# Optimal Reduced Order Modelling

Robert A. Milton, Solomon F. Brown

*Department of Chemical and Biological Engineering, University of Sheffield, Sheffield, S1 3JD, United Kingdom*

## Abstract

*Keywords:* Gaussian Process, Global Sensitivity Analysis, Sobol' Index, Surrogate Model

## 1. Introduction

A broad range of engineering topics ultimately concern the response of some noisy scalar quantity $y(\mathbf{x})$ to its $M$-dimensional input $\mathbf{x}$. It is extremely desirable to reduce $M$ as far as possible without materially impacting $y$ for a number of reasons.

Firstly, $y(\mathbf{x})$ is usually onerous to obtain, from pilot plant, laboratory or simulation. Therefore mitigating the input variables to be controlled and monitored is disproportionately advantageous.

Secondly, visualisation and analysis is far easier and more fruitful in a low dimensional space. There is epistemilogical value in model-order reduction for its own sake, and to inspire and guide advances in theory or modelling.

Thirdly, refined examination of the response throughout the input space demands an ensemble of results $y(\mathbf{x})$ whose numerosity grows exponentially with $M$. Worse still, the basic geometry of high-dimensional space contains some nasty surprises, collectively known as the curse of dimensionality. In particular, nearly all the hypervolume of a high-dimensional input space resides proximate to its boundary hypersurface. Taken without prejudice, the input hypervolume consists almost entirely of anomalies, flung to the outer reaches and extremal in one input dimension or other. The curse thus

---

manifests as a kind of engineering paranoia: failure (represented by extremal conditions) becomes almost inevitable as ever more inputs are measured or modelled.

The antidote to this perceived curse is to infer that monitoring more inputs must diminish their mutual independence or their relevance. The cure is to reduce the inputs to a set which is both mutually independent and highly relevant to the response. Topologically this corresponds to sweeping naively sampled input points away from the faces of a sampling hypercube. Most of the sample points get swept towards the heart of the input space, where operating conditions are normal and safe. The few remaining input samples get swept towards selected corners of the input hypercube, representing unsafe conditions in which several inputs record anomalous values simultaneously. Heuristically this captures failure due to inscrutably arcane causes, which is first detected in several measured variables at once, as well as failure due to unsafe levels in a single highly relevant input. Model order reduction basically evacuates then excises those regions of input space which are irrelevant to the response, or represent impossible combinations of dependent inputs. This can only by achieved by a principled analysis of the response $y(\mathbf{x})$ which detects both irrelevance and dependence. Which is dauntingly hard for all the same reasons it is so sorely needed.

### 1.1. Gaussian Process Surrogate

In the first instance, the difficulty and expense in obtaining and analyzing response data is mitigated by adopting a surrogate or emulator. Perhaps the most informative and general of these is the Gaussian Process (GP), widely investigated and employed over the past few decades. The response $y(\mathbf{x})$ to arbitrarily fixed input is modelled as the sum $f(\mathbf{x}) + e(\mathbf{x})$ of two Gaussian random variables encapsulating coherent signal and incoherent noise. The latter is characterised by a zero-mean distribution which is independent of the input:

$$e(\mathbf{x}) \sim \mathsf{N}\left[0, \sigma_{\mathbf{e}}^2\right]$$

The signal $f(\mathbf{x})$ is characterized by its covariance kernel $\sigma_{\mathbf{f}}^2 k(\mathbf{x}_n, \mathbf{x})$ which measures the similarity between inputs $\mathbf{x}_n$ and $\mathbf{x}$, and propagates any similarity to $y(\mathbf{x}_n)$ and $y(\mathbf{x})$. In the majoriy of applications, the kernel is naturally stationary, a function of $(\mathbf{x} - \mathbf{x}_n)$ alone. We shall further assume that the kernel is twice differentiable at its maximium ($\mathbf{x} = \mathbf{x}_n$). Once forced to exist, the Hessian at the maximum must be symmetric negative semidefinite and

therefore diagonalizes to

$$\partial_{\mathbf{xx}} \log k(\mathbf{x}, \mathbf{x}) =: -\Theta^\mathsf{T} \Lambda^{-2} \Theta$$

