# Simulation optimization via kriging: a sequential search using expected improvement with computing budget constraints

Ning Quan , Jun Yin , Szu Hui Ng & Loo Hay Lee

Published online: 10 Apr 2013.

Submit your article to this journal ✎

View related articles ✎

# Simulation optimization via kriging: a sequential search using expected improvement with computing budget constraints

NING QUAN, JUN YIN, SZU HUI NG* and LOO HAY LEE

*Department of Industrial and Systems Engineering, National University of Singapore, Singapore*
*E-mail: isensh@nus.edu.sg*

Metamodels are commonly used as fast surrogates for the objective function to facilitate the optimization of simulation models. Kriging (or the Gaussian process model) is a very popular metamodel form for deterministic and, recently, stochastic simulations. This article proposes a two-stage sequential framework for the optimization of stochastic simulations with heterogeneous variances under computing budget constraints. The proposed two-stage framework is based on the kriging model and incorporates optimal computing budget allocation techniques and the expected improvement function to drive and improve the estimation of the global optimum. Empirical results indicate that it is effective in obtaining optimal solutions and is more efficient than alternative metamodel-based techniques. The framework is also applied to a complex real ocean liner bunker fuel management problem with promising results.

Keywords: Simulation optimization, hetergeneous variances, kriging, optimal computing budget allocation, expected improvement, bunker fuel management

## 1. Introduction

Simulation is a widely applied tool to study and evaluate complex systems. A simulation model can imitate the behavior of an actual or anticipated artificial or physical system. It can capture the underlying mechanism and dynamics of a system, which enables decision makers to effectively manage daily operations and make long-term plans. It also provides a test bed to assess changes in operations and managerial policies (Greenwood *et al.*, 2005). In practice, however, in addition to assessing changes in policies, decision makers are often tasked with finding the best (optimal) decision/policy (Shim *et al.*, 2002). This requires combining simulation models with optimization techniques to find the best values of decision variables (inputs) that yield the optimal outputs of the system. Some existing optimization methods for stochastic simulation, including the sequential response surface methodology (see Angün *et al.* (2002)), the stochastic approximation method see (Kushner and Clark (1978)), the nested partitions method (see Shi and Olafsson (2000)), and other heuristic methods such as the genetic algorithm and simulated annealing approaches, are either insufficient for global optimization or computa-

tionally expensive. If we let $Y_0$ be the output of interest of the simulation, and $\theta$ the input decision variables, the general optimization problem has the form $\min_\theta f(\theta) = E(Y_0(\theta))$ where the objective is stochastic if $Y_0$ is random. Typically, replications are taken for each input setting, and the stochastic objective is estimated at that setting by averaging over the replications. Tekin and Sabuncuoglu (2004) and Fu (2007) provided a comprehensive review of the different approaches for simulation optimization.

Due to the stochastic and complex nature of most real-world systems, simulation models of these systems are themselves difficult to build and time consuming to execute. In many cases, decision makers cannot afford to explore a large area of the decision space or to conduct a lengthy search for the best set of decision variables. One feasible alternative is to build a metamodel, a simpler mathematical approximation of the objective, and use the metamodel as a surrogate to search for the optimal decision variable set. Since metamodels can run much faster than the underlying (computationally costly) simulation models, it is usually used as the objective function approximator of the simulation optimization problem (Nakayama *et al.*, 2002; Wan *et al.*, 2005). The most popular techniques used for metamodeling are based on parametric polynomial response surface approximations. Although polynomial response metamodels offer good approximations for simple models,

*Corresponding author

the main drawback of these models is their lack of flexibility to achieve a global fit. Furthermore, many of the complex simulation models and non-stationary models have highly non-linear responses, so linear polynomial approximations are not adequate. Various other types of metamodels such as multivariate adaptive regression splines, kriging, radial basis functions, artificial neural networks, and support vector regression have been proposed in recent years. A review of these metamodel performances and applications in engineering and decision support systems can be found in Simpson *et al.* (2001), Chen *et al.* (2006), and Li *et al.* (2010).

Among the different types of metamodels, the kriging model is an increasingly popular metamodel form as it is more adaptable than the regression-based models and not as complicated and time consuming as artificial intelligence techniques. In the performance study of various metamodels for stochastic computer models, Li *et al.* (2010) applied the previously mentioned metamodels to four popular test functions with different degrees of noise. The metamodels were judged on robustness, accuracy of fit, and efficiency, and kriging performed well in all three categories compared with the others. As metamodel-based optimization is a relatively efficient optimization method compared with other optimization techniques (Tekin and Sabuncuoglu, 2004), and the success of metamodel-based optimization depends critically on the fit of the metamodel, with its outstanding advantages in global fitting, especially under budget constraints, we focus here on kriging-based optimization.

Since its origins in geo-statistics (Matheron, 1963), kriging has been successfully applied in many deterministic computer experiments (Pham and Wagner, 1994; Gupta *et al.*, 2006; Roshan, 2006; Wu and Sun, 2007). Recently, there has also been an increasing interest in adopting kriging metamodels for stochastic simulations, including discrete event simulations. Ankenman *et al.* (2010) and Yin *et al.* (2011) proposed the Modified Nugget Effect Kriging (MNEK) and the stochastic kriging model, respectively, to address the more general heteroscedastic case.

In applying metamodels as a surrogate for optimizing the simulation model, a sequential approach is typically taken. Jones *et al.* (1998) proposed a sequential optimization method based on the kriging metamodel and the Bayesian global optimization approach. The proposed method applied the Expected Improvement (EI) function and the Efficient Global Optimization (EGO) algorithm to balance the local and global searches for the optimum of an unknown response surface as the solution to the global optimization of the corresponding deterministic simulation model. This method is a kriging metamodel-based optimization method developed from the Bayesian-based optimization methods in Mockus (1994). Kleijnen *et al.* (2011) extended this sequential optimization approach by introducing an improved estimator of the kriging variance through bootstrapping. As the originally proposed EI function and EGO algorithm are designed for deterministic scenarios, it considered the allocation of the design points as the only

design option and focused on balancing the search within the local area of the current optimum and the entire sample space. However, for stochastic simulations, the random variability of the stochastic response can considerably affect the metamodel fit (Yin *et al.*, 2009) and therefore the search for the optimum. In this situation, the experimental design is further affected by the stochastic noise in the simulation. Hence, in addition to reducing the spatial uncertainty by observing new design points, the experimenter must also consider the additional design option of adding replications to new and existing design points to reduce the random variability. Huang *et al.* (2006) adapted the EGO scheme for stochastic simulation models and proposed the Sequential Kriging Optimization (SKO) method for optimizing stochastic systems. With the MNEK model and augmented EI function, the SKO algorithm accounts for the influence of random noise. However, SKO only considers homoscedastic cases where the random noise function is assumed to have constant variances throughout the entire sample space. For the more general case with heterogeneous variances, SKO is unable to capture the behavior of the stochastic simulation model due to an incorrect assumption on the variances of the stochastic response. Hence, the estimated global optimum obtained by the SKO with augmented EI function can be far away from the true optimum due to the inadequate fit of the kriging model.

Picheny *et al.* (2010) extended the EI-based optimization algorithm to the case with normally distributed noise and non-constant variances. In addition, they proposed a more general quantile-based criterion, Expected Quantile Improvement (EQI) to take into account the user's risk tolerance. The higher the user sets the quantile, the more conservative will be the criterion and *vice versa*. Their algorithm accounts for limited computing budget and also considers the variance of the noise at unsampled locations when searching for a new point. This gives the algorithm a desirable characteristic of favoring exploration at the start where the available budget is high and becoming more conservative toward the end. However, it also requires the noise variance function be known, and the algorithm's computational complexity is greater compared with traditional EI. In addition, the algorithm proposed by Picheny *et al.* (2010) that features online allocation does not allow backtracking, meaning that once a point has been selected by the criterion and sampled until a condition is met, that point is never re-visited again. In an iterative algorithm where more and more information about the objective function is revealed as the algorithm progresses, this characteristic may not be ideal.

In this article, we develop a two-stage sequential framework that can be adopted for the optimization of stochastic simulation with heterogeneous variances. Similar to Jones (2001) and Kleijnen *et al.* (2011), we adopt the term *optimization* to mean a minimization or maximization, even when there are no constraints. The proposed two-stage framework is able to correctly account for the influence of

non-constant variances in the design space and hence balances the need to reduce spatial uncertainty and random variability with local and global search in a typical stochastic simulation scenario. This article differs from the work presented in Picheny *et al*. (2010) in that the noise variance function need not be known, and the proposed algorithm is less computationally demanding. This article is organized as follows. An ocean liner bunker fuel management optimization problem is introduced in the next section. Then, in Section 3, we review the kriging metamodel, original EI function proposed by Jones *et al*. (1998), and augmented EI function proposed by Huang *et al*. (2006). In Section 4, a two-stage sequential design framework is proposed for stochastic simulation optimization with heterogeneous variances. The proposed framework is applied to several test functions in Section 5 to illustrate its performance, and the ocean liner bunker fuel management problem is revisited in Section 6. Finally, conclusions are drawn and opportunities for further improvements are highlighted.

## 2. Motivating example

In recent years, increasing bunker prices have threatened shipping liner companies' accounting bottom line. To compensate for increasing prices, over 200 shipping companies have adopted the strategy of slow steaming, especially in long-haul loops such as Asia to Europe and Asia to North America. From empirical estimations, lowering the cargo vessel speed by 20% can save up to 50% of daily fuel consumption (Notteboom and Vernimmen, 2009). Although slow steaming results in having to increase the number of vessels in service to maintain a satisfactory service level, this additional capital cost is usually offset by the savings from fuel costs. Another challenge facing shipping companies is deciding where to refuel (bunker) along a given route and the amount to bunker at these selected ports.

A recent research study by Yao *et al*. (2012) modeled the shipping fuel management planning problem as a mixed-integer linear program. The linear program minimizes bunker fuel-related costs by determining optimal ship speeds between ports and selecting bunkering ports and bunkering amounts. Based on real data obtained from a shipping company, they empirically modeled the relationship between fuel consumption rate and ship speed for different ship sizes and expressed the daily fuel consumption rate as $F(V) = k_1 V^3 + k_2$, where $k_1$ and $k_2$ are two constants that depend on ship size and $V$ is ship speed in knots.

This previous work can be extended to consider the operational-level problem by considering more realistic operational conditions. This includes accounting for fuel consumption rate uncertainty and bunker price variation. With random fuel consumption, the ship's bunker fuel inventory safety level becomes an important system-level decision as a fuel stock-out scenario is possible and can be very costly to liner operations. In addition, as the distances between each port within the service can differ by thousands of nautical miles, safety levels at each leg should differ, each reflecting about an equal risk of a fuel stock-out within the leg. Due to the complexities of the real operational system and impact of the random components, a simulation optimization approach is applied to a rolling-horizon discrete event simulation model to determine the safety levels for a given shipping service,

In the rolling-horizon discrete event simulation, the following takes place upon the arrival of a ship at each port.

1. Observe the time taken and fuel consumed from the previous leg.
2. Update the bunker inventory level based on the observed consumption. Update bunker prices at all ports based on current market prices.
3. Solve a mathematical optimization problem to obtain the optimal fuel management decisions for the bunkering ports, bunkering amounts, and ship speeds.
4. Bunker fuel amounts at current port (if necessary) and set ship speed for next leg.

