

Tilburg University

Constrained Optimization in Random Simulation

Angun, Ebru; Kleijnen, Jack

Publication date:
2022

Document Version
Early version, also known as pre-print

[Link to publication in Tilburg University Research Portal](#)

Citation for published version (APA):

Angun, E., & Kleijnen, J. (2022). *Constrained Optimization in Random Simulation: Efficient Global Optimization and Karush-Kuhn-Tucker Conditions*. (CentER Discussion Paper; Vol. 2022-022). CentER, Center for Economic Research.

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

No. 2022-022

**CONSTRAINED OPTIMIZATION IN RANDOM SIMULATION:
EFFICIENT GLOBAL OPTIMIZATION AND
KARUSH-KUHN-TUCKER CONDITIONS**

By

Ebru Angün, Jack P.C. Kleijnen

30 August 2022

ISSN 0924-7815
ISSN 2213-9532

Constrained optimization in random simulation: efficient global optimization and Karush-Kuhn-Tucker conditions

Ebru Angün (a) and Jack P.C. Kleijnen (b)

(a) Industrial Engineering Department

Galatasaray University

Istanbul, Turkey

E-mail: eangun@gsu.edu.tr

Orcid: <https://orcid.org/0000-0002-8199-9746>

(b) Department of Management

Tilburg School of Economics and Management

Tilburg University

Postbox 90153, 5000 LE Tilburg, Netherlands

E-mail: kleijnen@tilburguniversity.edu

Orcid: <https://orcid.org/0000-0001-8413-2366>

Abstract

We develop a novel method for solving constrained optimization problems in random (or stochastic) simulation; i.e., our method minimizes the goal output subject to one or more output constraints and input constraints. Our method is indeed novel, as it combines the Karush-Kuhn-Tucker (KKT) conditions with the popular algorithm called "efficient global optimization" (EGO), which is also known as "Bayesian optimization" and is related to "active learning". Originally, EGO solves non-constrained optimization problems in deterministic simulation; EGO is a sequential algorithm that uses Kriging (or Gaussian process) metamodeling of the underlying simulation model, treating the simulation as a black box. Though there are many variants of EGO—for these non-constrained deterministic problems and for variants of these problems—none of these EGO-variants use the KKT conditions—even though these conditions are well-known (first-order necessary) optimality conditions in white-box problems. Because the simulation is random, we apply stochastic Kriging. Furthermore, we allow for variance heterogeneity

and apply a popular sample allocation rule to determine the number of replicated simulation outputs for selected combinations of simulation inputs. Moreover, our algorithm can take advantage of parallel computing. We numerically compare the performance of our algorithm and the popular proprietary OptQuest algorithm, in two familiar examples (namely, a mathematical toy example and a practical inventory system with a service-level constraint); we conclude that our algorithm is more efficient (requires fewer expensive simulation runs) and effective (gives better estimates of the true global optimum).

Keywords: Simulation, design of experiments; simulation, statistical analysis; artificial intelligence; computational-experiments; inventory-production

JEL: C0, C1, C9, C15, C44

1 Introduction

In this paper we solve constrained optimization problems in *random* or *stochastic* simulation, combining *efficient global optimization* (EGO) and the *Karush-Kuhn-Tucker* (KKT or briefly KT) conditions; Kleijnen et al. (2022) investigates a similar solution, in the simpler context of *deterministic* simulation. EGO is a popular sequential statistical algorithm that uses *Kriging*—or *Gaussian process* (GP) metamodeling—and was originally developed for non-constrained optimization in deterministic simulation; see Jones et al. (1998). EGO is also known as *Bayesian optimization* (BO), and is related to *active learning*; see the many references in Frazier (2018). Though there are many variants of EGO (for non-constrained optimization problems and for variants of these problems), none of these EGO-variants use the KKT conditions—even though these conditions are well-known first-order necessary optimality conditions in white-box optimization; see Boyd and Vandenberghe (2009).

To denote the multiple random *simulation outputs*, we use the symbol w_h ($h = 0, 1, \dots, t-1$) where $E(w_0)$ is the goal output to be minimized and $E(w_{h'})$ ($h' = 1, \dots, t-1$) are the $(t-1)$ constrained outputs such that the output constraints are $E(w_{h'}) \leq c_{h'}$ with prespecified $c_{h'}$. For example, in inventory management, the users may wish to minimize the expected cost such that the expected out-of-stock fraction does not exceed 10%. Like Kleijnen et al. (2022), we assume k continuous *simulation inputs* x_j ($j = 1, \dots, k$). Whereas Kleijnen et al. (2022) assumes that these x_j vary across a given *rectangular* experimental domain $l_j \leq x_j \leq u_j$, we assume—like Kleijnen et al. (2010)—more general input constraints; namely, $f_g(\mathbf{x}) \leq c_g$ ($g = 1, \dots, v$) such as $x_1 \leq x_2$. Furthermore, we assume that k is relatively small (say $1 \leq k \leq 20$ (it is well-known that Kriging requires small k : “curse of dimensionality”). For example, $k = 2$ in inventory management of a single item (or stock keeping unit, SKU) through the popular (s, S) system (with $S > s$). We define the

k -dimensional vector $\mathbf{x} = (x_1, \dots, x_k)'$.

Compared with deterministic simulation, our problem is more complicated because $E[w_h(\mathbf{x})]$ must be estimated. We estimate this $E[w_h(\mathbf{x})]$ from (say) m replications; by definition, replications give *independently and identically distributed* (IID) $w_{h;r}$ and $w_{h;r'}$ with $r \neq r'$ ($r, r' = 1, \dots, m$). So, replications use nonoverlapping (independent) streams of *pseudorandom numbers* (PRN). L’Ecuyer (2015) reviews PRN generators—in sequential and parallel computing—focussing on multiple independent streams and substreams, including software libraries (Obviously, the t output types are statistically dependent or correlated instead of IID.)

In our analysis we allow for *variance heterogeneity*; i.e., $\text{var}[w_h(\mathbf{x})]$ depends on \mathbf{x} , and denotes the so-called *intrinsic* variance caused by the PRN. Jalali et al. (2017) investigates variance heterogeneity in *unconstrained* optimization solved through several variants of EGO. Ankenman et al. (2010) also investigates variance heterogeneity and derives *stochastic Kriging* (SK), but focuses on prediction (instead of optimization) so the SK performance is measured through the *mean squared prediction error* (MSPE).

Moreover, we focus on random simulation with *independent random numbers* (IRN) and $m(\mathbf{x}_i)$ replications where \mathbf{x}_i denotes \mathbf{x} in the simulated input combination i where $i = 1, \dots, n$. We use the popular *sample allocation* rule $m(\mathbf{x}_i)/m(\mathbf{x}_{i'}) = \text{var}[w(\mathbf{x}_i)]/\text{var}[w(\mathbf{x}_{i'})]$.

Note that *common random numbers* (CRN) with $m(\mathbf{x}_i) = m$ is usually analyzed by SK requiring $m > n$. We also derive *ordinary Kriging* (OK) with CRN requiring $m > 1$. CRN are meant to reduce the variance of the estimated difference between two means. In constrained optimization, however, we need to estimate the magnitude of the slacks $c_{h'} - w_{h'}(\mathbf{x})$; CRN may give a relatively bad estimator of these slacks.

The literature on Kriging for constrained optimization in random simulation is extremely limited. Wang and Ierapetritou (2018) considers a single output constraint, whereas we consider the general case of $t \geq 1$ constraints; moreover, that publication assumes IRN with $m_i = m$ and independent w_0 (objective) and w_1 (constraint), whereas our (s, S) inventory with a service-level (or fill-rate) constraint gives correlated outputs. Whereas EGO originally used either the *expected improvement* (EI) or the *probability of feasibility* (PI) infill criterion (or acquisition function), that publication proposes “feasibility-enhanced Expected Improvement” that combines the well-known *augmented expected improvement* (AEI) and PI; we use neither AEI, nor PI. A survey of optimization including random simulation is Mahévas et al. (2019). A recent publication using PI is Guidetti et al. (2022).

Our algorithm can take advantage of *parallel computing*; i.e., we implement our algorithm on a personal computer (PC) with 16 cores, and use MATLAB’s default with 12 cores. Guidetti et al. (2022) develops a parallel BO procedure for constrained optimization, which uses PI. Note that Qing et al. (2022) develops a parallel procedure for constrained multi-objective optimization using a Pareto frontier and—instead of EI—an entropy infill criterion; Luo et al. (2015) develops a parallel procedure for so-called ranking-and-selection

(R&S); obviously, those problems are related to—yet different from—our type of problems.

We numerically investigate the performance of our algorithm, in a mathematical toy example inspired by Gramacy et al. (2016, p. 2) (and also used in Pourmohamad and Lee (2020) and Pourmohamad and Lee (2021)) and an (s, S) inventory system with a service-level constraint originally investigated in Bashyam and Fu (1998).

We organize the rest of this paper as follows. Section 2 formalizes our constrained optimization problem and the corresponding KT conditions; this section also interprets these conditions as a *least squares* (LS) model. Section 3 summarizes SK, derives SK estimators of the output gradients ∇_h , and discusses choosing the number of simulation replications m_i and input combinations n . Section 4 details a novel variant of EGO for optimization in random simulation with output and input constraints; i.e., this algorithm minimizes the well-known *modified EI* (or MEI) multiplied by a cosine that quantifies how well the KT conditions hold. Section 5 designs a *Monte Carlo* (MC) experiment with the toy example that controls the degree of variance heterogeneity; Subsection 5.1 compares the performance of our algorithm and OptQuest (which is a popular proprietary simulation optimization algorithm that combines several algorithms). Section 6 investigates an (s, S) inventory discrete-event simulation model with a service-level constraint (this constraint is used because practitioners find it hard to quantify out-of-stock costs, including loss-of-goodwill; this model is the major building block in practical inventory models with multiple SKUs); Subsection 6.1 compares the performance of our algorithm and OptQuest. Section 7 presents conclusions.

2 Constrained optimization and Kuhn-Tucker conditions

In Section 1, we mentioned that we focus on

$$\begin{aligned} & \min_{\mathbf{x}} E[w_0(\mathbf{x})] \text{ subject to} \\ & E[w_{h'}(\mathbf{x})] \leq c_{h'} \text{ with } h' = 1, \dots, t-1, \\ & f_g(\mathbf{x}) \leq c_g \text{ with } g = 1, \dots, v. \end{aligned} \tag{1}$$

Obviously, $f_g(\mathbf{x}) \leq c_g$ includes the usual *box constraints*

$$\mathbf{l} \leq \mathbf{x} \leq \mathbf{u} \quad (\mathbf{l} = (l_1, \dots, l_k)', \mathbf{u} = (u_1, \dots, u_k)'). \tag{2}$$

Note that we assume that the problem uses expected values instead of probabilities such as $PR[w_{h'}(\mathbf{x}) > c_{h'}] \leq 0.10$; classic Kriging also assumes expected values instead of probabilities (Kriging with probabilities

is referenced in Kleijnen (2015, p. 222)). In practice, the analysts and their clients cannot observe the true $E[w_h(\mathbf{x})]$; they can observe the actual averages during a finite period (or short term); in the long term, “we are all dead” (as the economist Keynes said) and the system and its environment will change so the analysts should find a new solution. However, for the toy problem detailed in Section 5, we do know the true means $E[w_h(\mathbf{x})]$; i.e., we do know whether the estimated optimal solution is feasible. In practice, post-processing may obtain so many replications of the (relatively expensive) simulation that the analysts are (say) 95% certain that the expected values are estimated within 1%.

We denote the *gradient* of $w_h(\mathbf{x})$ ($h = 0, 1, \dots, t - 1$) at \mathbf{x}_* by $\nabla_h(\mathbf{x}_*) = (\partial w_h / \partial x_1, \dots, \partial w_h / \partial x_k)' |_{\mathbf{x}_*}$. A constraint is *binding* (or “active”) at a specific \mathbf{x}_* if this \mathbf{x}_* lies at the boundary between the feasible area and the infeasible area; i.e., if this \mathbf{x}_* changes \leq into $=$ in (1), so the slack becomes exactly zero. Let $A_\lambda(\mathbf{x}_*)$ denote the *index set* with the indices h'' of the *binding output* constraints at \mathbf{x}_* (obviously, this set is empty if there are no binding output constraints at this \mathbf{x}_*).

Note that the box constraint in (2) implies the $2k$ *input constraints* $x_j \leq u_j$ and $-x_j \leq -l_j$. For these constraints we define $g = 1, \dots, v = 2k$, so ∇_g is a k -dimensional vector with 1 or -1 at position j , and 0 at the $k - 1$ remaining positions. In our toy example, all input constraints are box constraints; however, in our (s, S) inventory example there is the additional constraint $s \leq S$ or $s - S \leq 0$, so the corresponding gradient is $(1, -1)'$.

Analogous to $A_\lambda(\mathbf{x}_*)$, let $A_\mu(\mathbf{x}_*)$ denote the *index set* with the indices g' of the *binding input* constraints at \mathbf{x}_* . Denote the *Lagrangian* multiplier for binding output constraint h'' by $\lambda_{h''}$, and for binding input constraint g' by $\mu_{g'}$. Altogether, the *KT stationarity conditions* are

$$-\nabla_0(\mathbf{x}_*) = \sum_{h'' \in A_\lambda(\mathbf{x}_*)} \lambda_{h''}(\mathbf{x}_*) \nabla_{h''}(\mathbf{x}_*) + \sum_{g' \in A_\mu(\mathbf{x}_*)} \mu_{g'}(\mathbf{x}_*) \nabla_{g'}(\mathbf{x}_*)$$

with $\lambda_{h''}(\mathbf{x}_*) \geq 0$ and $\mu_{g'}(\mathbf{x}_*) \geq 0$. (3)

Obviously, (3) is an *orthogonal projection* of $-\nabla_0$ onto the linear space spanned by $\nabla_{h''}$ and $\nabla_{g'}$. Because LS is also an orthogonal projection, we interpret (3) as a LS model with the estimates $\tilde{\lambda}_{h''}$ and $\tilde{\mu}_{g'}$ where we use a tilde to denote LS (LS is a mathematical criterion; in the next sections, we shall use a hat to denote an estimator that uses a statistical criterion such as likelihood or bias resulting in a *maximum likelihood estimator* (MLE) or an unbiased estimator). We define the vectors $\boldsymbol{\lambda}_{h''} = (\lambda_{h''})$ and $\boldsymbol{\mu}_{g'} = (\mu_{g'})$ where the number of elements of these vectors is $\#A_\lambda(\mathbf{x}_*)$ and $\#A_\mu(\mathbf{x}_*)$. Next we define $\boldsymbol{\nu} = (\boldsymbol{\lambda}'_{h''}, \boldsymbol{\mu}'_{g'})'$, and $\boldsymbol{\Delta} = (\nabla_{h''}, \nabla_{g'})$. Using (3), we define the LS model with the explained (dependent) variable $-\nabla_0$ and the

explanatory (independent) variables $\mathbf{\Delta}$. This model gives the LS estimator $\tilde{\boldsymbol{\nu}} = (\tilde{\boldsymbol{\lambda}}'_{h''}, \tilde{\boldsymbol{\mu}}'_{g'})'$:

$$\tilde{\boldsymbol{\nu}} = (\mathbf{\Delta}'\mathbf{\Delta})^{-1}\mathbf{\Delta}'(-\nabla_0). \quad (4)$$

This $\tilde{\boldsymbol{\nu}}$ gives the LS model of the KT conditions defined in (3):

$$-\tilde{\nabla}_0(\mathbf{x}_*) = \sum_{h'' \in A_\lambda(\mathbf{x}_*)} \tilde{\lambda}_{h''}(\mathbf{x}_*) \nabla_{h''}(\mathbf{x}_*) + \sum_{g' \in A_\mu(\mathbf{x}_*)} \tilde{\mu}_{g'}(\mathbf{x}_*) \nabla_{g'}(\mathbf{x}_*)$$

with $\tilde{\lambda}_{h''}(\mathbf{x}_*) \geq 0$ and $\tilde{\mu}_{g'}(\mathbf{x}_*) \geq 0$. (5)

To quantify the *validity* of this model, we use $\cos[-\nabla_0(\mathbf{x}_*), -\tilde{\nabla}_0(\mathbf{x}_*)]$, which we abbreviate to $\widetilde{\cos}(\mathbf{x}_*)$. Ideally, $-\nabla_0(\mathbf{x}_*)$ and $-\tilde{\nabla}_0(\mathbf{x}_*)$ *coincide* so $\widetilde{\cos}(\mathbf{x}_*) = 1$. However, in our search for the optimum, we require only that these two vectors point into the “same” direction; i.e., \mathbf{x}_* lies “near” a stationary point. So, we require that $0 \leq \widetilde{\cos}(\mathbf{x}_*) \leq 1$. Obviously, $\cos[-\nabla_0(\mathbf{x}_*), -\tilde{\nabla}_0(\mathbf{x}_*)] = \cos[\nabla_0(\mathbf{x}_*), \tilde{\nabla}_0(\mathbf{x}_*)]$. Finally, we use the following well-known formula (see, e.g., Kolman and Hill (2008, p. 350)) where \bullet is the usual symbol for the inner product of two vectors:

$$\widetilde{\cos}(\mathbf{x}_*) = \frac{\nabla_0(\mathbf{x}_*) \bullet \tilde{\nabla}_0(\mathbf{x}_*)}{\|\nabla_0(\mathbf{x}_*)\| \times \|\tilde{\nabla}_0(\mathbf{x}_*)\|}. \quad (6)$$

Unfortunately, in simulation the functions $w_h(\mathbf{x})$ are unknown; simulation does give the *input/output* (I/O) data $(\mathbf{x}_i, w_{i,h;r})$. We can use these data to estimate *metamodels* that approximate the *black-box* functions $w_h(\mathbf{x})$. Actually we use Kriging metamodels, as we explain next.

3 Kriging in random simulation

Subsection 3.1 summarizes SK. Subsection 3.2 derives Kriging estimators of the output gradients ∇_h . Subsection 3.3 discusses choosing the number of simulation replications m_i and input combinations n .

