

# Multi-fidelity optimization via surrogate modelling

BY ALEXANDER I. J. FORRESTER\*, ANDRÁS SÓBESTER AND  
ANDY J. KEANE

*Computational Engineering and Design Group, School of Engineering Sciences,  
University of Southampton, Southampton SO17 1BJ, UK*

This paper demonstrates the application of correlated Gaussian process based approximations to optimization where multiple levels of analysis are available, using an extension to the geostatistical method of *co-kriging*. An exchange algorithm is used to choose which points of the search space to sample within each level of analysis. The derivation of the co-kriging equations is presented in an intuitive manner, along with a new variance estimator to account for varying degrees of computational ‘noise’ in the multiple levels of analysis. A multi-fidelity wing optimization is used to demonstrate the methodology.

**Keywords:** co-kriging; kriging; noise; subset selection; wing design

## 1. Introduction

Let  $f(\mathbf{x})$  denote an *objective function* that maps a vector  $\mathbf{x}$  of length  $k$  to some scalar measure of merit.  $\mathbf{x}^*$  is sought, such that  $f(\mathbf{x}^*)$  is the maximum or minimum of  $f(\mathbf{x})$ . Such *optimization* problems (of which this is the simplest form) are encountered in most branches of science. Of interest to us in this paper is the class of optimization problems where  $f(\mathbf{x})$  is expensive to obtain or we know its value only at a limited number of  $\mathbf{x}$  values. Thus, efficiency, in terms of the number of  $\mathbf{x}$  values that must be mapped to their objective function values before we get to within a reasonable distance of  $\mathbf{x}^*$ , is the key feature of the optimization algorithms considered.

While there is often no getting around the fact that a thorough search of a highly nonlinear, multi-dimensional  $f$  landscape will require sampling at a large number of sites, a useful shortcut is often employed. Based on a relatively small number of measurements, we can build a statistical approximation of the objective landscape, which, provided  $f$  is smooth and continuous and the measurements are reasonably uniformly spread, will be accurate enough to guide the search towards promising areas of the landscape.

If, at the same time, this *surrogate* of the expensive function is very cheap to evaluate, we can be as thorough as we like in searching it. Thus, most of the computational effort will be concentrated on what looks likely to be the

\* Author for correspondence (alexander.forrester@soton.ac.uk).

neighbourhood of the global optimum. An approximation technique that has received much attention in recent years is *kriging*<sup>1</sup>—we shall delve into the details of this later on, as our chosen multi-fidelity optimization technique is an extension of kriging. But what exactly do we mean by multi-fidelity optimization?

The objective function evaluation system may feature lower-fidelity, cheaper models (in addition to the main, high-fidelity calculation of  $f$ ). In most cases, the fast (but less trustworthy) and the slow (but more accurate) objective values can be obtained independently. Thus, we can learn more about the objective function by additionally measuring the cheap function(s) on a large number of  $\mathbf{x}$  values. This process is usually referred to as multi-fidelity optimization.

The use of secondary, correlated quantities to improve the accuracy of the model of the primary objective is not a new concept. For example, Hevesi *et al.* (1992) predict average annual precipitation values near a potential nuclear waste disposal site using a sparse set of precipitation measurements from the region along with the correlated and more easily obtainable elevation map of the area. Kennedy & O'Hagan (2000) approach the subject from the perspective of model building using objectives resulting from computational simulations of varying fidelities and costs. This is also the context of the work presented here; additionally, we introduce a method of optimization using correlated models of different costs and fidelities. We focus on combining co-kriging (the multi-response extension to kriging alluded to earlier) with a Bayesian model update criterion designed to balance the exploration of the search space and the quick exploitation of promising basins of attraction on the  $f$  landscape. The criterion is based on a newly derived error estimate that reflects the presence of noise in the observed data.

After a brief overview of kriging in §2, we discuss co-kriging in §3, followed by considerations on how to create sampling plans for multiple levels of analyses (§4). We put these pieces together in §5, where we describe an optimization strategy designed for multi-fidelity inputs. An aircraft wing design problem is used to demonstrate and compare the co-kriging approach and other multi-level strategies in §6. Section 7 discusses the remaining issues related to the use of correlated surrogates and summarizes our conclusions.

## 2. Kriging

We will consider co-kriging as a natural extension to the popular method of kriging and hence begin with a very brief introduction to the kriging method in order to introduce concepts and notation that follow through to co-kriging. Equations are shown without derivations, which are very similar to those presented later for co-kriging; the reader may wish to consult Jones (2001) or Forrester *et al.* (2006*b*) for more information on kriging.

As with all surrogate-based methods, to approximate a function  $f$  we start with a set of sample data—usually computed at a set of points in the domain of interest determined by a sampling plan (which will be discussed further in §4).

<sup>1</sup>Originally called *krigeage* by Matheron (1963) after D. G. Krige—a South African mining engineer who pioneered the method in the 1950s for determining ore-grade distributions based on core samples (Krige 1951).

The kriging prediction of  $f$  is built from a mean base term,  $\hat{\mu}$  (the circumflex denotes a maximum likelihood estimate, MLE) plus a stationary Gaussian process,  $Z(\mathbf{x})$ , representing the local features of  $f$  around the  $n$  sample points,  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}^T$ ,  $\mathbf{x} \in \mathbb{R}^k$ .  $Z(\mathbf{x})$  has zero mean and covariance

$$\text{cov}[Z(\mathbf{x}^{(n+1)}), Z(\mathbf{x}^{(i)})] = \sigma^2 \psi^{(i)}, \tag{2.1}$$

where  $\psi^{(i)}$  are correlations between a random variable  $Y(\mathbf{x})$  at the point to be predicted ( $\mathbf{x}^{(n+1)}$ ) and at the sample data points

$$\psi^{(i)} = \text{corr}[Y(\mathbf{x}^{(n+1)}), Y(\mathbf{x}^{(i)})] = \exp\left(-\sum_{j=1}^k \hat{\theta}_j \|x_j^{(n+1)} - x_j^{(i)}\|^{\hat{p}_j}\right). \tag{2.2}$$

The *hyper-parameter*,  $p_j$ , can be thought of as determining the smoothness of the function approximation. In many situations, we can assume that there will not be any discontinuities and use  $p_j=2$  rather than using an MLE. This means that the basis function is infinitely differentiable through a sample point (when  $\|\mathbf{x}^{(n+1)} - \mathbf{x}^{(i)}\| = 0$ ) and that the function is in the same form as a Gaussian distribution with variance  $1/\hat{\theta}_j$ .  $\hat{\theta}$  therefore can be thought of as determining how quickly the function changes as  $\mathbf{x}^{(n+1)}$  moves away from  $\mathbf{x}^{(i)}$ , with high and low  $\hat{\theta}_j$  values indicating an active or inactive function along dimension  $j$ , respectively. The kriging prediction is found as the value at  $\mathbf{x}^{(n+1)}$  that, maximizes the likelihood, given the sample data and MLEs of the hyper-parameters, and is given by

$$\hat{y}(\mathbf{x}^{(n+1)}) = \hat{\mu} + \sum_{i=1}^n b^{(i)} \psi^{(i)}(\mathbf{x}^{(n+1)}, \mathbf{x}^{(i)}). \tag{2.3}$$

The constants  $b_i$  are given by the column vector  $\mathbf{b} = \Psi^{-1}(\mathbf{y} - \mathbf{1}\hat{\mu})$ , where  $\Psi$  is an  $n \times n$  symmetric matrix of correlations between the sample data;  $\mathbf{y}$  is a column vector of responses,  $\{y(\mathbf{x}^{(1)}), \dots, y(\mathbf{x}^{(n)})\}^T$ ;  $\mathbf{1}$  is an  $n \times 1$  column vector of ones; and the MLE  $\hat{\mu} = \mathbf{1}^T \Psi^{-1} \mathbf{y} / \mathbf{1}^T \Psi^{-1} \mathbf{1}$ .

### 3. Co-kriging

We now consider how to build an approximation of a function that is expensive to evaluate which is enhanced by data from cheaper analyses of the function. Combining multiple sets of data naturally leads to a complex notation and we will try to simplify this by limiting ourselves to two datasets:<sup>2</sup> the most accurate expensive data with values  $\mathbf{y}_e$  at points  $\mathbf{X}_e$  and the less accurate cheap data  $\mathbf{y}_c$  at points  $\mathbf{X}_c$  ( $\mathbf{X}_e \subset \mathbf{X}_c$ ).<sup>3</sup> These data are concatenated to give the combined set of points

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_c \\ \mathbf{X}_e \end{pmatrix} = (\mathbf{x}_c^{(1)}, \dots, \mathbf{x}_c^{(n_c)}, \mathbf{x}_e^{(1)}, \dots, \mathbf{x}_e^{(n_e)})^T.$$

<sup>2</sup>Our methods can be extended to multiple code levels following the notation used by Kennedy & O’Hagan (2000).

<sup>3</sup>We require partially collocated points in our estimation of a scaling parameter between the data. It is possible to construct the presented co-kriging model with completely non-collocated points by using kriging estimates  $\hat{\mathbf{y}}_c(\mathbf{X}_e)$ .