This procedure yields the bunkering amounts, ports, and ship speeds one port at a time with each arrival in the simulation. Incorporating the mixed-integer linear program in the rolling-horizon discrete event simulation model facilitates a more realistic operational representation of the bunker fuel management problem, and enables better evaluations of system strategies (such as bunker safety levels).

Here we consider the Asia-Europe Express (AEX) service route of a shipping company (see Fig. 1). As each run of the simulation incorporates the rolling-horizon mathematical optimization model to determine the optimal fuel management strategy, each run is quite time consuming, taking around 5 minutes on an Intel Core 2 Duo 2.4 GHz computer. Consequently, one must carefully allocate the simulation budget to efficiently and effectively search the multidimensional decision space for the optimal bunker safety levels.

## 3. Background and basics

In this article, we develop a new methodology for performing simulation runs to explore and estimate a response model of the computer simulation model in order to determine the global optimal point. Here we adopt the kriging model as a surrogate for the simulation model, which helps provide an estimate of the response and also the noise and uncertainty. We also adopt a sequential search and allocation algorithm with a criterion to determine design points aimed at eventually identifying the optimal decision point.

**Fig. 1.** AEX service route (distances in nautical miles).

### 3.1. *Kriging basics*

The general kriging metamodel assumes that the stochastic simulation response can be modeled as realizations of a random process given as

$$Y(x) = Z(x) + \varepsilon(x), \tag{1}$$

where $Z(x)$ represents the deterministic mean function of the stochastic response (usually the expectation of a simulation's performance measure). $Z(x)$ can be further decomposed into the process mean function $\mu(x) = \mathbf{F}(x)\boldsymbol{\beta}$, where $\mathbf{F}(x)$ is a vector of $p$ functions of $x$ and $\boldsymbol{\beta}$ is the vector of model parameters, and a spatial process $\delta(x)$ with $E[\delta(x)] = 0$, $\text{var}[\delta(x)] = \sigma_z^2$ and $\text{cov}(\delta(x_i), \delta(x_j)) = \sigma_z^2 R_\theta(x_i, x_j)$. A popular choice of $R$ is the exponential family of correlation functions of the form $R(x_i, x_j) = \prod_{t=1}^{k} \exp(-\theta_t d_{ij,t}^p)$, where the correlation parameter $\theta_t$ is the sensitivity parameter that controls how fast the correlation decays with distance in the $t$th dimension, $d_{ij,t}$ is the distance between $x_i$ and $x_j$ in the $t$th dimension, $k$ is the dimension of input $x$, and $p$ controls the general smoothness of the response. $\varepsilon(x)$ is the random error function with zero mean and unknown variance $\sigma_\varepsilon^2(x)$. To model the dependence of $\sigma_\varepsilon^2(x)$ on location $x$, it is further assumed that $\delta(x)$ and $\sigma_\varepsilon^2(x)$ are independent. $\sigma_\varepsilon^2(x)$ can be modeled using non-parametric models (Opsomer *et al.*, 1999) or a spatial process with $E[\sigma_\varepsilon^2(x)] = 0$ and $\text{cov}[\sigma_\varepsilon^2(x_i), \sigma_\varepsilon^2(x_j)] = \sigma_\varepsilon^2 R_\varepsilon$ (Ankenman *et al.*, 2010; Ng and Yin, 2012). In the deterministic simu-

lation context, $\varepsilon(x)$ is not present (i.e., $\varepsilon(x) = 0$). Typically, the mean function $Z(x)$ is of experimenter's interest and the predictor for $Z(x)$ at any point $x_0$ can be expressed as

$$\hat{Z}(x_0) = \mathbf{F}(x_0)\boldsymbol{\beta} + \mathbf{c}(x_0)^{\mathrm{T}} \mathbf{R}^{-1}(\bar{\mathbf{Y}} - \mathbf{F}\boldsymbol{\beta}), \tag{2}$$

where $\mathbf{c}(x_0)$ is the spatial correlation vector of $x_0$ to all $m$ existing observation locations and $\bar{\mathbf{Y}} = \{\bar{Y}_1, \bar{Y}_1, \ldots, \bar{Y}_m\}^{\mathrm{T}}$ is the vector of the $m$ observed sample means. With the model form of Equation (1), $\mathbf{c}(x_0)$ and $\mathbf{R}$ are given as $\mathbf{c}(x_0) = \{corr_z(d_{01}), corr_z(d_{02}), \ldots, corr_z(d_{0m})\}^{\mathrm{T}}$ and

$$\mathbf{R} = \mathbf{R}_z + \mathbf{R}_\varepsilon = \begin{pmatrix} 1 & corr_z(d_{12}) & \ldots & corr_z(d_{1m}) \\ corr_z(d_{21}) & 1 & \ldots & corr_z(d_{2m}) \\ \vdots & \vdots & \ddots & \vdots \\ corr_z(d_{m1}) & corr_z(d_{m2}) & \ldots & 1 \end{pmatrix}$$
$$+ \begin{pmatrix} \eta_1 & 0 & \ldots & 0 \\ 0 & \eta_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \eta_m \end{pmatrix}$$

where $corr_z(d_{ij})$ is the correlation function representing the spatial correlation between $\delta(x_i)$ and $\delta(x_j)$ and $\eta_i = \sigma_\varepsilon^2/\sigma_Z^2$ is the ratio of the variances of the random error function $\varepsilon(x)$ to the mean function $Z(x)$. The predictive variance (or

mean squared prediction error) can be derived as

$$s^2(x_0) = \sigma_Z^2 \big[ 1 - \mathbf{c}(x_0)^{\mathrm{T}} \mathbf{R}^{-1} \mathbf{c}(x_0) + (\mathbf{F}(x_0) - \mathbf{F}^{\mathrm{T}} \mathbf{R}^{-1} \mathbf{c}(x_0))^{\mathrm{T}}$$
$$\times (\mathbf{F}^{\mathrm{T}} \mathbf{R}^{-1} \mathbf{F})^{-1} (\mathbf{F}(x_0) - \mathbf{F}^{\mathrm{T}} \mathbf{R}^{-1} \mathbf{c}(x_0)) \big] \quad (3)$$

This provides the prediction uncertainty at location $x_0$. Similar to the predictor mean, the predictor variance is dependent on both the spatial correlation and random variability. This modeling form is known as the MNEK model (Yin *et al*., 2011), and a similar form known as stochastic kriging can be found in Ankenman *et al*. (2010). In the absence of random noise $\varepsilon(x)$, the predictor mean and variance in Equations (2) and (3) reduce to the deterministic kriging predictor (see Santner *et al*. (2003, p. 55)):

$$\hat{Z}_z(x_0) = \mathbf{F}(x_0)\boldsymbol{\beta}_z + \mathbf{c}(x_0)^{\mathrm{T}} \mathbf{R}_z^{-1}(\mathbf{Y} - \mathbf{F}\boldsymbol{\beta}_z), \quad (4)$$

which is the best linear unbiased predictor with predictor variance:

$$s_z^2(x_0) = \sigma_z^2 \big[ 1 - \mathbf{c}(x_0)^{\mathrm{T}} \mathbf{R}_z^{-1} \mathbf{c}(x_0) + \big(\mathbf{F}(x_0) - \mathbf{F}^{\mathrm{T}} \mathbf{R}_z^{-1} \mathbf{c}(x_0)\big)^{\mathrm{T}}$$
$$\times \big(\mathbf{F}^{T} \mathbf{R}_z^{-1} \mathbf{F}\big)^{-1} \big(\mathbf{F}(x_0) - \mathbf{F}^{\mathrm{T}} \mathbf{R}_z^{-1} \mathbf{c}(x_0)\big) \big] \quad (5)$$

This deterministic kriging metamodel form treats the response surface as a Gaussian random process, interpolating exactly at observed points and representing the uncertainty about the surface at unobserved locations.

As $\boldsymbol{\beta}$ $\boldsymbol{\theta}$, $\sigma_z^2$, $\sigma_\varepsilon^2$ are typically unknown, the maximum likelihood method is commonly used to estimate these parameters (refer to Yin *et al*. (2011) for derived maximum likelihood estimators). Although substituting these estimates into Equations (3) and (5) results in biased estimators of the predictor variance (Yin *et al*., 2009; Kleijnen *et al*., 2011), Ankenmann *et al*. (2010) showed that the penalty for estimating these parameters (particularly $\mathbf{R}$) is typically small. Moreover, in its application in optimization, Kleijnen *et al*. (2011) also showed that using the estimated predictor variances with the EI is robust.

The kriging predictor in Equations (2) and (3) can be further decomposed in terms of the deterministic predictor in Equations (4) and (5), with

$$\hat{Z}(x_0) = \mathbf{F}(x_0)(\boldsymbol{\beta}_z + \boldsymbol{\beta}_\varepsilon) + \mathbf{c}(x_0)^{\mathrm{T}}(\mathbf{R}_z + \mathbf{R}_\varepsilon)^{-1}(\bar{\mathbf{Y}} - \mathbf{F}(\boldsymbol{\beta}_z + \boldsymbol{\beta}_\varepsilon))$$
$$= \mathbf{F}(x_0)\boldsymbol{\beta}_z + \mathbf{c}(x_0)^{\mathrm{T}}\mathbf{R}_z^{-1}(\mathbf{Y} - \mathbf{F}\boldsymbol{\beta}_z) + \mathbf{F}(x_0)\boldsymbol{\beta}_\varepsilon + \mathbf{c}(x_0)^{\mathrm{T}}\boldsymbol{\Xi}_\varepsilon$$
$$\quad (6)$$

where $\boldsymbol{\beta}_\varepsilon$ is a vector and $\boldsymbol{\Xi}_\varepsilon$ is a vector that depends on $\mathbf{R}_z$ and $\mathbf{R}_\varepsilon$. The predictive variance is given as

$$s^2(x_0) = \sigma_z^2 \big[ 1 - \mathbf{c}(x_0)^{T}(\mathbf{R}_z + \mathbf{R}_\varepsilon)^{-1}\mathbf{c}(x_0)$$
$$+ \big(\mathbf{F}(x_0) - \mathbf{F}^{T}(\mathbf{R}_z + \mathbf{R}_\varepsilon)^{-1}\mathbf{c}(x_0)\big)^{T}$$
$$\times \big(\mathbf{F}^{T}(\mathbf{R}_z + \mathbf{R}_\varepsilon)^{-1}\mathbf{F}\big)^{-1}$$
$$\times \big(\mathbf{F}(x_0) - \mathbf{F}^{T}(\mathbf{R}_z + \mathbf{R}_\varepsilon)^{-1}\mathbf{c}(x_0)\big) \big]$$

$$= \sigma_z^2 \big[ 1 - \mathbf{c}(x_0)^{T}\mathbf{R}_z^{-1}\mathbf{c}(x_0) + \big(\mathbf{F}(x_0) - \mathbf{F}^{T}\mathbf{R}_z^{-1}\mathbf{c}(x_0)\big)^{T}$$
$$\big(\mathbf{F}^{T}\mathbf{R}_z^{-1}\mathbf{F}\big)^{-1} \big(\mathbf{F}(x_0) - \mathbf{F}^{T}\mathbf{R}_z^{-1}\mathbf{c}(x_0)\big) \big] + \sigma_z^2 \boldsymbol{\Omega}_\varepsilon \quad (7)$$

where $\boldsymbol{\Omega}_\varepsilon$ is a function that depends on $\mathbf{R}_z$ and $\mathbf{R}_\varepsilon$. The terms $\boldsymbol{\beta}_\varepsilon$, $\boldsymbol{\Xi}_\varepsilon$, and $\boldsymbol{\Omega}_\varepsilon$ represent the influence of random noise on the estimators and will all equal to zero when $\varepsilon(x) = 0$.