3.1 Stochastic Kriging

We start with the formulas for the SK predictor and its variance derived in Ankenman et al. (2010). We let $e_{h;r}(\mathbf{x}_i)$ —or briefly $e_{i,h;r}$ —denote the *intrinsic noise* in replication r ($r = 1, \dots, m_i$) of simulation output h ($h = 0, 1, \dots, t - 1$) at simulated “old” input combination \mathbf{x}_i ($i = 1, \dots, n$). Kriging (or GP) assumes that $e_{i,h;r}$ has a Gaussian distribution with zero mean. Averaging the m_i replications gives the average simulation output $\bar{w}_h(\mathbf{x}_i)$ —or briefly $\bar{w}_{i,h}$ —and the average intrinsic noise $\bar{e}_h(\mathbf{x}_i)$ —or $\bar{e}_{i,h}$. Let $\boldsymbol{\Sigma}_{\bar{e},h}$

denote the *covariance matrix* of $\bar{e}_{i,h}$. IRN imply that this matrix is a diagonal matrix. We ignore the covariances between the t simulation outputs when fitting the t Kriging metamodels; i.e., we use *univariate Kriging* instead of multivariate Kriging or co-Kriging (Ankenman et al. (2010) considers only univariate SK). Actually, univariate Kriging drastically simplifies the analysis and may be superior, because multivariate Kriging requires the specification of an appropriate covariance matrix and the estimation of the individual covariances in that matrix; see Kleijnen and Mehdad (2014). Let $M_h(\mathbf{x})$ denote the *extrinsic noise* in the Kriging model for output h . Kriging assumes that e is independent of M . These assumptions give the following SK *predictor* for the *new* point \mathbf{x}_* (input combination not yet simulated), which uses the *n old* simulation I/O data $(\mathbf{X}, \bar{\mathbf{w}}_h)$ where $\mathbf{X} = \mathbf{X}_{n \times k}$ and $\bar{\mathbf{w}}_h = \bar{\mathbf{w}}_{h;n}$:

$$\hat{y}_h(\mathbf{x}_*) = \hat{\mu}_h + \hat{\boldsymbol{\sigma}}_{M;h}(\mathbf{x}_*)'(\hat{\boldsymbol{\Sigma}}_{M;h} + \hat{\boldsymbol{\Sigma}}_{\bar{e};h})^{-1}(\bar{\mathbf{w}}_h - \hat{\mu}_h \mathbf{1}_n), \quad (7)$$

which *plugs-in* the MLE $\hat{\boldsymbol{\psi}}_h = \hat{\boldsymbol{\psi}}(\mathbf{X}, \bar{\mathbf{w}}_h)$ of $\boldsymbol{\psi}_h$ where $\boldsymbol{\psi}_h$ denotes the vector with the SK parameters including the mean $\mu_h = E(y_h)$ and the variance $\tau_h^2 = Var(y_h)$ so the $n \times n$ correlation matrix $\mathbf{R}_h = (\rho_{i;i';h})$ equals $\tau_h^{-2} \boldsymbol{\Sigma}_{M;h}$ with the $n \times n$ covariance matrix $\boldsymbol{\Sigma}_{M;h} = (\text{Cov}(y_{h;i}, y_{h;i'}))$ where $y_{h;i} = y_h(\mathbf{x}_i)$ and $i' = 1, \dots, n$, and the n -dimensional correlation vector $\boldsymbol{\rho}_h(\mathbf{x}_*) = \tau_h^{-2} \boldsymbol{\sigma}_{M;h}(\mathbf{x}_*)$ with the n -dimensional covariance vector $\boldsymbol{\sigma}_{M;h}(\mathbf{x}_*) = (\sigma_{h;*}) = (\text{Cov}(y_{h;*}, y_{h;i}))$ with $y_{h;*} = y_h(\mathbf{x}_*)$; $\hat{\boldsymbol{\Sigma}}_{\bar{e};h}$ is the classic unbiased estimator (defined in (9)); $\mathbf{1}_n$ denotes the n -dimensional vector with all elements equal to 1; $M(\mathbf{x})$ denotes a zero-mean stationary GP (see the subscript M of $\boldsymbol{\Sigma}_{M;h}$).

The estimated MSPE of $\hat{y}_h(\mathbf{x}_*)$ is

$$\begin{aligned} s^2[\hat{y}_h(\mathbf{x}_*)] &= \hat{\tau}_h^2 - \hat{\tau}_h^4 \boldsymbol{\rho}_h'(\mathbf{x}_*) [\hat{\tau}_h^2 \mathbf{R}_h + \hat{\boldsymbol{\Sigma}}_{\bar{e};h}]^{-1} \boldsymbol{\rho}_h(\mathbf{x}_*) \\ &+ \hat{\delta}_h^2 [\mathbf{1}_n' (\hat{\tau}_h^2 \mathbf{R}_h + \hat{\boldsymbol{\Sigma}}_{\bar{e};h})^{-1} \mathbf{1}_n]^{-1} \\ &\text{with } \hat{\delta}_h = 1 - \mathbf{1}_n' (\hat{\tau}_h^2 \mathbf{R}_h + \hat{\boldsymbol{\Sigma}}_{\bar{e};h})^{-1} \boldsymbol{\rho}_h(\mathbf{x}_*) \hat{\tau}_h^2. \end{aligned} \quad (8)$$

Our equation (8) equals (27) in Ankenman et al. (2010)—apart from $\hat{\boldsymbol{\delta}}$ in the latter equation, which is actually a scalar $\hat{\delta}$. Our equation also resembles (9) in Jalali et al. (2017).

Assuming a *constant* intrinsic variance σ_h^2 is unrealistic, so we allow \mathbf{x} to affect both the mean and the variance of the distribution of $w_h(\mathbf{x})$; i.e., the intrinsic error has *heterogeneous* variances $\sigma_h^2(\mathbf{x}) = var[w_h(\mathbf{x})]$.

Following Ankenman et al. (2010), we estimate $\text{var}[\bar{w}_{i,h}]$ through

$$s^2(\bar{w}_{i,h}) = \frac{s^2(w_{i,h})}{m_i} \text{ with } s^2(w_{i,h}) = \frac{\sum_{r=1}^{m_i} (w_{i,h;r} - \bar{w}_{i,h})^2}{m_i - 1}$$

$$\text{with } \bar{w}_{i,h} = \frac{\sum_{r=1}^{m_i} w_{i,h;r}}{m_i}. \quad (9)$$

SK with IRN uses these $s^2(\bar{w}_{i,h})$ ($i = 1, \dots, n$) on the *diagonal* of $\widehat{\Sigma}_{\bar{e};h}$ in (7) and (8). These $s^2(\bar{w}_{i,h})$ require only $m_i > 1$ ($s^2(w_{i,h})$ uses the denominator $m_i - 1$, so it is unbiased, unlike the MLE with denominator m_i).

In Appendix 1 we discuss CRN and its analysis through either SK or OK.

3.2 Kriging estimators of output gradients

Because the KT conditions use the (unknown) gradients $\nabla_h(\mathbf{x}_*)$, we may use $\widehat{y}_h(\mathbf{x}_*)$ to estimate these $\nabla_h(\mathbf{x}_*)$. Inspired by the derivation of $\nabla[\widehat{y}_{\text{OK};h}(\mathbf{x}_*)]$ in Kleijnen et al. (2022) (also see Erickson et al. (2021)), we notice that in (7) all factors except $\widehat{\sigma}_{M;h}(\mathbf{x}_*)$ depend only on the old I/O simulation data $(\mathbf{X}, \bar{\mathbf{w}}_{h;n})$, which also determine $\widehat{\psi}_h = \widehat{\psi}(\mathbf{X}, \bar{\mathbf{w}}_{h;n})$. So, we may write

$$\nabla[\widehat{y}_h(\mathbf{x}_*)] = \nabla[\widehat{\sigma}_{M;h}(\mathbf{x}_*)' \mathbf{c}_h]$$

$$\text{with } \mathbf{c}_h = (\widehat{\Sigma}_{M;h} + \widehat{\Sigma}_{\bar{e};h})^{-1} (\bar{\mathbf{w}}_h - \widehat{\mu}_h \mathbf{1}_n). \quad (10)$$

Actually, we use the most popular type of correlation function in Kriging metamodels of simulation models; namely, the *anisotropic Gaussian correlation function* with parameter vector $\boldsymbol{\theta}_h = (\theta_{h;1}, \dots, \theta_{h;k})'$ where $\theta_{h;j} \geq 0$:

$$\rho(\boldsymbol{\theta}_h, \mathbf{x}, \mathbf{x}') = \prod_{j=1}^k \exp[-\theta_{h;j} (x_j - x'_j)^2] = \exp[-\sum_{j=1}^k \theta_{h;j} (x_j - x'_j)^2]. \quad (11)$$

This correlation function implies $\boldsymbol{\psi}_h = (\mu_h, \tau_h^2, \theta_{h;1}, \dots, \theta_{h;k})'$. This function and (10) give

$$\frac{\partial[\widehat{y}_h(\mathbf{x}_*)]}{\partial x_{*;j}} = -2\tau_h^2 \theta_{h;j} \{ \sum_{i=1}^n c_{h;i} (x_{*;j} - x_{i;j}) \exp[\sum_{j'=1}^k -\theta_{h;j'} (x_{*;j'} - x_{i;j'})^2] \}, \quad (12)$$

where $c_{h;i}$ is element i ($i = 1, \dots, n$) of \mathbf{c}_h in (10) so (12) defines element j ($j = 1, \dots, k$) of $\nabla[\widehat{y}_h(\mathbf{x}_*)]$.

3.3 Number of replications and input combinations

The number of replications m_i does not need to be a constant m , as we assume that the intrinsic noise varies with the simulation input \mathbf{x}_i ($i = 1, \dots, n$). Intuitively, a good *sample allocation rule* stipulates that m_i increases as $\sigma_h^2(\mathbf{x}_i)$ increases. Because $\widehat{y}_h(\mathbf{x}_*)$ uses the sample averages $\bar{w}_{h;i}$, our algorithm obtains so many

replications that these averages have the same estimated intrinsic variance. Let $m_{i;s}$ denote the available number of replications for input combination i in stage (or iteration) s . First we discuss $m_{i;s}$ for $s = 1, 2, \dots$; next we discuss the initial stage $s = 0$ (we assume that readers are not confused by s denoting the estimated standard deviation in other contexts in our paper).

Let $\hat{m}_{b;s}$ denote the estimated *desired* number of replications in stage s for input combination b —with $b = 1, \dots, B$ and $1 \leq B \leq n$ —“near” the estimated boundary (17). (Adding replications only to combinations near the boundary implies that our algorithm does not obtain expensive simulation replications for points “far away” from this boundary; consequently, the simulation observations not close to this boundary get small weights when predicting outputs near the boundary, as the extrinsic noise and the intrinsic noise are high).

If $B = 1$, then one replication is added to combination b ; next, our algorithm uses all available observations $m_{b;s}$ to recompute $\bar{w}_{b;h;s}$ and $s^2(\bar{w}_{b;h;s})$, etc. If $B > 1$, then we let $m_{\min;s}$ denote the smallest number of actually simulated replications among the B combinations near the boundary in stage s (obviously, all t output types use the same input). Let $b_{\min;s}$ denote the input combination corresponding with $m_{\min;s}$. We then compute the maximum of $\hat{m}_{b;s}$ over the t outputs, as follows:

$$\hat{m}_{b;s} = \max_h \left\lceil \frac{s^2(\bar{w}_{b;h;s})}{s^2(\bar{w}_{b_{\min;s};h;s})} m_{\min;s} \right\rceil \quad (13)$$

where $\lceil \cdot \rceil$ is the usual symbol for the ceiling function (the floor function $\lfloor \cdot \rfloor$ would imply that no replications are added if the B values of $s^2(\bar{w}_{b;h;s})$ differ only slightly, so our algorithm would get stuck). If (13) gives $\hat{m}_{b;s} > m_{b;s}$, then we add a single replication to combination b . This new observation on $w_{b;h;s}$ changes $s^2(w_{b;h;s})$; therefore, we recompute $\hat{m}_{b;s}$ through (13). If none of the B combinations requires adding an observation at this stage, then our algorithm uses all available observations $m_{b;s}$.

It is well known that the variance of a variance estimator (such as $s^2(w_{b;h;s})$) may be relatively high. Therefore we start with $m_{\min;s} = 10$ in the toy example, but in the inventory example we start with $m_{\min;s} = 2$ because $E[s^2(w_{b;h;s})]$ is small

Next we discuss the *preprocessing* stage, testing whether the initial design gives acceptable Kriging models for the t outputs. We start with a small initial number of replications m_{\min} and input combinations n_{\min} ; e.g., in the toy problem, we choose $m_i = m_{\min} = 10$ and $n = n_{\min} = 6$. (Ankenman et al. (2010) gives the rule-of-thumb $m_i \geq 10$, and Wang and Ierapetritou (2018) uses $m_i = m = 50$ as default; Pearce et al. (2022) starts with 5 replications). Kleijnen and Van Beers (2022) derives a *leave-one-out cross-validation* (LOO-CV) test for deterministic simulation; now we derive a similar test for random simulation. So, we delete input combination i (with $i = 1, \dots, n$) and its m_i replications for the t outputs $w_{i;h;r}$ (with $h = 0$,

..., $t - 1$ and $r = 1, \dots, m_i$), and compute the Kriging estimators $\hat{y}_{-i,h}$ and $s^2(\hat{y}_{-i,h})$ from the remaining I/O data $(\mathbf{X}_{-i}, \mathbf{W}_{-i})$; i.e., we use $(\mathbf{X}_{-i}, \mathbf{W}_{-i})$ instead of (\mathbf{X}, \mathbf{W}) in (7) and (8). Next, we restore the complete I/O data set (\mathbf{X}, \mathbf{W}) , and delete another I/O combination i' —until we have left out each combination once. We construct two-sided CIs with “per comparison” error rates. Because these CIs should hold for all t output types and all n points, we apply the “Bonferroni” inequality; i.e., we divide by $n \times t$ to obtain an “experimentwise” error rate α_E (e.g., $\alpha_E = 0.20$), and we use $z_{1-[\alpha_E/(2nt)]}$, which is the $1 - \alpha_E/(2nt)$ quantile of the standard normal variable z . Finally, if

$$\max_{i,h} \left[\frac{|\bar{w}_{i,h} - \hat{y}_{i,h}|}{\sqrt{s^2(\bar{w}_{i,h}) + s^2(\hat{y}_{i,h})}} \right] > z_{1-[\alpha_E/(2nt)]}$$

with $i = 1, \dots, n$ and $h = 0, \dots, t - 1$,

(14)

then we reject the t Kriging metamodels, and apply the allocation rule including (13). (Our preprocessing is inspired by Pedrielli et al. (2020), but differs in many details; e.g., that publication does not use (14).)

If \mathbf{x} is subject to box constraints only (see, e.g., the toy problem), then—in this initial stage—we use the most popular design type in Kriging—namely, *Latin hypercube sampling* (LHS). If \mathbf{x} is subject to an additional constraint such as $s \leq S$ in the (s, S) inventory simulation, then we still can use LHS—provided we *reparameterize* the problem; i.e., we use LHS with $x_1 = s$ and $x_2 = Q = S - s$ (a similar reparameterization is used in Law (2015, p. 636) for an (s, S) inventory simulation without a service-level constraint, analyzed through a local second-order polynomial instead of a global Kriging metamodel). In LHS we use *midpoints*, which is an option in MATLAB’s function *lhsdesign*. This option fixes the distances among the n projected values of a specific input j (with $j = 1, \dots, k$). We expect that this option gives better $\hat{\psi}_h$ in the anisotropic correlation function (11), because these projected values are not clustered; actually, these values lie apart from each other at distances that are multiples of $1/n$.

We start with $n = (k + 1)(k + 2)/2$ if $k \leq 6$ and $n = 5k$ if $k > 6$, following Tao et al. (2020)’s n for sequential optimization in deterministic simulation.

Note that Loeppky et al. (2009) proposes $n = 10k$, as a rule-of-thumb for sensitivity analysis via Kriging. This rule is also used in Wang and Ierapetritou (2018) for sequential optimization in random simulation. For $n = 5k$, Tao et al. (2020) refers to Liu et al. (2017). Pandita et al. (2019) states that $5k$ to $10k$ is used, and also refers to Sóbester et al. (2005) (the latter recommends n equal to 35% of the total computer budget, which is assumed to be fixed). For sequential optimization in deterministic simulation with $k = 2$, Pourmohamad and Lee (2021) uses $n = 20$ in the “modified Townsend” example, but uses $n = 10$ in their toy example that we shall change into a random simulation.

4 Efficient global optimization

Optimization should avoid convergence to a local optimum. Therefore, Jones et al. (1998)—the seminal publication on EGO for *unconstrained* optimization in *deterministic* simulation—derives an estimator of EI at \mathbf{x} , and uses this $\widehat{\text{EI}}(\mathbf{x})$ to select the next point to be simulated. For *unconstrained* optimization in *random* simulation with heterogeneous noise, Jalali et al. (2017) evaluates six variations of EI; actually, we define our estimator $\widehat{\text{MEI}}$ in (16) below. For *constrained* optimization in *deterministic* simulation, Kleijnen et al. (2022) uses an estimator of EI multiplied by an estimator of the cos defined in (6). Now we focus on *constrained* optimization in *random* simulation.

The simulation estimators $\bar{w}_{i:h}$ may be rather inaccurate, given $\sigma(\mathbf{x}_i)/m_i$. Therefore our algorithm uses the SK predictor \hat{y}_h (defined in (7)), which has a lower MSPE than $\bar{w}_{i:h}$ has—assuming there are “enough” old points “near” \mathbf{x}_i , and the bias of the Kriging predictor is “small”. However, our algorithm may give estimated optimal solutions (say) $\hat{\mathbf{x}}_{\min}$ that are actually infeasible; e.g., if \hat{y}_h has a symmetric distribution, then $PR[\hat{y}_{h'}(\hat{\mathbf{x}}_{\min}) > c_{h'}]$ may be as high as 50%. Therefore, our algorithm accepts $\hat{\mathbf{x}}_{\min}$ as the optimal input combination only if

$$\hat{y}_{h'}(\hat{\mathbf{x}}_{\min}) + z_{(1-\alpha)} s[\hat{y}_{h'}(\hat{\mathbf{x}}_{\min})] \leq c_{h'} \text{ with } \alpha = \alpha_{\text{infe}} \quad (15)$$

where α_{infe} is specified by the users to control the probability of an infeasible solution selected by our algorithm ($s(\hat{y}_{h'})$ is the square root of $s^2(\hat{y}_{h'})$ defined in (8)). For example, if α_{infe} is 50%, 10%, 1%, or 0.1%, then the “safety factor” $z_{(1-\alpha)}$ in (15) is 0.0000, 1.2816, 2.3263, or 3.0902.