When $|\mathbf{x} - \mathbf{x}_n|$ is large the kernel value is miniscule in any any relevant direction. The kernel details are therefore largely irrelevant to the response *any* time $|\mathbf{x} - \mathbf{x}_n|$ is large, advocating (if not justifying) the Taylor approximation

$$k(\mathbf{x}_n, \mathbf{x}) = \exp\left(-\frac{(\mathbf{x} - \mathbf{x}_n)^\mathsf{T} \Theta^\mathsf{T} \Lambda^{-2} \Theta (\mathbf{x} - \mathbf{x}_n)}{2}\left(1 + O(|\mathbf{x} - \mathbf{x}_n|)\right)\right)$$

This paper is exclusively concerned with kernels of this form. The differentiability we have imposed forces the power spectrum of the signal $f$ to decay rapidly. Modes of response oscillating rapidly with $\mathbf{x}$ are intepreted as noise by the GP, as the kernel smoothes $y$ into $f$. Such regularization is often, but not always, desirable, to avoid wildly unreliable interpolation of an overfit regression. We shall expand on this point later.

*1.2. Kernel Optimization*

The cure to the curse of dimensionality is to find orthogonal rotation matrix $\Theta$ and diagonal lengthscale matrix $\Lambda$ which best fit observed responses $y(\mathbf{X}^\mathsf{T})$. The largest lengthscales in $\Lambda$ mark the least relevant directions, to be pruned. However, the best fit must optimize $M(M+1)/2 + 2$ hyperparameters simultaneously to determine $\Theta, \Lambda, \sigma_\mathbf{f}^2$ and $\sigma_\mathbf{e}^2$. An exploitative, directed optimization of such numerosity is almost bound to get bogged down in local optima. Exploratory grid search is astronomically expensive $O(\exp(M(M+1)/2)$, likewise any random sampling which is not hopelessly sparse. Perhaps for these reasons, $\Theta$ has always been fixed as identity in the literature. The lengthscales comprising $\Lambda$ are also usually identical, furnishing a radial basis function (RBF) kernel. The few studies where $\Lambda$ is not identical speak of an automatic relevance determination (ARD) kernel, with model order reduction in mind. However, admitting no rotation $\Theta$ severely restricts the reductions available. If relevant inputs are mutually dependent, they must all be retained. Rotation, on the contrary, combines them into a low dimensional subspace.

*1.3. Global Sensitivity Analysis*

This paper proposes to achieve kernel optimization indirectly, via global sensitivity analysis (GSA). The surrogate expectation

$$\mathsf{E}_\Omega[y(\mathbf{x})] = \mathsf{E}_\Omega[f(\mathbf{x})] =: \bar{f}(\mathbf{x})$$

has a variation (over $\mathbf{x} \in \mathbb{R}^M$) which can be apportioned by Sobol' index

$$S_{\mathbf{m}}((\Theta)_{\mathbf{m} \times \mathbf{M}}) := \mathsf{Var_x}\left[\mathsf{E_x}\left[\bar{f}(\mathbf{x})|(\Theta\mathbf{x})_{\mathbf{m}}\right]\right] / \mathsf{Var_x}\left[\bar{f}(\mathbf{x})\right] \leq 1$$

to subspaces $(\Theta\mathbf{x})_{\mathbf{m}}$ of dimension $m \leq M$. These may be calculated analytically for the exponential quadratic kernel used here. To cure to the curse of dimensionality is to find $(\Theta)_{\mathbf{m} \times \mathbf{M}}$ such that $S_{\mathbf{m}} \approx 1$ for $m \ll M$. The rotation sub-matrix $(\Theta)_{\mathbf{m} \times \mathbf{M}}$ has a manageable number of elements if $m$ is small. This paper takes the most economical approach, maximizing $S_{\mathbf{m}}$ for $m = 1, \ldots, M$ in turn, to find the most relevant direction, then the second most relevant, and so on. Other approaches are considered in Section 4.