### 3.2. *Sequential optimization algorithms*

In metamodel-based optimization, a sequential approach is typically taken because a one-shot approach can be inefficient since it may focus sampling effort in regions of the design space that provide little information about the optimum. In this article, we adopt a sequential global optimization approach, allowing for more effective use of available computing resources to balance the allocation for the search and also replication numbers.

### 3.2.1. *Improvement functions*

In the selection of search points for sequential-based optimization algorithms, several criteria such as minimizing the response or minimizing a lower bound of the response, maximizing the probability of improvement, or maximizing the expected improvement, have been proposed (see Jones (2001)). Of these search criteria, the improvement functions are most popular due to their ability to trade-off exploitation and exploration in order to converge to the global minimum. Here we review the improvement function initially proposed for deterministic simulations and an extension for stochastic simulations.

Jones *et al*. (1998) defined the improvement function at any potential design point $x_0$ as $I(x_0) = \max[Z_{\min} - Z_{\mathrm{p}}(x_0), 0]$, where $Z_{\min}$ indicates the already observed (simulated) minimum of the mean function and $Z_{\mathrm{p}}(x_0)$ is the random variable that accounts for the predictive uncertainty of the kriging predictor at design point $x_0$. At any unsampled design points, as the mean response $Z_{\mathrm{p}}(x_0)$ is unknown, the improvement is averaged over the uncertainty in the response instead. The expected improvement is defined as

$$E\big[\max\big[Z_{\min} - Z_{\mathrm{p}}(x_0), 0\big]\big]. \quad (8)$$

This criterion is intuitive as it actively searches for design points that maximize the response improvement over the best, while considering also the uncertainty in the response at unobserved points. Based on the results from Jones *et al*. (1998), under the assumptions of the deterministic kriging model where $Z_{\mathrm{p}}(x_0) \sim N(\hat{Z}_{\mathrm{p}}(x_0), s_z^2(x_0))$, Equation (8) can
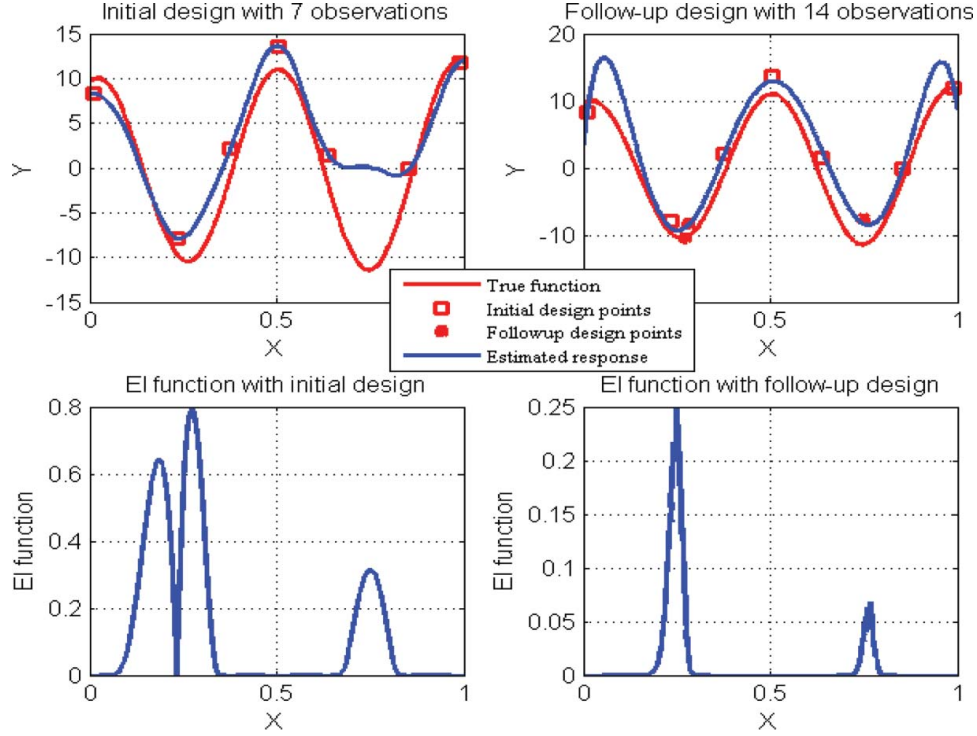
**Fig. 2.** EI function and response metamodel with noisy test function (white noise) (color figure provided online).

be computed by

$$
\begin{aligned}
E[\max[Z_{\min} - Z_{\mathrm{p}}(x_0), 0]] \\
= (Z_{\min} - \hat{Z}_z(x_0)) \Phi \left( \frac{Z_{\min} - \hat{Z}_z(x_0)}{S_z(x_0)} \right) \\
+ S_z(x_0) \varphi \left( \frac{Z_{\min} - \hat{Z}_z(x_0)}{S_z(x_0)} \right).
\end{aligned}
\tag{9}
$$

Huang *et al.* (2006) adapted the EI function to address stochastic simulation systems. They adopted the nugget effect kriging model (Cressie, 1993), where the random errors are assumed to be identical and independently distributed throughout the design space; i.e., $\eta_i = \eta$ for all $i$, to model the stochastic responses of the simulation and proposed the following Augmented EI (AEI) function to drive the sequential search for the design points:

$$
\begin{aligned}
AEI[I(x_0)] = E \left[ \max \left[ Z_{\min} - Z_{\mathrm{p}}(x_0), 0 \right] \right] \\
\times \left( 1 - \frac{\sigma_\varepsilon}{\sqrt{S^2(x_0) + \sigma_\varepsilon^2}} \right).
\end{aligned}
\tag{10}
$$

If $x_0$ is the current point with the best response, then the AEI is equal to the relative reduction in the prediction error at $x_0$ with the addition of another replication. The authors use this factor to represent the diminishing returns of additional replications at the current point of best response.

### 3.3. *Limitations of EGO and SKO in noisy heteroscedastic situations*

Although applications of EGO (with EI) and SKO (with AEI) have been successful in deterministic and homogeneously random situations respectively, direct application to heterogeneously random situations is not straightforward. To better understand the behavior of the algorithms and criteria, we first look at potential limitations of each in certain stochastic situations considered in this article.

We first illustrate the limitations of EGO in a simple random situation. Consider the following example where the noisy test function is given as

$$
Y(x) = Z(x) + \varepsilon(x) = (2x + 9.96) \cos(13x - 0.26) + \varepsilon(x),
\tag{11}
$$

where $\varepsilon(x)$ is a white noise function with zero mean, variance $\sigma_\varepsilon^2 = 2$, and $x \in [0, 1]$. The test function is sinusoidal with a slightly increasing gradient, with one local minimum at 0.2628 and a global minimum at 0.7460.

Starting with an initial seven-point Latin Hypercube Design (LHD) and five replications per point, the initial deterministic ordinary kriging model fit and EI function are given in the left two plots of Fig. 2. The main problem with fitting a deterministic kriging model to stochastic simulation outputs is that the kriging response interpolates the sample means, which can lead to a misleading estimate of the global minimum as seen in Fig. 2. We see that even after seven more input observations, the estimated response has mistakenly identified the wrong minimum point because of
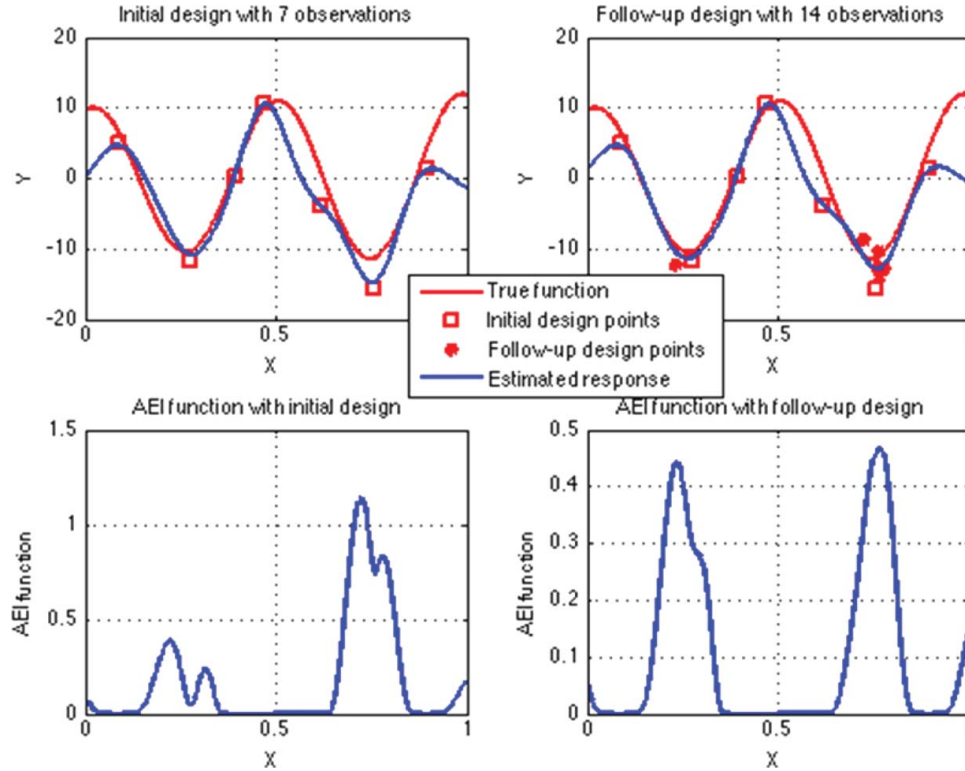
**Fig. 3.** AEI function and response metamodel with noisy test function (white noise) (color figure provided online).

a low observation near the local minimum on the left. This problem could be mitigated if there was a mechanism that could improve the precision of the sample means around the local minima. Since the EGO algorithm lacks such a feature to deal with random variability, the algorithm is heavily penalized by a misleading initial fit.

Applying the sequential SKO algorithm, we obtain a much better fit and results. Figure 3 displays the initial nugget effect model fit, AEI function, and final response fit. As can be seen, with the appropriate metamodel form and sequential criterion, the results are much improved.

However, in the more general case of stochastic simulations with heterogeneous noise, the AEI function cannot properly account for the influence of random noise with non-constant variances. One straightforward method is to substitute the constant value $\sigma_\varepsilon$ in Equation (10) with the function $\sigma_\varepsilon(x)$. We modify the test function in Equation (11) so that the variance of the random noise now is a location-dependent function with the linear form $\sigma_\varepsilon^2(x) = 1 + x, x \in [0, 1]$. Figure 4 displays the response model fit and the modified AEI function.