Note that $s^2(\hat{y}_{h'})$ underestimates the true MSPE, because it assumes known Kriging parameters ψ_h and an unbiased predictor $\hat{y}_{h'}$. Because (15) uses the Studentized slack $\{[\hat{y}_{h'}(\hat{\mathbf{x}}_{\min}) - c_{h'}]/s[\hat{y}_{h'}(\hat{\mathbf{x}}_{\min})]\}$, we would like to use the Student t -statistic with (say) v degrees of freedom t_v . However, we do not know the correct value of v , so we use z ; of course, $z_{(1-\alpha)} < t_{v;(1-\alpha)}$.

So, our algorithm estimates the best feasible goal output by $\hat{y}_{0; \min} = \min_{1 \leq i \leq n} \hat{y}_0(\mathbf{x}_i)$ if (15) holds (our initial design and PS satisfy all input constraints); else $\hat{y}_{0; \min} = \infty$. This $\hat{y}_{0; \min}$ gives

$$\begin{aligned} \widehat{\text{MEI}}(\mathbf{x}) &= (\hat{y}_{0; \min} - \hat{y}_0(\mathbf{x})) \Phi \left(\frac{\hat{y}_{0; \min} - \hat{y}_0(\mathbf{x})}{s[\hat{y}_0(\mathbf{x})]} \right) + \\ & s[\hat{y}_0(\mathbf{x})] \phi \left(\frac{\hat{y}_{0; \min} - \hat{y}_0(\mathbf{x})}{s[\hat{y}_0(\mathbf{x})]} \right). \end{aligned} \quad (16)$$

Obviously, $\widehat{\text{MEI}}(\mathbf{x})$ is positive, but it may become so small that we stop our search; actually, our algorithm stops as soon as $\widehat{\text{MEI}}(\mathbf{x}) < \varepsilon \times |\hat{y}_{0; \min}|$ where ε is some small positive number; e.g., 0.01.

To estimate how well the KT conditions hold at \mathbf{x} , we need to estimate which output constraints are *binding* at \mathbf{x} . Therefore we use the following decision rule:

$$\text{if } \frac{|\widehat{y}_{h'}(\mathbf{x}) - c_{h'}|}{s[\widehat{y}_{h'}(\mathbf{x})]} \leq z_{1-\alpha/2}, \text{ then } h' \in A_\lambda(\mathbf{x}). \quad (17)$$

Our algorithm searches for the solution (say) $\widehat{\mathbf{x}}_o$ of

$$\max_{\mathbf{x}} \widehat{\text{MEI}}(\mathbf{x}) \times \cos\{\widehat{\nabla}[\widehat{y}_0(\mathbf{x})], \widetilde{\nabla}[\widehat{y}_0(\mathbf{x})]\}, \quad (18)$$

while satisfying the $\#A_\lambda(\mathbf{x})$ binding output constraints in (17) and all v input constraints in (1); in (18)

$\cos\{\widehat{\nabla}[\widehat{y}_0(\mathbf{x})], \widetilde{\nabla}[\widehat{y}_0(\mathbf{x})]\}$ equals

$$\widetilde{\text{cos}}[\mathbf{x}] = \frac{\widehat{\nabla}[\widehat{y}_0(\mathbf{x})] \bullet \widetilde{\nabla}[\widehat{y}_0(\mathbf{x})]}{\|\widehat{\nabla}[\widehat{y}_0(\mathbf{x})]\| \times \|\widetilde{\nabla}[\widehat{y}_0(\mathbf{x})]\|} \quad (19)$$

where $\widetilde{\nabla}[\widehat{y}_0(\mathbf{x})]$ is the LS estimator

$$-\widetilde{\nabla}[\widehat{y}_0(\mathbf{x})] = \sum_{h'' \in A_\lambda(\mathbf{x})} \widetilde{\boldsymbol{\lambda}}_{h''}(\mathbf{x}) \widehat{\nabla}[\widehat{y}_{h''}(\mathbf{x})] + \sum_{g' \in A_\mu(\mathbf{x})} \widetilde{\boldsymbol{\mu}}_{g'}(\mathbf{x}) \nabla_{g'}(\mathbf{x})$$

with $\widetilde{\boldsymbol{\lambda}}_{h''}(\mathbf{x}) \geq 0$ and $\widetilde{\boldsymbol{\mu}}_{g'}(\mathbf{x}) \geq 0$. (20)

Note that deterministic simulation uses EI so EGO does not select old points \mathbf{x}_i (because $\widehat{\text{EI}}(\mathbf{x}_i) = 0$) for the next simulation, whereas random simulation may give $\widehat{\text{MEI}}(\mathbf{x}_i) > 0$ so our algorithm may revisit old points. More precisely, if $\mathbf{x} = \mathbf{x}_i$, then the first term in (16) is not positive and (because $\widehat{y}_0(\mathbf{x}_i)$ is not an exact interpolator) the second term is positive.

To find the solution $\widehat{\mathbf{x}}_o$ for the problem defined in (18) through (20), we use MATLAB's *pattern search* (PS). PS is detailed on

<https://nl.mathworks.com/help/gads/patternsearch.html>.

That website describes several types of optimization problems, including “Nonlinear Constraints” (we quote text on the PS website so the readers can easily retrieve the relevant text). PS has several “Options”, including “InitialMeshSize” and “Tolerance on the mesh size”. For both options we select values that depend on the specific example; i.e., for the toy problem with $0 \leq x_j \leq 1$ ($j = 1, 2$), we select the “InitialMeshSize” equal to 0.10 (or 10% of the input range), and for the inventory problem with $600 \leq s \leq 1,200$ and $600 \leq S \leq 1,800$ we select the value 20. During its search, PS decreases (or increases) this mesh size when the next iterate is worse (or better) than the current iterate (i.e., PS uses an “adaptive mesh”). “Tolerance on the

mesh size” stops the PS search as soon as the mesh size becomes smaller than this tolerance. For the toy problem with $0 \leq x_j \leq 1$, we select a tolerance equal to 0.001, and for the inventory problem we select a tolerance of 5. PS also has an “exitflag — Reason patternsearch stopped”; e.g., the flag value -2 means “No feasible point found”. (We do not use the option “UseParallel” in PS; in our algorithm we do use parallel computing.) PS requires an “Initial point”. To reduce the probability of PS getting *trapped* in a local optimum, PS may be *restarted* several times. Actually, we use $10k$ starting points that satisfy the v input constraints $f_g(\mathbf{x}) \leq c_g$ ($g = 1, \dots, v$), but may violate the output constraints (MATLAB advises ten feasible restart points). PS requires relatively little computer time in our algorithm, as PS uses \hat{y}_h instead of w_h (expensive simulation output). To select specific initial (so, $s = 0$) starting points, we use LHS. Obviously, PS also ends with $10k$ points, but some of these endpoints may coincide. To select *specific* starting points after the initial stage (so, $s = 1, 2, \dots$), we also use the new points that were selected by our algorithm after stage $s = 0$. However, in the first stages, there are not enough new points, so PS also uses LHS. (Wang and Ierapetritou (2018) also uses $10k$ starting points selected by LHS, but uses MATLAB’s `fmincon` to find infill points).

If PS cannot find a feasible solution (exitflag is -2), then we halve α and apply PS to the problem with this new α -value (obviously, this new value implies a wider CI around the estimated binding output constraints). We stop our algorithm as soon as the current α is lower than some small value (say) α_{\min} ; e.g., $\alpha_{\min} = 0.01$.

After PS has given $\hat{\mathbf{x}}_o$, our algorithm must decide on m_i . Our algorithm applies the sample allocation rule (13); i.e., it may add (expensive) *replications*, at old or new points. Next, it uses all available simulated I/O combinations, to re-estimate the SK parameters.

Whereas deterministic simulation gives the true outputs $w_{i,h}$, random simulation gives only the estimates $w_{i,h,r}$ ($r = 1, \dots, m_i$). Consequently, random simulation may require many more iterations, and still give an estimated global optimum that lies far away from the true optimum. Therefore our algorithm *restarts* with different initial designs (in general, restarting is popular in optimization). We use 12 restarts (actually, at the start of our experimentation we used a non-parallel computer; we may then add a single restart at a time—until we are satisfied with the result or we run out of time). Our algorithm uses PS with $n = 10k$ restarts, not in parallel. (In practice, the users may observe the development of the estimated optimum, and decide whether they wish to continue the search for the true optimum.)

We use the following hardware: an HPE ProLiant DL380 Gen10 server with two Intel(R) Xeon(R) Silver 4110 CPUs (16 cores) with 2.10 GHz, and 59.1 GB of usable RAM out of 64 GB. We use MATLAB R2019b; also see the Parallel Computing Toolbox on

<https://www.mathworks.com/help/parallel-computing/parallel-preferences.html>.

We use the PRN generator called *mrg32k3a*, which can create 2^{63} independent PRN substreams of length

2^{127} ; see

<https://www.mathworks.com/help/matlab/ref/randstream.randstream.create.html>.

This generator is also recommended in L’Ecuyer (2015). This generator enables us to control the PRN seeds such that restarts and macroreplications use non-overlapping substreams.

Note that initially we used a non-parallel implementation without restarting (or 1 “restart” instead of 12 restarts). In the toy example, this implementation gave many final values for $\hat{\mathbf{x}}_{\min}$ close to a local optimum instead of the global optimum. More precisely, in this example 60 macroreplications without restarts gave frequency diagrams for $\hat{\mathbf{y}}_{0; \min}$ with two “tops” (or modes), which we interpret as tops corresponding with the global and local optima.

In Algorithm 1, we present a pseudo-code for a single restart.

Algorithm 1

1. Select the algorithm’s control variables $\alpha_E, \alpha, \alpha_{\text{infe}}, \alpha_{\min}, \varepsilon, m, n$.

Comment: The variable α_E features in (14), α and α_{infe} in (15), α_{\min} after (20), ε below (16), and m and n in Section 3.3.

2. Sample an initial design matrix with n combinations of the k simulation inputs x_j that satisfies the input constraints $f_g(\mathbf{x}) \leq c_g$ with $g = 1, \dots, v$.

Comment: This Step gives the $n \times k$ matrix \mathbf{X} with the k -dimensional vectors \mathbf{x}_i where $i = 1, \dots, n$. We select $n = 6$ if $k = 2$.

3. Use \mathbf{X} as input for the given simulation model, to obtain m replications of the $n \times t$ matrix with (t -variate) simulation outputs $\mathbf{W}_{n \times t; r}$.

Comment: $\mathbf{W}_{n \times t; r}$ consists of the n -dimensional vectors $\mathbf{w}_{n; h; r}$ with $h = 0, \dots, t - 1$ and $r = 1, \dots, m$. We select the initial value $m = 10$ in the toy example and $m = 2$ in the inventory example.

4. Apply LOO-CV to $(\mathbf{X}, \mathbf{W}_{n \times t; r})$.

Comment: LOO-CV uses (14).

5. If LOO-CV rejects one or more SK models, then obtain a new replication for the combination with the smallest estimated variance among the t output types, and apply the allocation rule to compute the desired number of replications \hat{m}_i ; if necessary, then obtain additional replications, and update m_i ; return to Step 4.

Comment: To compute \hat{m}_i , use (13). To decide on additional replications, use the text below (13). Steps 4 and 5 constitute the preprocessing stage.

6. Use $(\mathbf{X}, \mathbf{W}_{n \times t; r})$ to compute $\hat{\psi}_h = \hat{\psi}(\mathbf{X}, \mathbf{w}_{n; h; r})$, $\hat{\Sigma}_{\bar{e}; h}$, $\hat{y}_h(\mathbf{x}_i)$, $s^2[\hat{y}_h(\mathbf{x}_i)]$.

Comment: $\hat{\Sigma}_{\bar{e}; h}$ is a diagonal matrix that uses (9).

7. Use $\hat{y}_h(\mathbf{x}_i)$ to compute $\hat{y}_{0; \min} = \min_i \hat{y}_0(\mathbf{x}_i)$ if $\hat{y}_{h'}(\mathbf{x}_i) + z_{(1-\alpha)} s[\hat{y}_{h'}(\mathbf{x}_i)] \leq c_{h'}$ (with $h' = 1, \dots, t-1$ and $\alpha = \alpha_{\text{infe}}$) and $f_g(\mathbf{x}) \leq c_g$ (with $g = 1, \dots, v$); else $\hat{y}_{0; \min} = \infty$.

Comment: This $\hat{y}_{0; \min}$ is needed to compute $\widehat{\text{MEI}}(\mathbf{x})$, in Step 9.

8. Use $\hat{y}_h(\mathbf{x}_i)$ to compute $\widehat{\nabla}[\hat{y}_h(\mathbf{x}_i)]$.

Comment: To compute $\widehat{\nabla}[\hat{y}_h(\mathbf{x}_i)]$, use (12). This $\widehat{\nabla}[\hat{y}_h(\mathbf{x}_i)]$ is needed to compute $\widetilde{\text{cos}}[\mathbf{x}]$, in Step 9.

9. Apply PS to estimate $\hat{\mathbf{x}}_o$, which maximizes $\widehat{\text{MEI}}(\mathbf{x}) \times \widetilde{\text{cos}}[\mathbf{x}]$ such that \mathbf{x} satisfies the $t-1$ CIs for the constrained outputs h' defined in (17) and the v input constraints defined in (1).

Comment: $\widehat{\text{MEI}}(\mathbf{x})$ is defined in (16) using \hat{y}_0 and $s^2(\hat{y}_0)$; $\widetilde{\text{cos}}[\mathbf{x}]$ is defined in (19) using the estimated gradients of the estimated binding output constraints $\widehat{\nabla}[\hat{y}_{h''}(\mathbf{x})]$ and the LS model with the explained variable $\widetilde{\nabla}[\hat{y}_0(\mathbf{x})]$ and the explanatory variables $\widehat{\nabla}[\hat{y}_{h''}(\mathbf{x})]$ and $\nabla_{g'}(\mathbf{x})$ —defined in (20).

10. If PS finds a feasible point $\hat{\mathbf{x}}_o$, then apply the allocation rule, obtain the required observations, and return to Step 6; else replace α by $\alpha/2$ in (17). If $\alpha \geq \alpha_{\min}$, then return to Step 9; else if $\alpha < \alpha_{\min}$ and $\widehat{\text{MEI}}(\mathbf{x}) < \varepsilon \times |\hat{y}_{0; \min}|$, then go to the final step (Step 11).

11. Present the estimated optimal goal output $\hat{y}_{0; \min}$.

After the multiple (namely, 12) *restarts*, our algorithm presents the minimum of the multiple estimated optimal goal outputs. Let $\hat{y}_{0; \min; q}$ denote $\hat{y}_{0; \min}$ of Step 11 in Algorithm 1 (single restart) of restart q (our PC with 12 cores implies $q = 1, \dots, 12$). Then all restarts together give the *overall* estimated optimal goal output; namely, $\hat{y}_{0; \min; q_o}$ with $q_o = \arg \min_q \hat{y}_{0; \min; q}$. Moreover, the algorithm determines information that corresponds with $\hat{y}_{0; \min; q_o}$; namely, $\hat{\mathbf{x}}_{o; q_o}$, $\hat{y}_{h'; q_o}(\hat{\mathbf{x}}_{o; q_o})$ ($h' = 1, \dots, t-1$), and the (say) 90% two-sided CI for $\hat{y}_{h; q_o}(\hat{\mathbf{x}}_{o; q_o})$ ($h = 0, \dots, t-1$). The confidence levels are supposed to apply per comparison; however, they are not exact, as the algorithm has many steps that use CIs. Because of the KT conditions, the estimated slacks $c_{h'} - \hat{y}_{h'; q_o}(\hat{\mathbf{x}}_{o; q_o})$ should be virtually zero for at least one output constraint.

Our algorithm uses 12 restarts (cores), so we select the following measure of *efficiency*. By definition, the *elapsed time* is the time between the starting time and the stopping time of our algorithm. That stopping time is determined by the *slowest* restart among the 12 restarts, which requires the maximum number of (expensive) simulation observations (the adapted EGO computations require negligible time). Restart q (with $q = 1, \dots, 12$) stops after T_q stages (so, $s = 0, 1, \dots, T_q$). Stage s has its own number of simulated points n_s , and each of these points has its own number of replications $m_{i; s}$. Let N_q denote the *total* number

of simulation observations of restart q , and let N denote N_q of the slowest restart:

$$N = \max_q N_q \text{ with } N_q = \sum_{s=0}^{T_q} \sum_{i=1}^{n_s} m_{i;s;q} \quad (21)$$

We use this N as a *pseudo-measure* for the elapsed time between the simultaneous start of the 12 restarts and the end of the slowest restart.

Note that we may add replications to either the B points near the boundary or all n points. We conjecture that both design options are equally effective, but that the former design option is more efficient. To test this conjecture, we compare these two options in the toy example, using one macroreplication with 12 restarts. Adding replications to all n points gives $\hat{\mathbf{x}}_{\min} = (0.1938, 0.4187)'$ with $\hat{y}_{0; \min} = 0.6159$, and $N = 101,671$. Adding replications to B points gives $\hat{\mathbf{x}}_{\min} = (0.1984, 0.4115)'$ with $\hat{y}_{0; \min} = 0.6122$ (nearly 1% better) and $N = 47,707$ (less than 50%).

Inspired by Tao et al. (2020, Table 1, line 32), we add an *extra Step* that finds the solution $\hat{\mathbf{x}}_{\min}$ for

$$\begin{aligned} & \min_{\mathbf{x}} \hat{y}_0(\mathbf{x}) \text{ subject to} \\ & \hat{y}_{h'}(\mathbf{x}) + z_{(1-\alpha)s}[\hat{y}_{h'}(\mathbf{x})] \leq c_{h'} \text{ with } \alpha = \alpha_{\text{infe}} \\ & f_g(\mathbf{x}) \leq c_g \quad g = 1, \dots, v \end{aligned} \quad (22)$$

where \hat{y}_h ($h = 0, \dots, t-1$) uses the simulation I/O data of all 12 restarts. To find this $\hat{\mathbf{x}}_{\min}$, our algorithm uses PS with 12 restarts in parallel. Unlike Tao et al., we check whether simulation with this input $\hat{\mathbf{x}}_{\min}$ gives a feasible solution and a smaller $\hat{y}_0(\hat{\mathbf{x}}_{\min})$; else $\hat{\mathbf{x}}_{\min}$ equals $\hat{\mathbf{x}}_{o; q_0}$ (estimate resulting from multiple restarts of our algorithm).

Note that these 12 restarts may give so many simulated combinations that a numerical problem arises (DACE gives the error message “Multiple design sites are not allowed”). To solve this problem, we combine I/O simulation data that lie “close together”; e.g., if the toy example (with $0 \leq x_{i;j} \leq 1$) gives $|x_{i;j} - x'_{i;j}| < 0.07$, then we replace these inputs and outputs by their averages $\bar{x}_{i;j}$ and \bar{w}_i .