As with kriging, the value at a point in  $\mathbf{X}$  is treated as if it were the realization of a Gaussian random variable. For co-kriging, we therefore have the random field

$$\mathbf{Y} = \begin{pmatrix} \mathbf{Y}_c(\mathbf{X}_c) \\ \mathbf{Y}_e(\mathbf{X}_e) \end{pmatrix} = (Y_c(\mathbf{x}_c^{(1)}), \dots, Y_c(\mathbf{x}_c^{(n_c)}), Y_e(\mathbf{x}_e^{(1)}), \dots, Y_e(\mathbf{x}_e^{(n_e)}))^T.$$

We follow the auto-regressive model of Kennedy & O’Hagan (2000), which assumes that  $\text{cov}\{Y_e(\mathbf{x}^{(i)}), Y_c(\mathbf{x})|Y_c(\mathbf{x}^{(i)})\} = 0, \forall \mathbf{x} \neq \mathbf{x}^{(i)}$ . This means that no more can be learnt about  $Y_e(\mathbf{x}^{(i)})$  from the cheaper code if the value of the expensive function at  $\mathbf{x}^{(i)}$  is known—a Markov property (i.e. we assume that the expensive simulation is correct and any inaccuracies lie wholly in the cheaper simulation). Gaussian processes  $Z_c(\cdot)$  and  $Z_e(\cdot)$  represent the local features of the cheap and expensive codes. Using the auto-regressive model, we are essentially approximating the expensive code as the cheap code multiplied by a scaling factor  $\rho$  plus a Gaussian process  $Z_d(\cdot)$  that represents the difference between  $\rho Z_c(\cdot)$  and  $Z_e(\cdot)$

$$Z_e(\mathbf{x}) = \rho Z_c(\mathbf{x}) + Z_d(\mathbf{x}). \tag{3.1}$$

Whereas in kriging we have a covariance matrix  $\text{cov}\{\mathbf{Y}(\mathbf{X}), \mathbf{Y}(\mathbf{X})\} = \sigma^2 \boldsymbol{\Psi}(\mathbf{X}, \mathbf{X})$ , we now have a more complex covariance matrix constructed as follows:

$$\begin{aligned} \text{cov}\{\mathbf{Y}_c(\mathbf{X}_c), \mathbf{Y}_c(\mathbf{X}_c)\} &= \text{cov}\{Z_c(\mathbf{X}_c), Z_c(\mathbf{X}_c)\} \\ &= \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_c), \\ \text{cov}\{\mathbf{Y}_e(\mathbf{X}_e), \mathbf{Y}_c(\mathbf{X}_c)\} &= \text{cov}\{\rho Z_c(\mathbf{X}_c) + Z_d(\mathbf{X}_c), Z_c(\mathbf{X}_e)\} \\ &= \rho \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_e), \\ \text{cov}\{\mathbf{Y}_e(\mathbf{X}_e), \mathbf{Y}_e(\mathbf{X}_e)\} &= \text{cov}\{\rho Z_c(\mathbf{X}_e) + Z_d(\mathbf{X}_e), \rho Z_c(\mathbf{X}_e) + Z_d(\mathbf{X}_e)\} \\ &= \rho^2 \text{cov}\{Z_c(\mathbf{X}_e), Z_c(\mathbf{X}_e)\} + \text{cov}\{Z_d(\mathbf{X}_e), Z_d(\mathbf{X}_e)\} \\ &= \rho^2 \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_e) + \sigma_d^2 \boldsymbol{\Psi}_d(\mathbf{X}_e, \mathbf{X}_e). \end{aligned}$$

The notation  $\boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_c)$ , for example, denotes a matrix of correlations of the form  $\psi_c$  between the data  $\mathbf{X}_e$  and  $\mathbf{X}_c$ . Our complete covariance matrix is thus

$$\mathbf{C} = \begin{pmatrix} \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_c) & \rho \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_e) \\ \rho \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_c) & \rho^2 \sigma_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_e) + \sigma_d^2 \boldsymbol{\Psi}_d(\mathbf{X}_e, \mathbf{X}_e) \end{pmatrix}. \tag{3.2}$$

The correlations are of the same form as equation (2.2), but there are two correlations,  $\psi_c$  and  $\psi_d$ , and we therefore have more hyper-parameters to estimate:  $\boldsymbol{\theta}_c$ ;  $\boldsymbol{\theta}_d$ ;  $\mathbf{p}_c$ ;  $\mathbf{p}_d$ ; and the scaling parameter  $\rho$ . Our cheap data are considered to be independent of the expensive data and we can find MLEs for  $\mu_c$ ,  $\sigma_c^2$ ,  $\boldsymbol{\theta}_c$  and  $\mathbf{p}_c$  by maximizing the ln-likelihood (ignoring constant terms)

$$-\frac{n_c}{2} \ln(\sigma_c^2) - \frac{1}{2} \ln|\det(\boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_c))| - \frac{(\mathbf{y}_c - \mathbf{1}\mu_c)^T \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_c)^{-1} (\mathbf{y}_c - \mathbf{1}\mu_c)}{2\sigma_c^2}. \tag{3.3}$$

By setting the derivatives of equation (3.3) w.r.t.  $\mu_c$  and  $\sigma_c^2$  to 0 and solving, we find MLEs of

$$\hat{\mu}_c = \mathbf{1}^T \Psi_c(\mathbf{X}_c, \mathbf{X}_c)^{-1} \mathbf{y}_c / \mathbf{1}^T \Psi_c(\mathbf{X}_c, \mathbf{X}_c)^{-1} \mathbf{1} \quad (3.4)$$

and

$$\hat{\sigma}_c^2 = (\mathbf{y}_c - \mathbf{1}\hat{\mu}_c)^T \Psi_c(\mathbf{X}_c, \mathbf{X}_c)^{-1} (\mathbf{y}_c - \mathbf{1}\hat{\mu}_c) / n_c. \quad (3.5)$$

Substituting equations (3.4) and (3.5) into (3.3) yields the concentrated ln-likelihood

$$-\frac{n_c}{2} \ln(\hat{\sigma}_c^2) - \frac{1}{2} \ln|\det(\Psi_c(\mathbf{X}_c, \mathbf{X}_c))| \quad (3.6)$$

and  $\hat{\theta}_c$  and  $\hat{\mathbf{p}}_c$  (if not set at  $\mathbf{2}$ ) are found by maximizing this equation.

To estimate the mean, variance, hyper-parameters and scaling parameter of the difference model ( $\mu_d$ ,  $\sigma_d^2$ ,  $\theta_d$ ,  $\mathbf{p}_d$  and  $\rho$ ), we first define

$$\mathbf{d} = \mathbf{y}_e - \rho \mathbf{y}_c(\mathbf{X}_e), \quad (3.7)$$

where  $\mathbf{y}_c(\mathbf{X}_e)$  are the values of  $\mathbf{y}_c$  at locations common to those of  $\mathbf{X}_e$  (the Markov property implies that we need to consider only this data). If  $\mathbf{y}_c$  is not available at  $\mathbf{X}_e$ , we may estimate  $\rho$  at little additional cost by using kriging estimates  $\hat{\mathbf{y}}_c(\mathbf{X}_e)$ , found from equation (2.3), using the already determined hyper-parameters  $\hat{\theta}_c$  and  $\hat{\mathbf{p}}_c$ . The ln-likelihood of  $\mathbf{d}$ , given  $\mathbf{y}_e$ , is now

$$\begin{aligned} & -\frac{n_e}{2} \ln(\sigma_d^2) - \frac{1}{2} \ln|\det(\Psi_d(\mathbf{X}_e, \mathbf{X}_e))| \\ & - \frac{(\mathbf{d} - \mathbf{1}\mu_d)^T \Psi_d(\mathbf{X}_e, \mathbf{X}_e)^{-1} (\mathbf{d} - \mathbf{1}\mu_d)}{2\sigma_d^2}, \end{aligned} \quad (3.8)$$

yielding MLEs of

$$\hat{\mu}_d = \mathbf{1}^T \Psi_d(\mathbf{X}_e, \mathbf{X}_e)^{-1} \mathbf{d} / \mathbf{1}^T \Psi_d(\mathbf{X}_e, \mathbf{X}_e)^{-1} \mathbf{1}$$

and

$$\hat{\sigma}_d^2 = (\mathbf{d} - \mathbf{1}\hat{\mu}_d)^T \Psi_d(\mathbf{X}_e, \mathbf{X}_e)^{-1} (\mathbf{d} - \mathbf{1}\hat{\mu}_d) / n_e, \quad (3.9)$$

with  $\hat{\theta}_d$ ,  $\hat{\mathbf{p}}_d$  (again, if not set at  $\mathbf{2}$ ) and  $\hat{\rho}$  found by maximizing

$$-\frac{n_e}{2} \ln(\hat{\sigma}_d^2) - \frac{1}{2} \ln|\det(\Psi_d(\mathbf{X}_e, \mathbf{X}_e))|. \quad (3.10)$$

Equations (3.6) and (3.10) must be maximized numerically using a suitable global search routine such as a genetic algorithm (GA). Depending upon the cost of evaluating the cheap and expensive functions  $f_c$  and  $f_e$ , for very high dimensional problems, the multiple matrix inversions involved in the likelihood maximization may render the use of the co-kriging model impractical (the size of the matrices depends directly on the quantities of data available, and the number of search steps needed in the MLE process is linked to the number of hyper-parameters being tuned). Typically, a statistical model used as a surrogate will be tuned many fewer times than the number of evaluations of  $f_e$  required by a direct search. The cost of tuning the model can therefore be allowed to exceed the cost of computing  $f_e$  and still provide significant speed-up. For large  $k$  and  $n$ ,

the time required to find MLEs can be reduced by using a constant  $\theta_{c,j}$  and  $\theta_{d,j}$  for all elements of  $\boldsymbol{\theta}_c$  and  $\boldsymbol{\theta}_d$  to simplify the maximization, though this may affect the accuracy of the approximation.