As seen in the upper left plot of Fig. 4, the estimated response surface of the initial design has two local optima: the current best on the left-hand side and a local sub-optimum on the right. Based on the functional form of the variance function, the area on the left-hand side has a lower variability than the area on the right-hand side. Hence, in the

follow-up sequential design, all seven new observations are located in the left region (area with lower noise and current best estimate). The AEI function still favors this left region after the follow-up design. This can be partially explained as the AEI function favors areas with lower variability due to the multiplicative effect of the augmented factor. Hence, it is possible that the algorithm will be trapped in the local area with low variability for a long time. This is not a desirable property for a global optimization method, especially when faced with limited budget constraints.

Another straightforward option is to simply use the original EI function based on the MNEK model (defined immediately after Equation (3)). The problem here is that the EI function is not very sensitive to decreases in predictor uncertainty $s(x)$ at the current best point. Figure 5 shows that at values of $Z_{\min} - \hat{Z}$ close to zero (as would be the case at the current best point), the EI function decreases slowly with respect to decreases in predictor uncertainty $s(x)$. This can lead to repeated selections of the current best design point at the expense of exploring other promising regions of the design space. The multiplicative factor in AEI was meant to address this issue by discouraging repeated selections of the current best point. However, applying the multiplicative factor to cases with heterogeneous noise would require the noise variance function be known, and the sampling issue mentioned before will still persist.
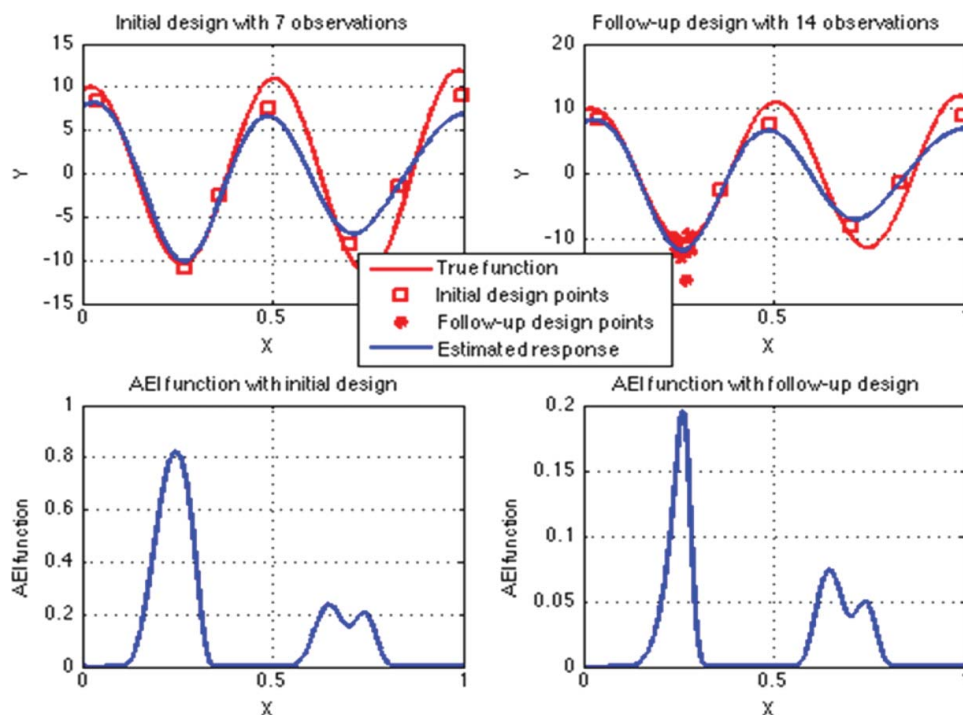
**Fig. 4.** Modified AEI function and response model with noisy test function (non-constant variance) (color figure provided online).

### 3.4. *Characteristics of good algorithms and criteria*

Based on the observations made above, we outline some challenges of adopting an EGO-type framework for simulations with heterogeneous noise.

1. The response model estimation is affected by both the location and computing effort at each design point (Yin *et al*., 2009; Ng and Yin, 2012). With heterogeneous variances and limited computing budget, the ideal distribution of computing effort is unlikely to be equal. An effective procedure should be able to make best use out of the limited computing budget in locating the global optimum.

2. The EGO algorithm has the desirable characteristic of balancing exploitation and exploration. In order to retain this desirable characteristic in a stochastic environment, a new procedure should be able to search globally without having to exhaustively search a local region.



**Fig. 5.** Contour plot of EI function of predictor mean difference and standard deviation using the MNEK model.

Good estimates $Z_{min}$ are also necessary, especially in situations where there are several optima with values close to the global optimum.

3. In situations with limited computing budget, it is useful for the algorithm to explore unexplored regions (improve chances of finding a better minimum) at the outset of the algorithm. However, as the budget is expended toward the end, the focus should be on fine-tuning in the current best area as there will likely be insufficient remaining budget to find and differentiate a lower minimum elsewhere.

In the next section, we will detail the development of our search and allocation procedure that considers these desirable characteristics.

## 4. Development of simulation optimization algorithm

As previously noted, the accuracy and fit of the kriging model can be drastically affected by random noise. It is therefore important to consider both the response and the noise levels carefully. As seen in Fig. 5, direct application of the MNEK model may not have an equitable effect on the EI criterion. The proposed methodology intends to address both considerations individually in a two-stage sequential approach. This two-stage iterative framework represents a "division of labor" approach toward the optimization of a stochastic objective function with heterogeneous noise variance. After the initial fit, each subsequent iteration of the algorithm is composed of a search stage followed by an allocation stage. In the first stage, denoted as the search stage, the Modified Expected Improvement (MEI) infill criterion is used to select a new point. In the second stage, denoted as the allocation stage, the Optimal Computing Budget Allocation (OCBA) technique (Chen *et al.*, 2000) is applied to distribute an additional number of simulation replications among sampled design points. Specifically, the algorithm's search stage is responsible for identifying potential global optimum points, whereas the allocation stage seeks to drive down uncertainty due to random variability at sampled points, with the goals of improving model fit at regions that contain local minimums and eventually correctly selecting the point of the global optimum.

The framework also features a division of allocation heuristic. Here we set the amount of computing budget per iteration as a constant, but the distribution of each iteration's budget between sampling and allocation stages changes as the algorithm progresses. At the start, most of the iteration's budget goes toward exploration (search stage); toward the end, the emphasis is shifted to identifying the point of best response (allocation stage).

Together, this iterative two-stage approach, along with the division of allocation heuristic, will sequentially search the design space for the global minimum while considering notably the fit of the surrogate model driving the search.

The two stages of the framework will be described in detail in the next subsections.

### 4.1. *The search stage*

In this search stage, the following modified EI function is proposed:

$$MEI = E\big(\max\big[Z_{min} - Z_p^*(x), 0\big]\big). \tag{12}$$

In this proposed criterion, $Z_{min}$ is the predicted response at the sampled point with the lowest sample mean, and $Z_p^*(x)$ is a normal random variable with mean given by the MNEK predictor at $x$ (Equation (2)) and variance given by the predictor's spatial prediction uncertainty $s_z^2$ (Equation (5)).

The use of the MNEK predictor for the response is a straightforward choice that provides an unbiased prediction given the heterogeneous nature of the noise. What is different from a straightforward adaptation of the EI criterion is the treatment of the predictor uncertainty. Since the allocation stage will by design address the stochastic noise, only $s_z^2$ is used in this search stage. This enables the search to focus on new points in promising regions with high predicted responses and new points that reduce the spatial uncertainty of the metamodel. In addition, by ignoring predictor uncertainty caused by random variability, $s_z^2$, the modified EI function assumes that the observations are made with infinite precision so the same point is never selected again. This allows the algorithm to quickly escape from a local optimum and brings the sampling behavior closer to that of the original EI criterion and its associated trade-off between exploration and exploitation.

### 4.2. *The allocation stage*

Since the search stage is dedicated to the selection of a new point, the allocation stage will have to manage random variability by intelligently allocating additional replications among sampled points. A related problem arose in the work on stochastic kriging performed by Ankenman *et al.* (2010). In that work, the ultimate goal was to improve the global fit of the stochastic kriging response surface; therefore, additional replications were distributed among sampled points with the objective of minimizing the integrated mean squared error. In the case here where the ultimate goal is global optimization, additional replications are distributed with the goal of maximizing the probability of the correct selection of the sampled point $x^{**}$ as the global optimum.

Alternative algorithms by Huang *et al.* (2006) and Picheny *et al.* (2010) identify the point with the best response (or global optimum) as the sampled point with the lowest predicted quantile response:

$$x^{**} = \underset{x=\{x_1, x_2, \dots, x_n\}}{\arg\min} \quad \hat{Z}(x) + cs(x).$$

where $c$ is user defined in accordance with the user's risk tolerance, with a higher $c$ corresponding to a lower level of risk tolerance and *vice versa*. For practical purposes, $x^{**}$ in these algorithms are also optimized over previously observed locations.

Compared with the above method, OCBA formulates the global optimum selection problem as an optimization exercise and thus provides a more rigorous way of identifying the sampled point with the best response. OCBA uses the sample mean and sample variance as response and response uncertainty estimators, respectively. As the allocation stage aims to improve the data used to fit the model, especially in promising regions, adopting these estimators at this stage is reasonable.

Assuming that we have $n$ sampled points, with each point $x_i$ having a sample mean given by $\bar{Y}_i$ sample variance $\hat{\sigma}_\varepsilon^2(x_i)$ and replication number $R_i$, then according to Theorem 1 provided by Chen *et al.* (2000), the Approximate Probability of Correct Selection (APCS) can be asymptotically maximized when available computing budget tends to infinity and when

$$\frac{N_i}{N_j} = \left( \frac{\hat{\sigma}_\varepsilon(x_i)/\Delta_{b,i}}{\hat{\sigma}_\varepsilon(x_j)/\Delta_{b,j}} \right)^2 i, j \in \{1, 2, \ldots, n\} \text{ and } i \neq j \neq b,$$ 

(13)

$$N_b = \hat{\sigma}_\varepsilon(x_b) \sqrt{\sum_{i=1, i \neq b}^{n} \frac{N_i^2}{\hat{\sigma}_\varepsilon^2(x_i)}},$$

(14)

where $N_i$ is the number of replications allocated to point $x_i$, $x_b$ is the point with the lowest sample mean and $\Delta_{b,i}$ is the difference between the lowest sample mean and the sample mean at point $x_i$. Adopting this allocation rule to maximize the APCS, at the end of the allocation stage, the sampled point with the lowest sample mean will be selected as the location of the best response. The estimate of the best response, $Z_{\min}$, which will be used in the MEI criterion will be given by the MNEK predictor.

An additional benefit of using the OCBA technique is that the majority of the additional replications will be allocated to points with low sample means and high sample variances in order to separate the best point from the contenders. Assuming that the lowest sample means lie close to local minimums, then OCBA can utilize the additional replications to tighten the metamodel fit at regions around the local minimums. This in turn can help the modified EI criterion make a better selection in the subsequent iteration.

