

Dice and Slice Simulation Optimization for High-Dimensional Discrete Problems

Harun Avci,^a Barry L. Nelson,^b Eunhye Song,^c Andreas Wächter^b

^aDepartment of Economics, Finance, and Quantitative Analysis, Kennesaw State University, Kennesaw, GA 30144, USA, havci@kennesaw.edu;

^bDepartment of Industrial Engineering and Management Sciences, Northwestern University, Evanston, IL 60208, USA,

nelsonb@northwestern.edu, andreas.waechter@northwestern.edu;

^cH. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA, esong32@gatech.edu

Although much progress has been made in simulation optimization, problems involving computationally expensive simulations having high-dimensional, *discrete* decision-variable spaces have been stubbornly resistant to solution. For this class of problems we propose Dice and Slice Simulation Optimization (DASSO). DASSO is a form of Bayesian optimization that represents the prior on the objective function implied by the simulation as a sum of low-dimensional Gaussian Markov random fields. This prior is consistent with the full-dimensional objective function, rather than assuming that it is actually separable. By working iteratively between posteriors on these low-dimensional “dice” and a full-dimensional “slice” of the decision-variable space, DASSO makes rapid progress with little algorithm overhead even on problems with more than a trillion feasible solutions. We achieve further computational savings by showing that we can find the best solution to simulate on each iteration without having to assess the potential of all solutions—as is traditionally done in Bayesian optimization—by identifying a small set of Pareto-optimal solutions in subsets of the dimensions. We prove that DASSO is asymptotically convergent to the optimal solution, while emphasizing that its most important feature is the ability to find good solutions quickly in problems beyond the capability of other methods.

Keywords: simulation; large-scale discrete simulation optimization; Gaussian Markov random fields; Bayesian optimization

1. Introduction

Simulation optimization (SO) is a technique for improving the performance of complex systems under uncertainty, where the design parameters of the simulation model represent the controllable decision variables of the system. Unlike deterministic optimization and stochastic programming problems, the objective function in SO has no explicit form and can only be estimated using simulation experiments. SO problems with discrete decision variables are known as discrete simulation optimization (DSO) problems and arise in many areas of operations research and management science.

For large-scale DSO problems, where only a small fraction of feasible solutions can reasonably be simulated, one can employ metaheuristics, adaptive random search, or inference-based algorithms. Metaheuristics typically provide no statistical performance guarantees. Adaptive random search algorithms, on the other hand, can be shown to achieve asymptotic convergence to the global or a local optimum. However, such convergence results have little practical meaning for large-scale problems if it is not possible to simulate even a modest fraction of the feasible solutions. See Fu et al. (2015) for a review of DSO algorithms.

Inference-based algorithms are closely related to Bayesian optimization (BO), a popular technique for optimizing black-box functions that is known for making rapid search progress. BO provides guidance for the search based on the conditional (posterior) distribution of a stochastic process representing the unknown objective function. A popular choice for the prior on the objective function is a Gaussian process (GP) due to its posterior being Gaussian if the simulation output distribution is also Gaussian. Although GPs with a

continuous domain have been adapted to solve DSO problems (see, e.g., Quan et al. 2013, Sun et al. 2014, Xie et al. 2016), Salemi et al. (2019) demonstrate empirically that continuous-decision-variable covariance functions may fail dramatically when used for guidance and optimality-gap inference in discrete problems. They also show that using a Gaussian Markov random field (GMRF) prior, as an alternative to a continuous GP, provides much better search guidance and stopping inference for DSO. Our new methods start from the foundation of the Gaussian Markov Improvement Algorithm (GMIA) of Salemi et al. (2019) but are dramatically different.

A GMRF is a multivariate Gaussian random vector associated with an undirected graph where the nodes represent feasible solutions and the edges determine correlation structure. Unlike GPs with a continuous domain, GMRFs can be naturally defined on a lattice, making it more suitable for DSO problems. The correlation structure is defined by a precision matrix, the inverse of the covariance matrix, where the nonzero elements correspond to edges in the defining, and ideally sparse, graph. The GMRF-based GMIA algorithm evaluates the complete expected improvement (CEI) of each feasible solution on each iteration to guide the search. CEI predicts each solution’s improvement in the objective function relative to the current sample-best solution based on the posterior distribution of the GMRF (Salemi et al. 2019, Semelhago et al. 2021).

Computing the full posterior distribution at every iteration of GMIA requires inversion of the precision matrix, a computationally expensive calculation. Fortunately, computing the CEI of each feasible solution requires only the diagonal and a single column of the covariance matrix (i.e., the inverse of the precision matrix) of the posterior. While the elements of the required column can be obtained by a direct back-solve, extracting the diagonal elements is computationally more challenging. To address this, Semelhago et al. (2017) propose an efficient way to extract the diagonal elements by exploiting the sparse structure of the precision matrix. They further speed up GMIA by employing a divide-and-conquer strategy to restrict the search within a small promising subset of feasible solutions for several iterations (Semelhago et al. 2021). Alternatively, Li and Song (2020) apply the Sherman-Morrison-Woodbury formula recursively, updating only the necessary elements of the covariance matrix until it is no longer cheaper than inversion.

These computational improvements greatly extend the reach of GMIA-based DSO. Nevertheless, even these variants encounter a limit on the problem size they can attack, for instance when the decision-variable dimension is high, since the computational cost of calculating the posterior distribution increases at least quadratically in the number of feasible solutions. To reach significantly beyond this limit, we propose the Dice and Slice Simulation Optimization (DASSO) algorithm that decomposes the prior distribution into a rigorously justified additive form via a functional analysis of variance (FANOVA). In this decomposition, the first-order terms represent independent GMRFs and the higher-order terms form a random effect with a chosen first-order term. The decomposition reduces the problem dimension to facilitate efficient posterior updates, moving high-dimensional DSO problems that are too large to solve with current technology into the realm of possibility. Our numerical analysis reveals that DASSO can obtain good feasible solutions rapidly on problems with more than a trillion feasible solutions, far beyond the reach of any other algorithm.

Of course, no single SO algorithm is appropriate for all problems. Our particular interest is in stochastic DSO problems with integer-ordered decision variables, very large numbers of feasible solutions, and a simulation that itself is so computationally expensive that one only expects to simulate a tiny fraction of the solution space, even given substantial time and high-performance computing. For such problems rapid progress toward better and better solutions, iteration by iteration, is essential. Although we do prove that DASSO is “convergent,” this is more of academic than practical interest for the class of problems on which we focus. Examples of such computationally prohibitive problems include the design of a new fuel injector production line described in Tongarlak et al. (2010), where 20 replications of a single feasible solution took 8 hours, and a “future mobility” problem that we worked on in collaboration with an automotive manufacturer, where each replication of a single feasible solution took a few hours. Computationally less expensive, yet still large-scale examples include the bike-sharing system modeled as a discrete-event simulation in Jian et al. (2016), and a multi-product inventory problem with an (s, S) policy for each product, where s and S are the reorder and order-up-to levels, respectively. We use this inventory problem for our numerical analysis and as a running example throughout the paper.

In summary, we propose the DASSO algorithm, a type of BO for solving large-scale, computationally expensive DSO problems whose feasible solutions are defined on a finite subset of a high-dimensional integer lattice. The key contributions are as follows: DASSO exploits an innovative representation that decomposes the prior on the objective function into an additive form, reducing the problem dimensionality to facilitate computationally efficient posterior updates, but without losing a full-dimensional representation as other additive decompositions do. Furthermore, DASSO achieves rapid search progress by identifying the best-CEI solution to simulate on each iteration while avoiding the computational overhead of actually computing the CEI of all solutions. In brief, DASSO makes large-scale DSO computationally efficient without sacrificing what makes BO so effective in small-scale problems, facilitating the solution of far larger problems than can be attacked by any competing BO algorithm.

The remainder of this article is organized as follows: Section 2 reviews the literature on high-dimensional BO and DSO. Section 3 briefly defines GMRFs and shows how they can be used to solve DSO problems. Section 4 explains how FANOVA justifies the prior distribution representation that underlies the DASSO algorithm. Empirical performance evaluations are in Section 5. Finally, conclusions are provided in Section 6. All proofs and many derivations are in the e-companion.

2. Literature Review

In their review paper, Binois and Wycoff (2022) classify the solution tactics for high-dimensional BO into three categories based on the structural model assumptions made to deal with the curse of dimensionality. Worth noting is that much of this literature addresses machine-learning problems or deterministic computer experiments, rather than stochastic simulation.

Under the assumption that only a subset of the decision variables are active, the first approach is to reduce the problem dimension by removing the decision variables that have little or no impact on the objective function. By doing so, a lower-dimensional problem with only influential decision variables can be solved instead. Such a lower-dimensional problem can be obtained by initially performing variable selection, if no expert knowledge is available. However, variable selection itself is more computationally challenging as the dimension increases. Also, when all decision variables have similar influence on the objective function, no such reduction is possible. Even when it can be done, the reduced-dimension problem may still be relatively high dimensional depending on the actual number of influential decision variables.

Assuming the existence of a lower-dimensional active subspace, the second approach is to find a projection of the decision variables to obtain an active subspace, then apply BO. Since such an active subspace, if it exists, is unknown, the projection must be estimated. One can parameterize the projection scheme (e.g., a projection matrix) and estimate the parameters (see, e.g., Garnett et al. 2014, Tripathy et al. 2016), obtain the projection from a sensitivity analysis (see, e.g., Djolonga et al. 2013), or sample it randomly (see, e.g., Wang et al. 2016). Yet it is critical (especially, when such an active subspace does not exist) to figure out which feasible solution to simulate next in an optimization search because that decision involves projecting back to the full space, which is needed to parameterize the simulation. Such a projection back may be inefficient when a large part of the space is infeasible. Moreover, the computational savings from a projection-based approach relies heavily on the size of the active subspace’s dimension (Mathesen et al. 2019).

Under the assumption of an additive structure for the objective function, the third approach is to decompose the objective function into a sum of functions, each with fewer decision variables. The motivation is to perform the optimization component-wise, making the search much more efficient, and thus addressing the computational burden of high-dimensional BO. Several papers in this stream assume that the objective function is a sum of independent low-dimensional functions with disjoint (separable) decision variable dimensions (see, e.g., Kandasamy et al. 2015, Gardner et al. 2017, Wang et al. 2018), while others allow possibly overlapping dimensions (see, e.g., Hoang et al. 2018, Rolland et al. 2018). They differ mainly in how they learn the additive structure. They all model the lower-dimensional functions as realizations of independent GPs. This implies that the prior on the objective function values is a GP with an additive kernel. A potential downside of the additivity assumption is that it implies that the interactions among the subsets of decision variables (which we refer to as *groups*) are assumed to be negligible. Further, the covariance matrix of the additive GP may not be invertible due to linear relationships among the solutions, which is more likely to occur when the feasible solution space is discrete (Durrande et al. 2010).

Our DASSO methodology is superficially related to this last approach in terms of decomposing the objective function, but it does not enforce a separability approximation. Instead, it employs discrete FANOVA to include *all* high-order interactions among the groups by iteratively employing two steps: The “dice” step avoids the computational burden of the high-order interactions by forming a random-effect that combines

the high-order terms with one of the lower-dimensional groups, and then selects a partial solution based on its posterior. Next, DASSO obtains a best-CEI solution from the posterior of the partial solution “slice,” a posterior that has the same (low) dimension as the hold-out group and that fully accounts for all higher-order interactions. Details are in Section 4.

Although FANOVA has been used mostly as a basis for global sensitivity analysis, some studies (see, e.g., Muehlenstaedt et al. 2012, Ginsbourger et al. 2016, Ulaganathan et al. 2016) focus on high-dimensional GP modeling on a continuous domain and provide uncertainty measures over predictions or more computationally efficient GP models. Our aim is different in that we solve SO problems over a discrete domain.

To our knowledge, the only other inference-based algorithm that is designed for high-dimensional DSO problems is the projected Gaussian Markov improvement algorithm (pGMIA) of Li and Song (2024); see also Mes et al. (2011). pGMIA partitions the dimensions into two groups, called “region” and “solution” layers, and projects the latter onto the former. The region layer is represented as a graph, where the objective function value at a node is the average of objective function values at the feasible solutions projected onto that node. pGMIA alternates between choosing a node in the region-layer graph and deciding the next feasible solution to simulate from the solutions projected onto the chosen node. For both decisions, it models the objective function values as a GMRF and adopts CEI as the criterion to choose which feasible solution to simulate next. Li and Song (2024) also describe pGMIA+, a multi-layer extension of pGMIA, and both algorithms outperform the state-of-the-art high-dimensional BO algorithms to which they compare. Therefore, pGMIA represents the state of the art against which we compare in Section 5.

