



COMPUTER MODEL CALIBRATION AS A METHOD OF DESIGN

Carl Ehrett

Comprehensive exam, Department of Mathematical Sciences

2018-08-20

Section 1

INTRODUCTION: COMPUTER MODEL CALIBRATION

Traditional computer model calibration

- In a problem of computer model calibration, we have:

$$y(\mathbf{x}) = f(\mathbf{x}) + \epsilon(\mathbf{x}) = \eta(\mathbf{x}, \boldsymbol{\theta}) + \delta(\mathbf{x}) + \epsilon(\mathbf{x}) \quad (1)$$

where $y(\mathbf{x})$ is the observed value at \mathbf{x} , $f(\mathbf{x})$ is the value of the true system at \mathbf{x} , $\epsilon(\mathbf{x})$ is the observation error (often iid Gaussian), $\eta(\mathbf{x}, \boldsymbol{\theta})$ is the value of the computer model and $\delta(\mathbf{x})$ captures the discrepancy between the computer model and the true system.

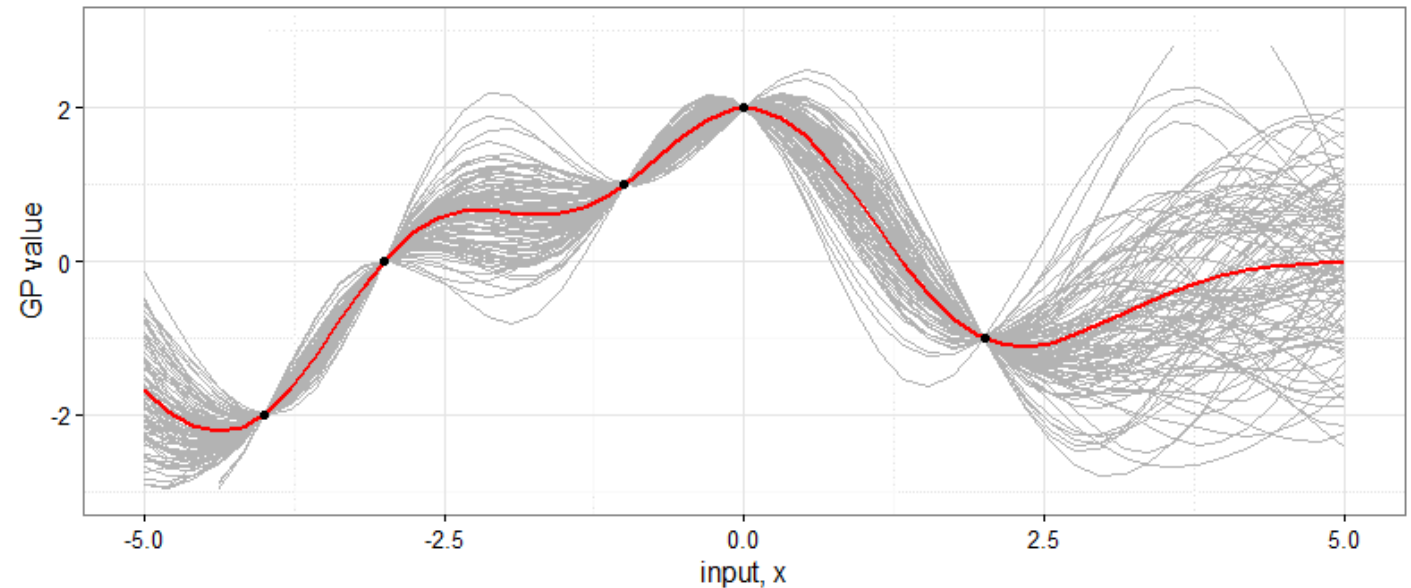
- Here, \mathbf{x} gives the known and/or controllable input settings for the system, and $\boldsymbol{\theta}$ gives the unknown parameters which must be given as inputs to the computer model.
- The purpose of computer model calibration is to use a set \mathbf{y} of observations to find the appropriate value of $\boldsymbol{\theta}$ so that $L(\delta)$ is minimized for loss $L(\cdot)$.

Calibration in the vein of Kennedy and O'Hagan (2001)

- Kennedy and O'Hagan propose a Bayesian method of calibration that results in a posterior distribution of $\theta|y$.
- A benefit of this approach is its inherent inclusion of uncertainty quantification.
- Despite identification difficulties between θ and $\delta(\cdot)$, this framework has become very influential, and has been the locus of much research.

Gaussian process emulators

- Kennedy and O'Hagan use Gaussian process (GP) emulators as code surrogates when calibrating computationally expensive models.



