{R}_w[:, i] (\boldsymbol{\Sigma}_i)^{-1} \mathbf{R}_w[:, i])) \quad (9)$$

where $\boldsymbol{\Sigma}_i = \hat{\sigma}_{y^F}^2 (\mathbf{B}^F)^T \mathbf{B}^F + \boldsymbol{\Sigma}_{wi}$ and $\boldsymbol{\Sigma}_{wi} = \text{diag}(\mathbf{S}_w[:, i])$.

2.3 Open Questions

2.3.1 Prediction Variance

Equation 7 describes the sampling distribution for the observed weights.

$$\mathbf{z}_i \sim N(\boldsymbol{\mu}_{wi}, \hat{\sigma}_y^2 ((\mathbf{B}^F)^T \mathbf{B})^{-1} + \boldsymbol{\Sigma}_{wi})$$

But, we could also sample from

$$\mathbf{Y}_i \sim N(\mathbf{B}^F \boldsymbol{\mu}_{wi}, \mathbf{B}^F \boldsymbol{\Sigma}_{wi} (\mathbf{B}^F)^T + \hat{\sigma}_y^2 \mathbf{I})$$

Which for some reason increases our predictive variability. I believe I understand now why the predictive variances are different, but I'm not sure what the solution is. Consider the following

$$Bz \sim BN(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) + \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (1)$$

$$z \sim \mathcal{N}(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) + (B^T B)^{-1} B^T \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (2)$$

$$Bz \sim BN(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) + B(B^T B)^{-1} B^T \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (3)$$

(1) and (3) are equivalent iff $B(B^T B)^{-1} B^T = \mathbf{I}$ which we might expect to be true as B came from SVD and is the product of an orthogonal matrix and a diagonal matrix. It turns out that $B(B^T B)^{-1} B^T = \mathbf{I}$ if we use all the columns of B , but not if we use only the first n_{PC} columns. For this reason we are getting different predictive variances.

2.3.2 Error Variance

For calibration we could compute σ_{yi}^2 for each experiment instead. We could use these individually for the likelihood, but average them at the end for the final error variance estimate. My thought is that this might improve calibration by reducing the chance that a good error estimate for one experiment is balanced out by a bad one for another experiment. However we are assuming that σ^2 is constant for all experiments. I tried this on the ball drop example (for only 1 random data set) and I got a worse estimate of \mathbf{t}^* so maybe its not a good idea.

3 Biased Calibration

3.1 Functional Response

We now consider the situation where the simulator is known to be a biased estimate of the observed data. Our statistical model is similar to Equation 7 although we must add another term to account for this bias.

$$\mathbf{y}^F(\mathbf{x}^F) = \eta(\mathbf{x}^F, \mathbf{t}^*) + \delta(\mathbf{x}^F) + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_{y^F}^2 \mathbf{I}_q) \quad (10)$$

where $\delta(\mathbf{x})$ accounts for the bias between the emulator and the observed data. We decompose our computer simulations as before and the full emulator is of the form $\eta(\mathbf{x}, \mathbf{t}) = \mathbf{B}\mathbf{w}(\mathbf{x}, \mathbf{t})$. We will similarly model the bias using a basis representation computed via SVD $\delta(\mathbf{x}) = \mathbf{D}\mathbf{v}(\mathbf{x})$, with a Gaussian Process on \mathbf{v} . A full GP as oppose to laGP is usually feasible for the bias as n is small to moderate. We will soon discuss how \mathbf{D} and \mathbf{v} are calculated.

We again calibrate using a maximum-likelihood approach where the likelihood is derived from equation 10.