With the hyper-parameters estimated, the co-kriging prediction of the expensive function is given by

$$\hat{y}_e(\mathbf{x}^{(n+1)}) = \hat{\boldsymbol{\mu}} + \mathbf{c}^T \mathbf{C}^{-1}(\mathbf{y} - \mathbf{1}\hat{\boldsymbol{\mu}}), \tag{3.11}$$

where

$$\mathbf{c} = \begin{pmatrix} \hat{\rho}\hat{\sigma}_c^2\psi_c(\mathbf{X}_c, \mathbf{x}^{(n+1)}) \\ \hat{\rho}^2\hat{\sigma}_c^2\psi_c(\mathbf{X}_e, \mathbf{x}^{(n+1)}) + \hat{\sigma}_d^2\psi_d(\mathbf{X}_e, \mathbf{x}^{(n+1)}) \end{pmatrix}$$

and  $\hat{\boldsymbol{\mu}} = \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} / \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}$ . The notation  $\psi_c(\mathbf{X}_e, \mathbf{x}^{(n+1)})$ , for example, denotes a column vector of correlations of the form  $\psi_c$  between the data  $\mathbf{X}_e$  and the new point  $\mathbf{x}^{(n+1)}$ . The derivation of equation (3.11) is given in appendix A. If we make a prediction at one of our expensive points,  $\mathbf{x}^{(n+1)} = \mathbf{x}_e^{(i)}$  and  $\mathbf{c}$  is the  $n_c + i$ th column of  $\mathbf{C}$ , then  $\mathbf{c}^T \mathbf{C}^{-1}$  is the  $n_c + i$ th unit vector and  $\hat{y}_e(\mathbf{x}_e^{(i)}) = \hat{\boldsymbol{\mu}} + \mathbf{y}_{(n_c+i)} - \hat{\boldsymbol{\mu}} = y_e^{(i)}$ . We see, therefore, that equation (3.11) is an interpolator of the expensive data. If we make a prediction at one of the cheap points,  $\mathbf{x}^{(n+1)} = \mathbf{x}_c^{(i)}$ ,  $\mathbf{c}$  is not a column of  $\mathbf{C}$  and the predictor will in some sense regress  $\mathbf{y}_c$  unless it coincides with  $\mathbf{y}_e$ .

The estimated mean-squared error in the predictor is given as follows:

$$s^2(\mathbf{x}) = \hat{\rho}^2\hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}. \tag{3.12}$$

Again, a derivation is given in appendix A. Since the co-kriging model is an interpolator of  $\mathbf{y}_e$ , we expect the error to be zero at the expensive sample points. For  $\mathbf{x}^{(n+1)} = \mathbf{x}_e^{(i)}$ ,  $\mathbf{c}^T \mathbf{C}^{-1}$  is the  $n_c + i$ th unit vector,  $\mathbf{c}^T \mathbf{C}^{-1} \mathbf{c} = \mathbf{c}^{(n_c+i)} = \rho_c^2\sigma_c^2 + \sigma_d^2$  and hence  $s^2(\mathbf{x})$  is indeed 0. For  $\mathbf{X}_c \setminus \mathbf{X}_e$ ,  $s^2(\mathbf{x}) \neq 0$  unless  $\mathbf{y}_e = \mathbf{y}_c(\mathbf{X}_e)$ . The error at these points is determined by the character of  $\mathbf{Y}_d$ . If the difference between  $\rho \mathbf{Y}_c(\mathbf{X}_e)$  and  $\mathbf{Y}_e(\mathbf{X}_e)$  is simple (characterized by low  $\theta_{d,j}$  values), the error will be low, whereas a more complex difference (high  $\theta_{d,j}$  values) will lead to high error estimates.

(a) *One-variable demonstration*

We will now look at how co-kriging behaves using an example of a simple one-variable function. Imagine that our expensive to compute data are calculated by the function  $f_e(x) = (6x - 2)^2 \sin(12x - 4)$ ,  $x \in \{0, 1\}$ , and a cheaper estimate of this data is given by  $f_c(x) = Af_e + B(x - 0.5) - C$ . We sample the design space extensively using the cheap function at  $\mathbf{X}_c = \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\}$ , but run the expensive function only at four of these points,  $\mathbf{X}_e = \{0, 0.4, 0.6, 1\}$ .

Figure 1 shows the functions  $f_e$  and  $f_c$  with  $A=0.5$ ,  $B=10$  and  $C=-5$ . A kriging prediction through  $\mathbf{y}_e$  gives a poor approximation to the deliberately deceptive function, but the co-kriging prediction lies very close to  $f_e$ , being better than both the standard kriging model and the cheap data. Despite the considerable differences between  $f_e$  and  $f_c$ , a simple relationship has been found between the expensive and the cheap data and the estimated error reduces almost to 0 at  $\mathbf{X}_c$  (figure 2).

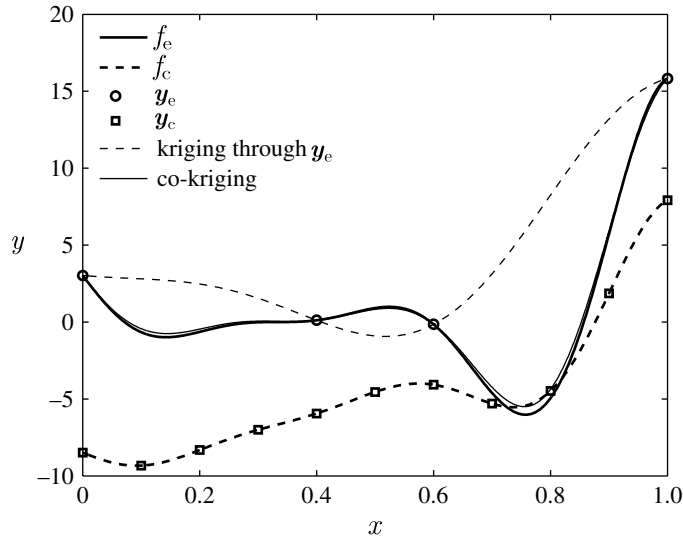


Figure 1. A one-variable co-kriging example. The kriging approximation using four expensive data points ( $\mathbf{y}_e$ ) has been significantly improved using extensive sampling from the cheap function ( $\mathbf{y}_c$ ).

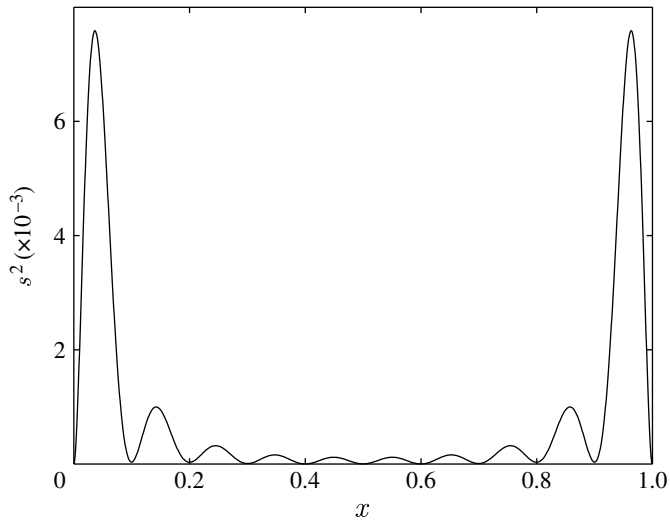


Figure 2. Estimated error in the co-kriging prediction in figure 1. The simple relationship between the data results in low error estimates at  $\mathbf{X}_c$  as well as  $\mathbf{X}_e$ .