- For a prior GP with mean $\mu(\cdot)$ and covariance $C(\cdot, \cdot)$, given a set of observations $\boldsymbol{\eta}$ at points \mathbf{X} , the posterior mean and covariance at a new set of points \mathbf{X}' are

$$\boldsymbol{\mu}_{\mathbf{X}'}^* = \boldsymbol{\mu}_{\mathbf{X}'} + \mathbf{C}_{\mathbf{X}', \mathbf{X}} \cdot \mathbf{C}_{\mathbf{X}, \mathbf{X}}^{-1} (\boldsymbol{\eta} - \boldsymbol{\mu}_{\mathbf{X}}), \quad \mathbf{C}_{\mathbf{X}', \mathbf{X}'}^* = \mathbf{C}_{\mathbf{X}', \mathbf{X}'} - \mathbf{C}_{\mathbf{X}', \mathbf{X}} \cdot \mathbf{C}_{\mathbf{X}, \mathbf{X}}^{-1} \cdot \mathbf{C}_{\mathbf{X}, \mathbf{X}'} \quad (2)$$
 where $\boldsymbol{\mu}_{\mathbf{A}}$ is a vector giving $\mu(a)$ for each $a \in \mathbf{A}$, and $\mathbf{C}_{\mathbf{A}, \mathbf{B}}$ is a matrix for which the i, j entry is $C(a_i, b_j)$ when $\mathbf{A} = (a_1, \dots, a_n)^T$, $\mathbf{B} = (b_1, \dots, b_m)^T$.

Benefits of GP emulators

Three benefits to GP priors in the context of computer model emulation:

1. They do not require detailed foreknowledge of the model's parametric form.
2. They easily interpolate the model output.
3. They facilitate uncertainty quantification.

The framework of Williams, Higdon, Gattiker et al. (2006)

- $\eta(\cdot, \cdot)$ is modeled with a mean-zero GP prior with covariance as the product power exponential function:

$$C((\mathbf{x}, \mathbf{t}), (\mathbf{x}', \mathbf{t}')) = \frac{1}{\lambda_\eta} \prod_{k=1}^p \rho_k^\eta (x_k - x'_k)^2 \times \prod_{k=1}^q \rho_{p+k}^\eta (t_k - t'_k)^2 \quad (3)$$

where $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^p$, $\mathbf{t} \in \mathcal{T} \subseteq \mathbb{R}^q$, with priors $\lambda_\eta \sim \text{Gamma}(5, 5)$, $\rho_k^\eta \sim \text{Beta}(0.1, 1)$.

- $\delta(\cdot)$ is similarly modeled with a mean-zero GP prior, with covariance

$$C_\delta(\mathbf{x}, \mathbf{x}') = \frac{1}{\lambda_\delta} \prod_{k=1}^p \rho_k^\delta (x_k - x'_k)^2 \quad (4)$$

with priors $\lambda_\delta \sim \text{Gamma}(5, 5)$, $\rho_k^\delta \sim \text{Beta}(0.3, 1)$.

The framework of Williams, Higdon, Gattiker et al. (2006)

- Let $\boldsymbol{\eta}$ be a vector of computer model outputs at inputs $((\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_n, \mathbf{t}_n))$ and \mathbf{y} be a vector of real world observations at $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m})$ with \mathbf{C}_y a matrix giving the (known) observation variance for \mathbf{y} . Then for $\mathcal{D} = (\boldsymbol{\eta}^T, \mathbf{y}^T)^T$ we have

$$\mathcal{D} | \boldsymbol{\theta}, \lambda_{\boldsymbol{\eta}}, \boldsymbol{\rho}^{\boldsymbol{\eta}}, \lambda_{\delta}, \boldsymbol{\rho}^{\delta}, \mathbf{C}_y \sim N(\mathbf{0}, \mathbf{C}_{\mathcal{D}}) \quad (5)$$

where $\mathbf{C}_{\mathcal{D}}$ is the sum of a matrix giving $C((\mathbf{x}_i, \mathbf{t}_i), (\mathbf{x}_j, \mathbf{t}_j))$ for the i, j entry with the matrix $\begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_y + \mathbf{C}_{\delta} \end{bmatrix}$. The matrix \mathbf{C}_{δ} has, for its i, j entry, $C_{\delta}(\mathbf{x}_{n+i}, \mathbf{x}_{n+j})$.

- The joint posterior under the model is

$$\pi(\boldsymbol{\theta}, \lambda_{\boldsymbol{\eta}}, \boldsymbol{\rho}^{\boldsymbol{\eta}}, \lambda_{\delta}, \boldsymbol{\rho}^{\delta}, \mathbf{C}_y | \mathcal{D}) \propto \pi(\mathcal{D} | \boldsymbol{\theta}, \lambda_{\boldsymbol{\eta}}, \boldsymbol{\rho}^{\boldsymbol{\eta}}, \lambda_{\delta}, \boldsymbol{\rho}^{\delta}, \mathbf{C}_y) \times \pi(\lambda_{\boldsymbol{\eta}}) \times \pi(\boldsymbol{\rho}^{\boldsymbol{\eta}}) \times \pi(\lambda_{\delta}) \times \pi(\boldsymbol{\rho}^{\delta}) \quad (6)$$

- This distribution can be explored via MCMC.

Section 2

ADAPTING THE CALIBRATION FRAMEWORK FOR OPTIMIZATION

Calibration to desired observations

- We can reconceptualize the Kennedy-O'Hagan calibration framework as a method for optimization.
- In the schema $y(\mathbf{x}) = f(\mathbf{x}) + \epsilon(\mathbf{x}) = \eta(\mathbf{x}, \boldsymbol{\theta}) + \delta(\mathbf{x}) + \epsilon(\mathbf{x})$, we now take $\boldsymbol{\theta}$ to be controllable inputs over which we wish to optimize model output.
- Instead of calibrating to a set of real-world observations \mathbf{y} , instead set \mathbf{y} to be a set of performance targets. Call performance targets used as manufactured data in such calibration *desired observations*.
- The result of calibration to desired observations (CDO) is a posterior distribution $\boldsymbol{\theta}|\mathbf{y}$. As under Kennedy-O'Hagan calibration, uncertainty quantification is included.