$$\begin{aligned} \mathbf{y}_i^F(\mathbf{x}_i^F) &= \eta(\mathbf{x}_i^F, \mathbf{t}^*) + \delta(\mathbf{x}_i^F) + \boldsymbol{\epsilon}_i \\ \implies \mathbf{y}_i^F(\mathbf{x}_i^F) &= \mathbf{B}^F \hat{\mathbf{w}}_i + \mathbf{D} \hat{\mathbf{v}}_i + \boldsymbol{\epsilon}_i \\ \mathbf{w}_i &\sim \mathcal{N}(\boldsymbol{\mu}_{wi}, \boldsymbol{\Sigma}_{wi}) \\ \mathbf{v}_i &\sim \mathcal{N}(\boldsymbol{\mu}_{vi}, \boldsymbol{\Sigma}_{vi}) \\ \boldsymbol{\epsilon}_i &\sim \mathcal{N}(\mathbf{0}, \sigma_{y^F}^2 \mathbf{I}_q) \\ \implies \mathbf{y}_i^F(\mathbf{x}_i^F) &\sim \mathbf{B}^F \mathcal{N}(\boldsymbol{\mu}_{wi}, \boldsymbol{\Sigma}_{wi}) + \mathbf{D} \mathcal{N}(\boldsymbol{\mu}_{vi}, \boldsymbol{\Sigma}_{vi}) + \mathcal{N}(\mathbf{0}, \sigma_{y^F}^2 \mathbf{I}_q) \\ \implies \mathbf{y}_i^F(\mathbf{x}_i^F) - \mathbf{B}^F \boldsymbol{\mu}_{wi} &\sim \mathcal{N}(\mathbf{D} \boldsymbol{\mu}_{vi}, \mathbf{B}^F \boldsymbol{\Sigma}_{wi} (\mathbf{B}^T)^F + \mathbf{D} \boldsymbol{\Sigma}_{vi} \mathbf{D}^T + \sigma_{y^F}^2 \mathbf{I}_q) \end{aligned} \quad (11)$$

Calibration is as follows

1. For a given \mathbf{t} , do step 1. from the calibration algorithm in Section 2.2 to get the matrices \mathbf{M}_w , \mathbf{S}_w

2. Compute bias matrix

$$\mathbf{R}_y = \mathbf{Y}^F - \mathbf{B}^F \mathbf{M}_w$$

3. Compute basis for bias

$$\text{SVD}(\mathbf{R}_y) = \mathbf{U}_b \mathbf{S}_b \mathbf{V}^T = \mathbf{D} \mathbf{V}^T$$

keeping only the first n_{pc}^* columns of \mathbf{D} and rows of \mathbf{V}^T where n_{pc}^* is chosen so that 95% of the variability in \mathbf{R}_y is accounted for.

4. Fit a separable GP to each component of the bias, estimating all parameters and predict at the same input locations.

$$\begin{aligned} \boldsymbol{\mu}_{vk}, \boldsymbol{\sigma}_{vk}^2 &= GP(X = [\mathbf{X}^F, \mathbf{t}], \\ Z &= \mathbf{V}^T[k, :], \\ XX &= [\mathbf{X}^F, \mathbf{t}]); \quad k = 1, \dots, n_{pc}^* \end{aligned} \quad (12)$$

5. Store predictive means and variances as before in matrices $\mathbf{M}_v = [\boldsymbol{\mu}_{v1} | \dots | \boldsymbol{\mu}_{vn_{pc}^*}]^T$, $\mathbf{S}_v = [\boldsymbol{\sigma}_{v1}^2 | \dots | \boldsymbol{\sigma}_{vn_{pc}^*}^2]^T$

6. Compute residuals

$$\mathbf{R}_\epsilon = \mathbf{R}_y - \mathbf{D}\mathbf{M}_v$$

and vectorize $\mathbf{r}_\epsilon = \text{Vec}(\mathbf{R}_\epsilon)$ to estimate

$$\hat{\sigma}_{y^F}^2 = \frac{1}{n * p} \mathbf{r}_\epsilon^T \mathbf{r}_\epsilon$$

7. Compute log-likelihood

$$ll \propto \sum_{i=1}^n -\frac{1}{2} \left(\log(|\Sigma_i|) + \mathbf{R}_\epsilon[:, i]^T \Sigma_i^{-1} \mathbf{R}_\epsilon[:, i] \right)$$

where $\Sigma_i = \mathbf{B}^F \Sigma_{wi} (\mathbf{B}^T)^F + \mathbf{D} \Sigma_{vi} \mathbf{D}^T + \hat{\sigma}_{y^F}^2 \mathbf{I}_q$ and $\Sigma_{vi} = \text{diag}(\mathbf{S}_v[:, i])$