3. Discrete Simulation Optimization with GMRFs

Our problem to minimize $y(\mathbf{x}) = \mathbb{E}[Y(\mathbf{x})]$ subject to $\mathbf{x} \in \mathcal{X}$, where \mathcal{X} is a finite subset of the d -dimensional integer lattice treated as a collection of (column) vectors. The objective, $y(\mathbf{x})$, can only be estimated via stochastic simulation. In particular, the output $Y_j(\mathbf{x}) = y(\mathbf{x}) + \epsilon_j(\mathbf{x})$ can be observed for \mathbf{x} on replication $j = 1, 2, \dots$, where $\{\epsilon_j(\mathbf{x})\}$ are independent and identically distributed with mean zero and finite (unknown) variance $\sigma^2(\mathbf{x})$ that may depend on \mathbf{x} . In the multi-product inventory problem, \mathbf{x} denotes the vector of reorder and order-up-to values for each product, $Y(\mathbf{x})$ represents the average cost per period obtained from a single replication, and $y(\mathbf{x})$ is the true expected average cost per period. We define GMRFs and explain how they have been used to solve the DSO problem in Sections 3.1–3.2, deferring our approach to Section 4.

3.1. Gaussian Markov Random Field (GMRF)

A GMRF is a multivariate Gaussian random vector $\mathbb{Y} = [\mathbb{Y}_1, \mathbb{Y}_2, \dots, \mathbb{Y}_n]^\top$ associated with an undirected labeled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the set of n nodes and \mathcal{E} is the set of edges (Rue and Held 2005). In a DSO setting, each node represents a solution \mathbf{x}_i associated with the corresponding element of \mathbb{Y} . Moreover, we impose the structure that solutions that differ by ± 1 in one coordinate are connected by edges. Letting

$\boldsymbol{\mu}$ and \mathbf{Q} be the mean vector and precision matrix of \mathbb{Y} , respectively, the probability density function of the GMRF can be written as $f(\mathbf{y}|\boldsymbol{\mu}, \mathbf{Q}) = (2\pi)^{-n/2} |\mathbf{Q}|^{1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{Q}(\mathbf{y} - \boldsymbol{\mu})\right)$. The precision matrix is the inverse of the covariance matrix and is positive-definite. The choice of \mathcal{E} makes \mathbf{Q} have no more than $2d + 1$ nonzero elements. The diagonal elements of \mathbf{Q} are the conditional precisions: $Q_{ii} = \text{Prec}(\mathbb{Y}_i | \mathbb{Y}_{\mathcal{V} \setminus \{i\}})$, where \mathbb{Y}_S denotes the subvector of \mathbb{Y} including only the nodes in any $S \subset \mathcal{V}$. The off-diagonal entries are proportional to conditional correlations; specifically, $Q_{ij} = -\text{Corr}(\mathbb{Y}_i, \mathbb{Y}_j | \mathbb{Y}_{\mathcal{V} \setminus \{i,j\}}) \sqrt{Q_{ii}Q_{jj}}$.

A GMRF \mathbb{Y} is Markovian in the sense that it has the local Markov property: $\mathbb{Y}_i \perp \mathbb{Y}_{\mathcal{V} \setminus (\mathcal{N}(i) \cup \{i\})} | \mathbb{Y}_{\mathcal{N}(i)}$, where $\mathcal{N}(i) = \{j: \{i, j\} \in \mathcal{E}\}$ is the set of neighbors of node i in \mathcal{G} . This implies that the random variable \mathbb{Y}_i , conditional on the values of its neighbors, is independent of the values at non-neighboring nodes. Consequently, $Q_{ij} \neq 0$ if and only if $\{i, j\} \in \mathcal{E}$.

3.2. Gaussian Markov Improvement Algorithm (GMIA)

Since the vector of objective function values, denoted by $\mathbf{y} = [y(\mathbf{x}_1), y(\mathbf{x}_2), \dots, y(\mathbf{x}_n)]^\top$, is unknown, it can be modeled as a realization of the GMRF $\mathbb{Y} = [\mathbb{Y}(\mathbf{x}_1), \mathbb{Y}(\mathbf{x}_2), \dots, \mathbb{Y}(\mathbf{x}_n)]^\top \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}^{-1})$ with mean vector $\boldsymbol{\mu}$ of size $n \times 1$ and precision matrix \mathbf{Q} of size $n \times n$; this is the prior.

Suppose that only the “design points,” a subset of the feasible solutions, have been simulated. Let $\mathcal{D} \subseteq \mathcal{X}$ denote the current subset of design points, and partition \mathcal{X} into the two disjoint sets \mathcal{D} and $\mathcal{U} = \mathcal{X} \setminus \mathcal{D}$. Notice that \mathcal{D} is the set of feasible solutions that have been simulated, and \mathcal{U} is the set of feasible solutions that have not (i.e., are unsimulated). Using these disjoint sets, the vector \mathbb{Y} can be partitioned into $\mathbb{Y}_{\mathcal{U}}$ and $\mathbb{Y}_{\mathcal{D}}$, which are the subvectors of \mathbb{Y} including the solutions in \mathcal{U} and \mathcal{D} , respectively, that is, $\mathbb{Y} = \begin{pmatrix} \mathbb{Y}_{\mathcal{U}} \\ \mathbb{Y}_{\mathcal{D}} \end{pmatrix}$. Similarly, the vector $\boldsymbol{\mu}$ can be partitioned into $\boldsymbol{\mu}_{\mathcal{U}}$ and $\boldsymbol{\mu}_{\mathcal{D}}$.

The sample mean vector of simulated values at the design points, denoted by $\bar{\mathbf{Y}}_{\mathcal{D}}$, can be represented as a realization of $\mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbb{Y}_{\mathcal{D}} + \boldsymbol{\epsilon}$, with $\mathbb{Y}_{\mathcal{D}}$ and $\boldsymbol{\epsilon}$ independent, and $\boldsymbol{\epsilon} \sim \mathcal{N}(\bar{\mathbf{0}}_{|\mathcal{D}|}, \boldsymbol{\Sigma}^\epsilon)$, where $\bar{\mathbf{0}}_{|\mathcal{D}|}$ is a $|\mathcal{D}|$ -dimensional vector of zeros and $\boldsymbol{\Sigma}^\epsilon$ is the intrinsic covariance matrix of the stochastic noise inherent to $\bar{\mathbf{Y}}_{\mathcal{D}}$. When the design points are simulated independently, $\boldsymbol{\Sigma}^\epsilon$ is a diagonal matrix, whereas it is a dense matrix when the design points are simulated with common random numbers; due to the expense of inverting dense matrices we opt for independent simulations.

The conditional distribution of \mathbb{Y} given $\mathbb{Y}_{\mathcal{D}}^\epsilon = \bar{\mathbf{Y}}_{\mathcal{D}}$ is

$$\mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu}_{\mathcal{U}} \\ \boldsymbol{\mu}_{\mathcal{D}} \end{pmatrix} + \bar{\mathbf{Q}}^{-1} \begin{pmatrix} \bar{\mathbf{0}}_{|\mathcal{U}|} \\ [\boldsymbol{\Sigma}^\epsilon]^{-1} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \end{pmatrix}, \bar{\mathbf{Q}}^{-1}\right), \text{ where } \bar{\mathbf{Q}} = \mathbf{Q} + \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & [\boldsymbol{\Sigma}^\epsilon]^{-1} \end{pmatrix} \quad (1)$$

is the conditional precision matrix and $\mathbf{0}_{a \times b}$ is the $a \times b$ matrix of zeros. The unsimulated solutions are first in the GMRF representation because we assume the solutions have been so ordered.

Let $\tilde{\mathbf{x}}$ denote the current sample-best solution (among the design points in \mathcal{D}) based on sample means, i.e., $\tilde{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{D}} \bar{\mathbf{Y}}_{\mathcal{D}}(\mathbf{x})$, where $\bar{\mathbf{Y}}_{\mathcal{D}}(\mathbf{x})$ is the element of $\bar{\mathbf{Y}}_{\mathcal{D}}$ associated with \mathbf{x} . Then, the CEI of $\mathbf{x} \in \mathcal{X}$

relative to $\tilde{\mathbf{x}}$ is $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) = \mathbb{E}[\max\{\mathbb{Y}(\tilde{\mathbf{x}}) - \mathbb{Y}(\mathbf{x}), 0\} \mid \mathbb{Y}_{\mathcal{D}}^{\epsilon} = \bar{\mathbf{Y}}_{\mathcal{D}}]$. Let $m(\mathbf{x})$ and $v(\mathbf{x})$ denote the conditional mean and conditional variance of $\mathbb{Y}(\mathbf{x})$, respectively, and $c(\tilde{\mathbf{x}}, \mathbf{x})$ denote the conditional covariance between $\mathbb{Y}(\tilde{\mathbf{x}})$ and $\mathbb{Y}(\mathbf{x})$. Then the conditional variance of the difference $\mathbb{Y}(\tilde{\mathbf{x}}) - \mathbb{Y}(\mathbf{x})$ is $v(\tilde{\mathbf{x}}, \mathbf{x}) = v(\tilde{\mathbf{x}}) + v(\mathbf{x}) - 2c(\tilde{\mathbf{x}}, \mathbf{x})$. Further, the CEI of $\mathbf{x} \in \mathcal{X} \setminus \{\tilde{\mathbf{x}}\}$ can be expressed as

$$\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) = (m(\tilde{\mathbf{x}}) - m(\mathbf{x})) \Phi\left(\frac{m(\tilde{\mathbf{x}}) - m(\mathbf{x})}{\sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})}}\right) + \sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})} \phi\left(\frac{m(\tilde{\mathbf{x}}) - m(\mathbf{x})}{\sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})}}\right),$$

where ϕ and Φ are the standard normal probability density and cumulative distribution functions, respectively (Salemi et al. 2019). At each iteration, GMIA simulates the sample-best solution, $\tilde{\mathbf{x}}$, and the solution with the largest CEI, then updates the posterior distribution and continues. CEI is an enhancement of the “expected improvement” criterion for optimizing deterministic black-box functions (Jones et al. 1998) that has been shown to be particularly effective in stochastic simulation (see also Chen and Ryzhov 2019). All else being equal, $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$ is decreasing in $m(\mathbf{x})$ and increasing in $v(\tilde{\mathbf{x}}, \mathbf{x})$, which is the property DASSO exploits to significantly improve computational efficiency as described in Section 4.6.

3.3. Intermezzo

A DSO algorithm for large-scale, expensive simulation optimization needs to be both search-effective and computationally feasible. The search-effectiveness of BO depends on the choice of GP prior, particularly its covariance structure. Hidden in the GMIA summary above is a sparse parameterization for \mathbf{Q} that has been shown empirically to be exceptional for guiding the search in DSO problems (Salemi et al. 2019, Semelhago et al. 2017, 2021, Li and Song 2020, 2024) as well as computationally advantageous; see Section EC.5 in the e-companion. Clearly this formulation is something we do not want to lose.

However, even though the form of \mathbf{Q} in GMIA is sparse, evaluating the posterior distribution, and even calculating the CEIs for every feasible solution on every iteration, becomes computationally prohibitive if the number of feasible solutions is (say) in the trillions. Therefore, GMIA and its variants are constrained by the size of the problem, particularly its dimension. To overcome this limitation, we propose the DASSO algorithm that has a search-effective yet computationally feasible covariance structure while avoiding the need to calculate the CEI of every feasible solution on each iteration, all without approximations.

4. Dice and Slice Simulation Optimization (DASSO)

Recall that the aim is to minimize $y(\mathbf{x})$ subject to $\mathbf{x} \in \mathcal{X}$, where \mathcal{X} is a finite subset of the d -dimensional integer lattice, and the objective can only be estimated via stochastic simulation. The set of feasible solutions can be represented as $\mathcal{X} = \times_{i=1}^d \mathcal{X}_i$, where \mathcal{X}_i is a finite subset of the one-dimensional integer lattice with $n_i = |\mathcal{X}_i|$. The number of feasible solutions is $n = \prod_{i=1}^d n_i$ which increases exponentially in d ; thus, only a small fraction of feasible solutions can be simulated for high-dimensional problems with expensive simulations even if $|\mathcal{X}_i|$ is small. In the following subsections, we show how DASSO tackles this problem.