Model shortcoming

- Given the schema $y(\mathbf{x}) = f(\mathbf{x}) + \epsilon(\mathbf{x}) = \eta(\mathbf{x}, \boldsymbol{\theta}) + \delta(\mathbf{x}) + \epsilon(\mathbf{x})$, “error” can be modeled as either $\delta(\mathbf{x}) + \epsilon(\mathbf{x})$ or just as $\epsilon(\mathbf{x})$ if a (nonzero) discrepancy is not included.
- Either works for achieving a distribution on $\boldsymbol{\theta}$ values which lead to output approximating desired observations \mathbf{y} .
- The benefit of using only $\epsilon(\cdot)$ is convenience in setting optimization priorities, and computational efficiency.
- The benefit of including $\delta(\cdot)$ is that the resulting error analysis is more accurate. The posterior distribution on $\delta(\cdot)$ describes the discrepancy between the desired observations and the optimal model output.

Hyperparameter estimation

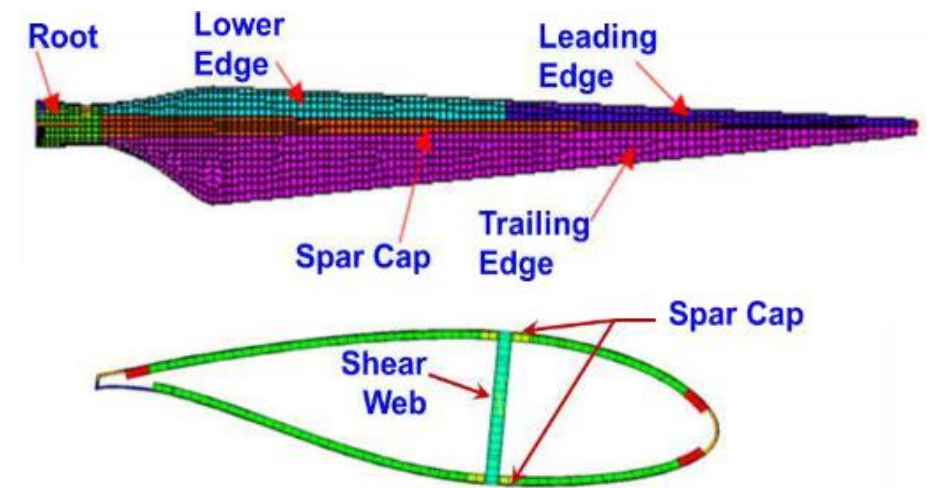
- Rather than estimate the hyperparameters λ_η, ρ^η as part of a full Bayesian analysis, we use their MLEs.
- This practice is common in calibration for convenience, but is crucial here for a different reason.
- These hyperparameters are part of our code surrogate. We want them to “learn” only from the code observations. We do **not** want them to “learn” from our (fake) desired observations \mathbf{y} .