## 2. Methodology

Let $\mathbf{X}$ be the $(N \times M)$ design matrix of observed inputs eliciting the $N$ response $y(\mathbf{X}^\intercal)$. The observations are standardized such that

$$(\mathbf{0})_{\mathbf{M}} = \mathsf{E}[\mathbf{x}_n] := \sum_{n=1}^{N}(\mathbf{X})_{n \times \mathbf{M}}^\intercal \quad ; \quad 1 = \mathsf{Var}[\mathbf{x}_n] = N^{-1}\sum_{n=1}^{N}(\mathbf{X})_{n \times \mathbf{M}}(\mathbf{X})_{n \times \mathbf{M}}^\intercal$$

$$0 = \mathsf{E}[y(\mathbf{x}_n)] := \sum_{n=1}^{N}(y(\mathbf{X}^\intercal))_n \quad ; \quad 1 = \mathsf{Var}[y(\mathbf{x}_n)] = N^{-1}y(\mathbf{X}^\intercal)^\intercal y(\mathbf{X}^\intercal)$$

where boldface subscripts refer to the multi-indices

$$\emptyset =: \mathbf{0} \subseteq \mathbf{m} := (1, \ldots, m) \subseteq \mathbf{M} \tag{1}$$

which always precede superscript operations (such as transposition or inversion). For brevity, we shall admit vector Gaussian probability densities $p((\mathbf{z})_{\mathbf{m}}; (\mathbf{Z})_{\mathbf{m} \times \mathbf{J}}, \Sigma_{\mathbf{z}})$ such that

$$(p((\mathbf{z})_{\mathbf{m}}; (\mathbf{Z})_{\mathbf{m} \times \mathbf{J}}, \Sigma_{\mathbf{z}}))_j$$
$$:= (2\pi)^{-M/2} |\Sigma_{\mathbf{z}}|^{-1/2} \exp\left(-\frac{(\mathbf{z} - (\mathbf{Z})_{\mathbf{m} \times j})^\intercal \Sigma_{\mathbf{z}}^{-1}(\mathbf{z} - (\mathbf{Z})_{\mathbf{m} \times j})}{2}\right) \tag{2}$$

naturally collapsing to the (scalar) normal multivariate density when $J = 1$.

## 2.1. Gaussian Process Surrogate

Non-parametric GP regression fits signal $f$ and noise $e$ Gaussian processes to

$$y(\mathbf{X^\intercal}) = f(\mathbf{X^\intercal}) + e(\mathbf{X^\intercal}) \tag{3}$$

This work exclusively employs objective Bayesian priors

$$f(\mathbf{X^\intercal}) \sim \mathsf{N}\big[(\mathbf{0})_\mathbf{N}, \sigma_\mathbf{f}^2 k(\mathbf{X^\intercal}, \mathbf{X^\intercal})\big]$$
$$e(\mathbf{X^\intercal}) \sim \mathsf{N}\big[(\mathbf{0})_\mathbf{N}, \sigma_\mathbf{e}^2 (\mathbf{1})_{\mathbf{N}\times\mathbf{N}}\big]$$

built on an ARD kernel [1, 2]

$$k(\mathbf{x}_n, \mathbf{x}) := (2\pi)^{M/2} \, |\Lambda| \, p\big(\mathbf{x}; \mathbf{x}_n, \Lambda^2\big) \tag{4}$$

with diagonal positive definite lengthscale matrix $\Lambda$. Bayesian conditioning ultimately furnishes the predictive process

$$y(\mathbf{x}) \sim \mathsf{N}\big[\bar{f}(\mathbf{x}), \Sigma_\mathbf{f}(\mathbf{x}) + \sigma_\mathbf{e}^2\big]$$