5 Mathematical toy example

We display the deterministic toy example of Gramacy et al. (2016) (also used in several other publications mentioned in Section 1) in Fig. 1, which includes the following special points: A is the global minimum; B is a local minimum; C is a local minimum with a binding input constraint; D and E are local maxima; A

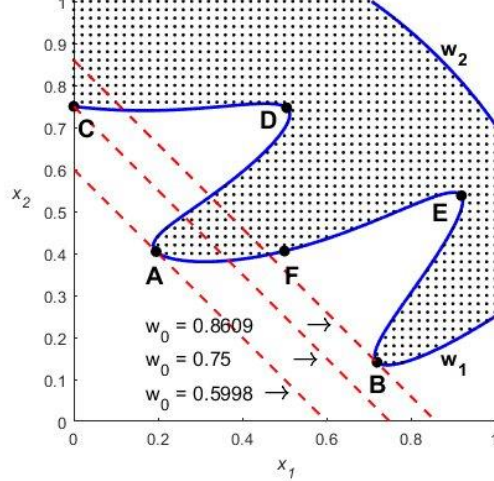


Figure 1: Mathematical toy example with inputs x_1 and x_2 , deterministic goal output w_0 and constraints w_1 and w_2 , and special points A through F

through E satisfy the KT conditions, whereas F does not; A through F are investigated in Kleijnen et al. (2022).

Now we modify this deterministic example such that it becomes relevant for random simulation; i.e., we introduce *intrinsic noise* $e_h(\mathbf{x})$ with $h = 0, 1, 2$, and replace the deterministic outputs $w_h(\mathbf{x})$ by their *expected values* $E[w_h(\mathbf{x})]$:

$$\begin{aligned}
 \min_{\mathbf{x}} E[w_0(\mathbf{x})] &= E[x_1 + x_2 + e_0(\mathbf{x})] \\
 E[w_1(\mathbf{x})] &= E\left\{\frac{3}{2} - x_1 - 2x_2 - \frac{1}{2} \sin[2\pi(x_1^2 - 2x_2)] + e_1(\mathbf{x})\right\} \leq 0 \\
 E[w_2(\mathbf{x})] &= E\left[-\frac{3}{2} + x_1^2 + x_2^2 + e_2(\mathbf{x})\right] \leq 0 \\
 0 &\leq x_j \leq 1 \quad j = 1, 2.
 \end{aligned} \tag{23}$$

We assume that e_h has *heterogeneous variances* that are determined by \mathbf{x} ; so, we let $\sigma_{e;h}^2(\mathbf{x})$ denote these variances. To select specific values for these $\sigma_{e;h}^2(\mathbf{x})$, we follow Jalali et al. (2017, p. 293) where the “best case” means that the standard deviation $\sigma_{e;h}(\mathbf{x})$ decreases as $E[w_h(\mathbf{x})]$ decreases, and the “worst” case means that $\sigma_{e;h}(\mathbf{x})$ increases as $E[w_h(\mathbf{x})]$ decreases; Jalali et al. (2017) defines these two cases, assuming no constraints when minimizing $E[w_0(\mathbf{x})]$. We, however, consider multiple outputs. For simplicity’s sake, we apply Jalali et al.’s approach to each of these outputs. We limit our experiments to the best case (so, as

$w_h(\mathbf{x})$ decreases, $\sigma_{e;h}(\mathbf{x})$ decreases) with relatively small noise:

$$\begin{aligned}\sigma_{e;0}(\mathbf{x}) &= 0.30 + 0.45E[w_0(\mathbf{x})] \\ \sigma_{e;1}(\mathbf{x}) &= 1.1507 + 0.45E[w_1(\mathbf{x})] \\ \sigma_{e;2}(\mathbf{x}) &= 0.975 + 0.45E[w_2(\mathbf{x})].\end{aligned}\tag{24}$$

These experiments still require many simulation observations; e.g., macroreplication # 1 requires $N = 47,707$. Appendix 2 details the derivation of formulas for small noise versus big noise in the best case versus the worst case.

For simplicity’s sake, we assume *zero cross-correlations* between different types of output in the toy example (the inventory example in Section 6 has two negatively correlated output types). To sample the univariate $e_h(\mathbf{x})$, we use MATLAB’s *normrnd* where $E[e_h(\mathbf{x})] = 0$ and $\sigma_{e;h}(\mathbf{x})$ follows from (24). Given these $e_h(\mathbf{x})$, we compute the corresponding w_h ($h = 0, 1, 2$) using (23). Appendix 4 shows how our algorithm proceeds through various stages.

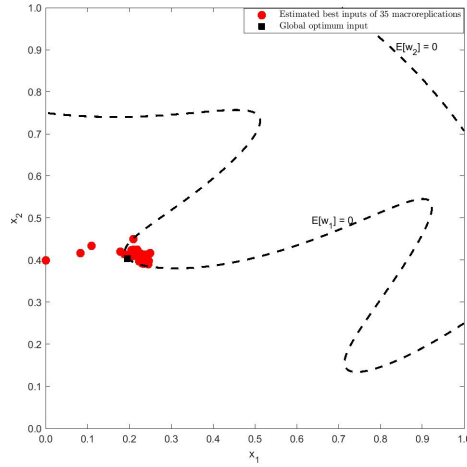


Figure 2: Final estimated optimal input combinations $\hat{\mathbf{x}}_{\min}$ in toy problem with 35 macroreplications, each with 12 restarts, $\alpha_{\text{infe}} = 10\%$, $m_{\min} = 10$, and 20 PS starting points

Fig. 2 presents $\hat{\mathbf{x}}_{\min}$, which denotes the estimated optimal input combinations after 12 restarts, in 35 macroreplications; each macroreplication uses $\alpha_{\text{infe}} = 10\%$, $m_{\min} = 10$, and 20 PS starting points in each iteration. We add that 5 of these 35 estimates result from the extra step at the end of our algorithm; see (22). This Figure clearly shows that 3 of these 35 macroreplications are actually infeasible—which agrees with $\alpha_{\text{infe}} = 10\%$. The other 32 macroreplications give estimates close to the true optimum.

The left-hand side of Fig. 3 presents the boxplot of the final estimated optimal goal outputs $\hat{y}_0(\hat{\mathbf{x}}_{\min})$ in

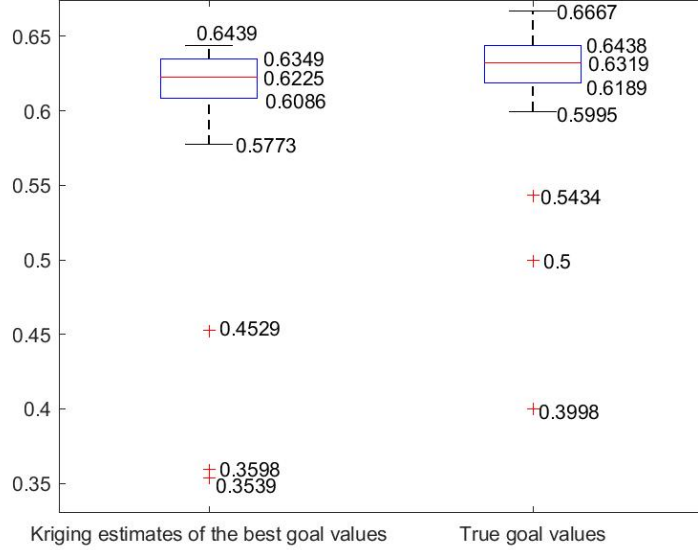


Figure 3: Boxplot of $\hat{y}_{\text{SK},0}(\hat{\mathbf{x}}_{\min})$ (left-hand side) and $E[w_0(\hat{\mathbf{x}}_{\min})]$ (right-hand side) in toy problem; 35 macroreplications, 12 restarts, $\alpha_{\text{infe}} = 10\%$, $m_{\min} = 10$, and 20 PS starting points

35 macroreplications (12 restarts, $\alpha_{\text{infe}} = 10\%$, $m_{\min} = 10$, 20 PS starting points). The three outliers below the lowest quartile correspond with the three infeasible input combinations in Fig. 2. The sample median is 0.6225 (the true optimum is 0.5998, so the estimate is nearly 4% higher). The right-hand side presents the boxplot of $E[w_0(\hat{\mathbf{x}}_{\min})]$ (we know this mean, because this example is only a mathematical toy problem).

Fig. 4 presents the boxplot for N in 35 macroreplications. Its median is 14,382; i.e., random simulation with large response variances may require much computer time.

Note that SK combined with CRN requires the estimation of the elements of the $n \times n$ covariance matrix $\Sigma_{\bar{e};h}$ —which requires $m > n$ (see Dykstra (1970)). In the toy example, we select $n = 6$ and $m = 7$ in the preprocessing stage. Based on five macroreplications—where α_{infe} is either 0.10 or 0.01—we conclude that SK with CRN performs worse than OK with CRN. We conjecture that the cause is the estimation of the extra Kriging parameter $\hat{\Sigma}_{\bar{e};h}$. Therefore, we focus on SK without CRN.

5.1 Our algorithm versus OptQuest in toy example

In this subsection, we numerically compare the performance of OptQuest and our algorithm. OptQuest is a black-box optimizer, as it is proprietary. Both OptQuest and our algorithm treat the simulation model as a black box; i.e., in our experiment, OptQuest receives the input $(\mathbf{x}_i, w_{i,h,r})$ (which is the simulation I/O with $i = 1, \dots, n$; $h = 0, 1, 2$; $r = 1, \dots, m_i$).

For OptQuest we select m_i through the sequential procedure in Law (2015, p. 505) for estimating the

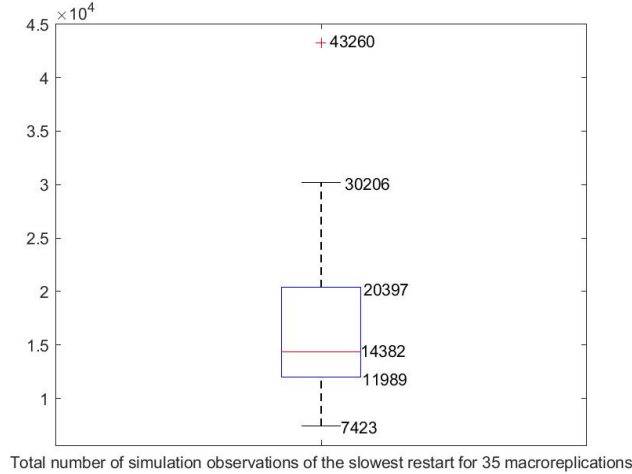


Figure 4: Elapsed times N in toy problem; 35 macroreplications, 12 restarts, $\alpha_{\text{infe}} = 10\%$, $m_{\text{min}} = 10$, and 20 PS starting points

mean $E(w)$ with a relative error of γ (with $0 < \gamma < 1$) and confidence level $1 - \alpha$. This procedure starts with an initial number of replications $m = m_0 \geq 10$ (we use our own symbols). Next, it computes the half-length of the two-sided $1 - \alpha$ CI:

$$l = t_{m-1; 1-\alpha/2} \frac{s(w)}{\sqrt{m}}. \quad (25)$$

If

$$\frac{l}{|\bar{w}|} < \frac{\gamma}{1 + \gamma}, \quad (26)$$

then stop adding replications; else, add one replication, and update (25) and (26), etc.—until (26) holds. (Alternative procedures are discussed in Kleijnen (2015, pp. 107–108).)

However, we point out that if (the denominator) $|\bar{w}|$ is “close” to zero, then this procedure gives a “high” \hat{m} (desired number of replications). To solve this problem, we return to the original optimization problem defined in (1), which includes the $(t - 1)$ output constraints $E[w_{h'}(\mathbf{x})] \leq c_{h'}$. Our toy example defined in (23), however, has the constraints $E[w_1(\mathbf{x})] = E\{3/2 - x_1 - \dots\} \leq 0$ and $E[w_2(\mathbf{x})] = E[-3/2 + x_1^2 + \dots] \leq 0$. To derive \hat{m} , we replace $w_{h'}(\mathbf{x})$ by $w_{h'}(\mathbf{x}) + c_{h'}$ (which has the same variance as $w_{h'}(\mathbf{x})$ has); i.e., in the toy problem we use $w_1(\mathbf{x}) - 3/2$ and $w_2(\mathbf{x}) + 3/2$ (we do not replace $w_0(\mathbf{x})$). We still use $w_{h'}(\mathbf{x})$ when we use OptQuest to search for the optimum.

Because we have t outputs, we wish that $m \geq \hat{m}$ holds for all t output; obviously, \mathbf{x}_i has its own $E[w_h(\mathbf{x}_i)]$ and $Var[w_h(\mathbf{x}_i)]$. Altogether, we define the estimated *desired* number of replications for \mathbf{x}_i

$$\hat{m}(\mathbf{x}_i) = \max_t [\hat{m}_t(\mathbf{x}_i)]. \quad (27)$$

We think that our procedure for determining m agrees with one of the OptQuest options that also uses γ and α (this option is part of the OptQuest black box).

Note that our own algorithm uses the allocation rule (13), and allows the search to return to old points and add replications. More precisely, our rule selects the *relative* number of replications for different input combinations, whereas Law’s rule selects the absolute number of replications for a given input combination; our rule uses the estimated variances of the simulation outputs, whereas Law’s rule uses the estimated standard deviations (so it dampens the effect of variance heterogeneity).

Because (27) may require many replications if γ is small (e.g., 5%), we may select a relatively high γ (e.g., 10%). The latter selection implies a relatively high $Var[w_h(\mathbf{x}_i)]$, so we may then select a relatively high safety factor $z_{(1-\alpha)}$ in (15). (which defines the problem solved by OptQuest and our procedure).

We start OptQuest with the same *initial stage* as we use for our algorithm; i.e., we start OptQuest with $(\mathbf{x}_i, w_{i;h;r})$ with $n = 6$ input combinations selected through LHS (with midpoints), and m_i ($i = 1, \dots, n$) replications where m_i follows from (27). This start is quite general; i.e., n follows from Tao et al. (2020) and these m_i are based on Law (2015), which is the most popular textbook on simulation. Our algorithm uses 12 cores (of our PC) and starts each core with its own initial $n \times k$ LHS design and initial m_i , per macroreplication. However, OptQuest does not use multiple cores (we do use our 12 cores to speed-up the simulation of the toy example). Moreover, we start each macroreplication of OptQuest with the $n \times k = 6 \times 2$ LHS design that gave the best results for our algorithm in the corresponding macroreplication (so we give OptQuest a good start in each macroreplication). Of course, we select the m_i in OptQuest using Law’s procedure defined in (25) etc.

After these *initial* I/O simulation data, OptQuest asks for simulation I/O data for a specific *new* input combination \mathbf{x}_{n+1} (which is a point “near” one or more points selected by the space-filling LHS-design). So, we provide $(\mathbf{x}_{n+1}, w_{n+1;h;r})$ to OptQuest, using (27) to decide on m_{n+1} .

We let OptQuest solve the analogue of (22), replacing $\hat{y}_h(\mathbf{x}_i)$ by $w_{h;r}(\mathbf{x}_i)$. The three outputs ($w_{0;r}(\mathbf{x}_i)$, $w_{1;r}(\mathbf{x}_i)$, $w_{2;r}(\mathbf{x}_i)$) are uncorrelated, in the toy example (in the inventory example the—possibly nonnormal—simulation outputs are statistically dependent). By definition, replications are IID. Furthermore, the simulation outputs of the n combinations are independent because the simulation does not use CRN). Each of the 12 cores of our PC works in parallel (independently); i.e., restart g ($g = 1, \dots, 12$) uses an initial $n \times k = 6 \times 2$ LHS design. After this initial stage (or iteration) $s = 0$, the next stages $s = 1, 2, \dots, T_g$ use independent simulation outputs with heterogeneous variances $\sigma_h^2(\mathbf{x}_i) = var[w_h(\mathbf{x}_i)]$. We can keep track of N (defined in (21)) per macroreplication in our algorithm, so we can allow OptQuest to run the same N .

Fig. 5 and Fig. 6 show $\hat{\mathbf{x}}_{\min}$ for OptQuest in 35 macroreplications. Fig. 5 shows $\hat{\mathbf{x}}_{\min}$ when we stop OptQuest manually as soon as the total number of simulation runs exceeds N of the corresponding

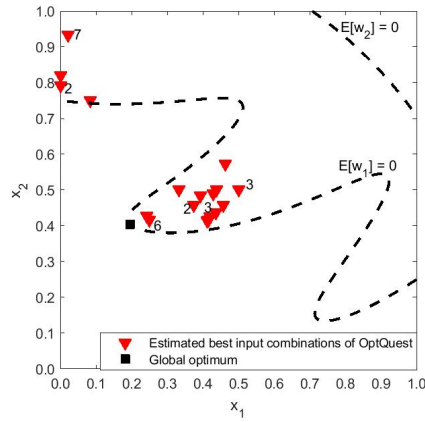


Figure 5: OptQuest's \hat{x}_{\min} with the stopping criterion N ; 35 macroreplications with $\alpha_{infe} = 10\%$ for the toy problem

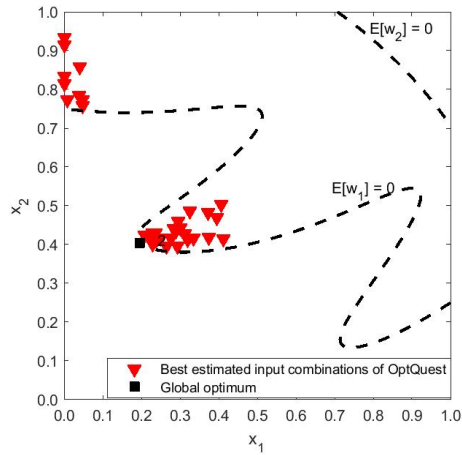


Figure 6: OptQuest's \hat{x}_{\min} with the stopping criterion n ; 35 macroreplications with $\alpha_{infe} = 10\%$ for the toy problem

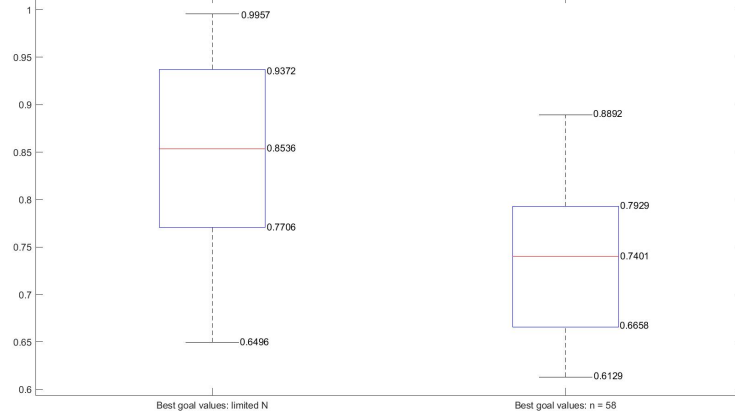


Figure 7: OptQuest’s $\hat{\mathbf{x}}_{\min}$ in 35 macroreplications with $\alpha_{infe} = 10\%$ for the toy problem: limited N on the left, and $n = 58$ on the right

macroreplication of our algorithm. Actually, several macroreplications give the same $\hat{\mathbf{x}}_{\min}$ (e.g., $\hat{\mathbf{x}}_{\min} = (0.022, 0.9337)'$ in seven macroreplications), as denoted by the numbers in this Figure. Whereas OptQuest gives no infeasible $\hat{\mathbf{x}}_{\min}$ -results, our algorithm gives 3 out-of-35 infeasible $\hat{\mathbf{x}}_{\min}$ -results (remember that $3/35$ is controlled by $\alpha_{infe} = 10\%$ in our algorithm). Furthermore, this Figure shows that 11 ($= 7 + 2 + 1 + 1$) of the 35 $\hat{\mathbf{x}}_{\min}$ -results are close to the local minimum $\mathbf{x}_C = (0, 0.7)'$, and the remaining 24 ($= 35 - 11$) $\hat{\mathbf{x}}_{\min}$ -results are close to the global minimum $\mathbf{x}_A \approx (0.1954, 0.4044)'$. Our algorithm, however, gives $\hat{\mathbf{x}}_{\min}$ -results close to \mathbf{x}_A in all 35 macroreplications; see again Fig. 2. Fig. 6 shows $\hat{\mathbf{x}}_{\min}$ -results when we follow the recommendation in OptQuest and set the total number of input combinations n between 58 and 65. Actually, even the smaller n (namely, 58) causes OptQuest to run longer than N does; yet, 9 of the 35 macroreplications give $\hat{\mathbf{x}}_{\min}$ -results near \mathbf{x}_C (and 26 $\hat{\mathbf{x}}_{\min}$ -results closer to \mathbf{x}_A). We add that 2 macroreplications give the same $\hat{\mathbf{x}}_{\min}$; namely, $\hat{\mathbf{x}}_{\min} = (0.2500, 0.4167)'$ (the number 2 is not displayed in Fig. 6, because this plot has too many triangles close together).