We have chosen the relationship between our low- and high-fidelity test functions in order to show how the hyper-parameters of the co-kriging model behave. The hyper-parameters pertaining to the cheap data,  $\boldsymbol{\theta}_c$  and  $\mathbf{p}_c$ , are affected only by this data and behave as described in §2. Moving on to the scaling parameter  $\rho$ , if our cheap model parameter  $A$  (the multiplying term linking the cheap and expensive functions) is varied such that  $1/A \in \{-10, 10\}$ , we obtain the values for  $\hat{\rho}$  shown in figure 3 and see that  $\hat{\rho} \approx 1/A$ . Similar trials show that the parameters  $B$  and  $C$  have no effect on  $\rho$ ; thus we see that  $\rho$  is purely a scaling parameter. Note that  $\hat{\rho}$  is only an *indicator* of the scaling, since this value is estimated based on the data

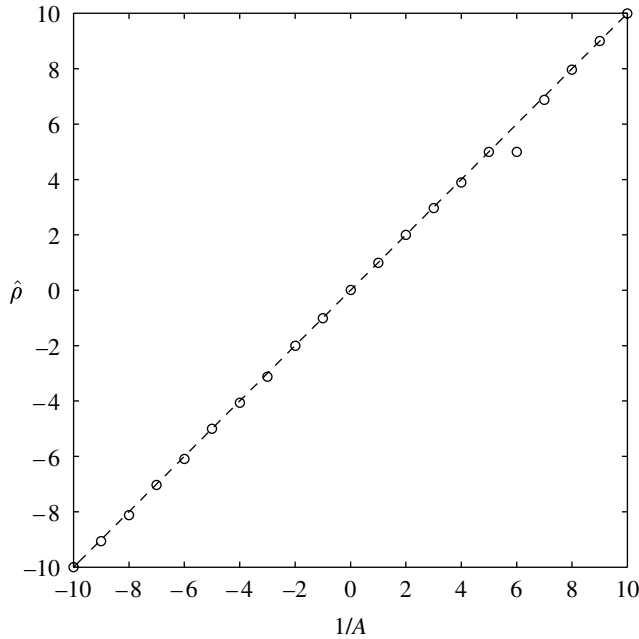


Figure 3. This plot of  $\hat{\rho}$  versus  $1/A$  shows that the MLE for  $\rho$  is a scaling factor between  $Z_c(\cdot)$  and  $Z_e(\cdot)$ , following the formulation in equation (3.1). There is a singularity at  $1/A=0$ , hence we have used  $1/A=0.01$  at this point.

available. For the data in figure 1,  $\hat{\rho} = 1.87$  (close to the true value of 2), but for small samples of  $\mathbf{y}_e$  the MLE may be misleading. The slight deviations of the data in figure 3 from  $\hat{\rho} = 1/A$  (shown as a dashed line), the most significant of which is at  $1/A=6$ , are where our GA search has not found the true MLE.

Recall that  $\mathbf{d} = \mathbf{y}_e - \rho \mathbf{y}_c(\mathbf{X}_e)$  (equation (3.7)) and hence, with  $\hat{\rho} \approx \mathbf{y}_e/\mathbf{y}_c$ ,  $\mathbf{d}$  represents the difference in trends between the cheap and the expensive data. Thus, for our one-variable example, when  $B, C=0$ ,  $\hat{\mu}_d, \hat{\sigma}_d^2 \rightarrow 0$  for all values of  $A$ , if  $\hat{\rho}$  is estimated accurately. Figure 4 shows how  $\hat{\sigma}_d^2$  varies for  $B \in \{-10, 10\}$  and we see that as  $B \rightarrow 0$  and therefore  $\mathbf{d} \rightarrow \mathbf{0}$ ,  $\hat{\sigma}_d^2$  also approaches 0. Note that  $\boldsymbol{\theta}_e$  and  $\hat{\boldsymbol{\rho}}_e$  will not be affected since the correlation in equation (2.2) is unaffected by the scaling of the objective data (it is, however, affected by the scaling of  $\mathbf{X}$ ).

Our choice of cheap function for the above example is somewhat contrived, but this has allowed us to show that the co-kriging model and its parameters are behaving as we would expect. For our test function, the correction process  $Z_d(\cdot)$  is linear. Co-kriging will work effectively for more complex correction processes with the proviso that  $Z_d(\cdot)$  must be simpler to model than  $Z_e(\cdot)$ . In §6 we demonstrate the benefits of co-kriging using an engineering design problem.

#### 4. Sampling plans

In §3 we have shown how to build the co-kriging model based on a set of sample data. We now consider how to choose the points  $\mathbf{X}_c$  and  $\mathbf{X}_e$  to give us the best prediction  $\hat{\mathbf{y}}_e$ .



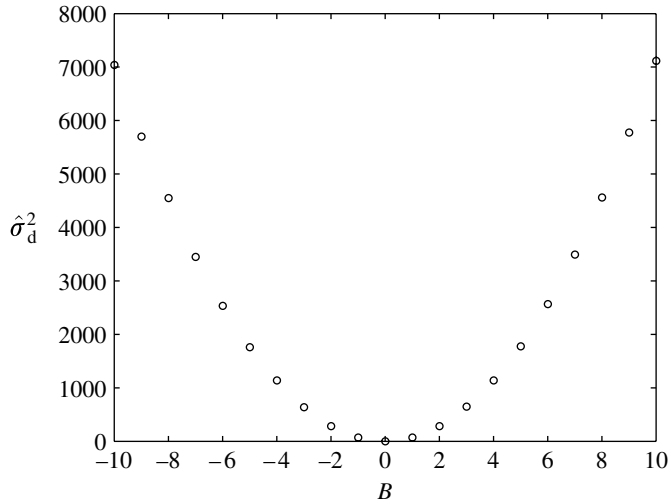


Figure 4. Variance  $\hat{\sigma}_d^2$  as the cheap function coefficient  $B$  is altered.  $\hat{\sigma}_d^2$  reduces to 0 as the difference between  $f_e$  and  $f_c$  can be modelled purely by the scaling parameter  $\rho$ .

The results of experiments, whether physical or computational, are corrupted by experimental error, that is, they deviate from the ‘true’ result. This can be due to human error, systematic deviations (due to a flaw or an inevitable physical or numerical limitation of the experiment) and, *only in physical experiments*, random error (usually linked to the limited accuracy of the instruments). When deciding on the sampling plan, that is, when choosing  $\mathbf{X}$ , there is nothing we can do about the first two error components and, while physical sampling plans often include replicates to reduce the random error, computer experiment plans are simply designed to cover the parameter space reasonably uniformly. The goal is, of course, to enable the model fitted to these points to give an accurate global approximation of the unknown objective function landscape.

The definition of ‘uniform coverage’ is by no means obvious and a substantial body of literature exists on the subject—here we work with the optimality criterion of Morris & Mitchell (1995). According to this, the plan with the best space-filling properties is that which maximizes the smallest distance between any pair of points within the sample (several additional ‘tie-breaker’ criteria are given in case of multiple optima). Additionally, we restrict our search to a class of plans known as Latin Hypercubes (LH; McKay *et al.* 1979), which are built to ensure uniformly spread projections of all points on all axes (figure 5a shows a random LH). We thus generate an initial Morris–Mitchell optimal LH plan  $\mathbf{X}_c$ .

A more interesting question is how do we select an  $n_e$ -element subset  $\mathbf{X}_e$  of  $\mathbf{X}_c$ , where the expensive simulations are to be run? Again, we wish to cover the parameter space evenly and hence turn to the Morris–Mitchell criterion, but this time we are dealing with a limited, discrete parameter space and thus the problem becomes a combinatorial one. Since selecting the subset that satisfies this is an *NP-complete* problem and an exhaustive search would have to examine  $n_c C_{n_e} = n_c! / n_e!(n_c - n_e)!$  subsets (clearly infeasible for all but very moderate cardinalities), here we use an exchange algorithm to select  $\mathbf{X}_e$  (Cook & Nachtsheim 1980).

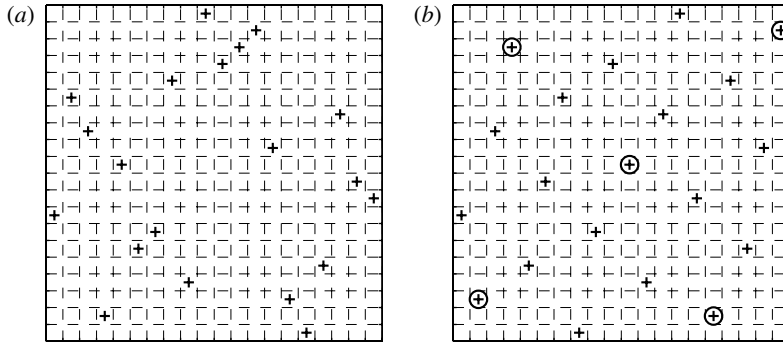


Figure 5. Two 20-point sampling plans: (a) a random LH and (b) a Morris–Mitchell optimal LH (plus symbol) with a 5-point subset found using the exchange algorithm (circle).

We start from a randomly selected subset  $\mathbf{X}_e$  and calculate the Morris–Mitchell criterion. We then exchange the first point  $\mathbf{x}_e^{(1)}$  with each of the remaining points in  $\mathbf{X}_c \setminus \mathbf{X}_e$  and retain the exchange that gives the best Morris–Mitchell criterion. This process is repeated for each remaining point  $\mathbf{x}_e^{(2)} \dots \mathbf{x}_e^{(n_e)}$ . A number of restarts from different initial subsets can be employed to avoid local optima. Figure 5b shows a Morris–Mitchell optimal LH with a subset chosen using this exchange algorithm.