4.1. Discrete FANOVA

The following Discrete FANOVA (DFA) representation of the objective function justifies the prior distribution at the heart of DASSO: Suppose we assign probabilities to each feasible solution for the purpose of measuring the global sensitivity in terms of variance of the output response with respect to solutions' component dimensions. Specifically, assume that the component dimensions of a feasible solution are independent and follow discrete uniform distributions. That is, letting $\mathbb{X} = [\mathbb{X}_1, \mathbb{X}_2, \dots, \mathbb{X}_d]^\top$ denote a random vector, we have $P\{\mathbb{X} = \mathbf{x}\} = \prod_{i=1}^d P\{\mathbb{X}_i = x_i\}$, $\forall \mathbf{x} = [x_1, x_2, \dots, x_d]^\top \in \mathcal{X}$, and $P\{\mathbb{X}_i = x\} = 1/n_i$, $\forall x \in \mathcal{X}_i$, for each component dimension \mathbb{X}_i of \mathbb{X} , where x_i denotes the i th component of \mathbf{x} . Therefore,

$$P\{\mathbb{X} = \mathbf{x}\} = \prod_{i=1}^d P\{\mathbb{X}_i = x_i\} = \frac{1}{\prod_{i=1}^d n_i} = \frac{1}{n}, \quad \forall \mathbf{x} = [x_1, x_2, \dots, x_d]^\top \in \mathcal{X}. \quad (2)$$

Thus, all feasible solutions are equally likely, or equally weighted in a non-probabilistic sense.

For any subset \mathbf{u} of component indices $\mathbf{d} = \{1, 2, \dots, d\}$, let $\mathbf{x}_{\mathbf{u}} = (x_i)_{i \in \mathbf{u}}$ denote the lower dimensional component of \mathbf{x} associated with \mathbf{u} . The number and set of lower dimensional components associated with \mathbf{u} are $n_{\mathbf{u}} = \prod_{i \in \mathbf{u}} n_i$ and $\mathcal{X}_{\mathbf{u}} = \times_{i \in \mathbf{u}} \mathcal{X}_i$, respectively. Also, let $-\mathbf{u}$ denote the complementary set $\mathbf{d} \setminus \mathbf{u}$. In the DFA decomposition, we first define the grand mean to be $y_{\emptyset}(\mathbf{x}) = \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{X}} y(\mathbf{x}) \triangleq \bar{y}$, $\forall \mathbf{x} \in \mathcal{X}$. Notice that $y_{\emptyset}(\mathbf{x})$ is a constant that does not depend on any component dimensions of \mathbf{x} . Then, for $\mathbf{u} \subseteq \mathbf{d}$, the effect of $\mathbf{x}_{\mathbf{u}}$ can be expressed as $y_{\mathbf{u}}(\mathbf{x}) = \frac{1}{n_{-\mathbf{u}}} \sum_{\mathbf{x}_{-\mathbf{u}} \in \mathcal{X}_{-\mathbf{u}}} (y(\mathbf{x}) - \sum_{\mathbf{v} \subsetneq \mathbf{u}} y_{\mathbf{v}}(\mathbf{x}))$, $\forall \mathbf{x} \in \mathcal{X}$. Notice that $y_{\mathbf{u}}(\mathbf{x})$ depends on \mathbf{x} only through $\mathbf{x}_{\mathbf{u}}$ as the effect of $\mathbf{x}_{-\mathbf{u}}$ is averaged out. Using these effects, $y(\mathbf{x})$ can be represented as a DFA decomposition: $y(\mathbf{x}) = \sum_{\mathbf{u} \subseteq \mathbf{d}} y_{\mathbf{u}}(\mathbf{x})$. Under the probability measure in (2), we can show that these functions have certain properties.

PROPOSITION 1. *The following properties hold if \mathbb{X} follows Distribution (2): (a) $E[y(\mathbb{X})] = \bar{y}$, (b) $E[y_{\mathbf{u}}(\mathbb{X})] = 0$, $\forall \mathbf{u} \subseteq \mathbf{d}$ with $\mathbf{u} \neq \emptyset$, (c) $E[y_{\mathbf{u}}(\mathbb{X})y_{\mathbf{v}}(\mathbb{X})] = 0$, $\forall \mathbf{u}, \mathbf{v} \subseteq \mathbf{d}$ with $\mathbf{u} \neq \mathbf{v}$, (d) $\text{Var}[y(\mathbb{X})] = \sum_{\mathbf{u} \subseteq \mathbf{d}} \text{Var}[y_{\mathbf{u}}(\mathbb{X})]$.*

Proposition 1 states that (a) the overall mean of $y(\mathbb{X})$ is \bar{y} ; (b) other than $y_{\emptyset}(\mathbf{x})$, each $y_{\mathbf{u}}(\mathbb{X})$ in the DFA decomposition has zero mean; (c) $y_{\mathbf{u}}(\mathbb{X})$ and $y_{\mathbf{v}}(\mathbb{X})$ for $\mathbf{u} \neq \mathbf{v}$ are uncorrelated; and (d) the overall variance of $y(\mathbb{X})$ decomposes into the sum of the variances of each individual function. Proposition 1 holds for any $y(\cdot)$ regardless of its functional form, and it motivates the DASSO prior distribution on \mathbb{Y} .

4.2. Decomposition

Partition \mathbf{d} into $g \geq 1$ disjoint non-empty sets $\mathbf{g}^{(1)}, \mathbf{g}^{(2)}, \dots, \mathbf{g}^{(g)}$, i.e., $\bigcup_{\rho \in \mathcal{G}} \mathbf{g}^{(\rho)} = \mathbf{d}$ and $\mathbf{g}^{(\rho)} \cap \mathbf{g}^{(\varrho)} = \emptyset$ for $\rho \neq \varrho$, where $\mathcal{G} = \{1, 2, \dots, g\}$. These sets represent “groups” indexed by the superscripts in parentheses. For instance, in the multi-product inventory problem, groups can represent subsets of products. The indexing of the groups is updated throughout the algorithm. Later we will single out group g for special treatment.

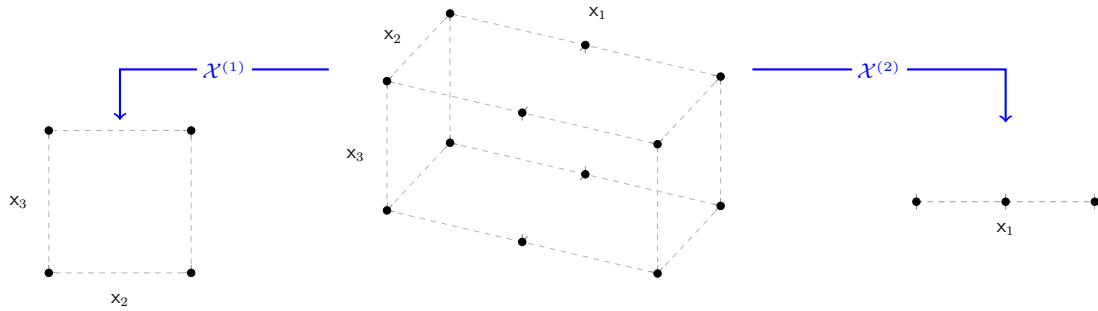


Figure 1 Illustration of a decomposition into two groups with $\mathbf{g}^{(1)} = \{2, 3\}$ and $\mathbf{g}^{(2)} = \{1\}$.

Throughout the paper, ρ and ϱ index groups. For each $\rho \in \mathcal{G}$, let $\mathcal{X}^{(\rho)}$ and $\mathbf{x}^{(\rho)}$ represent the corresponding components of \mathcal{X} and \mathbf{x} , respectively, indexed by $\mathbf{g}^{(\rho)}$, i.e., $\mathcal{X}^{(\rho)} = \times_{i \in \mathbf{g}^{(\rho)}} \mathcal{X}_i$ and $\mathbf{x}^{(\rho)} = (x_i)_{i \in \mathbf{g}^{(\rho)}}$. Notice that $\mathcal{X}^{(\rho)}$ is a finite subset of the $d^{(\rho)}$ -dimensional integer lattice with $n^{(\rho)}$ distinct points, where $d^{(\rho)} = |\mathbf{g}^{(\rho)}|$ and $n^{(\rho)} = |\mathcal{X}^{(\rho)}| = \prod_{i \in \mathbf{g}^{(\rho)}} n_i$. Further, notice that $\sum_{\rho \in \mathcal{G}} d^{(\rho)} = d$ and $n = \prod_{\rho \in \mathcal{G}} n^{(\rho)}$. Figure 1 illustrates a decomposition into two groups for a simple 3-dimensional example with $n = 12$ solutions.

For each group $\rho \in \mathcal{G}$, we define the *group function* $y^{(\rho)}(\mathbf{x}) = \sum_{\mathbf{u} \subseteq \mathbf{g}^{(\rho)}: \mathbf{u} \neq \emptyset} y_{\mathbf{u}}(\mathbf{x})$. That is, $y^{(\rho)}(\mathbf{x})$ is the summation of all functions in the DFA decomposition indexed by a subset of $\mathbf{g}^{(\rho)}$. Notice that $y^{(\rho)}(\mathbf{x})$ is a function of only $\mathbf{x}^{(\rho)}$, and thus can be rewritten as $y^{(\rho)}(\mathbf{x}^{(\rho)})$. Then, based on the DFA decomposition of $y(\mathbf{x})$, we can write $y(\mathbf{x}) = \beta_0 + \sum_{\rho \in \mathcal{G}} y^{(\rho)}(\mathbf{x}^{(\rho)}) + y^{(r)}(\mathbf{x})$, where $\beta_0 = \bar{y}$, and $y^{(r)}(\mathbf{x})$ represents the remainder. Namely, $y^{(r)}(\mathbf{x})$ is the summation of $y_{\mathbf{u}}(\mathbf{x})$ for $\mathbf{u} \subseteq \mathbf{d}$ such that $\mathbf{u} \neq \emptyset$ and $\mathbf{u} \not\subseteq \mathbf{g}^{(\rho)}$ for all $\rho \in \mathcal{G}$. Thus, $y^{(r)}(\mathbf{x})$ is a function of *all* component dimensions of \mathbf{x} and captures the interactions across the groups.

Let $\mathbf{y}^{(\rho)} = [y^{(\rho)}(\mathbf{x}_1^{(\rho)}), y^{(\rho)}(\mathbf{x}_2^{(\rho)}), \dots, y^{(\rho)}(\mathbf{x}_{n^{(\rho)}}^{(\rho)})]^\top$ be the vector of group ρ components of the objective function. Also, let $\mathbf{y}^{(r)}$ denote vector of remainder terms $[y^{(r)}(\mathbf{x}_1), y^{(r)}(\mathbf{x}_2), \dots, y^{(r)}(\mathbf{x}_n)]^\top$. Using the $\mathbf{y}^{(\rho)}$'s and $\mathbf{y}^{(r)}$, the vector of objective function values at all solutions in \mathcal{X} can be expressed as $\mathbf{y} = \beta_0 \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}} \mathbf{T}^{(\rho)} \mathbf{y}^{(\rho)} + \mathbf{y}^{(r)}$, where $\vec{\mathbf{1}}_n$ is an n -dimensional vector of ones and $\mathbf{T}^{(\rho)}$ is the transformation matrix associated with group ρ . The (k, l) th element of $\mathbf{T}^{(\rho)}$ is 1 if $\mathbf{x}_l^{(\rho)}$ is the corresponding lower dimensional component of solution \mathbf{x}_k , and 0 otherwise. Notice that $\mathbf{T}^{(\rho)}$ is an $n \times n^{(\rho)}$ matrix containing only one nonzero element with value 1 in each row. To assist optimization we impose a prior distribution on \mathbf{y} .