Section 3

APPLICATION: WIND TURBINE BLADE MATERIAL DESIGN

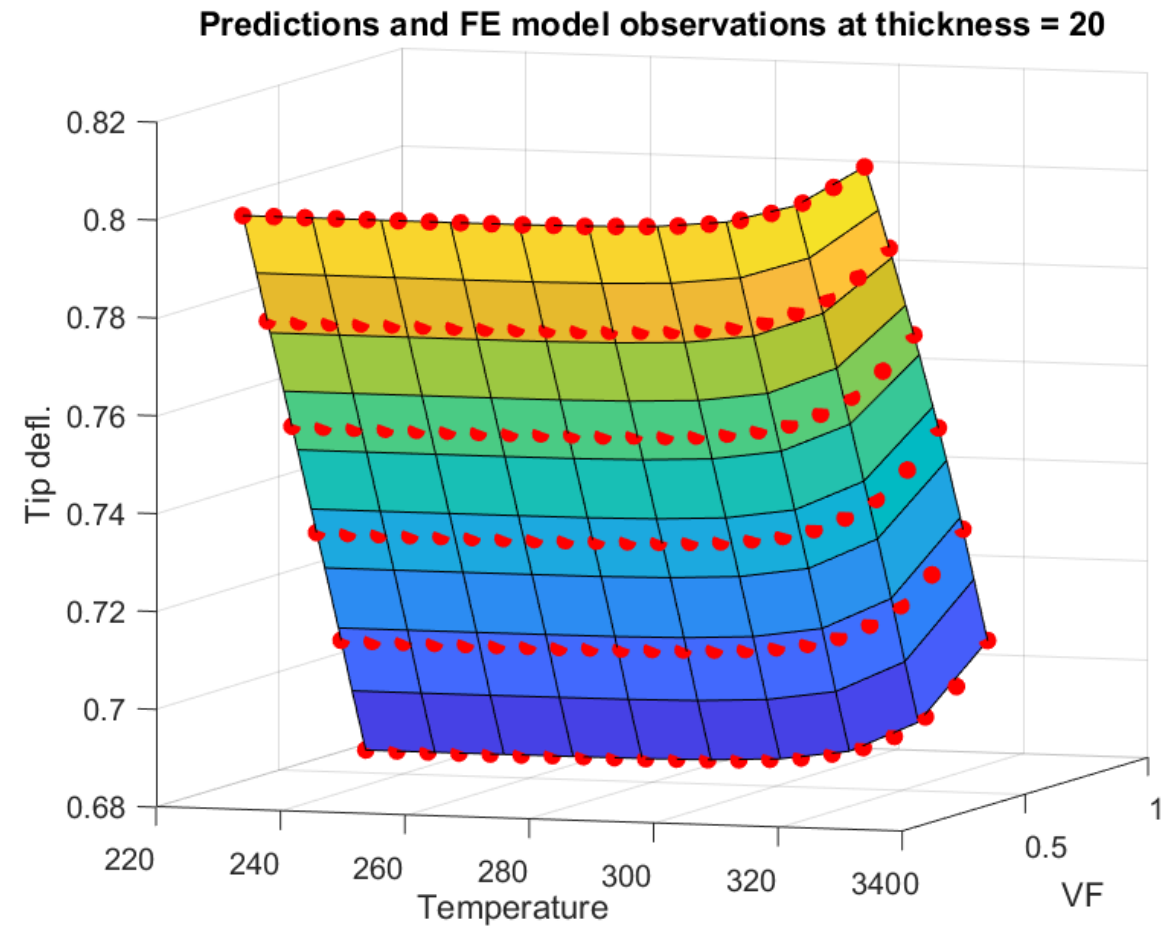
Wind turbine blade design problem

- Design problem: given a blade of fixed geometry, determine the optimal *volume fraction* $v \in (0,1)$ and *thickness* in mm k , across a range of temperatures h .
- We want to minimize three model outputs: tip deflection in meters d , rotation in radians r , and cost in USD c .



Wind turbine blade computer model

- The computer model is a finite element model implemented in ANSYS. It is too computationally expensive to deploy in MCMC.
- We collected 504 observations in a latin hypercube on the inputs (h, v, k).
- The output is trivariate. Dummy variables were added to convert the observations to 1512 univariate observations.
- Hyperparameters MLEs: $\widehat{\lambda}_{\eta} = 0.0152$,
 $\widehat{\rho}^{\eta} = (0.936, 0.651, 0.674, 0.480, 0.968)$



Model for wind turbine application

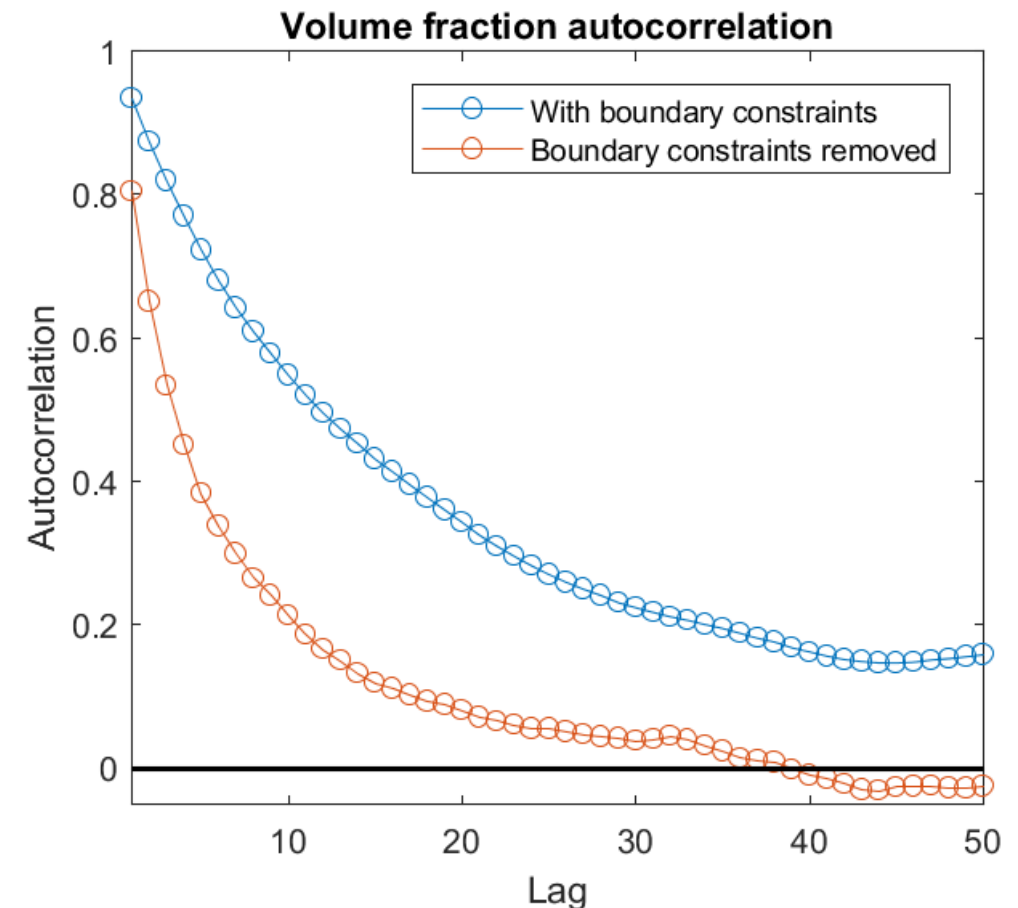
- Various options are available for the priors $\pi(\boldsymbol{\theta})$ and $\pi(\mathbf{C}_y)$. But in full generality, the model is given by:

$$\begin{aligned}\pi(\boldsymbol{\theta}, \mathbf{C}_y | \mathcal{D}) &\propto \pi(\mathcal{D} | \boldsymbol{\theta}, \mathbf{C}_y) \times \pi(\boldsymbol{\theta}) \times \pi(\mathbf{C}_y) \\ &\propto |\mathbf{C}_{\mathcal{D}}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathcal{D}^T \mathbf{C}_{\mathcal{D}}^{-1} \mathcal{D}\right) \times \pi(\boldsymbol{\theta}) \times \pi(\mathbf{C}_y)\end{aligned}\quad (7)$$

- This constitutes an expansion of the framework of Williams et al. (2006), to include estimation of \mathbf{C}_y .
- Here, \mathbf{C}_y is assumed to be diagonal. Options include homoscedasticity across model outputs ($\mathbf{C}_y = \sigma^2 \mathbf{I}$), or allowing σ_i^2 to vary for $i = 1, 2, 3$ corresponding to the three model outputs.

MCMC details

- Boundary constraints: $\theta \in [0.2, 0.6] \times [10, 25]$, $\sigma_i^2 > 0$.
- MCMC is improved by removing these boundary constraints by logit-transforming each element of θ , and taking the log of each σ_i^2 .
- I used normal proposal densities with adaptive covariance matrices to promote acceptance rates $\sim 23\%$.
- Adaptive algorithm: initialize $\tau = 1$. Every 100 draws during burn-in, let Σ be the sample covariance of all previous draws. If fewer than 20 of previous 100 draws were accepted, update $\tau = 0.75\tau$; if more than 25 accepted, update $\tau = 1.25\tau$. Set new proposal covariance to be $\tau\Sigma$.
- All computations are carried out on the log scale where possible, for computational efficiency.



Options for $\pi(\mathbf{C}_y)$

- Two independent sets of options:
 1. Specify predetermined values for σ_i^2 , $i = 1, 2, 3$, or else set a prior on them.
 2. Constrain $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$, or else allow them to take different values.
- Specifying predetermined values is useful for implementing priorities. A prior is useful when one has insufficient knowledge of the Pareto front to implement priorities.
- With $\pi(\boldsymbol{\theta})$ uniform, when specifying values of $\boldsymbol{\sigma}^2$ the full model is:

$$\pi(\boldsymbol{\theta}, \mathbf{C}_y | \mathcal{D}) \propto |\mathbf{C}_\mathcal{D}|^{-\frac{1}{2}} \exp(\mathcal{D}^T \mathbf{C}_\mathcal{D}^{-1} \mathcal{D}) \quad (8)$$

and with the reference prior $\pi(\mathbf{C}_y) = \prod_{k=1}^3 1/\sigma_k^2$, the full model is:

$$\pi(\boldsymbol{\theta}, \mathbf{C}_y | \mathcal{D}) \propto |\mathbf{C}_\mathcal{D}|^{-\frac{1}{2}} \exp(\mathcal{D}^T \mathbf{C}_\mathcal{D}^{-1} \mathcal{D}) \times \prod_{k=1}^3 \frac{1}{\sigma_k^2} \quad (9)$$