A rule of thumb for the number of points that should be used in the sampling plan is  $n = 10k$ . When using a particularly cheap analysis,  $n_c$  may be rather greater than this, allowing us to build a more accurate model, and if the relationship between  $f_c$  and  $f_e$  is simple,  $n_e$  may be somewhat fewer—the advantage of the co-kriging method. We will return to this subject in §7.

### 5. Co-kriging based optimization

In order to confirm and enhance the predictions of a surrogate model, it is usual to update the model with new evaluations in promising areas. Co-kriging (and kriging) treats the value of the function at  $\mathbf{x}$  as if it were the realization of a Gaussian random variable  $Y(\mathbf{x})$ , with a probability density function

$$\frac{1}{\sqrt{2\pi}s(\mathbf{x})} \exp - \frac{1}{2} \left( \frac{Y_e(\mathbf{x}) - \hat{y}_e(\mathbf{x})}{s(\mathbf{x})} \right)^2,$$

with mean given by the predictor,  $\hat{y}_e(\mathbf{x})$  (equation (3.11)) and variance,  $s^2(\mathbf{x})$  (equation (3.12)). This allows us to model our uncertainty about the predictions we make. The most plausible value at  $\mathbf{x}$  is  $\hat{y}_e(\mathbf{x})$ , with the probability decreasing as  $Y_e(\mathbf{x})$  moves away from  $\hat{y}_e(\mathbf{x})$ . Since there is an uncertainty in the value of  $\hat{y}_e(\mathbf{x})$  we can calculate the expectation of it being an improvement,  $I = \min\{\mathbf{y}_e\} - Y(\mathbf{x})$ , on the best value calculated so far as

$$E[I(\mathbf{x})] = \int_{-\infty}^{\infty} \max\{\min\{\mathbf{y}_e\} - Y_e(\mathbf{x}), 0\} \phi(Y_e(\mathbf{x})) dY_e$$

$$= \begin{cases} (\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})) \Phi \left( \frac{\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})}{s(\mathbf{x})} \right) + s(\mathbf{x}) \phi \left( \frac{\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})}{s(\mathbf{x})} \right) & \text{if } s > 0, \\ 0 & \text{if } s = 0, \end{cases}$$

where  $\Phi(\cdot)$  and  $\phi(\cdot)$  are the normal cumulative distribution function and probability density function, respectively. By maximizing  $E[I(\mathbf{x})]$  we can find the best new point at which to sample the design space. Note that  $E[I(\mathbf{x})]=0$  when  $s(\mathbf{x})=0$  so that there is no expectation of improvement at a point that has already been sampled and therefore no possibility of re-sampling, which is a necessary characteristic of an updating criterion when using deterministic computer experiments, and guarantees global convergence; given that there is no possibility of re-sampling, as the number of updates based on the maximum  $E[I(\mathbf{x})]$  increases, the design space will become densely populated and so the global optimum will be found (Locatelli 1997).

(a) Co-kriging regression

When using multi-fidelity analyses, we can modify the interpolating co-kriging formulation in §3 such that each analysis can be regressed appropriately to filter any noise present in the data. Different levels of filtering may be required for each analysis. For example, data from an empirical low-fidelity computer code may be smooth and require no regression, but could be coupled to a discretized high-fidelity analysis that displays noise that must be filtered using regression. Conversely, a low-fidelity model may be discretized on a much coarser mesh, requiring a higher degree of regression than the high-fidelity code. In our regressing co-kriging formulation, we therefore employ two regression constants,  $\lambda_c$  and  $\lambda_e$ . These are added to the leading diagonal of the correlation matrices to give the covariance matrix  $\mathbf{C} + \lambda$ ,

$$\begin{pmatrix} \sigma_c^2\{\Psi_c(\mathbf{X}_c, \mathbf{X}_c) + \mathbf{I}_{(n_c \times n_c)}\lambda_c\} & \rho\sigma_c^2\{\Psi_c(\mathbf{X}_c, \mathbf{X}_e) + \begin{pmatrix} \mathbf{0}_{(n_c - n_e \times n_e)} \\ \mathbf{I}_{(n_e \times n_e)} \end{pmatrix}\lambda_c\} \\ \rho\sigma_c^2\{\Psi_c(\mathbf{X}_e, \mathbf{X}_c)\} & \rho^2\sigma_c^2\{\Psi_c(\mathbf{X}_e, \mathbf{X}_e) + \mathbf{I}_{(n_e \times n_e)}\lambda_c\} \\ +(\mathbf{0}_{(n_e \times n_e - n_e)}\mathbf{I}_{(n_e \times n_e)})\lambda_c\} & +\sigma_d^2\{\Psi_d(\mathbf{X}_e, \mathbf{X}_e) + \mathbf{I}_{(n_e \times n_e)}\lambda_e\} \end{pmatrix}, \quad (5.1)$$

where  $\mathbf{I}$  is an identity matrix and  $\mathbf{0}$  a zero matrix. The values of  $\lambda_c$  and  $\lambda_e$  are found along with the other hyper-parameters by maximizing equations (3.3) and (3.8). The regressing predictor is now

$$\hat{y}_{e,r}(\mathbf{x}^{(n_e+1)}) = \hat{\mu} + \mathbf{c}^T(\mathbf{C} + \lambda)^{-1}(\mathbf{y} - \mathbf{1}\hat{\mu}), \quad (5.2)$$

(subscript r denotes regression) with an estimated mean-squared error

$$s_r^2(\mathbf{x}) = \hat{\rho}^2\hat{\sigma}_c^2\hat{\lambda}_c + \hat{\sigma}_d^2\hat{\lambda}_e - \mathbf{c}^T(\mathbf{C} + \lambda)^{-1}\mathbf{c}. \quad (5.3)$$

Unlike equation (3.12), the regression error does not fall to 0 at sample points, and therefore Locatelli’s proof of convergence is invalid because  $\max\{E[I(\mathbf{x})]\}$  can occur at a previously sampled point. This is satisfactory when  $f_c$  and/or  $f_e$  are evaluated through physical experiments, since repeated experiments may be required to reduce measurement error. If  $f_c$  and  $f_e$  are found using deterministic computer experiments, re-sampling should be avoided and the convergence properties of  $\max\{E[I(\mathbf{x})]\}$  (or other error-based sampling strategies) can be restored by modifying the estimated error to eliminate the error due to the noise filtering, since this is the error associated with the data rather than the quality of the model. This is achieved by fitting a kriging interpolation through regressed data, re-interpolation.

## (b) Co-kriging re-interpolation

We now show the derivation of a co-kriging error estimate that reflects the uncertainty in predicting the underlying trend of data, rather than the overall uncertainty, which includes any error in noisy data. The derivation, although naturally complicated by multiple sets of data and regression constants, follows the same theme as for kriging re-interpolation (Forrester *et al.* 2006b). As mentioned above, to eliminate the error in the noisy data, we need to calculate the error of an interpolation through a regressed set of data. The column vectors of regressed cheap and expensive data are found from equations (2.3) (with the addition of the regression constant) and (5.2), and can be expressed as

$$\hat{\mathbf{y}}_{c,r} = \mathbf{1}\hat{\mu}_c + \Psi_c(\Psi_c + \hat{\lambda}_c \mathbf{I})^{-1}(\mathbf{y}_c - \mathbf{1}\hat{\mu}_c)$$

and

$$\hat{\mathbf{y}}_{e,r} = \mathbf{1}_{(n_e \times 1)}\hat{\mu} + \{\mathbf{c}(\mathbf{x}_e^{(1)})^T, \dots, \mathbf{c}(\mathbf{x}_e^{(n_e)})^T\}^T(\mathbf{C} + \lambda)^{-1}(\mathbf{y} - \mathbf{1}_{(n \times 1)}\hat{\mu}),$$

where the MLEs for the hyper-parameters are those found for the co-kriging model.