with signal mean and variance

$$\begin{aligned}\bar{f}(\mathbf{x}) &:= \sigma_\mathbf{f}^2 k(\mathbf{x}, \mathbf{X^\intercal})\mathbf{K}^{-1} y(\mathbf{X^\intercal}) \\ \Sigma_\mathbf{f}(\mathbf{x}) &:= \sigma_\mathbf{f}^2 k(\mathbf{x}, \mathbf{x}) - \sigma_\mathbf{f}^2 k(\mathbf{x}, \mathbf{X^\intercal})\mathbf{K}^{-1}\sigma_\mathbf{f}^2 k(\mathbf{X^\intercal}, \mathbf{x})\end{aligned} \tag{5}$$

where

$$\mathbf{K} := \sigma_\mathbf{f}^2 k(\mathbf{X^\intercal}, \mathbf{X^\intercal}) + \sigma_\mathbf{e}^2 (\mathbf{1})_{\mathbf{N}\times\mathbf{N}} \tag{6}$$

The $M + 2$ hyperparameters constituting $\Lambda, \sigma_\mathbf{f}$ and $\sigma_\mathbf{e}$ are simultaneously optimized for maximimum marginal likelihood $\mathsf{p}[y|\mathbf{X^\intercal}]$, using the GPy software library.

## 2.2. Global Sensitivity Analysis

Imagine a sample datum $\mathbf{u}$ is drawn from a standardized normal test distribution

$$\mathbf{u} \sim \mathsf{N}[(\mathbf{0})_\mathbf{M}, (\mathbf{1})_{\mathbf{M}\times\mathbf{M}}] \tag{7}$$

The datum basis is rotated to

$$\mathbf{x} =: \Theta^\intercal \mathbf{u} \tag{8}$$

eliciting the conditional surrogate responses

$$f_\mathbf{m}((\mathbf{u})_\mathbf{m}) := \mathsf{E}\big[\bar{f}(\Theta^\intercal\mathbf{u})|(\mathbf{u})_\mathbf{m}\big] \tag{9}$$

Knowledge of $\mathbf{u}$ herein ranges from totally conditional $f_{\mathbf{M}}(\mathbf{u}) = \bar{f}(\mathbf{x})$ to unconditional ignorance $f_{\mathbf{0}} = \mathsf{E}\big[\bar{f}(\mathbf{x})\big]$. Equations (4) to (7) enable analytic integration yielding

$$f_{\mathbf{m}}((\mathbf{u})_{\mathbf{m}}) = \tilde{\mathbf{f}}^{\mathsf{T}} \, \frac{p\big((\mathbf{u})_{\mathbf{m}}; (\mathbf{T})_{\mathbf{N}\times\mathbf{m}}^{\mathsf{T}}, (\Sigma)_{\mathbf{m}\times\mathbf{m}}\big)}{p((\mathbf{u})_{\mathbf{m}}; (\mathbf{0})_{\mathbf{m}}, (\mathbf{1})_{\mathbf{m}\times\mathbf{m}})} \tag{10}$$

where $\tilde{\mathbf{f}}$ is the Hadamard (element-wise) product $\circ$ of two vectors

$$\tilde{\mathbf{f}} := (2\pi)^{M/2} \, |\Lambda| \, p\big(\mathbf{0}; \mathbf{X}^{\mathsf{T}}, \Lambda^2 + \mathbf{1}\big) \circ \big(\sigma_{\mathbf{f}}^2 \mathbf{K}^{-1} y(\mathbf{X}^{\mathsf{T}})\big) \tag{11}$$

and

$$\mathbf{T} := \mathbf{X} \left(\Lambda^2 + \mathbf{1}\right)^{-1} \Theta^{\mathsf{T}} \tag{12}$$

$$\Sigma := \Theta \left(\Lambda^{-2} + \mathbf{1}\right)^{-1} \Theta^{\mathsf{T}} \tag{13}$$

According to these formulae, the unconditional surrogate response is

$$f_{\mathbf{0}} = \mathsf{E}\big[\bar{f}(\mathbf{x})\big] = \tilde{\mathbf{f}}^{\mathsf{T}}(\mathbf{1})_{\mathbf{N}} \tag{14}$$

which does not depend on $\Theta$ of course. Standardization of $y(\mathbf{X}^{\mathsf{T}})$ instills an expectation of precisely zero here if $\mathbf{x}_n \sim \mathsf{N}[(\mathbf{0})_{\mathbf{M}}, (\mathbf{1})_{\mathbf{M}\times\mathbf{M}}]$ (which is often not exactly true).