Starting with the group function $\mathbf{y}^{(\rho)}$, we model it as a realization of the GMRF

$$\mathbb{Y}^{(\rho)} = [\mathbb{Y}^{(\rho)}(\mathbf{x}_1^{(\rho)}), \mathbb{Y}^{(\rho)}(\mathbf{x}_2^{(\rho)}), \dots, \mathbb{Y}^{(\rho)}(\mathbf{x}_{n^{(\rho)}}^{(\rho)})]^\top \sim \mathcal{N}(\vec{\mathbf{0}}_{n^{(\rho)}}, [\mathbf{Q}^{(\rho)}]^{-1})$$

with precision matrix $\mathbf{Q}^{(\rho)}$ of size $n^{(\rho)} \times n^{(\rho)}$. The prior mean of $\mathbb{Y}^{(\rho)}$ is zero as a consequence of the DFA decomposition; see Property (b) of Proposition 1. Similarly, the remainder $\mathbf{y}^{(r)}$ can be modeled as a realization of the GMRF $\mathbb{Y}^{(r)} \sim \mathcal{N}(\vec{\mathbf{0}}_n, [\mathbf{Q}^{(r)}]^{-1})$ with precision matrix $\mathbf{Q}^{(r)}$ of size $n \times n$.

What sets our model apart from other additive approximations proposed for high-dimensional BO (cf. Kandasamy et al. 2015, Gardner et al. 2017, and Wang et al. 2018) is that we retain the term $\mathbb{Y}^{(r)}$, representing the inadequacy of a fully separable model. Therefore, we do not assume $y(\cdot)$ separates into the sum

of functions of groups, unlike other approaches. All of our GMRFs are mutually independent by Property (c) of Proposition 1 because uncorrelated implies independence for Gaussian processes.

In summary, as justified by the DFA decomposition, the vector \mathbf{y} of objective function values is modeled as the realization of a constant plus a linear combination of mean-zero, independent GMRFs

$$\mathbb{Y} = \beta_0 \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}} \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} + \mathbb{Y}^{(r)}. \quad (3)$$

The goal of defining the prior distribution on the DFA decomposition is to achieve significant computational benefits when obtaining the posterior of \mathbb{Y} via the posteriors of the g lower-dimensional GMRFs. However, (3) does not reduce computation *unless* $\mathbf{Q}^{(r)}$ is *diagonal*, since it is full-dimensional ($n \times n$); diagonal implies no spatial correlation representing interaction. This is the first issue we address.

A second issue, whether $\mathbf{Q}^{(r)}$ is diagonal or not, is that if the values of its diagonal elements are very large, which may occur when the contribution of $y^{(r)}(\mathbf{x})$ to the objective function value $y(\mathbf{x})$ is negligible, then $y(\mathbf{x})$ becomes closer to a purely separable function. Purely separable \mathbb{Y} (without $\mathbb{Y}^{(r)}$) has linearly dependent elements due to the repetitive rows in the $\mathbf{T}^{(\rho)}$'s, leading to a singular covariance matrix for \mathbb{Y} . The singular covariance matrix may cause computational issues when updating the posterior distribution of \mathbb{Y} . More importantly, the linear dependence of groups causes inconsistencies between \mathbf{y} and \mathbb{Y} unless the former has the same additive structure as the latter. See Section EC.4 in the e-companion for a small example illustrating the linear dependence issue.

One of our key contributions is addressing these two issues in a way that preserves computational efficiency without losing a full-dimensional representation. To do this we divide and conquer.

Our first innovation is to obtain the posterior distributions of groups $1, 2, \dots, g-1$ by modeling group g plus the remainder term as a simple *random effect*, leading to a computationally tractable posterior distribution for \mathbb{Y} . Specifically, we replace $\mathbf{T}^{(g)} \mathbb{Y}^{(g)} + \mathbb{Y}^{(r)}$ with a GMRF with prior distribution $\mathbb{W}^{(g)} \sim \mathcal{N}(\vec{\mathbf{0}}_n, \sigma_g^2 \mathbf{I}_n)$ with a positive but finite σ_g^2 , where \mathbf{I}_n is the $n \times n$ identity matrix. Therefore, the GMRF prior for \mathbf{y} becomes

$$\mathbb{Y} = \beta_0 \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} + \mathbb{W}^{(g)}, \quad (4)$$

where $\mathcal{G}^{(-g)} = \mathcal{G} \setminus \{g\}$ eases notation. The identity of “group g ” may change in each iteration of DASSO so that no group is restricted to be only a part of $\mathbb{W}^{(g)}$.

Both (3) and (4) model the same problem but with different prior distributions; (4) provides weaker inference on the component dimensions in group g , but is computationally tractable and free of the inconsistency issue. Moreover, updating the identity of g lets (4) learn all groups’ effects equitably.

Under (4), the distribution of \mathbb{Y} is a convolution of low-dimensional distributions and an easy-to-evaluate full-dimensional distribution. Moreover, for fixed \mathcal{G} , \mathbb{Y} depends only on lower dimensional components $\mathbf{x}^{(-g)} = (x_i)_{i \in \mathcal{d} \setminus \mathbf{g}^{(g)}}$, which we refer to as the “first $g-1$ components.” Therefore, it infers the same posterior

for any feasible solution with the same first $g - 1$ components. Utilizing (4), DASSO performs a “dice” stage to determine the first $g - 1$ components of the solution to simulate by maximizing the CEI defined over $\mathcal{X}^{(-g)} = \times_{i \in \mathcal{d} \setminus \mathbf{g}^{(g)}} \mathcal{X}_i$.

Our second innovation is to further model the component dimensions in group g by a $n^{(g)}$ -dimensional GMRF for all solutions with the same first $g - 1$ components, as selected in the dice stage. DASSO then performs a “slice” stage to find the components in group g to simulate using the single-group GMRF.

In summary, DASSO alternates between using a simple random-effect GMRF to represent the complex remainder term when choosing values for the $g - 1$ groups in the dice stage, and a detailed GMRF of the lower-dimensional slice of the remaining group with those values fixed. As we show in Section 5 this leads to rapid search progress with low computational overhead.

Remark: We have argued that the prior in (4) is ideal for DASSO, and the form of the prior leads to a computationally feasible posterior computation, as shown below. However, to be search-effective the prior must have appropriate parameters. The DASSO prior has $2 + \sum_{\rho=1}^g (d^{(\rho)} + 1)$ parameters in total, $d^{(\rho)} + 1$ for each non-remainder group, and two for the overall mean and the random-effect variance. In Section EC.5 in the e-companion, we describe how we estimate these parameters via maximum likelihood so that they are consistent with the decomposition in (4). In brief, we apply maximum likelihood estimation to the *differences* in the simulation outputs at carefully chosen *pairs* of design points so that through cancellation the difference of a pair reflects only the effect of a single group in (4). This accounts for all but two of the parameters, whose estimation is straightforward.

4.3. Bayesian Inference for Dice Stage

In this section we derive the posterior distribution for the dice stage. Recall that $\mathcal{D} \subseteq \mathcal{X}$ is the set of feasible solutions that have been simulated, and $\mathcal{U} = \mathcal{X} \setminus \mathcal{D}$ is the set of feasible solutions that have not. Using these disjoint sets, the random-effect vector $\mathbb{W}^{(g)}$ can be partitioned into $\mathbb{W}_{\mathcal{U}}^{(g)}$ and $\mathbb{W}_{\mathcal{D}}^{(g)}$, which are the subvectors of $\mathbb{W}^{(g)}$ including the solutions in \mathcal{U} and \mathcal{D} , respectively.

For each group $\rho \in \mathcal{G}$, let $\mathcal{D}^{(\rho)} \subseteq \mathcal{X}^{(\rho)}$ denote the set of lower-dimensional components of the solutions in \mathcal{D} corresponding to group ρ , and partition $\mathcal{X}^{(\rho)}$ into the two disjoint sets $\mathcal{D}^{(\rho)}$ and $\mathcal{U}^{(\rho)} = \mathcal{X}^{(\rho)} \setminus \mathcal{D}^{(\rho)}$. Assuming that the elements of $\mathbb{Y}^{(\rho)}$ are reordered, we can partition $\mathbb{Y}^{(\rho)}$ and precision matrix $\mathbf{Q}^{(\rho)}$ as

$$\mathbb{Y}^{(\rho)} = \begin{pmatrix} \mathbb{Y}_{\mathcal{U}}^{(\rho)} \\ \mathbb{Y}_{\mathcal{D}}^{(\rho)} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \vec{0}_{|\mathcal{U}^{(\rho)}|} \\ \vec{0}_{|\mathcal{D}^{(\rho)}|} \end{pmatrix}, \begin{pmatrix} \mathbf{Q}_{\mathcal{UU}}^{(\rho)} & \mathbf{Q}_{\mathcal{UD}}^{(\rho)} \\ \mathbf{Q}_{\mathcal{DU}}^{(\rho)} & \mathbf{Q}_{\mathcal{DD}}^{(\rho)} \end{pmatrix}^{-1} \right),$$

where $\mathbb{Y}_{\mathcal{U}}^{(\rho)}$ and $\mathbb{Y}_{\mathcal{D}}^{(\rho)}$ are the subvectors of $\mathbb{Y}^{(\rho)}$ including the points in $\mathcal{D}^{(\rho)}$ and $\mathcal{U}^{(\rho)}$, respectively. The covariance matrix of $\mathbb{Y}_{\mathcal{D}}^{(\rho)}$, $\Sigma_{\mathcal{DD}}^{(\rho)} = \left[\mathbf{Q}_{\mathcal{DD}}^{(\rho)} - \mathbf{Q}_{\mathcal{DU}}^{(\rho)} [\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1} \mathbf{Q}_{\mathcal{UD}}^{(\rho)} \right]^{-1}$, can be derived from the block matrix inversion formula. Similarly, transformation matrix $\mathbf{T}^{(\rho)}$ can also be partitioned into

$$\begin{pmatrix} \mathbf{T}_{\mathcal{UU}}^{(\rho)} & \mathbf{T}_{\mathcal{UD}}^{(\rho)} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}^{(\rho)}|} & \mathbf{T}_{\mathcal{DD}}^{(\rho)} \end{pmatrix}.$$

The lower-left block is a zero matrix because $\mathbf{x}^{(\rho)} \in \mathcal{D}^{(\rho)}$ if $\mathbf{x} \in \mathcal{D}$ from the definition of $\mathcal{D}^{(\rho)}$. Then, $\mathbb{Y}_{\mathcal{D}} = \beta_0 \vec{\mathbf{1}}_{|\mathcal{D}|} + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^{(\rho)} + \mathbb{W}_{\mathcal{D}}^{(g)}$. The mean vector and covariance matrix of $\mathbb{Y}_{\mathcal{D}}$ are $\boldsymbol{\mu}_{\mathcal{D}} = \beta_0 \vec{\mathbf{1}}_{|\mathcal{D}|}$ and $\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} = \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top + \sigma_g^2 \mathbf{I}_{|\mathcal{D}|}$, respectively.

Theorems 1 and 2 state the posterior distributions of $\mathbb{Y}^{(\rho)}$ for $\rho \in \mathcal{G}^{(-g)}$ and $\mathbb{W}^{(g)}$, respectively.