Note that only the  $\sigma^2$  terms in equation (3.12) depend on  $\mathbf{y}$ . We therefore need to substitute  $\hat{\mathbf{y}}_{c,r}$  into equation (3.5) and  $\hat{\mathbf{y}}_{e,r}$  into equation (3.9). Beginning with  $\hat{\sigma}_{c,ri}^2$ ,

$$\hat{\sigma}_{c,ri}^2 = \frac{(\mathbf{y}_c - \mathbf{1}\hat{\mu}_c)(\Psi_c + \lambda_c \mathbf{I})^{-1}\Psi_c(\Psi_c + \lambda_c \mathbf{I})^{-1}(\mathbf{y}_c - \mathbf{1}\hat{\mu}_c)}{n_c} \quad (5.4)$$

(subscript ri denotes re-interpolation). Note that if  $\lambda_c=0$ , then  $\hat{\sigma}_{c,ri}^2 = \hat{\sigma}_c^2$ . To find  $\hat{\sigma}_{d,ri}^2$ , we substitute  $\hat{\mathbf{y}}_{e,r}$  and  $\hat{\mathbf{y}}_{c,r}$  into equation (3.7) to give

$$\begin{aligned} \mathbf{d}_r &= \mathbf{1}_{(n_e \times 1)}\hat{\mu} + \{\mathbf{c}(\mathbf{x}_e^{(1)})^T, \dots, \mathbf{c}(\mathbf{x}_e^{(n_e)})^T\}^T(\mathbf{C} + \lambda)^{-1}(\mathbf{y} - \mathbf{1}_{(n \times 1)}\hat{\mu}) \\ &\quad - \hat{\rho}[\mathbf{1}\hat{\mu}_c + \Psi_c(\mathbf{X}_e)(\Psi_c(\mathbf{X}_e) + \lambda_c \mathbf{I})^{-1}(\mathbf{y}_c(\mathbf{X}_e) - \mathbf{1}\hat{\mu}_c)]. \end{aligned}$$

This is substituted into equation (3.9) in place of  $\mathbf{d}$  (again, note that if  $\lambda_e=0$  then  $\hat{\sigma}_{d,ri}^2 = \hat{\sigma}_d^2$ ). The expressions for  $\hat{\sigma}_{c,ri}^2$  and  $\hat{\sigma}_{d,ri}^2$  can now be substituted into equation (3.12) to find the re-interpolation error estimate that reduces to 0 at  $\mathbf{X}_e$ . With the errors due to noisy data removed, the estimated error of the predicted smooth trend is typically very small when compared with an interpolating model. This can lead to problems with floating-point underflow when calculating  $E[I(\mathbf{x})]$ . A scheme for avoiding this problem is presented in appendix B.

## 6. Example problem

We will demonstrate the benefits of co-kriging, and in particular regressing co-kriging, through the optimization of a generic transonic civil aircraft wing. The wing is defined by 11 variables: area; aspect ratio; kink position; sweep; inboard and outboard taper ratios; root, kink and tip thickness to chord ratio; tip washout; and kink washout. Our cheap analysis is in the form of an empirical drag estimation code, *Tadpole* (Cousin & Metcalf 1990). This returns a drag value based on curve fits to previously analysed wings in approximately 0.6 s. Our ‘expensive’ code is the linearized potential method *VSaero* (Maskew 1982)

with viscous coupling (approx. 2 min per drag evaluation). We aim to minimize drag/dynamic pressure ( $D/q$ ) for a fixed lift (determined by a wing weight calculated by Tadpole).

To allow visualization of the design landscape, we limit the search to the four variables that have the most impact on drag: area,  $S$ ; aspect ratio,  $\mathcal{R}$ ; sweep,  $A$ ; and inboard taper ratio,  $T_{\text{in}}$  (these variables were shown to be the most dominant by Keane (2003)). The Tadpole and VSAero design landscapes are shown using hierarchical axis plots in figures 6 and 7. With the fast empirical drag estimation provided by Tadpole, we have been able to build a high-resolution plot from  $11^4=14\,641$  runs of the code in 2.3 hours, while using the slower physics-based VSAero we have produced a low-resolution plot using  $3^2 \times 11^2=1089$  evaluations in 34.5 hours.<sup>4</sup> Each tile of the plots shows  $D/q$  for  $S \in \{150, 250\} m^2$  and  $\mathcal{R} \in \{6, 12\}$  for a fixed  $A$  and  $T_{\text{in}}$ .  $A$  and  $T_{\text{in}}$  vary from tile to tile, with the value at the lower left corner of the tile representing the value for the entire tile. The blank portions of figure 7 are regions where VSAero has failed to return a result for unusual geometries that lead to extreme flow regimes.<sup>5</sup> It is seen that the two codes produce results that follow the same general trend of lower  $D/q$  for higher  $\mathcal{R}$  and  $T_{\text{in}}$ , but the global optimum of the VSAero landscape goes against this general trend, with the lowest  $D/q$  at  $T_{\text{in}}=0.55$ .  $Z_d$  will not, therefore, be a trivial function as in our one-variable example but, since the cheap and expensive landscapes have the same general trend, should be simpler than  $Z_c$ .

For this problem, we start with an initial Morris–Mitchell sampling plan for  $\mathbf{X}_c$  of 100 points (note  $n_c > 10k$ ), to which a kriging model is fitted. A comparison of 100 further Tadpole calculations with predictions from the kriging model yields a high correlation coefficient of 0.96, indicating that the cheap code is approximated well. A subset of  $n_e=20$  points (note  $n_e < 10k$ ) is selected (using the exchange algorithm) from which to build the initial co-kriging model.  $\mathbf{X}_c$  and  $\mathbf{X}_e$  are shown in figure 6. Despite the apparent sparseness of the  $\mathbf{X}_e$  data, a good initial co-kriging prediction of the VSAero landscape is obtained, with a correlation coefficient of 0.96 when compared with the  $11^2 \times 3^2$  VSAero dataset. We follow an iterative process of updating the co-kriging model with new VSAero and Tadpole data at  $\max\{E[I(\mathbf{x})]\}$  and retuning hyper-parameters until the optimum is found. Since  $n_c$  is large and the correlation with independent data is high, we do not expect the Tadpole landscape to be altered unduly by updates. The GA-based tuning of  $\theta_c$  and  $\lambda_c$  requires a large number of computationally intensive inversions of the large  $n_c \times n_c$  covariance matrix. We therefore save time by limiting our re-tuning of the cheap hyper-parameters  $\theta_c$  and  $\lambda_c$  to every 10 updates, while we re-tune  $\theta_e$  and  $\lambda_e$  at each step.

The co-kriging method is compared with a  $\max\{E[I(\mathbf{x})]\}$  search of a kriging model built using only the VSAero data. A  $D/q$  that improves on the minimum from the previously computed  $3^2 \times 11^2$  set of data plus one standard deviation of the noise exhibited by the VSAero results (found from 10 VSAero evaluations of the same design with small perturbations),  $D/q=2.57+0.027$ , is used as a stopping

<sup>4</sup>Generating such plots would, of course, not be possible in most cases due to the high computational cost. Here we have computed this large quantity of data for illustrative purposes.

<sup>5</sup>In our co-kriging and kriging-based searches, we have directed the search away from regions of failure by inputting penalized values of  $\hat{y}_e(\mathbf{x}_e) + s^2(\mathbf{x}_e)$  when VSAero fails to return a result (Forrester *et al.* 2006c).

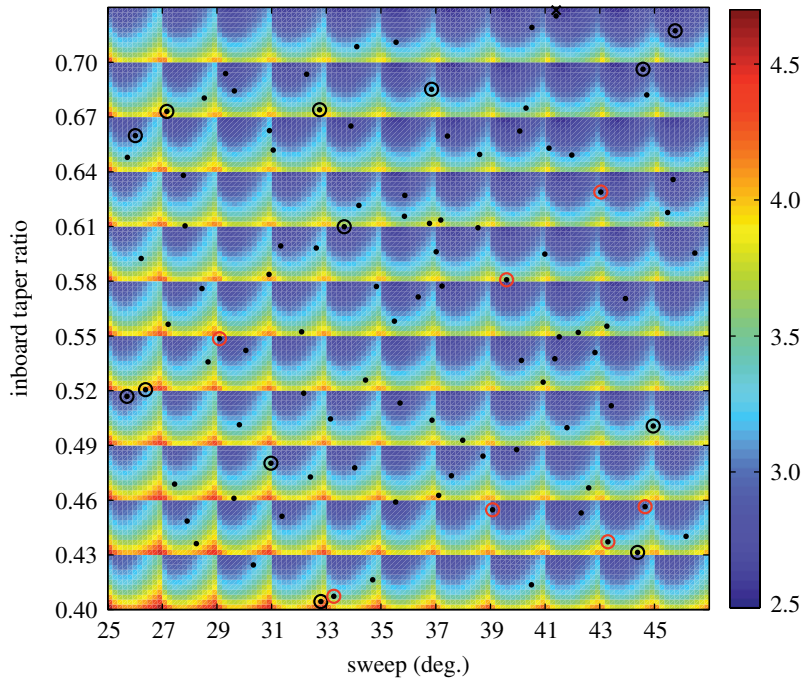


Figure 6. Tadpole calculated  $D/q$  with a typical sampling plan.  $\mathbf{X}_c$  (Tadpole evaluations) are shown as dots and are circled at locations of  $\mathbf{X}_e$  (VSaero evaluations). Red circles, failed VSaero simulations.