Conditional variances may now be calculated as

$$D_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}}) := \mathsf{Var}[f_{\mathbf{m}}((\mathbf{u})_{\mathbf{m}})] = \frac{\tilde{\mathbf{f}}^{\mathsf{T}} \, \mathbf{W}_{\mathbf{m}} \, \tilde{\mathbf{f}}}{\left|2(\Sigma)_{\mathbf{m}\times\mathbf{m}} - (\Sigma)_{\mathbf{m}\times\mathbf{m}}^2\right|^{1/2}} - f_{\mathbf{0}}^2 \tag{15}$$

where

$$(\mathbf{W}_{\mathbf{m}})_{n\times o} := \exp\left(\frac{-(\mathbf{T})_{n\times\mathbf{m}}(\Sigma)_{\mathbf{m}\times\mathbf{m}}^{-1}(\mathbf{T})_{n\times\mathbf{m}}^{\mathsf{T}} - (\mathbf{T})_{o\times\mathbf{m}}(\Sigma)_{\mathbf{m}\times\mathbf{m}}^{-1}(\mathbf{T})_{o\times\mathbf{m}}^{\mathsf{T}}}{2}\right)$$
$$\times \exp\left(\frac{+\left((\mathbf{T})_{n\times\mathbf{m}} + (\mathbf{T})_{o\times\mathbf{m}}\right)(\Psi)_{\mathbf{m}\times\mathbf{m}}^{-1}(\Sigma)_{\mathbf{m}\times\mathbf{m}}^{-1}\left((\mathbf{T})_{n\times\mathbf{m}}^{\mathsf{T}} + (\mathbf{T})_{o\times\mathbf{m}}^{\mathsf{T}}\right)}{2}\right) \tag{16}$$

and

$$\Psi := \Theta \left(\Lambda^{-2} + \mathbf{1}\right)^{-1} \left(2\Lambda^{-2} + \mathbf{1}\right) \Theta^{\mathsf{T}} \tag{17}$$

The proportion of response variance ascribable to the first $m$ basis directions of $\mathbf{u}$ is given by the Sobol' index

$$S_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}}) := D_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}})/D_{\mathbf{M}}(\Theta) \leq S_{\mathbf{M}}(\Theta) = 1 \tag{18}$$

Analytic expressions for $\partial_\Theta D_\mathbf{m}((\Theta)_{\mathbf{m} \times \mathbf{M}})$ have been obtained from Eq. (15) using standard formulae for differentating matrix inverses and determinants. As $D_\mathbf{m}$ projects $M$-dimensional $(\mathbf{x})_\mathbf{M}$ onto $m$-dimensional $(\mathbf{u})_\mathbf{M}$, the result is affected by just a few components of rotation:

$$(\partial_\Theta D_\mathbf{m}((\Theta)_{\mathbf{m} \times \mathbf{M}}))_{i \times j} \neq 0 \quad \Longrightarrow \quad i \leq m < M \tag{19}$$

In particular $D_\mathbf{M}(\Theta) = D_\mathbf{M}$ and $S_\mathbf{M}(\Theta) = 1$ are independent of $\Theta$, as there is no projection, only rotation, in transforming $(\mathbf{x})_\mathbf{M}$ into $(\mathbf{u})_\mathbf{M}$.

*2.3. Basis Optimization*

At this point in the analysis, everything has been fixed save the rotation

$$\mathbf{u} := \Theta \mathbf{x} \tag{20}$$

relating sampling distribution $\mathbf{u} \sim \mathsf{N}[(\mathbf{0})_\mathbf{M}, (\mathbf{1})_{\mathbf{M} \times \mathbf{M}}]$ to the input of the surrogate response $\bar{f}(\mathbf{x})$. This rotation will now be determined by maximizing the relevance – as measured by Sobol' index – of each $\mathbf{u}$-direction in turn. This means optimizing $\Theta$ in Eq. (20) row by row from top to bottom.