THEOREM 1. *For $\rho \in \mathcal{G}^{(-g)}$, the conditional distribution of $\mathbb{Y}^{(\rho)}$ given $\mathbb{Y}_{\mathcal{D}}^\epsilon = \bar{\mathbf{Y}}_{\mathcal{D}}$ is normal with mean vector $\mathbf{m}^{(\rho)}$ and covariance matrix $\bar{\boldsymbol{\Sigma}}^{(\rho)} = [\bar{\mathbf{Q}}^{(\rho)}]^{-1}$, where*

$$\mathbf{m}^{(\rho)} = [\bar{\mathbf{Q}}^{(\rho)}]^{-1} \begin{pmatrix} \vec{\mathbf{0}}_{|\mathcal{U}^{(\rho)}|} \\ [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \mathbf{E}^{(\rho)} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \end{pmatrix}, \quad \bar{\mathbf{Q}}^{(\rho)} = \mathbf{Q}^{(\rho)} + \begin{pmatrix} \mathbf{0}_{|\mathcal{U}^{(\rho)}| \times |\mathcal{U}^{(\rho)}|} & \mathbf{0}_{|\mathcal{U}^{(\rho)}| \times |\mathcal{D}^{(\rho)}|} \\ \mathbf{0}_{|\mathcal{D}^{(\rho)}| \times |\mathcal{U}^{(\rho)}|} & [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \mathbf{E}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \end{pmatrix},$$

$$\text{and } \mathbf{E}^{(\rho)} = [\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} - \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top + \boldsymbol{\Sigma}^\epsilon]^{-1}.$$

THEOREM 2. *The conditional distribution of $\mathbb{W}^{(g)}$ given $\mathbb{Y}_{\mathcal{D}}^\epsilon = \bar{\mathbf{Y}}_{\mathcal{D}}$ is normal with mean vector $\mathbf{m}^{(g)}$ and covariance matrix $\bar{\boldsymbol{\Sigma}}^{(g)} = [\bar{\mathbf{Q}}^{(g)}]^{-1}$, where*

$$\mathbf{m}^{(g)} = [\bar{\mathbf{Q}}^{(g)}]^{-1} \begin{pmatrix} \vec{\mathbf{0}}_{|\mathcal{U}|} \\ \mathbf{E}^{(g)} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \end{pmatrix}, \quad \bar{\mathbf{Q}}^{(g)} = \frac{1}{\sigma_g^2} \mathbf{I}_n + \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & \mathbf{E}^{(g)} \end{pmatrix},$$

$$\text{and } \mathbf{E}^{(g)} = [\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} - \sigma_g^2 \mathbf{I}_{|\mathcal{D}|} + \boldsymbol{\Sigma}^\epsilon]^{-1}.$$

Since $\bar{\boldsymbol{\Sigma}}^{(g)}$ is the inverse of a block-diagonal matrix, it has a block-diagonal structure. Although $\bar{\boldsymbol{\Sigma}}^{(g)}$ is an $n \times n$ matrix, its block-diagonal structure enables us to compute it by inverting only a $|\mathcal{D}| \times |\mathcal{D}|$ matrix. Corollary 1 below follows immediately from the fact that $\mathbb{Y} = \beta_0 \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} + \mathbb{W}^{(g)}$.

COROLLARY 1. *The conditional distribution of \mathbb{Y} given $\mathbb{Y}_{\mathcal{D}}^\epsilon = \bar{\mathbf{Y}}_{\mathcal{D}}$ is normal with the conditional mean vector and conditional covariance matrix*

$$\mathbf{m} = \beta_0 \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}^{(\rho)} \mathbf{m}^{(\rho)} + \mathbf{m}^{(g)} \text{ and } \bar{\boldsymbol{\Sigma}} = \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{T}^{(\rho)} \bar{\boldsymbol{\Sigma}}^{(\rho)} [\mathbf{T}^{(\rho)}]^\top + \bar{\boldsymbol{\Sigma}}^{(g)}.$$

Let $\mathbf{v} = \text{diag}(\bar{\boldsymbol{\Sigma}})$ and $\mathbf{c} = \bar{\boldsymbol{\Sigma}} \mathbf{e}_{\tilde{\mathbf{x}}}$, where $\text{diag}(\cdot)$ represents the diagonal of the corresponding matrix as a vector, and $\mathbf{e}_{\tilde{\mathbf{x}}}$ is the n -dimensional basis vector consisting of 0s and a single 1 in the position corresponding to $\tilde{\mathbf{x}}$. Then \mathbf{v} is the vector of conditional variances of \mathbb{Y} , and \mathbf{c} is the vector of conditional covariances between $\mathbb{Y}(\tilde{\mathbf{x}})$ and all components of \mathbb{Y} . For any \mathbf{x} , conditional mean $m(\mathbf{x})$, conditional variance $v(\mathbf{x})$ and conditional covariance $c(\tilde{\mathbf{x}}, \mathbf{x})$ are the respective components of \mathbf{m} , \mathbf{v} and \mathbf{c} associated with \mathbf{x} .

For group $\rho \in \mathcal{G}^{(-g)}$, let $m^{(\rho)}(\mathbf{x}^{(\rho)})$, $v^{(\rho)}(\mathbf{x}^{(\rho)})$ and $c^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \mathbf{x}^{(\rho)})$ be, respectively, the components of $\mathbf{m}^{(\rho)}$, $\mathbf{v}^{(\rho)}$ and $\mathbf{c}^{(\rho)}$, associated with lower dimensional component $\mathbf{x}^{(\rho)}$, where $\mathbf{v}^{(\rho)} = \text{diag}(\bar{\boldsymbol{\Sigma}}^{(\rho)})$ and $\mathbf{c}^{(\rho)} = \bar{\boldsymbol{\Sigma}}^{(\rho)} \mathbf{e}_{\tilde{\mathbf{x}}^{(\rho)}}$. Also, let $m^{(g)}(\mathbf{x})$, $v^{(g)}(\mathbf{x})$ and $c^{(g)}(\tilde{\mathbf{x}}, \mathbf{x})$ be, respectively, the components of $\mathbf{m}^{(g)}$, $\mathbf{v}^{(g)}$ and $\mathbf{c}^{(g)}$, associated with solution \mathbf{x} , where $\mathbf{v}^{(g)} = \text{diag}(\bar{\boldsymbol{\Sigma}}^{(g)})$ and $\mathbf{c}^{(g)} = \bar{\boldsymbol{\Sigma}}^{(g)} \mathbf{e}_{\tilde{\mathbf{x}}}$. Due to the special structure of $\mathbf{T}^{(\rho)}$, Corollary 1 implies that

$$m(\mathbf{x}) = \sum_{\rho \in \mathcal{G}^{(-g)}} m^{(\rho)}(\mathbf{x}^{(\rho)}) + m^{(g)}(\mathbf{x}), \quad v(\mathbf{x}) = \sum_{\rho \in \mathcal{G}^{(-g)}} v^{(\rho)}(\mathbf{x}^{(\rho)}) + v^{(g)}(\mathbf{x})$$

and

$$c(\tilde{\mathbf{x}}, \mathbf{x}) = \sum_{\rho \in \mathcal{G}^{(-g)}} c^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \mathbf{x}^{(\rho)}) + c^{(g)}(\tilde{\mathbf{x}}, \mathbf{x}).$$

Observe that $m(\mathbf{x})$, $v(\mathbf{x})$, and $c(\tilde{\mathbf{x}}, \mathbf{x})$ are obtained from terms extracted from low-dimensional posterior distributions and an easy-to-evaluate full-dimensional posterior distribution. This is as opposed to a *computationally impossible-to-evaluate* full-dimensional distribution. And because we can choose any way to group the component decision variables we can effectively manage the burden of the posterior calculations. These computational savings enable DASSO to deal with high-dimensional DSO problems.

From the terms $m(\mathbf{x})$, $v(\mathbf{x})$, and $c(\tilde{\mathbf{x}}, \mathbf{x})$, we can compute the CEI of \mathbf{x} . Recall that BO algorithms find good solutions rapidly because they are guided by evaluating—in our case simulating—the solution offering the most potential improvement based on inference from the posterior distribution, but this is a bottleneck if there are (say) trillions of solutions whose improvement must be assessed on each iteration. In Section 4.6, we make a significant computational improvement for the CEI calculation by exploiting the fact that $m(\mathbf{x})$, $v(\mathbf{x})$, and $c(\tilde{\mathbf{x}}, \mathbf{x})$ are written in additive forms, and thus we can easily determine if there is another solution with a larger CEI than a group of solutions without fully computing their CEIs. This provides a substantial computational saving because it suffices to calculate the CEIs for only a relatively small number of feasible solutions to find a CEI-maximizing solution; see Section 5.3.

Expressing $m(\mathbf{x})$, $v(\mathbf{x})$, and $c(\tilde{\mathbf{x}}, \mathbf{x})$ in additive forms also enables us to show that the CEIs of the feasible solutions in \mathcal{U} with the same first $g - 1$ components are identical; see Proposition 2 in Section 4.6. Therefore, there might be multiple solutions with the largest CEI. This is mainly a consequence of the random-effect group having a diagonal covariance matrix.

We think of the posteriors of the first $g - 1$ groups as individual *dice* because the overall posterior includes their summation. Once we “roll the dice” (select their CEI maximizing components), we can evaluate a *slice* of the full-dimensional posterior that we address in the next section.

4.4. Bayesian Inference for Slice Stage

Suppose that values of $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(g-1)}$, the first $g - 1$ components, are chosen in the dice stage. For simplicity, let \mathbf{z} denote these chosen components; notice that \mathbf{z} is not full-dimensional, i.e., its dimension is $d - d^{(g)}$. Then, the remaining *slice* of \mathbb{Y} , denoted by $\mathbb{Y}(\cdot|\mathbf{z})$, is only a function of the last component, $\mathbf{x}^{(g)}$, and thus lower dimensional. In the inventory problem this corresponds to fixing the reorder and order-up-to values of some products with the remaining yet to be chosen. Thus, $\mathbb{Y}(\cdot|\mathbf{z}) = \mathbf{S}_{\mathbf{z}}\mathbb{Y}$ is a random vector of size $n^{(g)} \times 1$, where $\mathbf{S}_{\mathbf{z}}$ is the transformation matrix to fix $\mathbf{x}^{(-g)}$ to \mathbf{z} . We refer to the distribution of $\mathbb{Y}(\cdot|\mathbf{z})$ as a “slice” rather than a “conditional” because it is a subvector of \mathbb{Y} and not conditioned on \mathbf{z} . The (k, l) th element of $\mathbf{S}_{\mathbf{z}}$ is 1 if solution \mathbf{x}_l is composed of $\mathbf{x}^{(-g)} = \mathbf{z}$ and $\mathbf{x}_k^{(g)}$, and 0 otherwise. Notice that $\mathbf{S}_{\mathbf{z}}$ is a $n^{(g)} \times n$ matrix containing one nonzero element with value 1 in each row. Under (3), $\mathbb{Y}(\cdot|\mathbf{z}) =$

$\beta_0 \mathbf{S}_z \vec{\mathbf{1}}_n + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbf{S}_z \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} + \mathbf{S}_z \mathbf{T}^{(g)} \mathbb{Y}^{(g)} + \mathbf{S}_z \mathbb{Y}^{(r)}$. The special structures of the transformation matrices \mathbf{S}_z and $\mathbf{T}^{(\rho)}$ imply that $\mathbf{S}_z \vec{\mathbf{1}}_n = \vec{\mathbf{1}}_{n(g)}$, $\mathbf{S}_z \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} = \mathbb{Y}^{(\rho)}(\mathbf{z}^{(\rho)}) \vec{\mathbf{1}}_{n(g)}$ for $\rho \in \mathcal{G}^{(-g)}$, and $\mathbf{S}_z \mathbf{T}^{(g)} = \mathbf{I}_{n(g)}$, where $\mathbf{z}^{(\rho)}$ represents the group ρ component of \mathbf{z} . Therefore,

$$\mathbb{Y}(\cdot|\mathbf{z}) = \left(\beta_0 + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbb{Y}^{(\rho)}(\mathbf{z}^{(\rho)}) \right) \vec{\mathbf{1}}_{n(g)} + \mathbb{Y}^{(g)} + \mathbf{S}_z \mathbb{Y}^{(r)}.$$

Ignoring the remainder part, $\mathbf{S}_z \mathbb{Y}^{(r)}$, for simplicity, we impose a GMRF prior with mean vector $\beta_z \vec{\mathbf{1}}_{n(g)}$ and precision matrix $\mathbf{Q}^{(g)}$ on $\mathbb{Y}(\cdot|\mathbf{z})$, where $\beta_z = \beta_0 + \sum_{\rho \in \mathcal{G}^{(-g)}} \mathbb{Y}^{(\rho)}(\mathbf{z}^{(\rho)})$. While $\mathbb{Y}(\cdot|\mathbf{z})$ could be modeled with any appropriate precision matrix, we use $\mathbf{Q}^{(g)}$ because of its convenience and simplicity; recall that the identity of g is updated over time and the parameters of $\mathbf{Q}^{(\rho)}$ are initially estimated for each group $\rho \in \mathcal{G}$.

Let $\mathcal{X}_z = \{\mathbf{x} \in \mathcal{X} : \mathbf{x}^{(-g)} = \mathbf{z}\}$ denote the set of feasible solutions whose first $g-1$ components are \mathbf{z} ; notice that $|\mathcal{X}_z| = n^{(g)}$. Also, let $\mathcal{D}_z = \mathcal{X}_z \cap \mathcal{D}$ denote the set of design points (i.e., feasible solutions that have been simulated) whose first $g-1$ components are \mathbf{z} . Partitioning $\mathbb{Y}(\cdot|\mathbf{z})$ into two subvectors accordingly, the conditional distribution of $\mathbb{Y}(\cdot|\mathbf{z})$ given $\mathbb{Y}_{\mathcal{D}}^\epsilon = \bar{\mathbf{Y}}_{\mathcal{D}}$ can be obtained from (1).