Figure 7 displays the boxplots for $\hat{y}_{0; \min}$, which correspond with $\hat{\mathbf{x}}_{\min}$ in Fig. 5 and Fig. 6. The median values of $\hat{y}_{0; \min}$ are 0.8536 and 0.7401, so both values are far from the true global optimal value 0.5998.

Finally, we obtain boxplots (not displayed) for the total number of simulation observations required in 35 macroreplications for OptQuest and our algorithm. The median value for OptQuest is 58,527, which is more than four times bigger than the median value 14,382 for our algorithm (see again Figure 4). The ranges of the two boxplots do not overlap. We conclude that our algorithm is more efficient (requires less elapsed time) than OptQuest, when applied to the toy example.

6 An (s, S) inventory simulation with a service-level constraint

As we have already discussed in Section 1, inventory simulations with a service-level constraint are of major practical relevance. Now we investigate an (s, S) inventory simulation with a service-level constraint that was originally investigated in Bashyam and Fu (1998). In this simulation, a company places an order (with its supplier) of size $S - I_p$ as soon as the *inventory position* I_p drops below the reorder level s —where I_p denotes the physical inventory plus outstanding orders minus “backorders”, which are customer orders that cannot be delivered immediately because the physical inventory is too small (backordering implies that there are no lost sales), and p denotes the review period $p = 1, 2, \dots, P$, where P is the end of the simulation run. Bashyam and Fu (1998) select $P = 30,000$, assuming this gives steady-state simulation outputs. (So, the order size is not a constant $Q = S - s$.)

Furthermore, Bashyam and Fu (1998) allows random demand D and lead times L , so the orders that the company places are not necessarily received in the order in which they are placed. This *order crossing* complicates the mathematical analysis, so analysts use simulation. Estimating the *optimal* inventory-control limits s_o and S_o is difficult, as the vast literature on inventory management shows.

Finally, Bashyam and Fu (1998) assumes that the goal variable— w_0 in our symbols—is the average total inventory *costs* excluding (the hard to quantify) cost of out-of-stock (causing back orders). The *disservice percentage* w_1 is the ratio of (i) total demand during P that is not satisfied from the physical stock, and (ii) total demand during that period. We assume that this percentage should not exceed 10%, so $E(w_1) \leq 0.10$ (Bashyam and Fu (1998) also investigates 1% and 5%, and Kleijnen and Wan (2007) investigates only 10%). Obviously, w_0 and w_1 are negatively correlated. (Our definition of w_1 implies that the service percentage is $1 - w_1$.)

Kleijnen and Wan (2007) assumes that D is exponentially distributed with mean 100, L is integer-valued Poisson distributed with mean 6, holding cost h is 1, variable ordering cost is 1, and fixed ordering cost K is 36 (we need the symbols h and K below; Bashyam and Fu (1998) investigates more cases.)

Table 1: Initial 6×2 design and its estimated outputs, in macroreplication 1; * denotes $\bar{w}_1 > 0.10$

s	S	\bar{w}_0	\bar{w}_1
750.0	1166.7	486.7022	0.2932 *
1183.3	1766.7	928.6275	0.0426
616.7	1100.0	453.7042	0.3549 *
816.7	1366.7	603.4781	0.1883 *
783.3	833.3	327.8253	0.5096 *
1116.7	1500.0	761.0117	0.0691

We use the following *input constraints* for s and S . For s we select the lower bound $E(D) \times E(L) =$

$100 \times 6 = 600$, which ignores the randomness of D and L , so it is too low to ensure proper service. For the upper bound we (rather arbitrarily) select $2 \times 600 = 1,200$. For S we (rather arbitrarily) select the same lower bound as for s (namely, 600), and for the upper bound we select $3 \times 600 = 1,800$. So, initially, we did not apply the reparametrization with $x_2 = Q = S - s$ (so, we did not apply LHS), but derived an alternative procedure that gives the first two columns in Table 1. For the same initial design we compare the performance of OptQuest and our algorithm. When we do parameterize this inventory problem, then we may vary Q around the popular *economic order quantity* (EOQ), which is $[2 \times K \times E(D)/h]^{1/2} = 85$; see Bashyam and Fu (1998). (Obviously, this EOQ ignores the randomness of D and L .)

Note that we must avoid numerical problems—when computing the Gaussian correlation function (11) and the corresponding gradient (12) (which contain $\exp[-\theta_{h,j'}(x_{*,j'} - x_{i,j'})^2]$)—so we *normalize* (or standardize) the simulation inputs s and S such that $0 \leq x_j \leq 1$. The input constraints in the (s, S) problem are $2k = 4$ box constraints plus the extra constraint $s \leq S$ (or $x_1 \leq x_2$).

Next we determine the number of *replications* m_i , using Law’s rule (26). It turns out that $m_0 = 10$ does not require any additional replications; i.e., the pilot sample is too big. Therefore we next select $m_0 = 2$, and find that m_i is very small; namely, 2, 3 or 4. We conclude that $P = 30,000$ gives very accurate estimates.

Running the simulation with the $n = 6$ initial input combinations of Table 1 and the m_i replications gives the outputs \bar{w}_0 and \bar{w}_1 in macroreplication 1; see the last two columns of this Table. Qualitative knowledge about (s, S) inventory systems implies that \bar{w}_1 (disservice fraction) decreases as s or S increases (because higher s or S keeps the system farther away from the “danger zone” where out-of-stocks occur). In this Table, we denote $\bar{w}_1 > 0.10$ by *; so, we estimate that only 2 of the 6 (s, S) combinations are feasible. Furthermore, this Table shows that the combination in the last row is the best initial solution; i.e., this combination gives an acceptable disservice and the lowest cost \bar{w}_0 .

Kleijnen and Wan (2007) compares several algorithms for estimating s_o and S_o , including OptQuest. We compare only our algorithm and the OptQuest algorithm (as we do for the toy example). For our comparison we apply both algorithms to the same (s, S) simulation model (there are different interpretations of (s, S) systems—e.g., not all interpretations use our definition of I_p —and programming codes—we follow the coding in Law (2015, pp. 48–61)).

6.1 Our algorithm versus OptQuest in inventory example

Fig. 8 displays the estimated optimal input combinations (\hat{s}_o, \hat{Q}_o) with $\hat{Q}_o = \hat{S}_o - \hat{s}_o$ in our (s, S) inventory simulation, in 22 macroreplications for our algorithm and OptQuest. We have 22 macroreplications because our OptQuest licence expired after six weeks, which enabled us to obtain 22 macroreplications for OptQuest

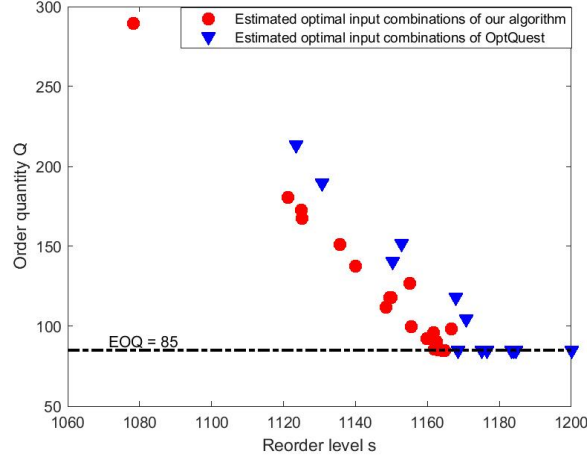


Figure 8: Estimated optimal input combination (\hat{s}_o, \hat{Q}_o) with $\hat{Q}_o = \hat{S}_o - \hat{s}_o$ in (s, S) inventory simulation, in 22 macroreplications for our algorithm and OptQuest

applied to our inventory example (after 35 macroreplications for the toy example). This Figure also displays the EOQ, which implies that all 22 (\hat{s}_o, \hat{Q}_o) -combinations—for our algorithm and OptQuest, respectively—give $\hat{Q}_o \geq \text{EOQ} = 85$. Moreover, our algorithm gives combinations with smaller \hat{s}_o and \hat{Q}_o than OptQuest gives, which decreases the costs but might endanger the service. Therefore we next construct Fig. 9.

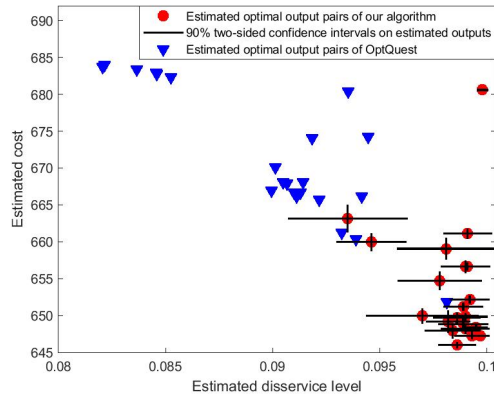


Figure 9: Scatterplot of estimated optimal disservice versus estimated optimal inventory cost, in 22 macroreplications; crosses denote 90% CIs

Fig. 9 displays a scatterplot of the estimated optimal disservice ($\bar{w}_1(\hat{\mathbf{x}}_{\min})$ in OptQuest, and $\hat{y}_1(\hat{\mathbf{x}}_{\min})$ in our algorithm) and the corresponding estimated optimal cost ($\bar{w}_{0;\min}$ and $\hat{y}_{0;\min}$) of our inventory simulation, in the same 22 macroreplications as in Fig. 8; the Figure also displays 90% CIs for each of these disservice and cost estimates. Projecting these 22 points for OptQuest and our algorithm onto the axes determines the boxplots. The median cost and disservice for our algorithm are 650 and 0.0990, and for OptQuest 668 and 0.0911. So, the percentage decrease in cost is 3%, while still satisfying the service-level threshold of 90%. In

practice, the cost saving is even higher when applying our algorithm instead of “experience” to optimize the (s, S) inventory system.

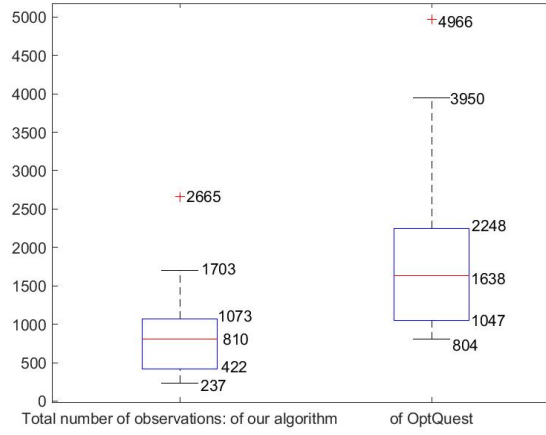


Figure 10: Boxplots of N (total number of simulation observations) of our algorithm (left) versus OptQuest (right) in 22 macroreplications of the inventory example

Fig. 10 displays boxplots of N (defined in (21) as a pseudo-measure for the elapsed times) of our algorithm and OptQuest. This Figure shows that the median N is 810 for our algorithm, and 1,638 for OptQuest; i.e., our algorithm requires much less time to get better estimates of the optimum.

7 Conclusions and future research

In this paper, we investigated constrained optimization problems in random simulation; i.e., we derived a novel algorithm, combining the KT conditions and EGO. More specifically, we account for both output constraints and input constraints, allow variance heterogeneity, and apply SK with a sample allocation rule. Our algorithm uses PS to estimate \mathbf{x} that maximizes the current SK estimate of EGO’s $\widehat{\text{MEI}}(\mathbf{x})$ multiplied by $\widetilde{\text{cos}}[\mathbf{x}]$, which quantifies how well the KT conditions hold at \mathbf{x} . Our algorithm uses multiple cores of our PC. Our algorithm turns out to be more efficient and effective than OptQuest is, in experiments with a mathematical example and an (s, S) inventory system.

Future research may investigate the following topics: (i) To improve the efficiency and effectiveness of our algorithm, we may customize PS or replace PS by a better optimizer. (ii) We may consider problems with $k \gg 2$ inputs. (iii) Besides continuous inputs we may investigate other types of input (e.g., integer inputs). (iii) Whereas we assume that the problem uses expected values, we may investigate probabilities such as $PR[w_{h'}(\mathbf{x}) > c_{h'}] \leq 0.10$; Kleijnen (2015, p. 222) gives references for Kriging with probabilities, and Chen et al. (2022) gives a more recent method and references.

Acknowledgement

We use the MATLAB code for stochastic Kriging that was developed by Sebastian Rojas Gonzalez (FWO Postdoctoral Fellow, Hasselt University and Ghent University, Belgium) while he was a Ph.D. student at KU Leuven (Belgium). To optimize our two examples through OptQuest (as an alternative for our own algorithm), we use a free 45-day trial version of OptQuest’s SimWrapper (<https://www.opttek.com/products/optquest/simwrapper>); we thank Ben Thengvall (OptTek) for his help with OptQuest and its interface with the MATLAB code for our examples. We thank Inneke van Nieuwenhuyse (Hasselt University) and Wim van Beers for their help with several prior versions of our paper

References

- Ankenman, B., B. Nelson, and J. Staum (2010) Stochastic kriging for simulation metamodeling. *Operations Research*, 58, no. 2, pp. 371–382
- Bashyam, S. and M.C. Fu (1998) Optimization of (s, S) inventory systems with random lead times and a service-level constraint. *Management Science*, 44, pp.243–256
- Boyd S. and L. Vandenberghe (2009), *Convex optimization; seventh printing with corrections*. Cambridge University Press, United Kingdom
- Chen, J. ,M.M. Aral, S-H. Kim, C. Park, and Y. Xie (2022): Constrained Bayesian optimization and spatio-temporal surveillance for sensor network design in the presence of measurement errors. *Engineering Optimization*, pp. 1–16
- Dykstra, R.L. (1970), Establishing the positive definiteness of the sample covariance matrix. *Annals of Mathematical Statistics*, 41, no.6, pp. 153–2154
- Erickson, C.B., B.E. Ankenman, M. Plumlee, and S.M. Sanchez (2021), Gradient-based criteria for sequential experiment design. *Quality and Reliability Engineering*, pp. 3084-3107
- Frazier, P.I. (2018), A tutorial on Bayesian optimization. <https://arxiv.org/pdf/1807.02811.pdf>
- Gramacy, R.B., G.A. Gray, S. Le Digabel, H.K.H. Lee, P. Ranjan, G. Wells, and S.M. Wild (2016) Modeling an augmented Lagrangian for blackbox constrained optimization, *Technometrics*, 58, no. 1, pp. 1–11 (Comments by several authors, and rejoinder by Gramacy et al. pp. 12–29)
- Guidetti, X., A. Rupenyan, L. Fassel, M.Nabavi, and J. Lygeros (2022), Advanced manufacturing configuration by sample-efficient batch Bayesian optimization. arXiv:2205.11827v1 [cs.LG] 24 May 2022