Row orthonormality leaves just $(M - m - 1)$ elements free in row $m$, which we encode as

$$(\Theta)_{\mathbf{m} \times \mathbf{M}} =: (\Xi)_{\mathbf{m} \times \mathbf{M}} \tilde{\Theta} \tag{21}$$

where $\Xi$ is orthogonal, and identical on the $(m - 1)$ rows already optimized

$$
\begin{aligned}
(\Xi)_{\mathbf{m} \backslash \{m\} \times \mathbf{M}} &= \mathbf{1}_{\mathbf{m} \backslash \{m\} \times \mathbf{M}} \\
(\Xi)_{m \times \mathbf{m} \backslash \{m\}} &= (\mathbf{0})_{1 \times \mathbf{m} \backslash \{m\}} \\
(\Xi)_{m \times m} &= \left(1 - \sum_{k=m+1}^{M} (\Xi)_{m \times k}^2 \right)^{1/2}
\end{aligned}
\tag{22}
$$

The last line induces a derivative adjustment

$$\frac{\partial}{\partial (\Xi)_{m \times k}} = \frac{\partial}{\partial (\Xi)_{m \times k}} - \frac{(\Xi)_{m \times k}}{(\Xi)_{m \times m}} \frac{\partial}{\partial (\Xi)_{m \times m}} \tag{23}$$

which should be exploited by the optimizer as a powerful repellant to orthonormality violations. This work uses a BFGS optimizer, fed an analytic Jacobian.

7

Given these constraints, row $m$ is optimally determined by

$$(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}} = \operatorname*{argmax}_{(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}}} S_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}}) = \operatorname*{argmax}_{(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}}} D_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}}) \qquad (24)$$

The optimal row $m$ is then incorporated in $\tilde{\Theta}$ and $\Xi$ according to

$$\tilde{\Theta} = \mathbf{Q}^{\mathsf{T}} \text{ where } \Theta^{\mathsf{T}} = \mathbf{Q}\mathbf{R} \text{ is the QR factorization of the update}$$
$$(\Xi)_{\mathbf{m}\times\mathbf{M}} = \mathbf{1}_{\mathbf{m}\times\mathbf{M}} \qquad (25)$$

ready to optimize row $m + 1$. Optimization followed by incorporation is performed for $m = 1, \ldots, M-1$ in turn to entirely optimize $\Theta$. The later rows could be left unoptimized, though they are successively cheaper to obtain.

*2.4. Summary*

The main loop of the optimizations reported in the next Section is:

---

1: **repeat**
2:     Fit GP surrogate to $y(\mathbf{X}^{\mathsf{T}})$, determining $\bar{f}(\mathbf{x})$ according to Section 2.1
3:     Set $\tilde{\Theta} \leftarrow \Theta \leftarrow \Theta_{\Pi} \leftarrow \mathbf{1}$
4:     **for** $m = 1$ **to** $M$ **do**
5:         According to Section 2.2, optimize

$$(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}} \leftarrow \operatorname*{argmax}_{(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}}} D_{\mathbf{m}}((\Theta)_{\mathbf{m}\times\mathbf{M}})$$

        where $(\Theta)_{\mathbf{m}\times\mathbf{M}} =: (\Xi)_{\mathbf{m}\times\mathbf{M}}\tilde{\Theta}$
6:         Update $\tilde{\Theta} \leftarrow \mathbf{Q}^{\mathsf{T}}$ where $\Theta^{\mathsf{T}} = \mathbf{Q}\mathbf{R}$
7:     **end for**
8:     Update the input basis to $\mathbf{X}^{\mathsf{T}} \leftarrow \Theta\mathbf{X}^{\mathsf{T}}$
9:     Update the overall rotation to $\Theta_{\Pi} \leftarrow \Theta\Theta_{\Pi}$
10: **until** $\Theta \approx \mathbf{1}$

---

The optimization in Step 5 is prone to fall into local optima, especially in the first iteration or two of the outermost loop. It is recommended that these early iterations explore the behaviour of $(\Xi)_{m\times\mathbf{M}\backslash\mathbf{m}}$ (by grid or randomized search) before attempting to exploit it (by gradient descent).