### 4.3. *An algorithm overview*

In the proposed framework, the procedure alternates between the search and allocation stages until the budget is consumed. As discussed, in such an approach, it is desirable to start off with more exploration and end by exploiting the best few minimums identified. Here, we adopt a simple division of allocation heuristic to divide the budget allocation
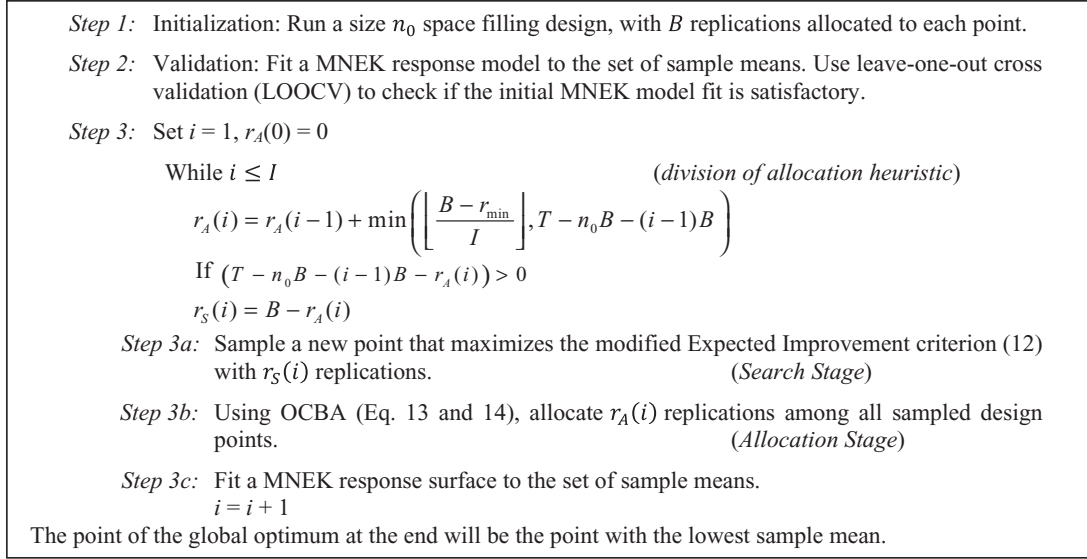
**Table 1.** Algorithm parameter definitions

| Parameter | Definition |
|---|---|
| $T$ | Total number of replications at the start |
| $B$ | Number of replications available for each iteration |
| $n_0$ | Size of initial space filling design. $n_0 B \leq T$ |
| $i$ | Current iteration. $i = 0, 1, 2, \ldots, I$, $I = \lceil (T - n_0 B)/B \rceil$ |
| $r_{\min}$ | Minimum number of replications for sampling a new point. $r_{\min} \leq B$ |
| $r_A(i)$ | Number of replications available for allocation stage at iteration $i$ |
| $r_S(i)$ | Number of replications available for search stage at iteration $i$ |

to the different stages on the proposed framework. Table 1 defines the parameters used in the algorithm.

Before the proposed algorithm can begin, the user needs to set the parameters $T$, $B$, $n_0$, and $r_{\min}$. $T$ will be constrained by practical considerations such as the total time available for the entire optimization exercise and the average length of a simulation run. The size of the initial design, $n_0$, can be set to $10k$ as suggested by Jones *et al.* (1998), where $k$ is the number of dimensions. This $10k$ recommendation has been extensively studied by Loeppky *et al.* (2009) and found to be a reasonable rule of thumb for initial designs. $r_{\min}$ and $B$ should be set such that there are sufficient replications available for the first allocation stage. As a rough guide, $B$ should be set such that the number of replications available for the first allocation stage is at least five times the number of design space dimensions. The algorithm and the computing resource allocation heuristic are described in Fig. 6.

Since the starting parameters that determine the number of iterations and the computing budget used per iteration are set prior to collecting any data, there is a possibility that the recommended parameter settings are unsuitable for the problem. In step 2, Leave-One-Out Cross-Validation (LOOCV) can provide feedback regarding the suitability of the initial parameters $T$, $B$, $n_0$, and $r_{\min}$. In LOOCV, the sample mean of a single design point in the initial design is used for validation while the remaining sample means are used to fit a MNEK response model. The MNEK response model fitted from $n_0 - 1$ sample means is then used to build a $(1 - \alpha)$ confidence interval for the expected response at the design point that was left out of the MNEK model. Ideally, the confidence interval should contain the sample mean obtained from the design point that was left out. This procedure is repeated until all initial design points have been validated. If a single design point or more fails the cross-validation test, it could mean that the current amount of computing budget is insufficient to deal with the noise in the response. Possible remedies include increasing $B$ or increasing the number of design points around the point(s) that fail the cross-validation test or applying a log

*Step 1:* Initialization: Run a size $n_0$ space filling design, with $B$ replications allocated to each point.

*Step 2:* Validation: Fit a MNEK response model to the set of sample means. Use leave-one-out cross validation (LOOCV) to check if the initial MNEK model fit is satisfactory.

*Step 3:* Set $i = 1$, $r_A(0) = 0$

While $i \leq I$        *(division of allocation heuristic)*

$$r_A(i) = r_A(i-1) + \min\left( \left\lfloor \frac{B - r_{\min}}{I} \right\rfloor, T - n_0 B - (i-1)B \right)$$

If $(T - n_0 B - (i-1)B - r_A(i)) > 0$

$r_S(i) = B - r_A(i)$

*Step 3a:* Sample a new point that maximizes the modified Expected Improvement criterion (12) with $r_S(i)$ replications.      *(Search Stage)*

*Step 3b:* Using OCBA (Eq. 13 and 14), allocate $r_A(i)$ replications among all sampled design points.      *(Allocation Stage)*

*Step 3c:* Fit a MNEK response surface to the set of sample means.
$i = i + 1$

The point of the global optimum at the end will be the point with the lowest sample mean.

**Fig. 6.** Sequential two-stage algorithm.

or inverse transformation to the response as suggested by Jones *et al.* (1998).

After successful validation, the computing budget set aside for the allocation stage ($r_A(i)$) increases by a block of $\lfloor (B - r_{\min})/I \rfloor$ replications with every iteration while it decreases by the same number for the search stage. This simple computing budget distribution heuristic gives the algorithm a desirable characteristic of placing more emphasis on exploration at the start and focusing more on exploitation at the end.

In summary, the proposed algorithm relies on two techniques to handle the tasks of selecting new points and managing the heterogeneous nature of the noise. From a modeling viewpoint, we can relate the algorithm back to the kriging model form and its overall iterative goal of optimization. The selection of a new point in the search stage can be seen as improving the fit of predictor $\hat{Z}(x_0)$ or the subsequent iterations by improving the estimation of $\hat{\beta}_z$ (see Equation (6)) and reducing the spatial uncertainty in the prediction, $\sigma_z^2[1 - \mathbf{c}(x_0)^{\mathrm{T}}\mathbf{R}_z^{-1}\mathbf{c}(x_0) + ((1 - \mathbf{F}^{\mathrm{T}}\mathbf{R}_z^{-1}\mathbf{c}(x_0))/\mathbf{F}^{\mathrm{T}}\mathbf{R}_z^{-1}\mathbf{F})^2]$ (or the first term in Equation (7)), especially for promising regions of $x_0$. The purpose of the allocation stage is twofold. This stage first improves the fit of $\hat{Z}(x_0)$ by improving the estimation of $\hat{\beta}_\varepsilon$ (see Equation (6)) and reducing the random uncertainty in the prediction, $\sigma_z^2 \mathbf{\Omega}_z$ (or the last term in Equation (7)). Second, the resource allocation goal of maximizing the probability of correct selection also improves the estimation of $Z_{\min}$ for the search stage in the next iteration.

### 4.4. *Further remarks on the algorithm*

First, by ignoring predictor uncertainty attributed by random noise ($s_\varepsilon^2$) in the modified EI criterion, noisy regions

of the design space may be penalized by the proposed criterion. The MNEK fit in highly noisy regions tends to revert closer to the mean, which results in a smaller $Z_{\min} - Z_p^*(x)$ value. This disadvantage is compounded when the prediction uncertainty at these regions is smaller than it actually is due to the absence of $s_\varepsilon^2$. Although this problem can be mitigated with the cross-validation test and allocation stage, a more comprehensive solution may be to dynamically allocate computing budget to a point until uncertainty due to noise is reduced to an acceptable level before moving on to a new design point.

Second, despite the fact that the algorithm is designed to address simulation systems that are stochastic with variances that usually change significantly across the design space the proposed approach can also be directly applied to the homogeneous case. In these situations where the estimated variances throughout the design space are approximately equal, the allocation criteria (Equations (13) and (14)) will naturally adjust to focus on the most promising points with lower sample means. The division of allocation heuristic can also be modified to focus more on the search in promising regions at the start and converge to the allocation stage toward the end to select the optimal $x^{**}$.

Third, as mentioned earlier, the allocation stage uses the sample means and variances as estimators in the OCBA method. An alternative approach for the allocation stage is to adapt OCBA to work with the kriging predictors. This will include accounting for the correlations between the outputs of the observed points.

Finally, although the proposed algorithm applies an EI-based criterion for the search stage, modifications can be applied to this stage of the algorithm to focus on a quantile-based criterion similar to the EQI. This can be done by replacing the EI with the analytical form of the EQI, Equation (8) of Picheny *et al.* (2010). The allocation

**Table 2.** Starting parameters

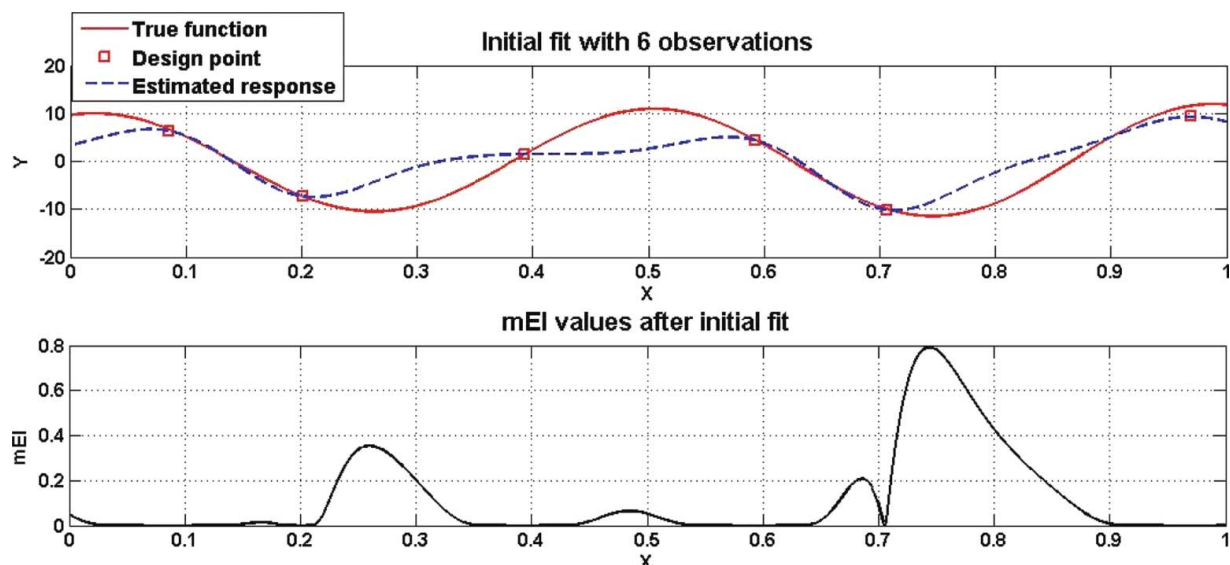| Parameter | Value |
|---|---|
| $T$ | 360 |
| $B$ | 40 |
| $n_0$ | 6 |
| $r_{\min}$ | 10 |

**Table 3.** Computing budget distribution

| Iteration | Search stage budget (number of replications allocated) | Allocation stage budget (number of replications allocated) |
|---|---|---|
| 1 | 30 | 10 |
| 2 | 20 | 20 |
| 3 | 10 | 30 |

stage will be similarly adjusted for the sample quantiles in place of the means.