- Jalali, H., I. van Nieuwenhuysse, and V. Picheny (2017), Comparison of Kriging-based algorithms for simulation optimization with heterogeneous noise. *European Journal of Operational Research*, 261, pp. 279–301
- Jones, D., M. Schonlau, and W. Welch (1998). Efficient global optimization of expensive blackbox functions. *Journal of Global Optimization*, 13, pp. 455–492
- Kleijnen, J.P.C. (2015), *Design and analysis of simulation experiments; second edition*. Springer
- Kleijnen, J.P.C. and E. Mehdad (2014), Multivariate versus univariate Kriging metamodels for multi-response simulation models. *European Journal of Operational Research*, 236, pp. 573–582
- Kleijnen, J. P. C. and W. van Beers (2022), Statistical tests for cross-validation of Kriging models. *INFORMS Journal on Computing*, 34, no. 1, pp. 607-621
- Kleijnen J.P.C., W. van Beers, and I. van Nieuwenhuysse (2010) Constrained optimization in simulation: a novel approach. *European Journal of Operational Research*, 202, pp. 164–174
- Kleijnen J.P.C., I. van Nieuwenhuysse, and W.C.M van Beers (2022), Constrained optimization in simulation: efficient global optimization and Karush-Kuhn-Tucker conditions. CentER Discussion Paper Series, No. 2022-015
- Kleijnen, J. P. C. and J. Wan (2007), Optimization of simulated systems: OptQuest and alternatives. *Simulation Modelling Practice and Theory*, 15, pp. 354–362
- Kolman, B. and D.R. Hill (2008), *Elementary Linear Algebra with Applications; 9th edition*. Pearson International Edition, Upper Saddle, New Jersey
- Law, A.M. (2015), *Simulation modeling and analysis; fifth edition*. McGraw-Hill, Boston
- L’Ecuyer, P. (2015), Random number generation with multiple streams for sequential and parallel computing. *Proceedings of the 2015 Winter Simulation Conference*, edited by L. Yilmaz, W. K. V. Chan, I. Moon, T. M. K. Roeder, C. Macal, and M. D. Rossetti
- Liu, H., S. Xu, X. Chen, X. Wang, and Q. Ma (2017), Constrained global optimization via a DIRECT-type constraint-handling technique and an adaptive metamodeling strategy. *Structural and Multidisciplinary Optimization*, 55, no. 1, pp. 155–177
- Loeppky, J.L., J. Sacks, and W. Welch (2009) Choosing the sample size of a computer experiment: a practical guide. *Technometrics*, 51, no. 4, pp. 366–376
- Luo, J., L.J. Hong, B.L. Nelson, and Y. Wu (2015), Fully sequential procedures for large-scale ranking-and-selection problems in parallel computing environments. *Operations Research*, 63, no. 5, pp. 1177–1194
- Mahévas, S., V. Picheny, P. Lambert, N. Dumoulin, L. Rouan, J. Soulié, D. Brockhoff, S. Lehuta, R. Le Riche, R. Faivre, H. Drouineau (2019), A practical guide for conducting calibration and decision-making optimisation with complex ecological models. *Preprints*, 2019120249 (<https://doi.org/10.20944/preprints201912.0249.v1>)

- Pandita, P., N. Awalgaonkar, I. Bilonis, and J. Panchal (2019), Learning arbitrary quantities of interest from expensive black-box functions through Bayesian sequential optimal design. arXiv:1912.07366v1 [stat.ML] 16 Dec 2019
- Pearce, M.A.L. , M. Poloczek, and J. Branke (2022), Bayesian optimization allowing for common random numbers. *Operations Research*, preprint
- Pedrielli, G., S. Wang, S.H. Ng (2020), An extended two stage sequential optimization approach: properties and performance. *European Journal of Operational Research*, 287, pp. 929–945
- Pourmohamad, T. and H.K.H. Lee (2020), The statistical filter approach to constrained optimization. *Technometrics*, 62, no. 3, pp. 303–312
- Pourmohamad T. and H.K.H. Lee (2021), Bayesian optimization via barrier functions. *Journal of Computational and Graphical Statistics*, accepted
- Qing, J., H.B. Moss, T. Dhaene, and I. Couckuyt (2022), $\{PF\}^2ES$: Parallel feasible Pareto frontier entropy search for multi-objective Bayesian optimization under unknown constraints. arXiv preprint arXiv:2204.05411, 2022
- Sóbester, A., S. J. Leary, and A. J. Keane (2005), On the design of optimization strategies based on global response surface approximation models. *Journal of Global Optimization*, 33, no. 1, pp. 31–59
- Tao, I., G. Zhao, and S. Ren (2020), An efficient Kriging-based constrained optimization algorithm by global and local sampling in feasible region. *Journal of Mechanical Design*, 142, no. 5, pp. 051401-1 – 051401-15
- Wang, Z. and M. Ierapetritou (2018), Constrained optimization of black-box stochastic systems using a novel feasibility enhanced Kriging-based method. *Computers and Chemical Engineering*, 118, pp. 210–223

Appendix 1: Common random numbers

In this appendix we investigate CRN with the usual assumption; namely, the number of replications per input combination is constant so $m(\mathbf{x}_i) = m$. We investigate two analysis methods; namely, SK and OK. CRN analyzed by SK requires the estimates $s[w_h(\mathbf{x}_i), w_h(\mathbf{x}_{i'})]$ with $i, i' = 1, \dots, n$. These $s[w_h(\mathbf{x}_i), w_h(\mathbf{x}_{i'})]$ give the $n \times n$ matrix $\widehat{\Sigma}_{\bar{e};h}$, which is *nonsingular* only if $m > n$, as Dykstra (1970) proves. Our initial design uses $n = (k + 1)(k + 2)/2$ if $k \leq 6$ and $n = 5k$ if $k > 6$, following Tao et al. (2020)'s n for sequential optimization in deterministic simulation. Whichever n is selected, a high m may be required to obtain a non-singular $\widehat{\Sigma}_{\bar{e};h}$.

The literature on SK with CRN is very limited. Chen et al. (2012) does investigate SK with CRN, reporting a higher MSPE and a lower variance for the estimated gradient $\widehat{\nabla}(\mathbf{x})$. We focus on constrained optimization, and try to explain why this optimization implies that CRN do not perform well. Obviously, CRN with $m(\mathbf{x}_i) = m$ conflicts with the popular *allocation* rules that require more replications for input combinations with higher output variability (quantified by $\text{var}[w(\mathbf{x}_i)]$ or $\text{var}[w(\mathbf{x}_i)]^{1/2}$).

Note that Pearce et al. (2022) allows CRN, introduces a novel GP metamodel, and uses the Knowledge Gradient (KG)—in a Bayesian framework—to solve unconstrained optimization. We, however, use a GP that is conceptually simpler and requires the estimation of fewer parameters; moreover, we use MEI (which is simpler and therefore more popular than KG) in constrained (instead of unconstrained) optimization. (Pearce et al. (2022) also reviews the literature on unconstrained optimization with or without CRN.)

CRN with $m_i = m$ gives m IID vectors $\mathbf{w}_{h;r}$, each with the n components $w_{i;h;r}$ with mean $E(w_{i;h})$ and variance $\sigma_{i;h}^2$, while $w_{i;h}$ and $w_{i';h}$ have covariance $\sigma_{i;i';h}$. Like Chen et al. (2012), we estimate this $\sigma_{i;i';h}$ through

$$s(\bar{w}_{i;h}, \bar{w}_{i';h}) = \frac{\sum_{r=1}^m (w_{i;h;r} - \bar{w}_{i;h})(w_{i';h;r} - \bar{w}_{i';h})}{(m-1)m}. \quad (28)$$

Now we derive OK with CRN, which requires neither $s^2[w_h(\mathbf{x}_i)]$ nor $s[w_h(\mathbf{x}_i), w_h(\mathbf{x}_{i'})]$; moreover, $m > 1$ suffices. Originally, OK was derived for deterministic simulation, which implies that $\widehat{\Sigma}_{\bar{e};h}$ vanishes; i.e., the OK predictor (say) $\widehat{y}_{\text{OK};h}$ and its estimated MSPE simplify drastically, compared with the SK predictor and its estimated MSPE. If $\widehat{\Sigma}_{M;h}$ is *nearly singular* (or ill-conditioned), then we may add a small *nugget* effect; i.e., $\widehat{\Sigma}_{\bar{e};h}$ becomes $d_h \mathbf{I}_{n \times n}$ where d_h is a small number that is determined by the CPU accuracy; see Gramacy and Lee (2012), and Lophaven et al. (2002, p. 11). However, the original interpretation of the nugget effect is found in geostatistics; namely, the noise is determined by measurement errors etc.; this noise is denoted by d_h —or better σ_h^2 ; see Cressie (1993).

Now we introduce a new OK estimator for random simulation with CRN and $m_i = m$. We point out that $\mathbf{w}_{h;r}$ is the simulation output h in replication r with a fixed seed, so \mathbf{w}_r is *not random* anymore (Pearce et

al. (2022) also emphasizes that the simulation output is deterministic, once the seed is fixed). Therefore, we apply OK to this $\mathbf{w}_{h;r}$ and obtain $\hat{y}_{\text{OK};h;r}(\mathbf{x}|\mathbf{X}, \mathbf{W}_{n \times t})$ or briefly $\hat{y}_{\text{OK};h;r}(\mathbf{x})$, which uses $\hat{\boldsymbol{\psi}}(\mathbf{X}, \mathbf{w}_{h;r})$:

$$\hat{y}_{\text{OK};h;r}(\mathbf{x}) = \hat{\mu}_{h;r} + \hat{\boldsymbol{\sigma}}_{M;h;r}(\mathbf{x})' \hat{\boldsymbol{\Sigma}}_{M;h;r}^{-1}(\mathbf{w}_{h;r} - \hat{\mu}_{h;r} \mathbf{1}_n) \quad (29)$$

This $\hat{y}_{\text{OK};h;r}(\mathbf{x})$ would be the *best linear unbiased predictor* (BLUP) if the Kriging parameters $\boldsymbol{\psi}_h$ were known and the Kriging metamodel were perfectly valid so $E[\hat{y}_{\text{OK};h;r}(\mathbf{x})] = E[w_h(\mathbf{x})]$. Furthermore, we assume that the (nonlinear) estimator $\hat{y}_{\text{OK};h;r}(\mathbf{x})$ is normally distributed. Then the best estimator of $E[w_h(\mathbf{x})]$ is the OK predictor averaged over the m replications:

$$\bar{y}_{\text{OK};h}(\mathbf{x}) = \frac{\sum_{r=1}^m \hat{y}_{\text{OK};h;r}(\mathbf{x})}{m}. \quad (30)$$

The $\mathbf{w}_{h;r}$ give m IID $\hat{y}_{\text{OK};h;r}(\mathbf{x})$, so

$$\text{var}[\bar{y}_{\text{OK};h}(\mathbf{x})] = \frac{\text{var}[\hat{y}_{\text{OK};h}(\mathbf{x})]}{m}. \quad (31)$$

So, $m \uparrow \infty$ implies $\text{var}[\bar{y}_{\text{OK};h}(\mathbf{x})] \downarrow 0$; i.e., $\bar{y}_{\text{OK};h}(\mathbf{x})$ converges to $E[w_h(\mathbf{x})]$.

If the predictor is unbiased, then the MSPE reduces to the *extrinsic variance*; also see the subscript M in $\hat{y}_h(\mathbf{x}_*)$ (obviously, $\hat{y}_{\text{OK};h}(\mathbf{x})$ is a special case of $\hat{y}_h(\mathbf{x})$; i.e., $\hat{\boldsymbol{\Sigma}}_{\bar{e};h}$ vanishes in $\hat{y}_h(\mathbf{x})$ and $s^2[\hat{y}_h(\mathbf{x})]$). We follow the usual practice, and estimate this variance by the plug-in estimator

$$s_M^2[\hat{y}_{\text{OK};h;r}(\mathbf{x})] = \hat{\tau}_{h;r}^2 - \hat{\boldsymbol{\sigma}}_{M;h;r}(\mathbf{x})' \hat{\boldsymbol{\Sigma}}_{M;h;r}^{-1} \hat{\boldsymbol{\sigma}}_{M;h;r}(\mathbf{x}) + \frac{[1 - \mathbf{1}'_n \hat{\boldsymbol{\Sigma}}_{M;h;r}^{-1} \hat{\boldsymbol{\sigma}}_{M;h;r}(\mathbf{x})]^2}{\mathbf{1}'_n \hat{\boldsymbol{\Sigma}}_{M;h;r}^{-1} \mathbf{1}_n}. \quad (32)$$

Though this variance estimator is not normally distributed, we simply use the average (instead of a robust estimator such as the sample median):

$$\overline{s_M^2}[\hat{y}_{\text{OK};h}(\mathbf{x})] = \frac{\sum_{r=1}^m s_M^2[\hat{y}_{\text{OK};h;r}(\mathbf{x})]}{m}. \quad (33)$$

Notice that (33) and (32) have the same expected value (namely, the extrinsic variance, given \mathbf{X} and $\mathbf{W}_{n \times t}$),

but (33) has a smaller variance:

$$\text{var}\{\overline{s_M^2}[\hat{y}_{\text{OK};h}(\mathbf{x})]\} = \frac{\text{var}\{s_M^2[\hat{y}_{\text{OK};h}(\mathbf{x})]\}}{m}.$$

So, $m \uparrow \infty$ implies $\text{var}\{\overline{s_M^2}[\widehat{y}_{\text{OK};h}(\mathbf{x})]\} \downarrow 0$; i.e., $\overline{s_M^2}[\widehat{y}_{\text{OK};h}(\mathbf{x})]$ converges to the true extrinsic variance.

It is well known that if $\mathbf{x} = \mathbf{x}_i$, then $s_M^2[\widehat{y}_{\text{OK};h;r}(\mathbf{x})]$ (defined in (32)) is exactly zero (OK is a perfect interpolator), for all r -values. Hence, its average $\overline{s_M^2}[\widehat{y}_{\text{OK};h}(\mathbf{x})]$ is also zero. However, if $\mathbf{x} = \mathbf{x}_i$, then there is still intrinsic noise; i.e., $\widehat{y}_{\text{OK};h;r}(\mathbf{x}_i) = w_{i,h;r}$ so

$$s^2[\widehat{y}_{\text{OK};h;r}(\mathbf{x}_i)] = s^2(w_{i,h;r}) = \frac{\sum_{r=1}^m (w_{i,h;r} - \bar{w}_{i,h})^2}{(m-1)}.$$

This gives

$$s^2[\widehat{y}_{\text{OK};h;r}(\mathbf{x}_i)] = \frac{s^2[\widehat{y}_{\text{OK};h;r}(\mathbf{x}_i)]}{m}.$$

For a general \mathbf{x} , we use

$$s_{\bar{e}}^2[\widehat{y}_{\text{OK};h}(\mathbf{x})] = \frac{\sum_{r=1}^m [\widehat{y}_{\text{OK};h;r}(\mathbf{x}) - \bar{y}_{\text{OK};h}(\mathbf{x})]^2}{(m-1)m}; \quad (34)$$

also see the subscript \bar{e} in $\widehat{y}_h(\mathbf{x}_*)$. So, $m \uparrow \infty$ gives $s_{\bar{e}}^2[\widehat{y}_{\text{OK};h}(\mathbf{x})] \downarrow 0$; also see (31). Altogether, we decide to use

$$s^2[\widehat{y}_{\text{OK};h}(\mathbf{x})] = \overline{s_M^2}[\widehat{y}_{\text{OK};h}(\mathbf{x})] + s_{\bar{e}}^2[\widehat{y}_{\text{OK};h}(\mathbf{x})]. \quad (35)$$

We use the average of $\widehat{\nabla}[\widehat{y}_{\text{OK};h;r}(\mathbf{x})]$ (estimated gradient of $\widehat{y}_{\text{OK};h}$ at \mathbf{x} , given \mathbf{X} and $\mathbf{W}_{n \times t, r}$ —for CRN with $r = 1, \dots, m$):

$$\overline{\widehat{\nabla}}[\widehat{y}_{\text{OK};h}(\mathbf{x})] = \frac{\sum_{r=1}^m \widehat{\nabla}[\widehat{y}_{\text{OK};h;r}(\mathbf{x})]}{m} \text{ with } h = 0, \dots, t-1. \quad (36)$$

Note that our method is inspired by Kleijnen (2015, pp. 109–110)’s sensitivity analysis (not optimization) through *linear regression* analysis in random simulation with CRN and variance heterogeneity; i.e., instead of *generalized LS* (GLS) with estimated $\widehat{\Sigma}_e$, this analysis may use *ordinary LS* (OLS). Furthermore, CRN may be combined with *antithetic random numbers* (ARN); e.g., Chih (2013) uses such a combination in a second-order polynomial metamodel.

Note that our algorithm uses only the point estimators $s^2[\widehat{y}_{\text{OK};h}(\mathbf{x})]$ and $\overline{\widehat{\nabla}}[\widehat{y}_{\text{OK};h}(\mathbf{x})]$, and not their estimated variances.

Note that LOO-CV with CRN (analyzed through SK or OK) implies that $\bar{w}_{i,h}$ and $\widehat{y}_{i,h}$ or $\widehat{y}_{\text{OK};i,h}$ (in the numerator) are not independent: $\widehat{y}_{i,h}$ or $\widehat{y}_{\text{OK};i,h}$ depend on $w_{i',h}$ ($i' \neq i$), and $\bar{w}_{i,h}$ also depends on $w_{i',h}$. However, we ignore this dependence.

If our algorithm uses OK with CRN, then it replaces \widehat{y} by \widehat{y}_{OK} .

The condition $m_i = m$ implies that our allocation rule does not apply. If PS revisits an old point, then we add one replication to each of the B points near the estimated boundary; if PS gives a new point, then

we add m replications to that point (so, $\widehat{\Sigma}_{\bar{e};h}$ remains well defined, in SK; OK does not use $\widehat{\Sigma}_{\bar{e};h}$). If $B \geq 5$, then we re-estimate all Kriging parameters; the condition $B \geq 5$ is meant to give acceptable MLEs of the parameters ψ_h .

The next algorithm is a pseudo-code for CRN analyzed by SK. This algorithm closely follows the algorithm for IRN analyzed by SK, presented in the main text. However, the condition $m_i = m$ implies that the algorithm does not apply the sample-allocation rule. We delete all comments on the algorithm in the main text.

Algorithm 2

1. Select the algorithm's control variables $\alpha_E, \alpha, \alpha_{\text{infe}}, \alpha_{\text{min}}, \varepsilon, m, n$.
2. Sample the initial design matrix \mathbf{X} with n combinations of the k simulation inputs x_j that satisfies the input constraints $f_g(\mathbf{x}) \leq c_g$ with $g = 1, \dots, v$.
3. Use \mathbf{X} as input for the simulation model, to obtain m replications using CRN, which gives $\mathbf{W}_{n \times t; r}$.
4. Apply LOO-CV to $(\mathbf{X}, \mathbf{W}_{n \times t; r})$.
5. If LOO-CV rejects the SK models, then obtain a new replication for all n combinations; replace m by $m + 1$; return to Step 4.
6. Use $(\mathbf{X}, \mathbf{W}_{n \times t; r})$ to compute $\widehat{\psi}_h, \widehat{\Sigma}_{\bar{e};h}, \widehat{y}_h(\mathbf{x}_i)$, and $s^2[\widehat{y}_h(\mathbf{x}_i)]$.
7. Use $\widehat{y}_h(\mathbf{x}_i)$ to compute $\widehat{y}_{0; \text{min}} = \min_i \widehat{y}_0(\mathbf{x}_i)$ if $\widehat{y}_{h'}(\mathbf{x}_i) + z_{(1-\alpha)} s[\widehat{y}_{h'}(\mathbf{x}_i)] \leq c_{h'}$ with $h' = 1, \dots, t - 1$ and $\alpha = \alpha_{\text{infe}}$; else $\widehat{y}_{0; \text{min}} = \infty$.
8. Use $\widehat{y}_h(\mathbf{x}_i)$ to compute $\widehat{\nabla}_h[\widehat{y}_h(\mathbf{x}_i)]$.
9. Apply PS to estimate $\widehat{\mathbf{x}}_o$, which maximizes $\widehat{\text{MEI}}(\mathbf{x}) \times \widehat{\text{cos}}[\mathbf{x}]$ such that \mathbf{x} satisfies the $t - 1$ CIs for the constrained outputs h' and the v input constraints.
10. If PS finds a feasible point $\widehat{\mathbf{x}}_o$ that is an old point, then add one replication to all B points near the estimated boundary; else if PS finds a feasible point $\widehat{\mathbf{x}}_o$ that is a new point, then add m replications to this point; return to Step 6.
11. If PS does not find a feasible point $\widehat{\mathbf{x}}_o$, then replace α by $\alpha/2$. If $\alpha \geq \alpha_{\text{min}}$, then return to Step 9; else if $\alpha < \alpha_{\text{min}}$ and $\widehat{\text{MEI}}(\mathbf{x}) < \varepsilon \times |\widehat{y}_{0; \text{min}}|$, then go to the final Step.
12. Present the estimated optimal goal output $\widehat{y}_{0; \text{min}}$.