As the input basis is updated at each step, $\mathbf{X} = \mathbf{U}$ ultimately. The key output of the algorithm is the overall rotation $\Theta_{\Pi}$ of the original basis for $\mathbf{x}$ to the optimal basis for $\mathbf{u}$.

## 3. Results

The algorithm described in Section 2.4 has been tested on a suite of test functions, using $N \in 100, 200, 400, 800, 1600$ data of $M = 5$ input dimensions. All quoted results are the mean over two folds (each with $N$ training data and $N$ test data). In each case an $N \times M$ design matrix $\mathbf{X}$ is sampled from a standard normal distribution. The input to the test function $f \colon [x_-, x_+]^M \to \mathbb{R}$ is generally constructed as

$$\hat{\mathbf{X}}^\intercal = (x_+ - x_-)c(\Phi\mathbf{X}^\intercal) + x_-(\mathbf{1})_{\mathbf{M} \times \mathbf{N}} \tag{26}$$

where $c \colon \mathbb{R}^M \to \mathbb{R}^M$ is the cumulative density function for $M$ independent standard normal random variables, and $\Phi$ is a test rotation matrix. The corresponding optimal input rotation from Section 2.4 is

$$\Theta_\Pi = \begin{cases} \Theta_1 & \text{if } \Phi \text{ is identity matrix } \mathbf{1} \\ \Theta_\mathbf{R} & \text{if } \Phi \text{ is a random rotation matrix } \Phi_\mathbf{R} \end{cases} \tag{27}$$

which should recover the random rotation as

$$\Theta_\mathbf{R} \cong \Theta_1 \Phi_\mathbf{R} \tag{28}$$

However, this is congruence, not equality. For each test function, the intial GP fit is assessed by test statistics from independent data (from the other fold), together with errors in the calculated Sobol' indices. The latter are important as they at the heart of subsequent calculations. The input basis is then optimized, calculating $\Theta_1$. A reduced dimensionality $\underline{M}$ for the optimized basis is determined as

$$\min \left\{ \underline{M} \leq M \mid S_{\underline{\mathbf{M}}} \geq 0.90 \right\} \tag{29}$$

A GP is fit to this reduced input, and its test statistics compared with the initial GP.

The whole procedure is then repeated (with entirely fresh data) to which a random input rotation $\Phi_\mathbf{R}$ is applied. The input basis is optimized, calculating $\Theta_\mathbf{R}$, whereas the reduced dimensionality $\underline{M}$ is kept from the unrotated analysis.

Finally the rotated and unrotated versions are compared for congruence, essentially we wish

$$|\epsilon_\Theta| = \left\|(\Phi_\mathbf{R}\Theta_\mathbf{1}\Theta_\mathbf{R}^\intercal - \mathbf{1})_{\underline{\mathbf{M}}\times\underline{\mathbf{M}}}\right\| / \underline{M} \tag{30}$$

where $\|\cdot\|$ is the usual Frobenius matrix norm and $\underline{\mathbf{M}}$ includes no irrelevant input directions (which might be rotated versions of each other). All inputs and outputs are standardized to a mean of zero and variance of 1. All functions are tested with and without output noise $e \sim \mathsf{N}[0, 0.001]$ added to $f(\hat{\mathbf{X}}^\intercal)$. In the noiseless cases $e = 0$ we floor the Gaussian process noise variance as $\epsilon_\mathbf{e} \geq 10^{-6}$ in order to prevent numerical instability in the calculation of Sobol' indices which is otherwise observed.