## 5. Two test functions

To illustrate the proposed sequential two-stage approach, numerical examples on the one-dimensional function and a two-dimensional tetramodal function are conducted. In the first example, the algorithm is illustrated step by step to show its workings and properties. In the second example, the performance of the proposed algorithm is compared with the EQI approach of Picheny *et al*. (2010).

### 5.1. *One-dimensional test function (illustration of algorithm)*

We applied the proposed algorithm to the one-dimensional example in Section 3.3 with the noise variance function $3(1 + x)^2$ to further illustrate how the algorithm works. An LHD was used for the initial fit and the starting parameters adopted are given in Table 2.

With these parameters' settings, three iterations of the algorithm were executed, with each iteration's computing budget distribution between the sampling and allocation stages as shown in Table 3.

Figure 7 shows the estimated response after the initial fit and the modified EI function values. Although the initial fit's estimation error is high, the estimated response correctly identifies the presence of two local minima, which is also reflected in the location of the two highest peaks in the modified EI plot.

At the end of iteration 1, the point with the highest modified EI value was sampled with 30 replications and 10 replications were given to the two points near the true global minimum as shown in the top plot of Fig. 8. With two points placed near the global minimum, spatial uncertainty in that region was reduced drastically, and as the bottom plot of Fig. 8 shows, the modified EI value now had a main peak near the other local minimum. This highlights the advantage of using the "division of labor" approach. The modified EI criterion is quick to move on from a local minimum and into another, while the task of refining the estimates at points with the best responses is handled by the OCBA technique.

In the subsequent iterations, the search swings from the local minimum on the left back to the global minimum in the final iteration (as shown in Fig. 9). This example also highlights the fact that as the algorithm progresses, the number of points "competing" for the best response will likely increase. Therefore, the amount of computing



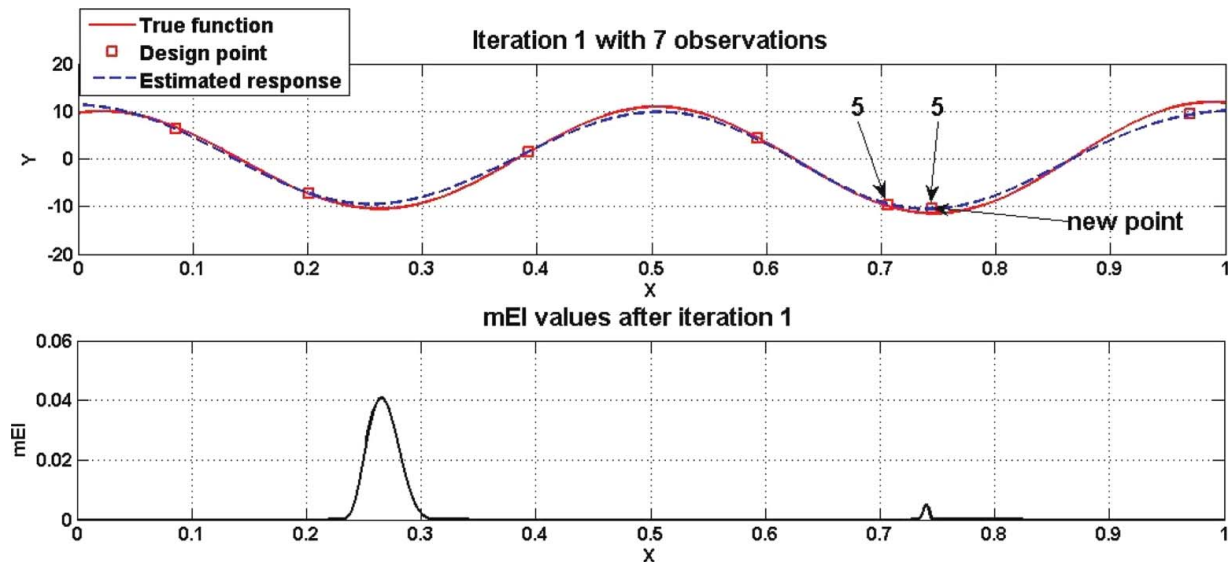**Fig. 7.** Initial fit (color figure provided online).

**Fig. 8.** Iteration 1 (color figure provided online).

budget set aside for the allocation stage needs to increase per iteration to keep up with the increasing number of sampled points.

A sensitivity analysis on the initial parameters $B$, $n_0$, and $r_{min}$ was subsequently conducted. $T$ was not varied as it was assumed to be constrained by the budget. We varied the three parameters by 25, 33, and 33% respectively. The results show that at the same noise level, the ability of the algorithm to locate a minimum that is within $\pm 0.02$ of the true optimum is robust to changes in both $B$ and $r_{min}$. However, reducing the initial design size $n_0$ reduces the chance of locating the global optimum. In this case, a smaller $B$ and $r_{min}$ is preferred, providing more runs for the search. When the variability increases to $10(1 - x)^2$, the number of times it mistakenly identifies the local minimum on the left increases across all settings of the parameters. In this higher variance case, a larger $B$ is preferred as it improves the initial fit, providing clearer distinctions between the local optimums. This increases the chance of locating a global minimum that is within $\pm 0.02$ of the true optimum. A larger $B$ also increases the weight on allocation at the start of the algorithm, providing more replications to drive down the noise and improve the initial fit.

Overall, in this example, we observe the following. Increasing $n_0$ improves the initial fit of the model, but this has to be traded off with the noise level and $B$. When the noise is low, smaller values of $B$ and $r_{min}$ are reasonable. When the variability is high, a higher $B$ is preferred. In
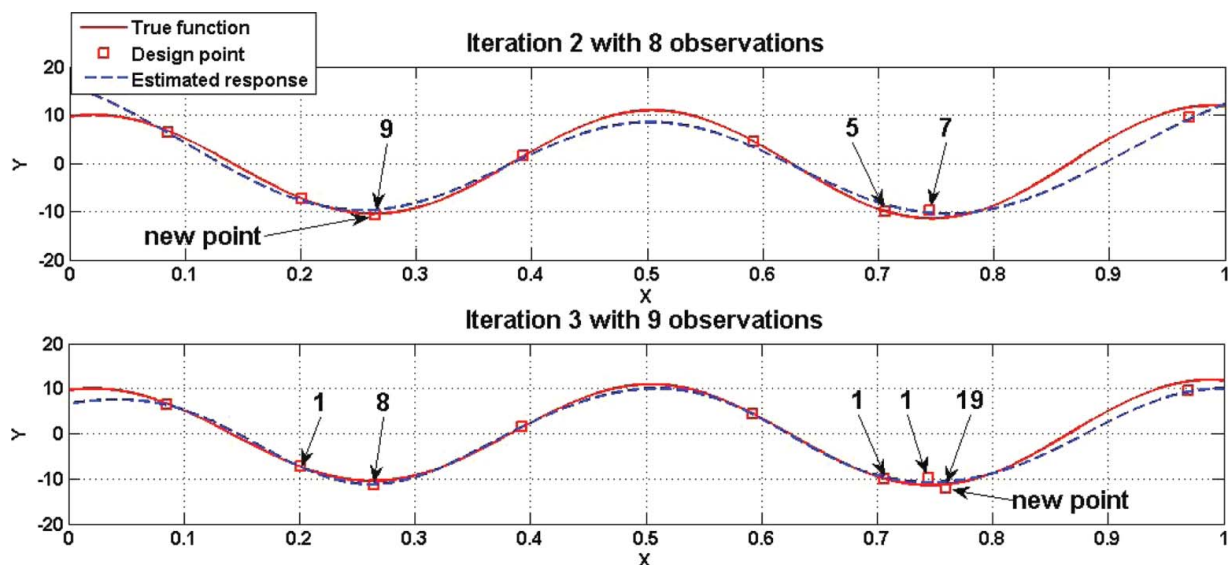


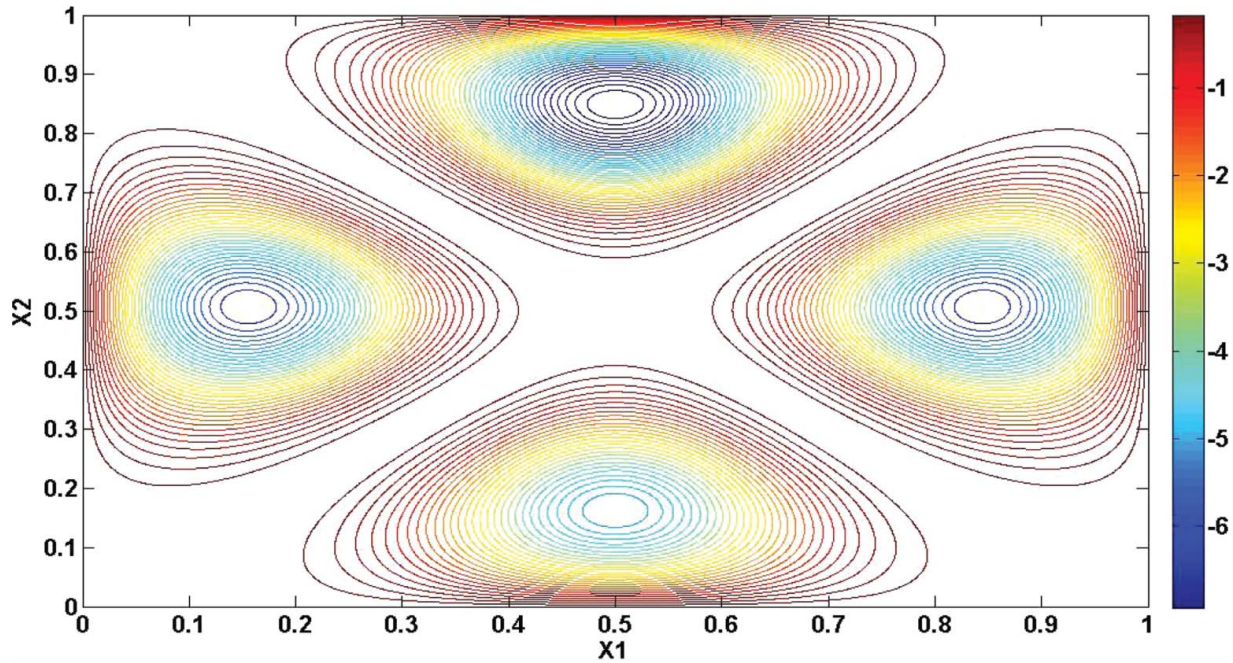**Fig. 9.** Iterations 2 and 3 (color figure provided online).