4.5. The DASSO Algorithm

DASSO is presented in Algorithm 1. At a high level, DASSO employs CEI to choose the values of the first $g-1$ components of \mathbf{x} in the dice stage, selects the remaining components of group g in the slice stage, simulates the selected solutions and then repeats. A more detailed description follows.

Algorithm 1 first initializes the design set \mathcal{D} by choosing a subset of feasible solutions. After simulating the solutions in \mathcal{D} , the GMRF parameters of the prior are estimated. DASSO then starts the dice stage: the identity of the last group is chosen and the parameter β_0 is re-estimated accordingly by using the simulation outputs from the solutions in \mathcal{D} . Re-estimation of β_0 is not necessary, but doing so helps the inference since the estimated value depends on the identity of the last group. Unlike the other GMRF parameters, estimation of the constant term is straightforward; see Section EC.5 in the e-companion. DASSO next finds a solution with the largest CEI, $\hat{\mathbf{x}} \in \arg \max_{\mathbf{x} \in \mathcal{X} \setminus \{\tilde{\mathbf{x}}\}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$, and fixes the first $g-1$ components to $\hat{\mathbf{x}}^{(-g)}$, i.e., $\mathbf{z} = \hat{\mathbf{x}}^{(-g)}$. The set of candidate solutions is restricted to $\mathcal{X}_z = \{\mathbf{x} \in \mathcal{X} : \mathbf{x}^{(-g)} = \mathbf{z}\}$ to form a slice. Finally, the current sample-best solution, $\tilde{\mathbf{x}}$, is simulated. Figure 2 illustrates the dice stage for a simple 3-dimensional example with $n = 64$ solutions.

In the slice stage, the parameter β_z is estimated by using the simulation outputs of the solutions in \mathcal{D}_z . If $\mathcal{D}_z = \mathcal{X}_z \cap \mathcal{D} = \emptyset$, i.e., if none of the solutions in \mathcal{D}_z has been simulated before, then a set of solutions in \mathcal{X}_z is chosen to simulate before the estimation. Until the stopping criterion of the slice stage is satisfied, the algorithm simulates $\arg \max_{\mathbf{x} \in \mathcal{X}_z} \text{CEI}(\tilde{\mathbf{x}}|\mathbf{z}, \mathbf{x})$ and $\tilde{\mathbf{x}}|\mathbf{z} = \arg \min_{\mathbf{x} \in \mathcal{D}_z} \bar{\mathbf{Y}}_{\mathcal{D}}(\mathbf{x})$, and updates the slice posterior distribution. Notice that the slice stage computes the CEI of each solution in \mathcal{D}_z with respect to the sample best within \mathcal{D}_z , but this is a small set. Once the stopping criterion of the slice stage is satisfied, the algorithm returns to the dice stage unless its stopping criterion is also satisfied.

Algorithm 1 DASSO

-
- 1: Initialize \mathcal{D} by choosing a subset of feasible solutions. Simulate the solutions in \mathcal{D} to estimate the parameters (i.e., $\mathbf{Q}^{(\rho)}$ and σ_ρ^2 for $\rho \in \mathcal{G}$, and β_0).
 - 2: **while** stopping condition of the *dice stage* not satisfied **do**
 - 3: Choose a group to be the last group g .
 - 4: Re-estimate β_0 by using the simulation outputs of the solutions in \mathcal{D} .
 - 5: Using the posterior distribution of \mathbb{Y} , calculate the CEIs of the solutions to find $\hat{\mathbf{x}}$ with the largest CEI, i.e., $\hat{\mathbf{x}} \in \arg \max_{\mathbf{x} \in \mathcal{X} \setminus \{\hat{\mathbf{x}}\}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$.
 - 6: Restrict the set of candidate solutions to $\mathcal{X}_z = \{\mathbf{x} \in \mathcal{X} : \mathbf{x}^{(-g)} = \mathbf{z}\}$, where $\mathbf{z} = \hat{\mathbf{x}}^{(-g)}$.
 - 7: Simulate the current sample-best solution $\tilde{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{D}} \bar{Y}_{\mathcal{D}}(\mathbf{x})$.
 - 8: Simulate some solutions in \mathcal{X}_z if $\mathcal{D}_z = \mathcal{X}_z \cap \mathcal{D} = \emptyset$ so that β_z can be estimated.
 - 9: **while** stopping condition of the *slice stage* not satisfied **do**
 - 10: For the solutions in \mathcal{X}_z , let $\mathbb{Y}(\cdot | \mathbf{z}) \sim \mathcal{N}(\beta_z \vec{\mathbf{1}}_{n(g)}, [\mathbf{Q}^{(g)}]^{-1})$ be a GMRF.
 - 11: Estimate β_z by using the simulation outputs of the solutions in \mathcal{D}_z .
 - 12: Let $\tilde{\mathbf{x}} | \mathbf{z} = \arg \min_{\mathbf{x} \in \mathcal{D}_z} \bar{Y}_{\mathcal{D}}(\mathbf{x})$ be the sample-best solution in \mathcal{D}_z .
 - 13: Using the posterior distribution of $\mathbb{Y}(\cdot | \mathbf{z})$, calculate the CEIs of the solutions in \mathcal{X}_z to find the solution with the largest CEI, i.e., $\arg \max_{\mathbf{x} \in \mathcal{X}_z} \text{CEI}(\tilde{\mathbf{x}} | \mathbf{z}, \mathbf{x})$.
 - 14: Simulate $\arg \max_{\mathbf{x} \in \mathcal{X}_z} \text{CEI}(\tilde{\mathbf{x}} | \mathbf{z}, \mathbf{x})$ and $\tilde{\mathbf{x}} | \mathbf{z}$.
 - 15: **end while**
 - 16: **end while**
-

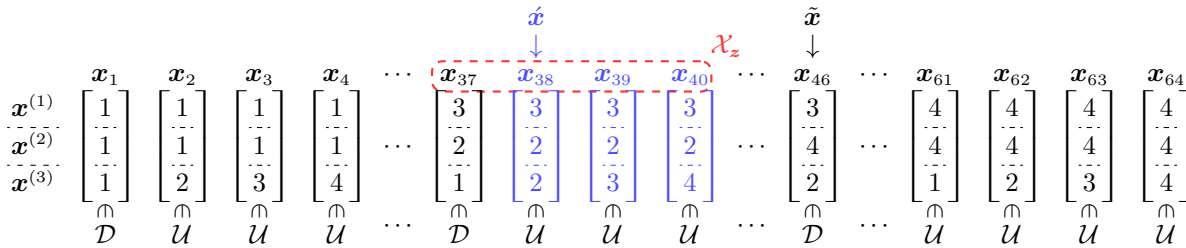


Figure 2 Illustration of the dice stage on a solution space $\mathcal{X} = \{1, 2, 3, 4\}^3$ with three 1-dimensional groups. Letting $g = 3$, suppose that solutions \mathbf{x}_{38} , \mathbf{x}_{39} and \mathbf{x}_{40} (in blue) have the largest CEI in the dice stage in Step 5. Then, the restricted set of candidate solutions is $\mathcal{X}_z = \{\mathbf{x}_{37}, \mathbf{x}_{38}, \mathbf{x}_{39}, \mathbf{x}_{40}\}$ (indicated with red dashed line).

DASSO requires user inputs for the decomposition and some criteria in Steps 1, 2, 3, 8, and 9 as well as an acquisition function such as CEI to guide the search. Below, we discuss the choices adopted in our numerical experiments in Section 5. First, notice that the user-defined decomposition (i.e., the groups) remains the same throughout the algorithm. Although DASSO could allow the decomposition to be updated, the GMRF parameters would need to be estimated for each new decomposition. To estimate the parameters for the user-defined decomposition in Step 1, we employ a parameter estimation method described in Section EC.5

of the e-companion. This method simulates specially constructed pairs of design points. Thus, to estimate the parameters for a new decomposition, additional design points that do not directly guide the search would need to be simulated to complete the pairs, wasting simulation effort. Therefore, we do not update the decomposition in our numerical experiments. Nevertheless, we empirically evaluate the performance of DASSO with different fixed decompositions.

Algorithm 1 contains two stopping conditions, one for the dice stage in Step 2 and another for the slice stage in Step 9. For the former, we adopt a fixed computation budget in terms of total number of replications, i.e., the algorithm terminates after a certain number of replications. For the latter, we similarly adopt a fixed budget in a sense that the number of slice-stage iterations per dice stage is one; in other words, Steps 9 and 15 are removed from the algorithm. We use fixed-budget stopping conditions for the sake of fair comparison in our numerical experiments. An alternative is using CEI for stopping, as proposed in Salemi et al. (2019), where the algorithm terminates once the largest CEI value drops below a specified threshold.

To update the identity of the group g in Step 3, we can impose a discrete distribution with support \mathcal{G} to randomly choose it. Alternatively this discrete distribution could adapt to some metrics measuring the impact of groups on the objective function values and changed from time to time during the search. For the sake of simplicity in our numerical experiments, we impose a discrete uniform distribution; that is, we choose group g randomly at each change. While intuitively it seems clear that each group should take turns as group g , the best way to assign it is an open question. Similarly, to choose some solutions from \mathcal{X}_z to simulate when $\mathcal{D}_z = \emptyset$ in Step 8, a discrete distribution with support $\mathcal{X}^{(g)}$ can be applied; we again adopt a discrete uniform distribution in our numerical experiments.

DASSO employs CEI as the acquisition function in both the dice and slice stages to choose the feasible solution to simulate next. Alternatively, acquisition functions such as the knowledge gradient (Frazier et al. 2009) and information-theoretic approaches (see, e.g., Hernández-Lobato et al. 2014, Wang and Jegelka 2017, Hvarfner et al. 2022) could be applied in the slice stage because the candidate solution set is relatively small. However, implementing them in the dice stage requires more careful consideration in large-scale problems. Specifically, identifying the solution that maximizes the acquisition function among a very large number of feasible solutions on each iteration by enumeration imposes too much computational overhead to be practical. An advantage of CEI, as described in Section 4.6, is that the structural properties of CEI support a strategy that identifies the best-CEI solution while only evaluating a small fraction of them.

Before providing the empirical evaluation to demonstrate the excellent finite-sample performance of DASSO in Section 5, we note that we establish the global convergence of DASSO under the very mild conditions and a small tweak to the algorithm; see Section EC.3 in the e-companion. We confess that this convergence result is strictly of academic interest, since we do not expect to approach anything like convergence in the class of problems we consider. The real test for DASSO and any competitor is whether it makes rapid improvement by exploiting limited simulation, unencumbered by algorithm overhead.

4.6. Computational Efficiency of Dice Stage

Although we have focused on the computational expense of the posterior distribution update, all expected-improvement type algorithms suffer the additional computational burden of needing to evaluate the potential of *every* feasible solution on *every* iteration. Choosing a solution heuristically or focusing on only a small subset of feasible solutions risks making a poor decision, which is fatal when simulations are so expensive. Fortunately, DASSO facilitates an elegant solution to this problem.

To make a decision in the dice stage, the CEI of each solution needs to be computed; see Step 5 of Algorithm 1. Aiming to reduce the bottleneck, we derive the following two relationships between the CEIs of solutions that have not been simulated:

PROPOSITION 2. *For two solutions \mathbf{x} and $\check{\mathbf{x}}$ in \mathcal{U} :*

- (R1) $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) = \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$ if $\mathbf{x}^{(\rho)} = \check{\mathbf{x}}^{(\rho)}$ for all $\rho \in \mathcal{G}^{(-g)}$.
- (R2) $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$ if for some group $\rho \in \mathcal{G}^{(-g)}$, all of the following conditions hold:
 - (C2.1) $\mathbf{x}^{(\varrho)} = \check{\mathbf{x}}^{(\varrho)}$ for all $\varrho \in \mathcal{G} \setminus \{\rho, g\}$,
 - (C2.2) $m^{(\rho)}(\mathbf{x}^{(\rho)}) \geq m^{(\rho)}(\check{\mathbf{x}}^{(\rho)})$, and
 - (C2.3) $v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \mathbf{x}^{(\rho)}) \leq v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \check{\mathbf{x}}^{(\rho)})$, where $v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \cdot) = v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}) + v^{(\rho)}(\cdot) - 2c^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \cdot)$.