Section 4

PARETO FRONT ESTIMATION

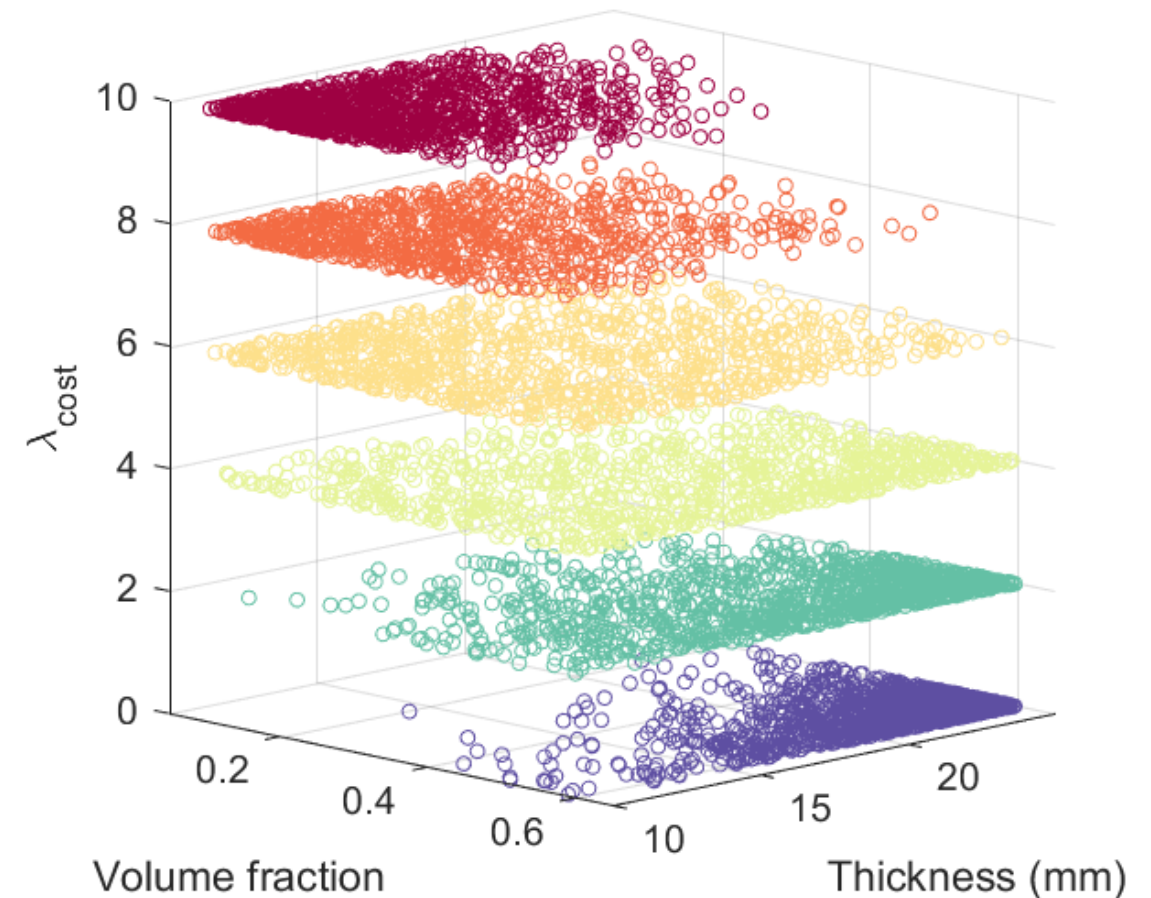
Exploring the Pareto front

- Sometimes priorities are not specific enough to be represented by a single set of desired observations.
- One might wish to explore optimal outputs over (e.g.) a range of costs.
- Two methods are described here for doing so: (1) exploiting known relationships in the model through a prior on θ , and (2) calibrating over a grid of “known” values of a subset of model outputs.

Exploiting known input/output relationships

- High volume fraction and high thickness are each known to correlate to high cost.
- Therefore, rather than include cost, we can use a prior $\pi(\boldsymbol{\theta}) = \exp -\lambda_c \|\boldsymbol{\theta}\|^2$ to penalize high-cost regions of the parameter space. The severity of the penalty depends upon $\lambda_c > 0$.
- Removing cost from the model brings computational benefits. We have now only 1008 observations rather than 1512.
- The effect of increasing the value of λ_c is to push $\boldsymbol{\theta}$ out of the upper (high-cost) region of its support toward the lower (low-cost) region.

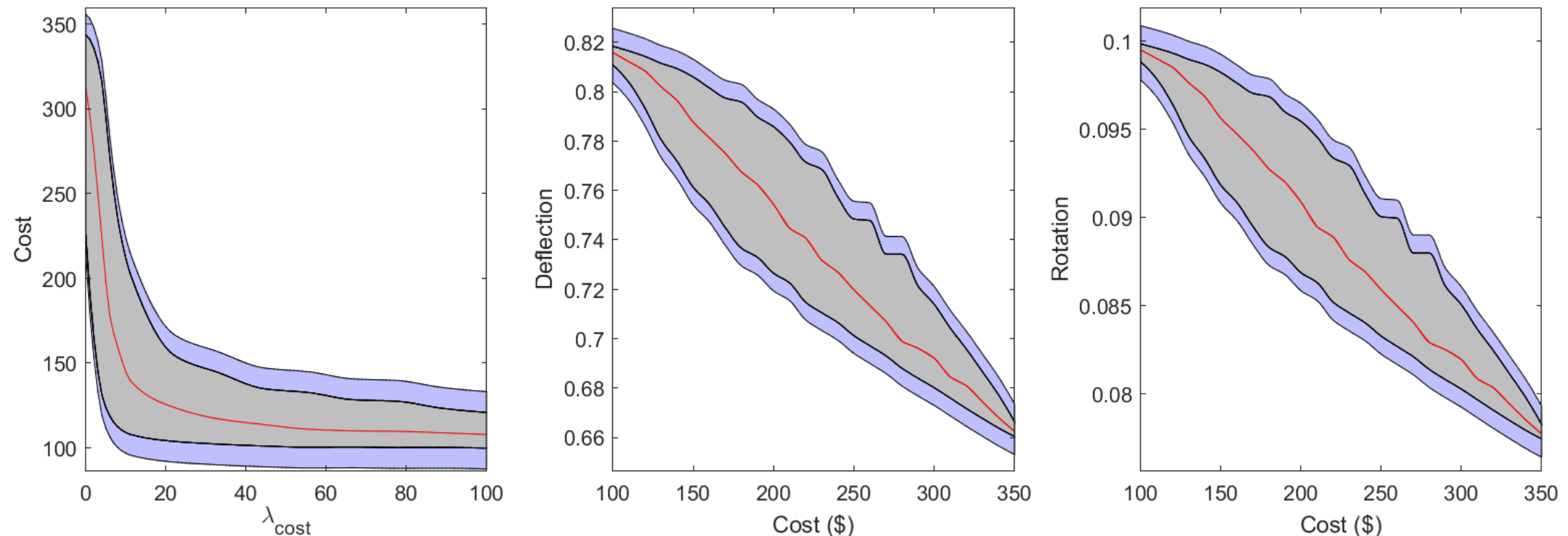
Posterior distribution of volume fraction and thickness