The next algorithm uses OK to analyze CRN with $m_i = m$ replications.

Algorithm 3

1. Select the algorithm's control variables $\alpha_E, \alpha, \alpha_{\text{infe}}, \alpha_{\text{min}}, \varepsilon, m, n$.
2. Sample the initial design matrix $\mathbf{X}_{n \times k}$ with n combinations of the k simulation inputs x_j that satisfies the input constraints $f_g(\mathbf{x}) \leq c_g$ with $g = 1, \dots, v$.
3. Use $\mathbf{X}_{n \times k}$ as input for the simulation model, to obtain m replications using CRN, which gives $\mathbf{W}_{n \times t; r}$.
4. Apply LOO-CV to $(\mathbf{X}_{n \times k}, \mathbf{W}_{n \times t; r})$.
5. If LOO-CV rejects the OK models, then obtain a new replication for all n combinations; replace m by $m + 1$; return to Step 4.
6. Use $(\mathbf{X}_{n \times k}, \mathbf{W}_{n \times t; r})$ to compute $\hat{\psi}_{h; r}, \hat{y}_{\text{OK}; h; r}(\mathbf{x}_i), s^2[\hat{y}_{\text{OK}; h; r}(\mathbf{x}_i)]$.
7. Use $\hat{y}_{\text{OK}; h}(\mathbf{x}_i)$ to compute $\bar{y}_{\text{OK}; 0; \text{min}} = \min_i \bar{y}_{\text{OK}; 0}(\mathbf{x}_i)$ if $\bar{y}_{\text{OK}; h'}(\hat{\mathbf{x}}_i) + z_{(1-\alpha)} s[\bar{y}_{\text{OK}; h'}(\hat{\mathbf{x}}_i)] \leq c_{h'}$ with $h' = 1, \dots, t - 1$ and $\alpha = \alpha_{\text{infe}}$; else $\bar{y}_{\text{OK}; 0; \text{min}} = \infty$.
8. Use $\hat{y}_{\text{OK}; h; r}(\mathbf{x}_i)$ to compute $\hat{V}[\hat{y}_{\text{OK}; h; r}(\mathbf{x}_i)]$ and its average $\bar{V}[\hat{y}_{\text{OK}; h}(\mathbf{x}_i)]$.
9. Apply PS to estimate $\hat{\mathbf{x}}_o$, which maximizes $\widehat{\text{MEI}}(\mathbf{x}) \times \bar{c}\text{os}[\mathbf{x}]$ such that \mathbf{x} satisfies the $t - 1$ CIs for the constrained outputs h' and the v input constraints.
10. If PS finds a feasible point $\hat{\mathbf{x}}_o$ that is an old point, then add one replication to all B points near the estimated boundary; else if PS finds a feasible point $\hat{\mathbf{x}}_o$ that is a new point, then add m replications to this point; return to Step 6.
11. If PS does not find a feasible point $\hat{\mathbf{x}}_o$, then replace α by $\alpha/2$. If $\alpha \geq \alpha_{\text{min}}$, then return to Step 9; else if $\alpha < \alpha_{\text{min}}$ and $\widehat{\text{MEI}}(\mathbf{x}) < \varepsilon \times |\bar{y}_{\text{OK}; 0; \text{min}}|$, then go to the final Step.
12. Present the estimated optimal goal output $\bar{y}_{\text{OK}; 0; \text{min}}$.

Now we discuss CRN in the *toy example*. We assume that CRN indeed “work”; i.e., CRN create *positive* correlations between $e_h(\mathbf{x})$ and $e_h(\mathbf{x}')$ (outputs of the same type h , at points \mathbf{x} and \mathbf{x}'). Furthermore, we assume that $\rho_{i; i'; h}$ increases, as the simulated input combinations \mathbf{x}_i and $\mathbf{x}_{i'}$ resemble each other more; i.e., we use the Gaussian correlation function, which uses $(x_{i; j} - x_{i'; j})^2$ (so, our assumption resembles the assumption for the *extrinsic* noise in a GP). Initially we used LHS with $n = 20$ (instead of 6) midpoints. Combining this LHS and the Gaussian correlation function implies that the toy example has $\min_i (x_{i; j} - x_{i'; j})^2$

$= (1/20)^2 = 0.0025$ and $\max_i (x_{i;j} - x_{i';j})^2 = 0.950^2 = 0.9025$. Obviously, if $(x_{i;j} - x_{i';j}) \downarrow 0$, then $\rho_{i;i';h} \uparrow 1$; if $(x_{i;j} - x_{i';j}) \uparrow \infty$, then $\rho_{i;i';h} \downarrow 0$. However, $n = 20$ implies $\max_i (x_{i;j} - x_{i';j}) = 0.950 (\ll \infty)$; we decide which value $\rho_{i;i';h}$ should have at $x_{i';j} = 0.975$ (biggest midpoint). Obviously, this value—together with $\rho_{i;i';h} = 1$ at $x_{i;j} = 0$ —determines θ_j , which features in the Gaussian correlation function. Furthermore, we define a $\rho_{i;i';h}$ between 0.5 and 1.0 as *high*, and between 0.0 and 0.5 as *low*. To select specific values for $\rho_{i;i';h}$ in these intervals, we use uniform sampling. Furthermore, we assume that $\theta_1 = \theta_2$, which is perfect for e_0 and e_2 ; for simplicity's sake, we also assume $\theta_1 = \theta_2$ for e_1 . For example, $\rho_{i;i';h} = 0.5$ gives $\theta_1 = \theta_2 = 0.3128$, and $\rho_{i;i';h} = 0.1$ gives $\theta_1 = \theta_2 = 2.0781$.

Altogether we obtain *many* covariance matrixes; namely, one matrix per output type, for small noise and big noise, per initial 20×20 design \mathbf{X} , which is sampled per macroreplication and increases as new points are added. To illustrate these matrixes, we display the upper-left 5×5 submatrix of the 20×20 (symmetric) covariance matrix for e_1 —for small noise and high CRN effect—with $\sigma_{e;1}(\mathbf{x})$ (instead of $\sigma_{e;1}^2(\mathbf{x})$) on the main diagonal, and the (scale-free) $\rho_{i;i';1}$ (so $\text{cov}[e_1(\mathbf{x}_i), e_1(\mathbf{x}_{i'})] = \rho_{i;i';1} \times \sigma_{e;1}(\mathbf{x}_i) \times \sigma_{e;1}(\mathbf{x}_{i'})$):

$$\begin{bmatrix} 1.21 & 0.97 & 0.94 & 0.94 & 0.98 \\ & 1.32 & 0.93 & 0.91 & 0.93 \\ & & 1.18 & 0.99 & 0.87 \\ & & & 0.81 & 0.87 \\ & & & & 0.76 \end{bmatrix} \quad (37)$$

The values $\rho_{1;2;1} = 0.97$ and $\rho_{1;3;1} = 0.94$ in this matrix imply that \mathbf{x}_1 and \mathbf{x}_2 are closer than \mathbf{x}_1 and \mathbf{x}_3 are. We must check whether $\Sigma_{e;20 \times 20}$ is *positive definite* (PD). MATLAB offers two methods: (i) “attempt Cholesky factorization”, and (ii) “check eigenvalues”; we use method (ii).

We assume *zero cross-correlations* $\rho_{i;0;1}$, $\rho_{i;0;2}$, and $\rho_{i;1;2}$ in the toy example, whereas a (random) *discrete-event* simulation model with multiple output types (so, $t > 1$) gives non-zero cross-correlations because the t outputs are generated by the same simulation model (also see the inventory example in Section 6). Unfortunately, we discovered that the toy example with non-zero cross-correlations gives serious numerical problems (we must then use a 60×60 covariance matrix that uses 20 3×3 submatrixes $\Sigma_e(\mathbf{x}_i, \mathbf{x}_{i'})$ instead of the preceding three 20×20 matrixes).

Once we have confirmed that $\Sigma_{h;n \times n}$ —with initial value $n = 20$ —is PD, we sample the n -variate $\mathbf{e}_{h;n}$ from the *multi-variate normal density* (MND) with all n means equal to zero and the non-diagonal covariance matrix $\Sigma_{h;n \times n}$:

$$\mathbf{e}_{h;n} \sim N_n(\mathbf{0}_n, \Sigma_{h;n \times n}) \quad (38)$$

where our toy example initially uses $n = 10k = 20$, so $\Sigma_{h;n \times n}$ is 20×20 and $h = 0, 1, 2$; we call the corresponding \mathbf{X} the *initial* or *pilot* design. To sample from a MND, we use the MATLAB function *mvnrnd*. Given $e_{h;n}$ sampled through (38), we compute the corresponding w_h ($h = 0, 1, 2$).

Next, our algorithm finds an estimated optimal point $\hat{\mathbf{x}}_o$. There are two possibilities: (i) This $\hat{\mathbf{x}}_o$ is an *old* point, so a replication is added for both this point and all other old points “close” to the boundary, using CRN; i.e., m increases by 1 and the augmented I/O simulation data becomes $(\mathbf{X}, w_{h,r}(\hat{\mathbf{x}}_{o,r}))$ (with $h = 0, 1, 2$ and $r = 1, \dots, m$) where $w_{h,r}(\hat{\mathbf{x}}_{o,r})$ uses $e_{h;n}$ sampled through (38). (ii) This $\hat{\mathbf{x}}_o$ is a *new* point (near the boundary), so m replications are obtained for this point using the same PRNs as the old points used for the corresponding replications; i.e., m does not increase but n does increase by 1 and the augmented I/O simulation data becomes $(\mathbf{X}, w_{h,r}(\hat{\mathbf{x}}_o))$ where $w_{h,r}(\hat{\mathbf{x}}_o)$ is sampled through (40)—explained next. If $\hat{\mathbf{x}}_o$ is “very close” to one of the old points, then DACE gives the error message “Multiple design sites are not allowed”. In our toy example (with $0 \leq x_j \leq 1$) we consider $\hat{\mathbf{x}}_o$ to be an old point if the absolute values of the differences in both input dimensions are less than 0.001 (so m becomes $m + 1$). Note that (by definition) each *iteration* after the initial design increases either m or n (i.e., adds either a column or a row to \mathbf{W}_h).

If $\hat{\mathbf{x}}_o$ is a new point, then $\hat{\mathbf{x}}_o$ should give outputs that are correlated with the outputs of the old outputs (because of CRN); e.g., if CRN happens to give relatively high old outputs (i.e., many sampled values of $e_{h;n}$ exceed $\mathbf{0}_n$), then the new outputs should also be relatively high (i.e., $e_{h;n+1}(\hat{\mathbf{x}}_o) > 0$). We define $e_{h;n+1} = (e_{h;n}, e_{h;o})'$ where $e_{h;o} = e_h(\hat{\mathbf{x}}_o)$. Obviously, $E(e_{h;n}) = \mathbf{0}_n$ and $\Sigma_{h;n \times n}$ is the $n \times n$ covariance matrix in (38). Furthermore, $\sigma_{o,h}^2$ denotes the variance of $e_{o,h}$ at $\hat{\mathbf{x}}_o$. Finally, $\sigma_{n;o}$ denotes the n -dimensional vector with the covariances between the n old intrinsic noises and the new intrinsic noise at $\hat{\mathbf{x}}_o$, for response h . Altogether, we define the following $(n + 1) \times (n + 1)$ covariance matrix

$$\Sigma_{h;(n+1) \times (n+1)} = \begin{bmatrix} \Sigma_{h;n \times n} & \sigma_{n;o} \\ \sigma'_{n;o} & \sigma_{o,h}^2 \end{bmatrix}. \quad (39)$$

Obviously, this matrix should be PD. If the sampled values for $\rho_{i;o,h}$ —together with the n old points \mathbf{x}_i ($i = 1, \dots, n$) and the new point \mathbf{x}_o —do not give a $\sigma_{n;o}$ that gives a PD matrix, then we resample $\rho_{i;o,h}$ and compute the corresponding $\theta_1 = \theta_2$ —until the matrix is PD. To sample from the *conditional* distribution of $e_h(\hat{\mathbf{x}}_o)$ given e_n , we change the mean vector and the covariance matrix in (38) (e.g., Mittelhammer (1996)):

$$e_h(\hat{\mathbf{x}}_o) \mid e_{h;n} \sim N(\sigma'_{n;o} \Sigma_{h;n \times n}^{-1} e_{h;n}, \sigma_{o,h}^2 - \sigma'_{n;o} \Sigma_{h;n \times n}^{-1} \sigma_{n;o}) \quad (40)$$

We have already seen that $\Sigma_{h;n \times n}$ is PD (see (37)), but now (40) uses $\Sigma_{h;n \times n}^{-1}$ —which might be *nearly*

singular (MATLAB then gives the error message: "the matrix is close to be singular" and we add a nugget effect). Fortunately, our toy example does not give these numerical problems.

Each *macroreplication* samples a new initial LHS design, independent of the designs in the other macroreplications. This new design implies a new (20×20) $\Sigma_{h;n \times n}$.

Appendix 2: The heterogeneous variances of the intrinsic noise in the toy problem

Jalali et al. (2017) gives the following two equations:

$$\begin{aligned} [\max_{\mathbf{x}} w_h(\mathbf{x}) + b_h] \times a_h &= \max_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}) \\ [\min_{\mathbf{x}} w_h(\mathbf{x}) + b_h] \times a_h &= \min_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}), \end{aligned} \quad (41)$$

and defines the *range* of $w_h(x)$ as $R_h = \max_{\mathbf{x}} E[w_h(\mathbf{x})] - \min_{\mathbf{x}} E[w_h(\mathbf{x})]$ and distinguishes between small noise $0.6R_h$ or $0.15R_h$ and big noise $6.0R_h$ or $1.5R_h$. Because that publication estimates the global minimum of an unconstrained problem, the *best* case means that $\sigma_{e;h}(\mathbf{x})$ decreases linearly as $w_h(\mathbf{x})$ decreases. So, the right-hand sides of (41) use

$$\begin{aligned} \max_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}) &= 0.60R_h \text{ or } 6.00R_h \\ \min_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}) &= 0.15R_h \text{ or } 1.50R_h. \end{aligned} \quad (42)$$

Combining (41) and (42), we can determine a_h and b_h in (41); this gives the desired linear equation

$$\sigma_{e;h}(\mathbf{x}) = a_h b_h + a_h E[w_h(\mathbf{x})]. \quad (43)$$

The *worst* case means that $\sigma_{e;h}(\mathbf{x})$ increases linearly as $E[w_h(\mathbf{x})]$ decreases. Combining the analogue of (41) with (42) gives

$$\begin{aligned} \{b_h + \min_{\mathbf{x}} E[w_h(\mathbf{x})]\} \times a_h &= \max_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}) \\ \{b_h + \max_{\mathbf{x}} E[w_h(\mathbf{x})]\} \times a_h &= \min_{\mathbf{x}} \sigma_{e;h}(\mathbf{x}). \end{aligned} \quad (44)$$

Using the data in the toy example, we compute $\min_{\mathbf{x}} E[w_0(\mathbf{x})] = 0$ and $\max_{\mathbf{x}} E[w_0(\mathbf{x})] = 2$ so $R_0 = 2$;

hence, (43) (best case with smallest noise) gives

$$\sigma_{e;0}(\mathbf{x}) = 0.30 + 0.45E[w_0(\mathbf{x})]. \quad (45)$$

Furthermore, $\min_x E[w_1(\mathbf{x})] = -1.5$ and $\max_x E[w_1(\mathbf{x})] = 1.6711$ so $R_1 = 3.1711$, and $\min_x E[w_2(\mathbf{x})] = -1.5$ and $\max_x E[w_2(\mathbf{x})] = 0.5$ so $R_2 = 2$. The main text gives $\sigma_{e;1}(\mathbf{x})$ and $\sigma_{e;2}(\mathbf{x})$ for the best case with smallest noise; furthermore, $\sigma_{e;h}(\mathbf{x})$ for the worst case with biggest noise is

$$\begin{aligned} \sigma_{e;0}(\mathbf{x}) &= 12.0 - 4.5E[w_0(\mathbf{x})] \\ \sigma_{e;1}(\mathbf{x}) &= 12.2766 - 4.5E[w_1(\mathbf{x})] \\ \sigma_{e;2}(\mathbf{x}) &= 5.25 - 4.5E[w_2(\mathbf{x})] \end{aligned} \quad (46)$$

Appendix 3: EI versus MEI

Jones et al. (1998) derives the following estimator of EI where $w_{\min} = \min_{1 \leq i \leq n} w(x_i)$ and we use our symbol $\hat{y}_{\text{OK}}(\mathbf{x})$:

$$\begin{aligned} \widehat{\text{EI}}(\mathbf{x}) &= (w_{\min} - \hat{y}_{\text{OK}}(\mathbf{x})) \Phi \left(\frac{w_{\min} - \hat{y}_{\text{OK}}(\mathbf{x})}{s[\hat{y}_{\text{OK}}(\mathbf{x})]} \right) \\ &\quad + s[\hat{y}_{\text{OK}}(\mathbf{x})] \phi \left(\frac{w_{\min} - \hat{y}_{\text{OK}}(\mathbf{x})}{s[\hat{y}_{\text{OK}}(\mathbf{x})]} \right). \end{aligned} \quad (47)$$

Inspired by Jalali et al. (2017), we define the following estimator of MEI where \mathbf{x}_{\min} is the point that gives w_{\min} , $\hat{y}(\mathbf{x}_{\min})$ replaces w_{\min} :

$$\begin{aligned} \widehat{\text{MEI}}(\mathbf{x}) &= (\hat{y}(\mathbf{x}_{\min}) - \hat{y}(\mathbf{x})) \Phi \left(\frac{\hat{y}(\mathbf{x}_{\min}) - \hat{y}(\mathbf{x})}{s[\hat{y}(\mathbf{x})]} \right) \\ &\quad + s[\hat{y}(\mathbf{x})] \phi \left(\frac{\hat{y}(\mathbf{x}_{\min}) - \hat{y}(\mathbf{x})}{s[\hat{y}(\mathbf{x})]} \right). \end{aligned} \quad (48)$$

Appendix 4: Details for toy problem

This Appendix shows how our algorithm proceeds through various stages (or iterations)—before it reaches its final estimated optimum.