**Fig. 10.** Contour plot of test function (color figure provided online).

general, the recommendations provided in Section 4.3 are conservative and can perform well in both noise situations.

## 5.2. *Two-dimensional function (comparative study)*

Like the proposed algorithm, EQI was developed for optimization of simulations with heterogeneous noise. Out of the two resource allocation schemes in EQI, we selected EQI with online allocation for comparison because it is less computationally demanding than the alternative. Furthermore, EQI with online allocation also has a two-stage approach that first selects a new point that maximizes the EQI criterion, followed by the allocation of additional computing resources to the point until a criterion is satisfied.

The test function used for the comparison is a modified version of the tetramodal function used by Keys and Rees (2004):

$$Z(x_1, x_2) = -5(1 - (2x_1 - 1)^2)(1 - (2x_2 - 1)^2)(4 + 2x_1 - 1)$$
$$\times \left(0.05^{(2x_1 - 1)^2} - 0.05^{(2x_2 - 1)^2}\right)^2. \quad (15)$$

Both dimensions of the test functions, $x_1$ and $x_2$, are scaled to [0, 1]. The test function's contour plot is shown in Fig. 10. The values of the local minima decrease along the $x_1$ axis. The global minimum is located at [0.85, 0.5] and has the value $-7.098$. The other three local minima are located at [0.5, 0.15], [0.5, 0.85], and [0.15, 0.5] with values $-6.041$, $-6.041$, and $-4.984$ respectively. The standard deviation of the noise is set as $1.2x_1$, which adds higher levels of noise to local minima with lower values.

The comparisons between the two algorithms were carried out for 100 macro-replications. The variance of the noise in EQI was taken from the actual function while the proposed algorithm was run under two settings—one with the noise variance known and the other with the noise variance estimated using sample variances. In every comparison, both algorithms started off with the same size 20 LHS design with 40 replications at each point and the same set of initial sample means. The 50th percentile was used in the EQI criterion. Other algorithm parameters are shown in Table 4.

In EQI, $T_{20}$ refers to the post initial fit budget and $\gamma$ is a user-defined parameter with value from zero to one that controls how much available budget is allocated to the newly sampled point. $t_e$ refers to the "step size," which is the number of replications allocated to an existing point or used to sample a new point with every EQI action. $\gamma$ was set to a value quoted by Picheny *et al.* (2010) while a step size of 10 was chosen to match $r_{min}$ in the proposed algorithm.

The algorithms were first evaluated on how close they were able to get to the true location and value of the global minimum. Table 5 shows the results.

**Table 4.** Algorithms' settings

| *Proposed algorithm* | | *EQI* | |
|---|---|---|---|
| *Parameter* | *Setting* | *Parameter* | *Setting* |
| $T$-$n_0 B$ | 200 | $T_{20}$ | 200 |
| $B$ | 40 | $t_e$ | 10 |
| $r_{min}$ | 10 | $\gamma$ | 0.5 |

**Table 5.** Performance of algorithms

| Evaluation criterion | EQI (noise variance known) | | Proposed algorithm (noise variance known) | | Proposed algorithm (noise variance estimated) | |
|---|---|---|---|---|---|---|
| | Average | Standard deviation | Average | Standard deviation | Average | Standard deviation |
| $\left|x^{**}_{\text{pred}} - x^{**}_{\text{true}}\right|$ | 0.318 | 0.253 | 0.318 | 0.250 | 0.312 | 0.252 |
| $\left|\hat{Z}(x^{**}_{\text{pred}}) - Z(x^{**}_{\text{true}})\right|$ | 1.669 | 0.998 | 1.652 | 0.998 | 1.669 | 0.991 |

Statistically, there was no significant difference between the three sets of results. The proposed algorithm was able to match EQI's performance in both evaluation criteria. For this example, it showed that the proposed algorithm was able to achieve a level of accuracy similar to EQI. Furthermore, the proposed algorithm achieved it with less computational effort than EQI. Taking the number of matrix inversions as a measure of computation effort required, we saw that the proposed algorithm required five inversions in total per macro-replication—one for the initial fit and four more for the subsequent iterations. As for EQI, there were a total of 200 001 inversions per macro-replication. With the design space discretized into a 100 by 100 grid, the EQI value at every one of these grid points would have to be evaluated which would involve the inversion of a $n + 1$ by $n + 1$ matrix, where $n$ is the number of sampled points.

Another advantage of the proposed algorithm over EQI is that the noise variance function need not be known. In cases where the noise variance function is unknown, EQI requires that the noise variance function be estimated. Unlike EQI, the proposed algorithm has no such requirement as it uses sample variances as estimates of the variance of the noise at sampled points.

On the other hand, the proposed algorithm's division of computing resource heuristic is not as flexible as EQI's online allocation scheme. Although the proposed algorithm is able to adaptively allocate replications among sampled points, the proposed algorithm's computing budget for every iteration and the total number of iterations are fixed prior to the start of the algorithm. In EQI's case, the computing resource distribution between its sampling and replication allocation stages is "online" rather than predetermined, so the resource distribution scheme makes use of increasing information as the EQI algorithm progresses while accounting for the diminishing computing budget. Other allocation schemes where the budget for each iteration $B(i)$ is dynamically determined can be explored for our algorithm and this is an area for future research.

## 6. Ocean liner case study revisited

We applied our proposed algorithm to the fuel management problem described in Section 2 to determine the bunker inventory safety levels for the AEX service. Similar to Yao

*et al.* (2012), the overall objective for this problem is to minimize the total bunker fuel-related costs. We adopt the cost model in Yao *et al.* (2012) with an additional penalty function to describe the impact on the liner operation if a fuel-out situation occurs at sea. Here we set the penalty function to be the product of the bunker capacity, bunker price at the next port, and a penalty factor. We set the penalty factor value at two. This reflects a high approximate cost of sending out an additional vessel to fill up at sea and can be adjusted based on individual company's costs and penalties. For example, if a ship with a bunker capacity of 3000 tons ran out of fuel before the port Yantian (where bunker costs \$460 per metric ton), the penalty cost will amount to $3000 \times 460 \times 2 = \$2760\,000$.

To account for the actual consumption rate variation, a noise term is added to the original empirical model: $F' = k_1 V^3 + k_2 + \varepsilon(V)$. Here we also assume that noise $\varepsilon$ is a function of the ship speed. In addition, based on real data obtained from a shipping company, we observe the following three characteristics.

1. The noise term follows a normal distribution with a zero mean and a standard deviation that increases with ship speed.
2. The coefficient of variation, $CV$, the ratio of standard deviation to mean daily bunker consumption rate, is a constant.
3. This constant is different for different sized ships, where the bigger the ship, the smaller the number.

Hence the noise function is a normal distribution, $N(0, (CV \times F(V))^2)$.

Bunker prices at each port were obtained from Bloomberg for the period of January 2009 to July 2009. The prices were estimated to follow a normal distribution with means given in Table 6. The variances were also estimated from the same set of data, and as they were approximately similar for all ports, the same variance, $\sigma^2 = 15$, was used for all ports.

First, we compared our proposed approach with the results of a brute-force search. A small computing budget of 80 was used, and we assumed a single bunker inventory safety level for all 15 legs of the route to facilitate the feasibility of the search. Table 7 provides relevant parameters for the AEX service obtained from the shipping company, some of which were estimated from the company's data.

**Table 6.** Bunker price parameters for ports in AEX route

| Ports | Mean price ($/ton) |
|---|---|
| Hakata | 456 |
| Kwangyang | 458 |
| Pusan | 465 |
| Shanghai | 463 |
| Kaoshiung | 465 |
| Hong Kong | 459 |
| Yantian | 460 |
| Singapore | 471 |
| Rotterdam | 463 |
| Hamburg | 468 |
| Thamesport | 471 |
| Colombo | 469 |
| Singapore | 459 |
| Hong Kong | 462 |
| Kaoshiung | 465 |

**Table 7.** Parameters for ports in AEX route

| Parameters | Values |
|---|---|
| Bunker capacity | 3000 tons |
| Number of port calls | 15 |
| Coefficient of variation | 0.07 |
| Variance of bunker price | 15 $\$^2$/ton$^2$ |
| $k_1$ | 0.007 297 tons/(day × knots$^3$) |
| $k_2$ | 71.4 tons/day |
| Ship size | 5000 TEU |

Here the bunker inventory safety level was defined as the percentage of the total bunker capacity (3000 metric tons). For the brute-force search, we observed 16 evenly distributed points from 0.5 to 10%, with five replications at each point. Table 8 provides the optimal safety levels obtained by the two methods, and the expected optimal costs were estimated using an additional 50 replications conducted at the optimal levels. From the table, we see that with the same computing budget, our proposed approach is able to obtain a lower cost solution. *T*-test results confirm that the expected costs obtained from the proposed approach is statistically lower (at $\alpha = 0.05$). To conduct a more comprehensive brute-force search, we observed 51 evenly distributed points in the same region, each with five replications. The last row of Table 8 shows the optimal results with the expected optimal costs computed based on 50 replications. Comparing our approach with the more comprehensive brute-force search, we see that our proposed approach was able to obtain a solution close to the solution from the comprehensive search (although statistically different at $\alpha = 0.05$) with less than 35% of the expense.

Overall, the results suggest that the proposed two-stage approach can more effectively distribute the computing budget to promising regions, finding better solutions with a fixed budget or finding equivalent results with a much smaller budget. This can be very useful for the problems with limited computing budgets.

To more realistically solve this problem, we relaxed the assumption on the bunker inventory safety level and considered different safety levels for different legs. This is more reasonable as the distances for each leg within the service can differ by thousands of nautical miles. With 15 different legs within the service, we considered this as a 15-dimension decision problem. A computing budget of 3000 replications was used, which is a reasonable amount reflecting the typical limited time that decision makers have in the light of changing prices and short dynamic time horizon. Two thousand replications were spread evenly among a $n_0 = 200$ (~10k) initial space-filling LHD, and the remaining were allocated according to our proposed approach described in Section 4.3. Table 9 shows the results.