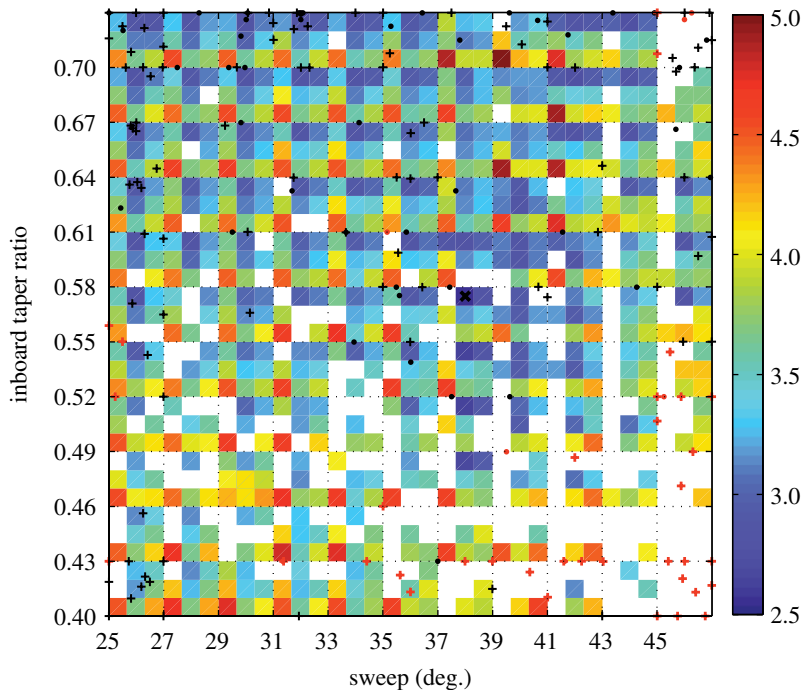


Figure 7. VSaero calculated  $D/q$  with co-kriging (dots) and kriging (crosses) updates from a typical search (the search that yielded values closest to the mean results in table 1).

Table 1. Performance comparison for the four-variable transonic wing problem. (Averaged over five searches.)

	initial model		best $D/q$ (m <sup>2</sup> )	function evaluations		
	$r^2$	RMSE	VSaero	Tadpole	succeeded VSaero	failed VSaero
kriging	0.949	0.155	2.621	0	119	40
co-kriging	0.962	0.108	2.556	185	79	27

criterion in both cases, i.e. similar quality answers are achieved from both the searches, allowing a direct comparison of the effort involved. We also limit the size of the kriging and co-kriging models to a maximum of 200 VSaero calculations. Beyond this, the time required for tuning the model becomes equivalent to the time required to run VSaero. Results, averaged over five searches from sampling plans produced using different random number seeds, are shown in [table 1](#).

The co-kriging based search consistently outperforms the kriging-based search; finding better optima for reduced numbers of VSaero evaluations and with fewer failed simulations. The initial prediction of the kriging model is almost as accurate as the co-kriging model (see correlation and RMSE in [table 1](#)), but the greater coverage of data in the co-kriging model leads to a better selection of successful update points in promising regions, as seen in [figure 7](#), which shows the distribution of update evaluations for a typical search. Moreover, the co-kriging updates are concentrated in regions of good designs, while the kriging updates are more widespread because, as there is a sparsity of data, there is high error and therefore high expectations of improvement in many areas (i.e. there is an emphasis on exploration over exploitation).

Only discrete values of  $\mathcal{A}$  and  $T_{\text{in}}$  can be displayed in [figure 7](#) and these design variables are rounded to the nearest  $2^\circ$  and 0.03 respectively, for the purpose of visualization. It should be borne in mind that there is, in fact, a distribution of designs between tiles. The reader should also appreciate that, although the major axis is much larger than the axes of the individual tiles, they both represent the same variation across a variable. Thus, although there seems to be a wide distribution of update points across  $\mathcal{A}$  and  $T_{\text{in}}$ , for the co-kriging updates these points represent tight clusters in regions of optimal feasible designs.

In §3 we discussed the use of two regression constants  $\lambda_c$  and  $\lambda_d$  in the co-kriging formulation. The final MLEs of these hyper-parameters for the search in [figure 7](#) were  $\hat{\lambda}_c = 1.2 \times 10^{-6}$  and  $\hat{\lambda}_d = 6.5 \times 10^{-3}$ . These values agree with our intuition that data from an empirical code (Tadpole) will be smooth and require little regression, i.e. a very small  $\hat{\lambda}_c$ , and data from a discretized physics-based code (VSaero) will be noisy and require smoothing, i.e. a larger  $\hat{\lambda}_d$ .

## 7. Discussion

Our results demonstrate that correlating analyses at multiple levels of fidelity can enhance the accuracy of a surrogate model of the highest level of analysis. This correlated model can be used to find optimal solutions more quickly. We have



presented a global optimization strategy using a co-kriging based method that allows for varying levels of noise filtering across multi-fidelity analyses and which converges towards the global optimum using expected improvement maximization.

We have shown a co-kriging based optimization using two levels of CFD code fidelity. The methods are extendable to multiple levels and there are many other examples of analyses that may be coupled. Here, both codes are predicting the same quantity (drag) but, as long as there is a useful correlation, the analyses may be related to completely different quantities as is the case with the work of Hevesi *et al.* (1992), which correlated precipitation with elevation: two different but obviously correlated quantities. Often only one form of analysis is available, although this may come in varying levels of fidelity. For example, Kennedy & O' Hagan (2000) used different finite-element mesh resolutions to produce cheap and expensive results. Forrester *et al.* (2006a) showed that partially converged CFD simulations correlate well with their converged counterparts. The convergence could be stopped at any number of points, thus providing many levels of fidelity.

In §4 we discussed the choice of  $n_c$  and  $n_e$  based on using more or fewer points than an  $n = 10k$  rule of thumb. In the example problem, we checked that the model built using our choices of  $n_c$  and  $n_e$  was accurate by assessing the correlation with a second set of data. Naturally the number of points required is problem dependent and additional validation data may be too expensive to compute. A possible way to reduce computational effort in situations when the cost of  $f_c$  is considerable is to start with a very low  $n_c$  and add data at points where the estimated error is maximized until an accurate cheap model is obtained. Model accuracy may be assessed using a leave-one-out cross-validation procedure. This method will produce a space-filling design with  $n_c$  controlled by the desired model accuracy. Maximum expected improvement updates *could* then proceed starting from an initial sample as small as  $n_e = 2$ , though it may be advantageous to follow a maximum error update strategy until an accurate model is produced.

In our example problem we limited ourselves to four variables. Increasing the number of variables naturally results in a more difficult optimization problem, requiring larger initial sampling plans and more update evaluations, and so the benefits of the co-kriging method are likely to increase. However, since we have noted that co-kriging relies on the relationship between analyses being simpler to model than the analyses themselves, we can assume that for additional variables which have little impact on the design objective the co-kriging model will not offer a significant speed-up. Thus, while we expect co-kriging to perform well in high dimensions, we should still try to isolate the most important variables, either through prior knowledge or through the analysis of elementary effects (Morris 1991).

We limited the total sample size (including updates) to  $n_e = 200$  in the example problem due to the expense of tuning the hyper-parameters for large quantities of data overtaking the expense of the CFD simulations. This effectively puts an upper limit on the number of variables that can be considered, since more variables require more points. However, in such cases where  $f_e$  is quick to evaluate when compared with the hyper-parameter tuning we can choose to simplify the correlation, and hence the tuning time, by using the same  $\theta_{c,j}$  and  $\theta_{e,j}$  across all dimensions. The extension of surrogate-based techniques to higher dimensions is a particularly active area of research and further developments in this area would be welcome.



This work has been supported by EPSRC grant code GR/T19209/01.

### Appendix A

To derive the co-kriging predictor, we follow a method similar to that for ordinary kriging. The basis of this method is that we wish our prediction of a new expensive point to be consistent with the observed data and the MLEs for the hyper-parameters. We therefore augment the observed data with a predicted value and maximize the likelihood of this augmented dataset by varying our prediction while keeping the hyper-parameters fixed. This gives us an MLE  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$ .

The augmented dataset is defined as  $\tilde{\mathbf{X}} = \{\mathbf{X}_c^T \ \mathbf{X}_e^T \ \mathbf{x}^{(n_e+1)T}\}^T$  and  $\tilde{\mathbf{y}} = \{\mathbf{y}_c^T \ \mathbf{y}_e^T \ y^{(n_e+1)}\}^T$ , with covariance matrix  $\tilde{\mathbf{C}}$  given by

$$\begin{pmatrix} \hat{\sigma}_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_c) & \rho \hat{\sigma}_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_c, \mathbf{X}_e) & \rho \hat{\sigma}_c^2 \psi_c(\mathbf{X}_c, \mathbf{x}^{(n_e+1)}) \\ \rho \hat{\sigma}_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_c) & \rho_c^2 \hat{\sigma}_c^2 \boldsymbol{\Psi}_c(\mathbf{X}_e, \mathbf{X}_e) + \hat{\sigma}_d^2 \boldsymbol{\Psi}_d(\mathbf{X}_e, \mathbf{X}_e) & \rho \hat{\sigma}_c^2 \psi_c(\mathbf{X}_e, \mathbf{x}^{(n_e+1)}) + \hat{\sigma}_d^2 \psi_d(\mathbf{X}_e, \mathbf{x}^{(n_e+1)}) \\ \left( \rho \hat{\sigma}_c^2 \psi_c(\mathbf{X}_c, \mathbf{x}^{(n_e+1)})^T \ \rho_c \hat{\sigma}_c^2 \psi_c(\mathbf{X}_e, \mathbf{x}^{(n_e+1)})^T + \hat{\sigma}_d^2 \psi_d(\mathbf{X}_e, \mathbf{x}^{(n_e+1)})^T \right) & & \rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 \end{pmatrix}$$

which, defining  $\mathbf{c}$  as a column vector of the covariance of  $\mathbf{X}$  and  $\mathbf{x}^{(n_e+1)}$ , can be expressed as

$$\tilde{\mathbf{C}} = \begin{pmatrix} \mathbf{C} & \mathbf{c} \\ \mathbf{c}^T & \rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 \end{pmatrix}. \tag{A1}$$