Furthermore, for two solutions \mathbf{x} and $\check{\mathbf{x}}$ in \mathcal{D} :

- (R3) $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$ if $\mathbf{x}^{(\rho)} = \check{\mathbf{x}}^{(\rho)}$ for all $\rho \in \mathcal{G}^{(-g)}$, $m^{(g)}(\mathbf{x}) \geq m^{(g)}(\check{\mathbf{x}})$, and $v^{(g)}(\tilde{\mathbf{x}}, \mathbf{x}) \leq v^{(g)}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$, where $v^{(g)}(\tilde{\mathbf{x}}, \cdot) = v^{(g)}(\tilde{\mathbf{x}}) + v^{(g)}(\cdot) - 2c^{(g)}(\tilde{\mathbf{x}}, \cdot)$.

Relation (R1) states that the CEIs of the feasible solutions in \mathcal{U} (unsimulated) with the same first $g - 1$ components are identical. If a solution \mathbf{x} is dominated by some other solution $\check{\mathbf{x}}$ in terms of CEI, then \mathbf{x} can be excluded from the candidate solutions to simulate in the next iteration. Recall that $\text{CEI}(\tilde{\mathbf{x}}, \cdot)$ decreases in $m(\cdot)$ and increases in $v(\tilde{\mathbf{x}}, \cdot)$. From this property, Relation (R2) stipulates that the dominance relationship between CEIs at solutions \mathbf{x} and $\check{\mathbf{x}}$ can be established without actually computing them by identifying the Pareto-efficient points for each group based on conditional means and variances. Therefore, we can identify solutions in \mathcal{U} whose CEIs cannot be the largest if one of their lower dimensional components is dominated in mean and variance; in the inventory problem this corresponds to reorder and order-up-to points that cannot have the largest expected improvement on the next iteration. Since these relations enable us to compare the CEIs of two solutions in \mathcal{U} without actually calculating them, the computational burden can be *drastically* reduced, especially when the number of feasible solutions is large. Relation (R3) similarly stipulates the dominance relationship between CEIs at design points (i.e., simulated solutions in \mathcal{D}) with the same first $g - 1$ components. Specifically, we can identify solutions in \mathcal{D} that cannot be the largest if their random-effect component is dominated in mean and variance. Relation (R3) does not reduce the computational burden as much as Relation (R2) does since the cardinality of \mathcal{D} is much smaller than \mathcal{U} , but it could help in a DSO problem for which a larger number of distinct solutions can be simulated.

To further elaborate on (R2), let

$$\mathcal{F}^{(\rho)}(\mathbf{x}^{(\rho)}) = \{\check{\mathbf{x}}^{(\rho)} \in \mathcal{X}^{(\rho)} \setminus \{\mathbf{x}^{(\rho)}\} : m^{(\rho)}(\mathbf{x}^{(\rho)}) \geq m^{(\rho)}(\check{\mathbf{x}}^{(\rho)}) \text{ and } v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \mathbf{x}^{(\rho)}) \leq v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \check{\mathbf{x}}^{(\rho)})\}$$

denote the set of points that Pareto-dominate $\mathbf{x}^{(\rho)} \in \mathcal{X}^{(\rho)}$ for $\rho \in \mathcal{G}^{(-g)}$. Moreover, $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$ for solutions \mathbf{x} and $\check{\mathbf{x}}$ in \mathcal{U} with $\mathbf{x}^{(\varrho)} = \check{\mathbf{x}}^{(\varrho)}$ for all $\varrho \in \mathcal{G} \setminus \{\rho, g\}$ and $\check{\mathbf{x}}^{(\rho)} \in \mathcal{F}^{(\rho)}(\mathbf{x}^{(\rho)})$. We define a point $\mathbf{x}^{(\rho)} \in \mathcal{X}^{(\rho)}$ as Pareto-efficient in group $\rho \in \mathcal{G}^{(-g)}$ if $\mathcal{F}^{(\rho)}(\mathbf{x}^{(\rho)}) = \emptyset$, i.e., if it is not Pareto-dominated by any other point. Then, the Pareto frontier, which consists of the Pareto-efficient points, of group $\rho \in \mathcal{G}^{(-g)}$ is $\mathcal{F}^{(\rho)} = \{\mathbf{x}^{(\rho)} \in \mathcal{X}^{(\rho)} : \mathcal{F}^{(\rho)}(\mathbf{x}^{(\rho)}) = \emptyset\}$. We use these Pareto frontiers to construct $\mathcal{F} = \{\mathbf{x} \in \mathcal{U} : \mathbf{x}^{(\rho)} \in \mathcal{F}^{(\rho)}, \forall \rho \in \mathcal{G}^{(-g)}\}$, which is the set of solutions in \mathcal{U} with Pareto-efficient points for all lower dimensional components except that of the last group. As a result of Proposition 2, the following corollary presents the conditions under which it suffices to calculate only the CEIs of solutions in \mathcal{F} to identify the solution with the largest CEI, in \mathcal{U} .

COROLLARY 2. *Let $\hat{\mathbf{x}} \in \arg \max_{\mathbf{x} \in \mathcal{U} \setminus \mathcal{F}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$. If there exist some solution $\check{\mathbf{x}} \in \mathcal{U}$ such that $\hat{\mathbf{x}}^{(\varrho)} = \check{\mathbf{x}}^{(\varrho)}$ for all $\varrho \in \mathcal{G} \setminus \{\rho, g\}$ and $\check{\mathbf{x}}^{(\rho)} \in \mathcal{F}^{(\rho)}(\hat{\mathbf{x}}^{(\rho)})$ for some group $\rho \in \mathcal{G}^{(-g)}$, then*

$$\max_{\mathbf{x} \in \mathcal{U} \setminus \mathcal{F}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) = \text{CEI}(\tilde{\mathbf{x}}, \hat{\mathbf{x}}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}}) \leq \max_{\mathbf{x} \in \mathcal{F}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}).$$

Furthermore, $\check{\mathbf{x}} \in \mathcal{F}$ because $\hat{\mathbf{x}} \notin \arg \max_{\mathbf{x} \in \mathcal{U} \setminus \mathcal{F}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$, otherwise.

When the number of feasible solutions is large, the condition in Corollary 2 holds almost all the time because considering the large size of \mathcal{U} , for all $\hat{\mathbf{x}} \in \arg \max_{\mathbf{x} \in \mathcal{U} \setminus \mathcal{F}} \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})$, it is very unlikely to simulate every solution $\check{\mathbf{x}}$ with $\hat{\mathbf{x}}^{(\varrho)} = \check{\mathbf{x}}^{(\varrho)}$ for all $\varrho \in \mathcal{G} \setminus \{\rho, g\}$ and $\check{\mathbf{x}}^{(\rho)} \in \mathcal{F}^{(\rho)}(\hat{\mathbf{x}}^{(\rho)})$, for some group $\rho \in \mathcal{G}^{(-g)}$. Therefore, Corollary 2 implies that it suffices to calculate the CEIs of solutions in $\mathcal{D} \cup \mathcal{F}$ to find a solution with the largest CEI in \mathcal{X} because for some solution with the largest CEI in $\mathcal{X} \setminus (\mathcal{D} \cup \mathcal{F}) = \mathcal{U} \setminus \mathcal{F}$, there exists a solution in \mathcal{F} with a larger CEI. Thus, a solution with the largest CEI in \mathcal{F} has a relatively large—typically the largest—CEI in \mathcal{U} . Further, following Relation (R1), the CEIs of solutions in $\mathcal{F} \subseteq \mathcal{U}$ with the same first $g - 1$ components are identical. Therefore, it is sufficient to calculate the CEIs of solutions in $\mathcal{D} \cup \tilde{\mathcal{F}}$ to obtain promising first $g - 1$ components, where $\tilde{\mathcal{F}} \subseteq \mathcal{F}$ includes only one solution for each first $g - 1$ components. Since $\mathcal{D} \cup \tilde{\mathcal{F}}$ is much smaller than \mathcal{X} , the reduction in computational overhead is huge, especially when the number of feasible solutions is very large; computational analysis conducted in Section 5 shows this tremendous benefit.

Despite the significant computational savings from using Pareto frontiers, DASSO stills encounters a limit on the problem size it can handle because the size of $\tilde{\mathcal{F}}$ increases exponentially in the number of groups. To extend DASSO beyond this limit, a heuristic approach can keep the number of CEI calculations computationally feasible. In particular, \mathcal{F} can be alternatively constructed as $\{\mathbf{x} \in \mathcal{U} : \mathbf{x}^{(\rho)} \in \tilde{\mathcal{F}}^{(\rho)}, \forall \rho \in \mathcal{G}^{(-g)}\}$, where $\tilde{\mathcal{F}}^{(\rho)}$ is some subset of $\mathcal{X}^{(\rho)}$ chosen heuristically at a desired size. We employ this approach in our numerical experiments in Section 5 when $\{\mathbf{x} \in \mathcal{U} : \mathbf{x}^{(\rho)} \in \mathcal{F}^{(\rho)}, \forall \rho \in \mathcal{G}^{(-g)}\}$ is too large to handle directly.

Computing CEIs requires the diagonal elements of $[\bar{\mathbf{Q}}^{(\rho)}]^{-1}$, the column of $[\bar{\mathbf{Q}}^{(\rho)}]^{-1}$ corresponding to $\tilde{\mathbf{x}}$, and $\mathbf{m}^{(\rho)}$, for each $\rho \in \mathcal{G}$. Instead of inverting $\bar{\mathbf{Q}}^{(\rho)}$ to obtain these elements, we use the block matrix inversion formula (Lemma EC.2 in Section EC.1 in the e-companion) to compute the elements more efficiently. Letting $\check{\Sigma}_{\mathcal{D}\mathcal{D}} = \Sigma_{\mathcal{D}\mathcal{D}} + \Sigma^\epsilon$, we have

$$[\bar{\mathbf{Q}}^{(\rho)}]^{-1} = \begin{pmatrix} [\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} + [\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} \mathbf{Q}_{\mathcal{U}\mathcal{D}}^{(\rho)} \mathbf{S}^{(\rho)} \mathbf{Q}_{\mathcal{D}\mathcal{U}}^{(\rho)} [\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} & -[\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} \mathbf{Q}_{\mathcal{U}\mathcal{D}}^{(\rho)} \mathbf{S}^{(\rho)} \\ -\mathbf{S}^{(\rho)} \mathbf{Q}_{\mathcal{D}\mathcal{U}}^{(\rho)} [\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} & \mathbf{S}^{(\rho)} \end{pmatrix},$$

and thus

$$\mathbf{m}^{(\rho)} = \begin{pmatrix} -[\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}]^{-1} \mathbf{Q}_{\mathcal{U}\mathcal{D}}^{(\rho)} \mathbf{S}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \mathbf{E}^{(\rho)} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \\ \mathbf{S}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \mathbf{E}^{(\rho)} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \end{pmatrix},$$

where $\mathbf{S}^{(\rho)} = \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)} - \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)}$ is the inverse of the Schur complement of $\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(\rho)}$ in $[\bar{\mathbf{Q}}^{(\rho)}]^{-1}$, for $\rho \in \mathcal{G}^{(-g)}$. Also, we have

$$[\bar{\mathbf{Q}}^{(g)}]^{-1} = \begin{pmatrix} \sigma_g^2 \mathbf{I}_{|\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & \mathbf{S}^{(g)} \end{pmatrix}, \text{ and thus } \mathbf{m}^{(g)} = \begin{pmatrix} \vec{\mathbf{0}}_{|\mathcal{U}|} \\ \mathbf{S}^{(g)} \mathbf{E}^{(g)} (\bar{\mathbf{Y}}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{D}}) \end{pmatrix},$$

where $\mathbf{S}^{(g)} = \sigma_g^2 \mathbf{I}_{|\mathcal{D}|} - (\sigma_g^2)^2 \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1}$ is the inverse of the Schur complement of $\mathbf{Q}_{\mathcal{U}\mathcal{U}}^{(g)}$ in $[\bar{\mathbf{Q}}^{(g)}]^{-1}$, making $\mathbf{S}^{(g)} \mathbf{E}^{(g)} = \mathbf{I}_{|\mathcal{D}|} - \frac{1}{\sigma_g^2} \mathbf{S}^{(g)}$. Further, to compute $\mathbf{E}^{(\rho)}$ more efficiently for $\rho \in \mathcal{G}^{(-g)}$, we use the Woodbury matrix identity (Lemma EC.1 in Section EC.1 in the e-companion):

$$\mathbf{E}^{(\rho)} = \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1} + \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \left([\Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)}]^{-1} - [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \right)^{-1} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top \check{\Sigma}_{\mathcal{D}\mathcal{D}}^{-1}.$$

This is more efficient than inverting $\Sigma_{\mathcal{D}\mathcal{D}} - \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top + \Sigma^\epsilon$ of size $|\mathcal{D}| \times |\mathcal{D}|$ for each $\rho \in \mathcal{G}^{(-g)}$ because it requires inverting two matrices of size $|\mathcal{D}^{(\rho)}| \times |\mathcal{D}^{(\rho)}|$ for each $\rho \in \mathcal{G}^{(-g)}$ in addition to inverting $\check{\Sigma}_{\mathcal{D}\mathcal{D}}$ of size $|\mathcal{D}| \times |\mathcal{D}|$ only once. Since $|\mathcal{D}^{(\rho)}|$ is bounded by $n^{(\rho)}$ and $|\mathcal{D}|$, and the increase on $|\mathcal{D}^{(\rho)}|$ is more slowly than the increase on $|\mathcal{D}|$ as the algorithm iterates, the benefit is greater in the later stages.