As expected, bunker safety inventory levels increased with increasing distances between port calls. With varying safety levels, the optimal total cost was lower than the single safety level at $1.013\,41 \times 10^8$. An interesting observation is that the increase in safety levels was not linear—the increments in safety levels decreased as distances increased. This was mainly due to bunker fuel management decisions in each rolling-horizon period. As observed in Yao *et al.* (2012), the time windows of the shipping schedule have a major impact on the speeds of each leg and hence the bunkering decisions. As observed in the simulation, the resulting bunkering ports were usually ports before a long leg. As the optimal bunkering amounts at these ports were generally high (due to the long leg ahead and quantity-based discounts considered), the influence of the safety levels was reduced. Moreover, ships speed up (or slow down) based on the time windows of the schedule, thereby increasing (decreasing) consumption. It was also observed that legs with tight time windows had on average higher travel speeds, thus requiring higher safety levels. The results here clearly illustrate the interrelatedness of

**Table 8.** Optimal results for single bunker inventory level

|  | Optimal minimum bunker inventory safety level (%) | Total replications required | Expected optimal cost ($ ×10$^{08}$) | Standard error of optimal cost (×10$^{03}$) |
|---|---|---|---|---|
| Two-stage approach | 5.92 | 80 | 1.14294 | 8.43 |
| Brute-force search (16 point design) | 5.57 | 80 | 1.14302 | 9.89 |
| Brute-force search (51 point design) | 5.85 | 255 | 1.14290 | 6.22 |

**Table 9.** Optimal results for 15 different legs

| Port legs | Distance *(nautical miles)* | *Optimal bunker inventory safety level (%)* |
|---|---|---|
| Hakata–Kwangyang | 171 | 1.3 |
| Kwangyang–Pusan | 86 | 1.1 |
| Pusan–Shanghai | 475 | 2.4 |
| Shanghai–Kaoshiung | 639 | 2.6 |
| Kaoshiung–Hong Kong | 350 | 2.2 |
| Hong Kong–Yantian | 45 | 1.0 |
| Yantian–Singapore | 1500 | 3.3 |
| Singapore–Rotterdam | 8339 | 7.4 |
| Rotterdam–Hamburg | 223 | 1.6 |
| Hamburg–Thamesport | 325 | 2.1 |
| Thamesport–Colombo | 6727 | 6.8 |
| Colombo–Singapore | 1567 | 3.4 |
| Singapore–Hong Kong | 1460 | 3.1 |
| Hong Kong–Kaoshiung | 350 | 2.1 |
| Kaoshiung–Hakata | 904 | 2.8 |

operational-level bunker fuel management decisions and system-level and safety-level decisions and the necessity to consider both together in a stochastic simulation environment to ensure the robustness of the decisions.

## 7. Conclusions

In this article, we proposed a two-stage sequential framework for the optimization of stochastic simulations with heterogeneous variances. The proposed two-stage framework is based on the kriging model and iteratively incorporates the OCBA technique and the modified EI function to drive and improve the estimation of the global optimum. We first illustrated our approach with several numerical examples. The empirical results indicated that the proposed approach is effective in obtaining the optimal solutions and required less computing time than other kriging-based optimization techniques proposed to address stochastic simulations. We also applied the approach to a real complex ocean liner simulation model to determine the optimal bunker inventory safety levels for a fuel management problem. The results from this problem provided invaluable insights on the inventory levels on a service route and clearly illustrated the interrelatedness of the bunker fuel management and safety levels. Future research includes developing adaptive schemes that dynamically distributes the budget for each iteration, studying in detail the convergence results of the algorithm and comparisons with other simulation optimization approaches.

The MATLAB codes to implement the proposed framework are available at http://www.ise.nus.edu.sg/staff/ngsh/download/matlabdocs/index.php/SOK.

## References

Angün, M.E., Gürkan, G., Den Hertog, D. and Kleijnen, J.P.C. (2002) Response surface methodology revised, in *Proceedings of the 2002 Winter Simulation Conference*, IEEE Press, Piscataway, NJ, pp. 377–383.

Ankenman, B., Nelson, B.L. and Staum, J. (2010) Stochastic kriging for simulation metamodeling. *Operations Research*, **58**, 371–382.

Chen, C.H., Lin, J., Yücesan, E. and Chick, S.E. (2000) Simulation budget allocation for further enhancing the efficiency of ordinal optimization. *Journal of Discrete Event Dynamic Systems: Theory and Applications*, **10**, 251–270.

Chen, C.P., Tsui, K.L., Barton, R.R. and Meckesheimer, M. (2006) A review on design, modeling and applications of computer experiments. *IIE Transactions*, **38**(4), 237–291.

Cressie, N. (1993) *Statistics for Spatial Data*, Wiley, Chichester, UK.

Fu, M.C. (2007) Are we there yet? The marriage between simulation and optimization. *OR/MS Today*, **34**(3), 16–17.

Greenwood, A.G., Vanguri, S., Eksioglu, B., Jain, P., Hill, T.W., Miller, J.W. and Walden, C.T. (2005) Simulation optimization decision support system for ship panel shop operations, in *Proceedings of the 2005 Winter Simulation Conference,* Volume 1, IEEE Press, Piscataway, NJ, pp. 2078–2086.

Gupta, A., Yu, D., Xu, L. and Reinikainen, T. (2006) Optimal parameter selection for electronic packaging using sequential computer simulations. *Transactions of the ASME, Journal of Manufacturing Science and Engineering*, **128**, 705–715.

Huang, D., Allen, T.T., Notz, W.I. and Zeng, N. (2006) Global optimization of stochastic black-box systems via sequential kriging metamodels. *Journal of Global Optimization*, **34**(3), 441–466.

Jones, D.R. (2001) A taxonomy of global optimization methods based on response surfaces. *Journal of Global Optimization*, **21**, 345–383.

Jones, D.R., Schonlau, M. and Welch, W.J. (1998) Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, **13**, 455–492.

Keys, A.C. and Rees, L.P. (2004) A sequential-design metamodeling strategy for simulation optimization. *Computers & Operations Research*, **31**, 1911–1932.

Kleijnen, J.P.C., Van Beers, W. and Van Nieuwenhuyse, I. (2011) Expected improvement in efficient global optimization through bootstrapping. *Journal of Global Optimization*, **54**(1), 59–73.

Kushner, H. and Clark, D. (1978) *Stochastic Approximation Methods for Constrained and Unconstrained Systems*, Springer, New York, NY.

Li, Y., Ng, S.H., Xie, M. and Goh, T.N. (2010) A systematic comparison of metamodeling techniques for simulation optimization in decision support systems. *Applied Soft Computing*, **10**(4), 1257–1273.

Loeppky, J.L., Sacks, J. and Welch, W.J. (2009) Choosing the sample size of a computer experiment: a practical guide. *Technometrics*, **51**(4), 366–376.

Matheron, G. (1963) Principles of geostatistics. *Economic Geology*, **58**, 1246–1266.

Mockus, J. (1994) Application of Bayesian approach to numerical methods of global and stochastic optimization. *Journal of Global Optimization*, **4**, 347–365.

Nakayama, H., Arakawa, M. and Sasaki, R. (2002) Simulation-based optimization using computational intelligence. *Optimization and Engineering*, **3**(2), 201–214.

Ng, S.H. and Yin, J. (2012) Bayesian kriging analysis and design for stochastic simulations. *ACM Transactions on Modeling and Computer Simulation*, **22**(3), Article 17.

Notteboom, T.E. and Vernimmen, B. (2009) The effect of high fuel costs on liner service configuration in container shipping. *Journal of Transport Geography*, **17**(5), 325–337.

Opsomer, J.D., Ruppert, D., Wand, W.P., Holst, U. and Hossler, O. (1999) Kriging with nonparameteric variance functionestimation. *Biometrics*, **55**, 704–710.

Pham, T. and Wagner, M. (1994) Filtering noisy images using kriging, in *Proceedings of the Fifth International Symposium on Signal Processing and its Applications*, Tew, J.D., Manivannan, S., Sadowski, D.A. and Seila, A.F. (eds), IEEE Press, Piscataway, NJ, pp. 184–191.

Picheny, V., Ginsbourger, D. and Richet, Y. (2010) Noisy expected improvement and on-line computation time allocation for the optimization of simulators with tunable precision, in *Proceedings of the Second International Conference on Engineering Optimization*, pp. 1–10.

Roshan, J.V. (2006) Limit kriging. *Technometrics*, **48**, 458–466.

Santner, T.J., Williams, B.J. and Notz, W.I. (2003) *The Design and Analysis of Computer Experiments*, Springer, New York, NY.

Shi, L. and Olafsson, S. (2000) Nested partitions method for global optimization. *Operations Research*, **48**, 390–407.

Shim, J.P., Warkentin, M., Courtney, J.F., Power, D.J., Sharda, R. and Carlsson, C. (2002) Past, present, and future of decision support technology. *Decision Support Systems*, **33**(2), 111–126.

Simpson, T.W., Peplinski, J., Koch, P.N. and Allen, J.K. (2001) Metamodels for computer-based engineering design: survey and recommendations. *Engineering with Computers*, **17**, 129–150.

Tekin, E. and Sabuncuoglu, I. (2004) Simulation optimization: a comprehensive review on theory and applications. *IIE Transactions*, **36**(11), 1067–1081.

Wan, X.T., Pekny, J.F. and Reklaitis, G.V. (2005) Simulation-based optimization with surrogate models—application to supply chain management. *Computers & Chemical Engineering*, **29**, 1317–1328.

Wu, H. and Sun, F. (2007) Adaptive kriging control of discrete-time nonlinear systems. *Control Theory & Applications*, **1**(3), 646–656.

Yao, Z., Ng, S.H. and Lee, L.H. (2012) A study on bunker fuel management for shipping liner services. *Computers & Operations Research*, **39**, 1160–1172.

Yin, J., Ng, S.H. and Ng, K.M. (2009) A study on the effects of parameter estimation on kriging model's prediction error in stochastic simulations, in *Proceedings of the Winter Simulation Conference*, Rossetti, M.D., Hill, R.R., Johansson, B., Dunkin, A. and Ingalls, R.G. (eds), IEEE Press, Piscataway, NJ, pp. 674–685.

Yin, J., Ng, S.H. and Ng, K.M. (2011) Kriging metamodel with modified nugget-effect: the heteroscedastic variance case. *Computers & Industrial Engineering*, **61**, 760–777.

## Biographies

Ning Quan earned a B.Eng. degree in Industrial and Systems Engineering at the National University of Singapore in 2012 and is currently pursuing a Ph.D. degree in Industrial Engineering at the University of Illinois at Urbana–Champaign. As an undergraduate, he did research on the use of kriging metamodels in simulation optimization. He is currently working on the topic of multi-disciplinary optimization, with applications in renewable energy systems.

Jun Yin received a B.E. degree from the University of Science and Technology of China in 2006. He is currently a Ph.D. candidate in the Department of Industrial & Systems Engineering at the National University of Singapore. He is currently also a research fellow in the Kuang-Chi Institute of Advanced Technology, Shenzhen, China. His research interests include design of experiments, stochastic simulation, and metamodeling. He is also involved with the design of complex electro-magnetic devices and their application.

Szu Hui Ng is an Associate Professor in the Department of Industrial and Systems Engineering at the National University of Singapore. She holds B.S., M.S., and Ph.D. degrees in Industrial and Operations Engineering from the University of Michigan. Her research interests include computer simulation modeling and analysis, design of experiments, and quality and reliability engineering. She is a member of IEEE and INFORMS and a senior member of IIE.

Loo Hay Lee is an Associate Professor and Deputy Head (Graduate Studies and Research) in the Department of Industrial and Systems Engineering, National University of Singapore. He received his B.S. (Electrical Engineering) degree from the National Taiwan University in 1992 and his S.M. and Ph.D. degrees in 1994 and 1997 from Harvard University. He is currently a senior member of IEEE and a member of INFORMS and ORSS. He is the Associate Editor for *IIE Transactions*, *IEEE Transactions on Automatic Control*, *Flexible Services and Manufacturing Journal*, and the *Asia Pacific Journal of Operational Research*. He is currently the co-editor for *Journal of Simulation* and is a member in the advisory board for *OR Spectrum*. His research interests include simulation-based optimization, maritime logistics and port operations, and supply chain modeling.