Fig. 11 shows dashed curves to denote the *white-box* constraints $E(w_{h'}) = 0$ ($h' = 1, 2$) (which are the analogues of the solid curves for the deterministic constraints). These curves—combined with the (non-displayed) goal function $E[w_0(\mathbf{x})] = x_1 + x_2$ —give the (black) square that denotes $\mathbf{x}_A \approx (0.1954, 0.4044)'$.

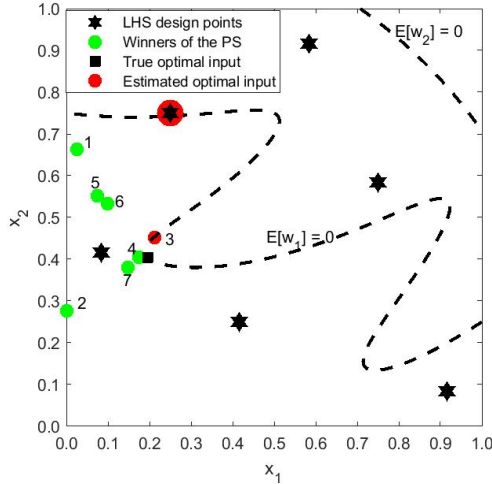


Figure 11: Toy example: input combinations in initial design and next iterations, for a restart ending near A in macroreplication 57 with $\alpha_{\text{infe}} = 10\%$, $m_{\text{min}} = 10$, and 20 PS starting points

Each *restart* begins in iteration 0 using a new initial LHS design. We discuss two of the twelve restarts, in macroreplication 57 (there are 60 macroreplications); namely, one restart that ends near \mathbf{x}_A and has relatively few iterations after iteration 0 (namely, seven so our discussion is relatively short), and one restart that ends near $\mathbf{x}_C = (0.00, 0.75)$ and has relatively many iterations including revisits of old input combinations. All restarts and macroreplications use different LHS designs, but use the same α_{infe} (namely, 10%), m_{min} (= 10), and number of PS starting points (20).

Table 2: Iterations of a restart ending near point A in macroreplication 57

Iter.	$\hat{\mathbf{x}}_{\text{min}}$	$\hat{y}_{0;\text{min}}$	$E[w_0(\hat{\mathbf{x}}_{\text{min}})]$	$\hat{\mathbf{x}}_o$	\widehat{MEI}	$\widetilde{\text{cos}}$
0	(0.2500, 0.7500)	0.9106	1.0000	(0.0250, 0.6625)	0.3712	0.9980
1	(0.2500, 0.7500)	0.9935	1.0000	(0.0006, 0.2769)	0.7276	1.0000
2	(0.2500, 0.7500)	0.9699	1.0000	(0.2125, 0.4515)	0.4198	0.9993
3	(0.2500, 0.7500)	0.9778	1.0000	(0.1740, 0.4047)	0.4632	0.9996
4	(0.2125, 0.4515)	0.5680	0.6640	(0.0750, 0.5515)	0.0138	0.6307
5	(0.2125, 0.4515)	0.5777	0.6640	(0.0989, 0.5325)	0.0125	0.3748
6	(0.2125, 0.4515)	0.5726	0.6640	(0.1490, 0.3797)	0.1033	0.9989
7	(0.2125, 0.4515)	0.6068	0.6640			

The (black) stars denote the initial $n = 6$ input combinations selected through LHS with midpoints. The preprocessing (with LOO-CV) obtains the required number of replications for these combinations. The star within a (red) ball denotes $\hat{\mathbf{x}}_{\text{min}} = (0.2500, 0.7500)'$, in iteration 0; this $\hat{\mathbf{x}}_{\text{min}}$ is also displayed in Table 2, in the row for iteration—or “Iter.”—0 (because we know the true I/O functions, we know that this $\hat{\mathbf{x}}_{\text{min}}$ is indeed feasible). The algorithm uses PS, to compute $\hat{\mathbf{x}}_o$ that maximizes the product of $\widehat{MEI}(\mathbf{x})$ and $\widetilde{\text{cos}}(\mathbf{x})$;

see Table 2, last three columns.

In the Figure, the (green) solid balls denote $\hat{\boldsymbol{x}}_o$ in iteration 1, ..., 7. The dark (red instead of light green) ball denotes the final estimate of $\hat{\boldsymbol{x}}_{\min} = (0.2125, 0.4515)'$, which equals $\hat{\boldsymbol{x}}_o$ obtained in iteration 2.

Furthermore, this Table shows that iteration 0 gives $\hat{y}_{0;\min} = 0.9106$. We (not the algorithm) also know the corresponding $E[w_0(\hat{\boldsymbol{x}}_{\min})] = 0.2500 + 0.7500 = 1.0000$. The winner of the PS is $\hat{\boldsymbol{x}}_o = (0.0250, 0.6625)'$. This $\hat{\boldsymbol{x}}_o$ maximizes the product of $\widehat{\text{MEI}}(\boldsymbol{x})$ and $\widetilde{\text{cos}}(\boldsymbol{x})$; actually, $\widehat{\text{MEI}}(\hat{\boldsymbol{x}}_o) = 0.3712$ (which is a sizable improvement of $\hat{y}_{0;\min} = 0.9106$) and $\widetilde{\text{cos}}(\boldsymbol{x}) = 0.9980$ (which is close to the maximum value 1). At the end of an iteration, the algorithm runs the simulation model with $\hat{\boldsymbol{x}}_o$ as input, obtains (not displayed) $w_{h;r}(\hat{\boldsymbol{x}}_o)$, and uses the augmented I/O data to re-estimate (not displayed) \hat{y}_h , $s[\hat{y}_h]$, and $\widehat{\nabla}_h$. Iteration 1 uses these \hat{y}_h and $s[\hat{y}_h]$, and estimates that $\hat{\boldsymbol{x}}_o = (0.0250, 0.6625)'$ (result of iteration 0) is infeasible; so, $\hat{\boldsymbol{x}}_{\min}$ remains the same as in iteration 0. However, $\hat{y}_{0;\min}$ becomes 0.9935. Of course, if $\hat{\boldsymbol{x}}_{\min}$ remains the same, then $E[w_0(\hat{\boldsymbol{x}}_{\min})]$ remains the same; namely, 1.0000. The re-estimated \hat{y}_h , $s[\hat{y}_h]$, and $\widehat{\nabla}_h$ give a new $\hat{\boldsymbol{x}}_o = (0.0006, 0.2769)'$. The algorithm obtains $w_{h;r}(\hat{\boldsymbol{x}}_o)$, and re-estimates \hat{y}_h , $s[\hat{y}_h]$, and $\widehat{\nabla}_h$.

Iteration 2 estimates $\hat{\boldsymbol{x}}_o = (0.0006, 0.2769)'$ to be infeasible; so, $\hat{\boldsymbol{x}}_{\min}$ does not change. The new $\hat{\boldsymbol{x}}_o$ becomes $(0.2125, 0.4515)'$.

Iteration 3 estimates this $\hat{\boldsymbol{x}}_o$ to be infeasible, so this $\hat{\boldsymbol{x}}_o$ does not become the new $\hat{\boldsymbol{x}}_{\min}$. Actually, the algorithm evaluates all old simulated combinations, and finds that $\hat{\boldsymbol{x}}_{\min}$ remains the same as in iteration 0. The new $\hat{\boldsymbol{x}}_o$ becomes $(0.1740, 0.4047)'$.

Iteration 4 re-estimates \hat{y}_h and $s[\hat{y}_h]$, and estimates $\hat{\boldsymbol{x}}_o = (0.1740, 0.4047)'$ to be infeasible. Moreover, the new $\hat{\boldsymbol{x}}_{\min}$ becomes an old $\hat{\boldsymbol{x}}_o$; namely, $(0.2125, 0.4515)'$ (see iteration 2). And so the algorithm proceeds.

Iteration 7 gives the final $\hat{\boldsymbol{x}}_{\min} = (0.2125, 0.4515)'$ with $\hat{y}_{0;\min} = 0.6068$ (iteration 4 has already given the same $\hat{\boldsymbol{x}}_{\min}$, but a different $\hat{y}_{0;\min} = 0.5680$). The algorithm stops because $\alpha = 0.0063$ falls below its lower bound 0.01.

Fig. 12 displays $\hat{y}_{0;\min}$ as a function of N_q for restart q that requires preprocessing with 237 runs (so the average number of replications per combination is $237/6 = 39.5$) and the total N is 483. The total number of replications ($\sum_{i=1}^n m_{i;s;q}$) in iteration $s = 7$ is much smaller than in iteration 1 (the distance between two consecutive balls projected onto the horizontal axis is smaller in iteration 7). This Figure also shows + signs that denote symmetric 90% CIs around $\hat{y}_{0;\min}$. Obviously, the lengths of these CIs tend to decrease as N increases. (We do not display the Figure for $E[w_0(\hat{\boldsymbol{x}}_{\min})]$, because that Figure resembles Fig. 12.)

Fig. 13 is the analogue of Fig. 11, but finishes near C instead of A.

Fig. 14 is the analogue of Fig. 12. It shows that this restart requires many more simulation runs in the preprocessing; namely, 1,497. The final value of N is 9,883. Its final solution is farther away from the true global optimum.

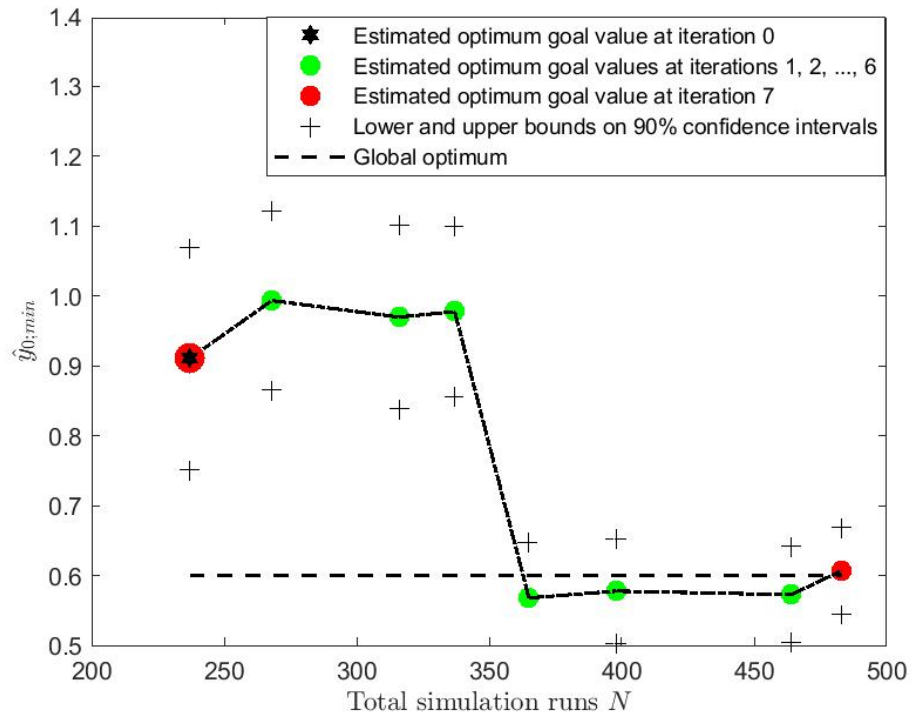


Figure 12: Toy example: estimated optimal goal value when simulation outputs are added in iterations 0 through 7, for a restart ending near A in macroreplication 57 with $\alpha_{\text{infe}} = 10\%$, $m_{\text{min}} = 10$, and 20 PS starting points

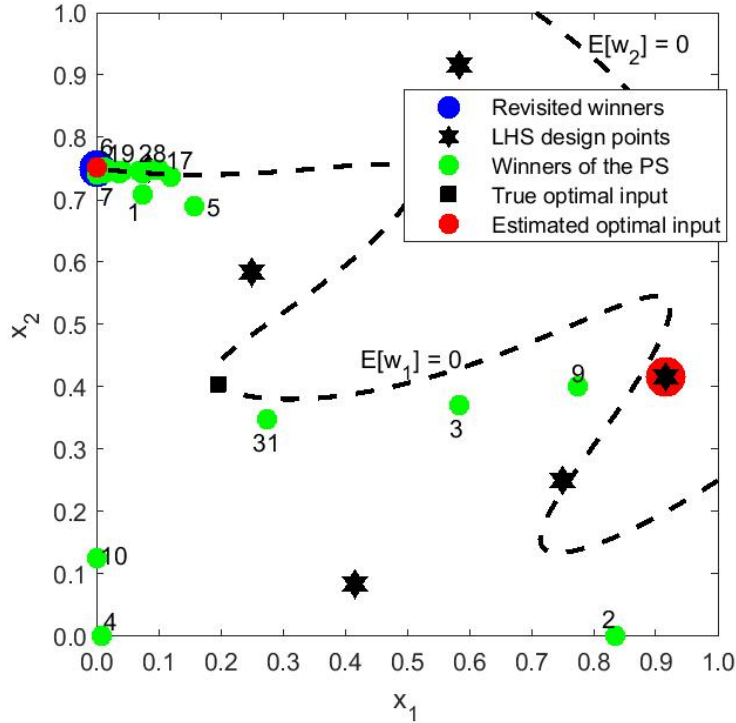


Figure 13: Toy example: input combinations in initial design and next iterations, for a restart ending near C in macroreplication 57 with $\alpha_{\text{infe}} = 10\%$, $m_{\text{min}} = 10$, and 20 PS starting points

References

- Chen, X., B. Ankenman, and B.L. Nelson (2012), The effects of common random numbers on stochastic kriging metamodels. *ACM Transactions on Modeling and Computer Simulation*, 22, no. 2, Article 7
- Chih M (2013), A more accurate second-order polynomial metamodel using a pseudo-random number assignment strategy. *Journal of the Operational Research Society*, 64:198–207
- Cressie, N.A.C. (1993), *Statistics for spatial data. (Revised edition)*. Wiley, New York
- Dykstra, R.L. (1970), Establishing the positive definiteness of the sample covariance matrix. *Annals of Mathematical Statistics*, 41, no.6, pp. 153–2154
- Gramacy, R.B. and H.K.H. Lee (2012), Cases for the nugget in modeling computer experiments. *Statistics and Computing*, 22, no 3, pp. 713–722
- Jalali, H., I. v an Nieuwenhuysse, and V. Picheny (2017), Comparison of Kriging-based algorithms for simulation optimization with heterogeneous noise. *European Journal of Operational Research*, 261, pp. 279–301
- Jones, D., M. Schonlau, and W. Welch (1998). Efficient global optimization of expensive blackbox functions. *Journal of Global Optimization*, 13, pp. 455–492

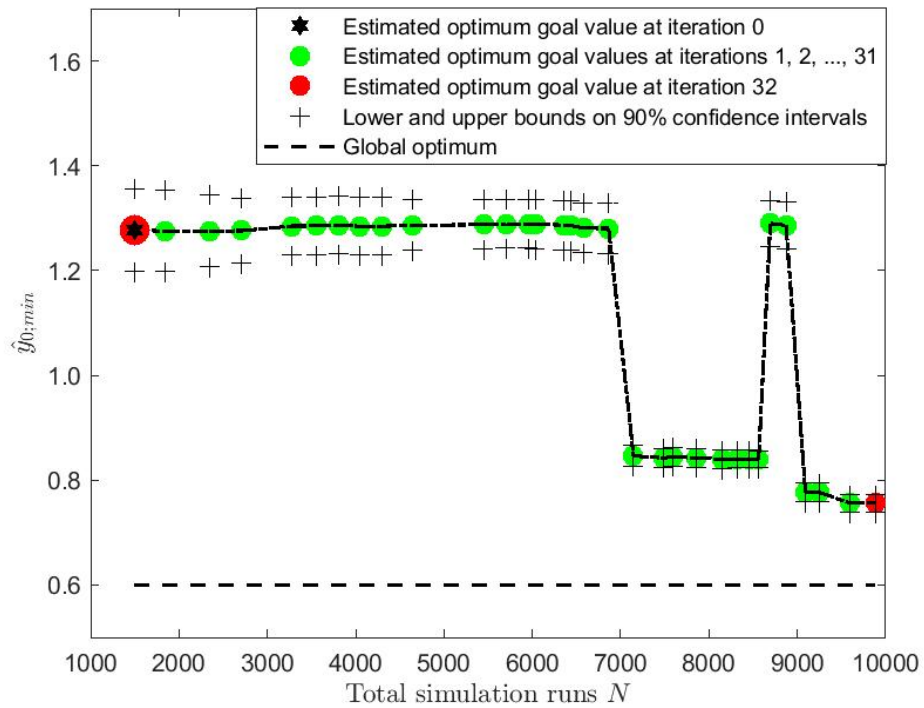


Figure 14: Toy example: estimated optimal goal value when simulation outputs are added in iterations 0, 1, ..., for a restart ending near C in macroreplication 57 with $\alpha_{\text{infe}} = 10\%$, $m_{\text{min}} = 10$, and 20 PS starting points

- Kleijnen, J.P.C. (2015), *Design and analysis of simulation experiments; second edition*. Springer Lophaven S.N., H.B. Nielsen, J. Sondergaard (2002), DACE: A Matlab Kriging Toolbox, Version 2.0 IMM Technical University of Denmark, Kongens Lyngby, Denmark
- Mittelhammer R.C. (1996). *Mathematical Statistics for Economics and Business*. Springer: New York
- Pearce, M.A.L. , M. Poloczek, and J. Branke (2022), Bayesian optimization allowing for common random numbers. *Operations Research*, preprint
- Tao, I., G. Zhao, and S. Ren (2020), An efficient Kriging-based constrained optimization algorithm by global and local sampling in feasible region. *Journal of Mechanical Design*, 142, no. 5, pp. 051401-1 – 051401-15