*3.1. Sine Function*

$$f(\hat{\mathbf{x}}) := \sin(\hat{\mathbf{x}}_1) \tag{31}$$

$$[x_-, x_+] := [-\pi, +\pi]$$
$$S_\mathbf{1} := S_{(1)} = 1$$

Fitting an initial GP recovers the exact Sobol' indices to within an accuracy of 0.005 (precision actually decreasing with the number of data). Optimizing the input basis improves this to $10^{-8}$, which can only be due to repeating GP regression as 4 iterations were run to optimize $\Theta$, even though convergence is immediate. Applying a rand[h]om rotation $\Phi_\mathbf{R}$, the exact Sobol' indices are recovered to within $10^{-8}$ after 3 iterations. The 1D active subspace measures are recorded in

| Noise | $N$ | $|\underline{\mathbf{u}}_1|$ |
|---|---|---|
| 0.0000 | 100 | 0.9993 |
| 0.0010 | 100 | 0.9961 |
| 0.0000 | 200 | 0.7310 |
| 0.0010 | 200 | 0.9997 |
| 0.0000 | 400 | 1.0000 |
| 0.0010 | 400 | 0.9999 |
| 0.0000 | 800 | 0.8919 |
| 0.0010 | 800 | 0.3323 |
| 0.0000 | 1600 | 1.0000 |
| 0.0010 | 1600 | 0.9999 |

| Noise | N | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|---|
| 0.0000 | 100 | 0.9993 | 0.9993 | 1.0000 | 1.0000 | 1.0000 |
| 0.0010 | 100 | 0.9961 | 0.9961 | 1.0000 | 1.0000 | 1.0000 |
| 0.0000 | 200 | 0.7310 | 0.7310 | 1.0000 | 1.0000 | 1.0000 |
| 0.0010 | 200 | 0.9997 | 0.9997 | 1.0000 | 1.0000 | 1.0000 |
| 0.0000 | 400 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 0.0010 | 400 | 0.9999 | 0.9999 | 1.0000 | 1.0000 | 1.0000 |
| 0.0000 | 800 | 0.8919 | 0.8919 | 1.0000 | 1.0000 | 1.0000 |
| 0.0010 | 800 | 0.3323 | 0.3323 | 1.0000 | 1.0000 | 1.0000 |
| 0.0000 | 1600 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 0.0010 | 1600 | 0.9999 | 0.9999 | 1.0000 | 1.0000 | 1.0000 |

*3.2. Decoupled Ishigami Function*

$$f(\mathbf{x}) := \left(1 + b\mathbf{x}_3^4\right)\sin(\mathbf{x}_1) + a\sin^2(\mathbf{x}_2) \tag{32}$$

$$a = 2.0 \quad ; \quad b = 0$$
$$S_1 = 0.5 \quad ; \quad S_2 = 1$$

*3.3. Ishigami Function*

$$f(\mathbf{x}) := \left(1 + b\mathbf{x}_3^4\right)\sin(\mathbf{x}_1) + a\sin^2(\mathbf{x}_2) \tag{33}$$

$$a = 7.0 \quad ; \quad b = 0.1$$
$$S_1 = 0.3139 \quad ; \quad S_2 = 0.7563 \quad ; \quad S_3 = 1$$

*3.4. Sobol' G Function*

$$f(\mathbf{x}) := \prod_{i=1}^{D} \frac{|4\mathbf{x}_i - 2| + \mathbf{a}_i}{1 + \mathbf{a}_i} \tag{34}$$

$$\mathbf{a}_i = (i-1)/2$$
$$S_1 = 0.4107 \quad ; \quad S_2 = 0.6541 \quad ; \quad S_3 = 0.8113 \quad ; \quad S_4 = 0.9203 \quad ; \quad S_5 = 1$$

## 4. Discussion

## References

[1] D. Wipf, S. Nagarajan, A new view of automatic relevance determination, in: Proceedings of the 20th International Conference on Neural Information Processing Systems, NIPS'07, Curran Associates Inc., USA, 2007, pp. 1625–1632.
URL http://dl.acm.org/citation.cfm?id=2981562.2981766

[2] R. M. Neal, Bayesian Learning for Neural Networks, Springer New York, 1996. doi:10.1007/978-1-4612-0745-0.