Pareto bands using λ_c

Calibrating over each of a range of λ_c , we achieve a comprehensive picture of how optimal performance varies with cost, with included uncertainty bands.

Cost vs. λ_{cost} , and performance metrics vs. cost with 90% credible interval including code uncertainty



Specifying known cost

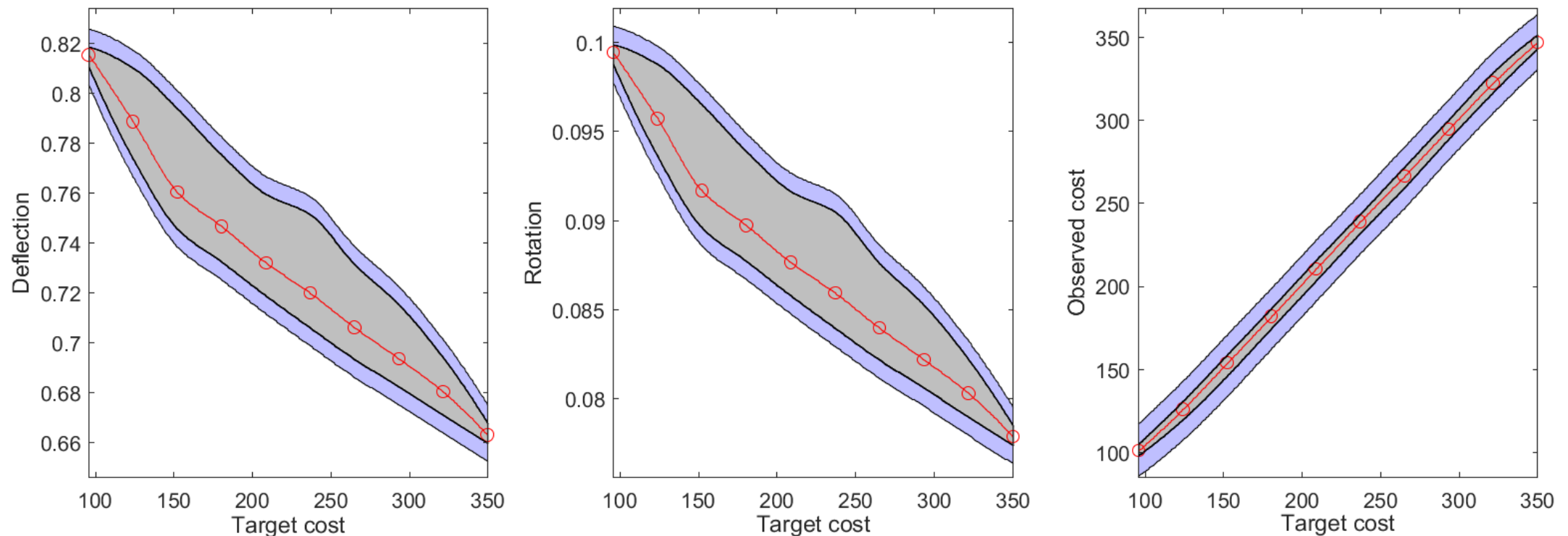
- Alternatively, one can achieve a similar result (even when no exploitable input/output relationship is known) by performing CDO over a grid of “known” costs.
- Set a grid $\{c_i\}_{i=1}^p$ of achievable costs, and for each grid point c_i , calibrate to desired observation (e.g.) $[0,0, c_i]$.
- In order to specify cost as known up to slight deviations, set σ_3^2 equal to a low value (here, 0.05 on standardized scale).
- When using the reference prior for the variance of deflection and rotation, and uniform prior on $\boldsymbol{\theta}$, the full posterior is:

$$\pi(\boldsymbol{\theta}, \mathbf{C}_y | \mathcal{D}) \propto |\mathbf{C}_{\mathcal{D}}|^{-\frac{1}{2}} \exp(\mathcal{D}^T \mathbf{C}_{\mathcal{D}}^{-1} \mathcal{D}) \times \prod_{i=1}^2 \frac{1}{\sigma_i^2} \quad (10)$$

Pareto bands using cost grid

Calibrating over each of a range of values of cost, we achieve an estimate of optimal performance with uncertainty bands across costs.