In equations (3.3) and (3.8) it is seen that only the last term of the ln-likelihood contains the sample data, and hence to find an MLE  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  we need to maximize

$$-\frac{1}{2}(\tilde{\mathbf{y}} - \mathbf{1}\boldsymbol{\mu})^T \tilde{\mathbf{C}}^{-1}(\tilde{\mathbf{y}} - \mathbf{1}\boldsymbol{\mu}),$$

which may be expressed as

$$-\frac{1}{2} \begin{pmatrix} \mathbf{y} - \mathbf{1}\hat{\boldsymbol{\mu}} \\ \hat{y}_e(\mathbf{x}^{(n_e+1)}) - \hat{\mu} \end{pmatrix}^T \begin{pmatrix} \mathbf{C} & \mathbf{c} \\ \mathbf{c}^T & \rho_c^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y} - \mathbf{1}\hat{\boldsymbol{\mu}} \\ \hat{y}_e(\mathbf{x}^{(n_e+1)}) - \hat{\mu} \end{pmatrix}. \tag{A2}$$

The inverse of the augmented covariance matrix  $\tilde{\mathbf{C}}^{-1}$  is found using the partitioned inverse formula (Theil 1971)

$$\begin{pmatrix} \mathbf{C}^{-1} + \mathbf{C}^{-1} \mathbf{c}(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})^{-1} \mathbf{c}^T \mathbf{C}^{-1} & -\mathbf{C}^{-1} \mathbf{c}(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})^{-1} \\ -(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})^{-1} \mathbf{c}^T \mathbf{C}^{-1} & (\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})^{-1} \end{pmatrix}. \tag{A3}$$

Substituting (A 3) into (A 2) and ignoring terms without  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  we obtain

$$\begin{aligned} & \left( \frac{-1}{2(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) (\hat{y}_e(\mathbf{x}^{(n_e+1)}) - \hat{\mu})^2 \\ & + \left( \frac{\mathbf{c}^T \mathbf{C}^{-1}(\mathbf{y} - \mathbf{1}\hat{\boldsymbol{\mu}})}{(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) (\hat{y}_e(\mathbf{x}^{(n_e+1)}) - \hat{\mu}). \end{aligned} \tag{A4}$$

This expression is maximized by taking the derivative with respect to  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  and setting it to 0,

$$\left( \frac{-1}{\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}} \right) (\hat{y}_e(\mathbf{x}^{(n_e+1)}) - \hat{\mu}) + \left( \frac{\mathbf{c}^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{1}\hat{\mu})}{(\rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) = 0. \tag{A 5}$$

Solving for  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  now gives

$$\hat{y}_e(\mathbf{x}^{(n_e+1)}) = \hat{\mu} + \mathbf{c}^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{1}\hat{\mu}). \tag{A 6}$$

For the estimated mean-squared error in this prediction, we again turn to Jones (2001) for our method of derivation. It is argued that our MLE  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  is more accurate if the likelihood has a definite maximum, i.e. a high curvature. Conversely, if different values of  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$  have similar likelihoods the MLE is less accurate. The error in the predictor is therefore inversely related to the curvature (the second derivative with respect to  $\hat{y}_e(\mathbf{x}^{(n_e+1)})$ ) of the augmented likelihood. We have already found the first derivative in equation (A 5) and taking the derivative and inverse of this equation we obtain

$$s^2(\mathbf{x}) \approx \rho^2 \hat{\sigma}_c^2 + \hat{\sigma}_d^2 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}. \tag{A 7}$$

This equation is not precisely the classically derived formula, which includes an extra term  $(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c})^2 / \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}$ . However, this extra term is so small that it can safely be neglected.

### Appendix B

The expected improvement (equation (5.1)) is typically calculated as

$$\begin{aligned} & (\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})) \left[ \frac{1}{2} + \frac{1}{2} \operatorname{erf} \left( \frac{\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})}{\sqrt{2}s(\mathbf{x})} \right) \right] \\ & + \frac{s(\mathbf{x})}{\sqrt{2\pi}} \exp \left[ - \frac{(\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x}))^2}{2s^2(\mathbf{x})} \right]. \end{aligned} \tag{B 1}$$

When using re-interpolation,  $s(\mathbf{x})$  is typically small and so  $\operatorname{erf}(\cdot) \rightarrow -1$  and  $\exp(\cdot) \rightarrow 0$ . This often leads to floating-point underflow and  $E[I(\mathbf{x})] = 0$ . Using the substitution  $a = (\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})) / \sqrt{2}s(\mathbf{x})$ , when  $a \ll -1$ ,  $\operatorname{erf}(a)$  can be expressed using a Maclaurin series expansion and the first term of (B 1) becomes

$$(\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})) \left[ \frac{1}{2\sqrt{\pi}} \exp(-a^2) \sum_{n=0}^{\infty} \frac{(-1)^n (2n-1)!!}{2^n} a^{-(2n+1)} \right].$$

Note that  $\exp(-a^2)$  appears in the second term in (B 1) and therefore  $E[I(\mathbf{x})]$  can now be expressed as

$$\begin{aligned} & \left[ (\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x})) \frac{1}{2\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{(-1)^n (2n-1)!!}{2^n} a^{-(2n+1)} + \frac{s(\mathbf{x})}{\sqrt{2\pi}} \right] \\ & \times \exp \left[ - \frac{(\min\{\mathbf{y}_e\} - \hat{y}_e(\mathbf{x}))^2}{2s^2(\mathbf{x})} \right]. \end{aligned} \tag{B 2}$$

We can now take natural logarithms and  $\ln E[I(\mathbf{x})]$  can be searched by the optimizer without problems with floating-point underflow.

## References

- Cook, R. D. & Nachtsheim, C. J. 1980 A comparison of algorithms for constructing exact D-optimal designs. *Technometrics* **22**, 315–324. (doi:10.2307/1268315)
- Cousin, J. & Metcalf M. 1990 The BAE Ltd. transport aircraft synthesis and optimization program. In *AHS, and ASEE, Aircraft Design, Systems and Operations Conference, Dayton, OH, 17–19 September 1990*.
- Forrester, A. I. J., Bressloff, N. W. & Keane, A. J. 2006a Optimization using surrogate models and partially converged computational fluid dynamics simulations. *Proc. R. Soc. A* **462**, 2177–2204. (doi:10.1098/rspa.2006.1679)
- Forrester, A. I. J., Keane, A. J. & Bressloff, N. W. 2006b Design and analysis of ‘noisy’ computer experiments. *AIAA J.* **44**, 2331–2339. (doi:10.2514/1.20068)
- Forrester, A. I. J., Sóbester, A. & Keane, A. J. 2006c Optimization with missing data. *Proc. R. Soc. A* **462**, 935–945. (doi:10.1098/rspa.2005.1608)
- Hevesi, J., Flint, A. & Istok, J. 1992 Precipitation estimation in mountainous terrain using multivariate geostatistics. Part II: isohyetal maps. *J. Appl. Meteorol.* **31**, 677–688. (doi:10.1175/1520-0450(1992)031<0677:PEIMTU>2.0.CO;2)
- Jones, D. R. 2001 A taxonomy of global optimization methods based on response surfaces. *J. Global Optim.* **21**, 345–383. (doi:10.1023/A:1012771025575)
- Keane, A. J. 2003 Wing optimization using design of experiment, response surface, and data fusion methods. *J. Aircraft* **40**, 741–750.
- Kennedy, M. C. & O’Hagan, A. 2000 Predicting the output from complex computer code when fast approximations are available. *Biometrika* **87**, 1–13. (doi:10.1093/biomet/87.1.1)
- Krige, D. G. 1951 A statistical approach to some basic mine valuation problems on the Witwatersrand. *J. Chem. Metallurg. Min. Eng. Soc. S. Afr.* **52**, 119–139.
- Locatelli, M. 1997 Bayesian algorithms for one-dimensional global optimization. *J. Global Optim.* **10**, 57–76. (doi:10.1023/A:1008294716304)
- Maskew, B. 1982 Prediction of subsonic aerodynamic characteristics: a case for low-order panel methods. *J. Aircraft* **19**, 157–163.
- Matheron, G. 1963 Principles of geostatistics. *Econ. Geol.* **58**, 1246–1266.
- McKay, M. D., Beckman, R. J. & Conover, W. J. 1979 A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* **21**, 239–245. (doi:10.2307/1268522)
- Morris, M. D. 1991 Factorial sampling plans for preliminary computational experiments. *Technometrics* **33**, 161–174. (doi:10.2307/1269043)
- Morris, M. D. & Mitchell, T. J. 1995 Exploratory designs for computer experiments. *J. Stat. Plan. Infer.* **43**, 381–402. (doi:10.1016/0378-3758(94)00035-T)
- Theil, H. 1971 *Principles of econometrics*. New York, NY: John Wiley.