Performance metrics vs. (known) target cost, with 90% credible interval including code uncertainty



Appendix

MCMC ALGORITHM

Algorithm 1:

Step A: Draw $\boldsymbol{\theta}^{(0)} \sim \text{Unif}([0, 1]^2)$, $(\sigma_j^2)^{(0)} \sim \text{Gamma}(2, 2)$ for $j = 1, 2, \dots, K$. Set \mathbf{C}_y using $(\sigma_j^2)^{(0)}$, $j = 1, 2, \dots, K$. Set \mathbf{C}_D using $\boldsymbol{\theta}^{(0)}$ and \mathbf{C}_y . Set $s^2 = 1$. Specify $\lambda_{cost} \geq 0$. Specify total number of iterations M and burn-in b . Set mult = 2. Set $\Sigma = I_2$. Set $A = 0$ and $A_\sigma = 0$. Set $i = 0$.

Step B: Repeat M times.

1. Update $i = i + 1$.
 2. Draw $\boldsymbol{\theta}^* \sim \text{logit}^{-1} \left(\text{MVN}(\text{logit}(\boldsymbol{\theta}^{(i-1)}), \Sigma) \right)$ and set \mathbf{C}_D^* using $\boldsymbol{\theta}^*$ and \mathbf{C}_y .
 3. Find:

$$\log \alpha = \log \left(|\mathbf{C}_D^*|^{-1/2} \exp(\mathcal{D}^T \mathbf{C}_D^{*-1} \mathcal{D}) \right) - \lambda_{cost} \cdot \|\boldsymbol{\theta}^*\|^2 - \log \left(|\mathbf{C}_D|^{-1/2} \exp(\mathcal{D}^T \mathbf{C}_D^{-1} \mathcal{D}) \right) \\ + \lambda_{cost} \cdot \|\boldsymbol{\theta}^{(i-1)}\|^2 + \log \left(\prod_{j=1}^2 (\theta_j^* (1 - \theta_j^*)) \right) - \log \left(\prod_{j=1}^2 (\theta_j^{(i-1)} (1 - \theta_j^{(i-1)})) \right)$$
 4. Draw $a \sim \text{Unif}(0, 1)$. If $a < \alpha$, set $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^*$, $\mathbf{C}_D = \mathbf{C}_D^*$, and increment $A = A + 1$; otherwise, $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)}$.
 5. Draw $[\sigma_1^{2*}, \sigma_2^{2*}, \sigma_3^{2*}]^T \sim \exp \left(\text{MVN} \left(\log \left[\sigma_1^{2(i-1)}, \sigma_2^{2(i-1)}, \sigma_3^{2(i-1)} \right]^T, s^2 I_3 \right) \right)$. Set \mathbf{C}_y^* using $[\sigma_1^{2*}, \sigma_2^{2*}, \sigma_3^{2*}]^T$. Set \mathbf{C}_D^* using $\boldsymbol{\theta}^{(i)}$ and \mathbf{C}_y^* .
 6. Find $\log \alpha = \log (|\mathbf{C}_D^*|^{-1/2} \exp(\mathcal{D}^T \mathbf{C}_D^{*-1} \mathcal{D})) - \log (|\mathbf{C}_D|^{-1/2} \exp(\mathcal{D}^T \mathbf{C}_D^{-1} \mathcal{D}))$. (Notice that the log ratio of priors $\log \left(\frac{\pi(\sigma_i^{2'})}{\pi(\sigma_i^2)} \right) = \log \frac{\sigma_i^2}{\sigma_i^{2'}} = \log \sigma_i^2 - \log \sigma_i^{2'}$ is cancelled out by the log Metropolis-Hastings correction for the asymmetrical proposal density: $\log \left(\frac{q(\sigma_i^2 | \sigma_i^{2'})}{q(\sigma_i^{2'} | \sigma_i^2)} \right) = \log \left(\frac{\sigma_i^{2'}}{\sigma_i^2} \right) = \log \sigma_i^{2'} - \log \sigma_i^2$.)
 7. Draw $a \sim \text{Unif}(0, 1)$. If $a < \alpha$, set $[\sigma_1^{2(i)}, \sigma_2^{2(i)}, \sigma_3^{2(i)}]^T = [\sigma_1^{2*}, \sigma_2^{2*}, \sigma_3^{2*}]^T$, $\mathbf{C}_y = \mathbf{C}_y^*$, $\mathbf{C}_D = \mathbf{C}_D^*$, and increment $A_\sigma = A_\sigma + 1$; otherwise, set $[\sigma_1^{2(i)}, \sigma_2^{2(i)}, \sigma_3^{2(i)}]^T = [\sigma_1^{2(i-1)}, \sigma_2^{2(i-1)}, \sigma_3^{2(i-1)}]^T$.
 8. If $i \leq b$ and $i \pmod{100} = 0$: Update mult = $1.5 \cdot \text{mult} \cdot \mathbf{1}_{A > 30} + 0.75 \cdot \text{mult} \cdot \mathbf{1}_{A < 20}$ and set $A = 0$. Update $\Sigma = \text{mult} \cdot \text{Cov}(\Theta)$, where $\Theta = [\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(i)}]^T$. Update $s^2 = 1.5 \cdot s^2 \cdot \mathbf{1}_{A_\sigma > 30} + 0.75 \cdot s^2 \cdot \mathbf{1}_{A_\sigma < 20}$ and set $A_\sigma = 0$.
-