5. Numerical Experiments

In this section, we empirically evaluate the performance of DASSO, investigate its sensitivity to properties of the objective function, and explore the importance of the chosen decomposition. A key feature of the DASSO prior is that it is consistent even with non-separable objective functions, but one might expect that a separable objective function, if one knows the proper separation, would be favorable.

For the numerical experiments we consider four examples: The first two examples are optimization of the high-dimensional Zakharov and Styblinski-Tang test functions with added stochastic noise. The third is a multi-product inventory problem which can be considered as a practical DSO problem with knowledge about a natural decomposition. The last is a stylized problem allowing control of how separable the objective function is. Recall that lack of separability—i.e., the contribution of the remainder term in the decomposition—is addressed in the slice stage, so these experiments assess the effectiveness of slicing. While these problems are all high-dimensional with large numbers of feasible solutions, they are not computationally expensive simulations. This allows us to run many long macroreplications that show convergence

behaviors, which we would not expect to do in practice. Therefore, our computation-time analysis is effectively an assessment of algorithm overhead, not simulation time.

In all experiments the parameters of the GMRF are estimated from s initial design points chosen by Latin hypercube sampling and sg additional design points selected to augment these s initial design points; this is the minimum number of additional points needed to parameterize the g groups plus remainder prior. See Section EC.5 in the e-companion for details of parameter estimation. All $s(g + 1)$ design points are simulated for r_0 replications. However, for better comparison in different settings, the DASSO search is initialized from only the s initial design points, which are common with the same simulation outputs for all settings. Each time a solution is selected, it is simulated r_d replications if it has been simulated previously and r_u replications otherwise. For the experiments, $s = 100$ and $r_0 = r_d = r_u = 10$ in Section 5.1, $s = 200$ and $r_0 = r_d = r_u = 10$ in Section 5.2, and $s = 15$, $r_0 = 20$, $r_d = 4$, and $r_u = 10$ in Sections 5.3 and 5.4. These are largely arbitrary choices and there is likely room for a more data-driven selections, a topic we reserve for future work. All computations are executed on a desktop computer with a Windows 10 operating system, a 2.9 GHz Intel Core i7 CPU, 32 GB of RAM, 8 cores and 16 logical processors.

5.1. Well-Known Test Functions

The k -dimensional Zakharov function $y(x_1, x_2, \dots, x_k) = \sum_{i=1}^k x_i^2 + \left(\sum_{i=1}^k 0.5ix_i\right)^2 + \left(\sum_{i=1}^k 0.5ix_i\right)^4$ is often used as a test case in BO as it presents a very challenging problem. The function is minimized at $\mathbf{x} = (0, 0, \dots, 0)^\top$. Letting $\mathcal{X} = \{-2, -1, 0, 1, 2\}^{10}$ be the set of feasible solutions, the 10-dimensional Zakharov function takes values between $[0, 9.15 \times 10^6]$ and has around 9.8 million feasible solutions. To make the problem stochastic, we add zero-mean normally distributed noise with variance 1.8^2 to objective function values. We also considered different levels of stochastic noise with variances 2.6^2 , 3.9^2 , and 5.2^2 , but chose not to exhibit the results because they show similar performance.

The k -dimensional Styblinski-Tang function $y(x_1, x_2, \dots, x_k) = \frac{1}{20} \sum_{i=1}^k (x_i^4 - 16x_i^2 + 5x_i)$ is another standard test case. Letting $\mathcal{X} = \{-6, -3, 0, 3, 6\}^{10}$ be the set of feasible solutions, the 10-dimensional Styblinski-Tang function takes values between $[-39, 375]$ and is minimized at $\mathbf{x} = (-3, -3, \dots, -3)^\top$. We add zero-mean normally distributed noise with variance 3^2 to objective function values.

We use these two problems to compare the performance of DASSO to pGMIA proposed by Li and Song (2024), the current state-of-the-art for high-dimensional DSO problems, and Bounce (Bayesian optimization using increasingly high-dimensional combinatorial and continuous embeddings) proposed by Papenmeier et al. (2023), a high-dimensional BO algorithm optimizing over combinatorial, continuous, or mixed spaces.

pGMIA was shown empirically to outperform the multi-resolution GMIA of Salemi et al. (2019) and four state-of-the-art high-dimensional BO algorithms: Random EMbedding Bayesian Optimization (Wang et al. 2016), Sparse Axis-Aligned Subspace Bayesian Optimization (Eriksson and Jankowiak 2021), and High-Dimensional Bayesian Optimization and High-Dimensional Batch Bayesian Optimization (Wang et al. 2017). Thus, comparing DASSO with pGMIA indirectly assesses DASSO against these algorithms.

Unlike DASSO, which decomposes the prior distribution into an additive form, pGMIA batches the dimensions into two layers and hierarchically optimizes each layer by projecting one layer onto the other. We note that pGMIA is implemented in MATLAB while DASSO is implemented in Python. Since pGMIA with the random projection criterion (pGMIA-R) is observed to empirically outperform other benchmarks, we compare DASSO to that version. In their numerical experiments, Li and Song (2024) set the numbers of region- and solution-layer dimensions to 5 for pGMIA. The closest equivalence for DASSO is a decomposition consisting of two 5-dimensional groups. Let \mathcal{G}_5 denote such a decomposition with the sets of component indices being $\mathbf{g}^{(1)} = \{1, 2, 3, 4, 5\}$ and $\mathbf{g}^{(2)} = \{6, 7, 8, 9, 10\}$. This assignment of component indices to groups is arbitrary, and when we tried other assignments we observed no difference in performance.

Papenmeier et al. (2023) demonstrate that Bounce empirically outperforms five state-of-the-art algorithms designed for combinatorial, continuous, or mixed input domains: BO with Dictionaries (Deshwal et al. 2023), CAtegorical Spaces, or Mixed, OPTimisatiOn with Local-trust-regIons & TAilored Non-parametric (Wan et al. 2021), COMbinatorial Bayesian Optimization (Oh et al. 2019), Sequential Model-based Algorithm Configuration (Hutter et al. 2011), and Random Decomposition Upper-Confidence Bound (Ziomek and Ammar 2023). To address the computational burden of high-dimensional BO, Bounce defines a GP surrogate in an iteratively refined lower-dimensional subspace called the target space, and maximizes the acquisition function within promising regions of the target space by using a trust-region-based method. Bounce is also implemented in Python.

Bounce is not designed to handle stochastic output, unlike DASSO and pGMIA that are specifically created to be effective on such problems. Therefore, whenever a solution is queried by Bounce, we run 10 replications and treat their average as deterministic response. To align Bounce as closely as possible to DASSO with decomposition \mathcal{G}_5 , we set its initial target dimensionality to 5.

Applications of large-scale, high-dimensional DSO typically exist between two extremes: Simulations that execute so slowly that even substantial optimization algorithm overhead is negligible relative to simulation time, and simulations that execute so fast that algorithm overhead is the bottleneck to computational feasibility. Tongarlak et al. (2010) is an example of the former (20 replications of a single feasible solution took 8 hours), and the Zakharov and Styblinski-Tang functions are examples of the latter. DASSO is designed to be effective when the number of feasible solutions or replications to simulate is computationally limited (i.e., the simulation is very slow), while imposing much less algorithm overhead than other methods.

To see this, Figures 3 and 4 show the mean optimality gap of the 10-dimensional Zakharov and Styblinski-Tang functions, respectively, as the distance from the optimal value versus the total number of simulation replications (left) and versus the wall-clock time (right), averaged over 100 macro-replications of DASSO, pGMIA, and Bounce. When each simulation replication takes significant time it is more critical to make rapid improvement in the objective value with fewer replications. On the other hand, when each simulation replication is very fast, the computational overhead of the algorithm matters. The figures show

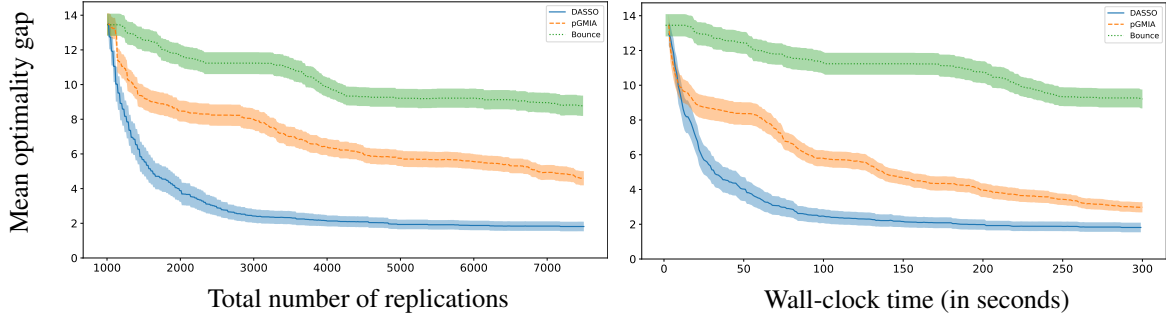


Figure 3 Mean optimality gap of the Zakharov function vs. total number of replications and wall-clock time across 100 macro-replications for DASSO (blue solid line) with decomposition \mathcal{G}_5 , pGMIA (orange dashed line), and Bounce (green dotted line). The shaded area around each curve shows point-wise ± 2 standard error of the average.

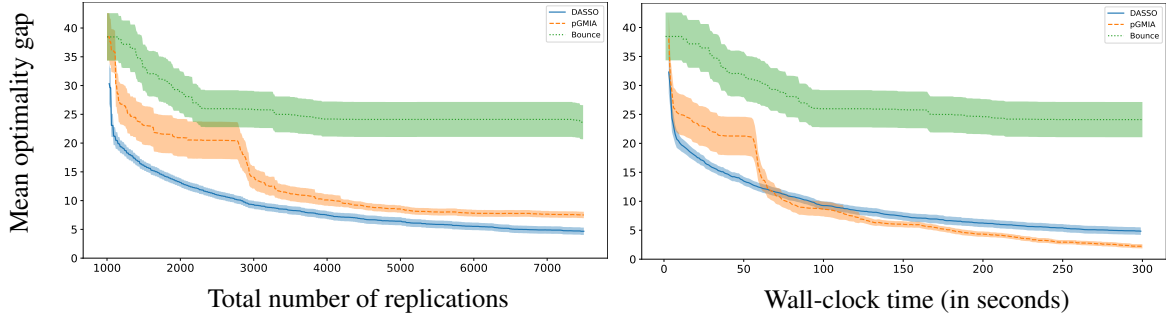


Figure 4 Mean optimality gap of the Styblinski-Tang function vs. total number of replications and wall-clock time across 100 macro-replications for DASSO (blue solid line) with decomposition \mathcal{G}_5 , pGMIA (orange dashed line), and Bounce (green dotted line). The shaded area around each curve shows point-wise ± 2 standard error of the average.

that DASSO outperforms pGMIA and Bounce in terms of progress per replication for both Zakharov and Styblinski-Tang. In terms of progress per second, DASSO is the best for Zakarov and slightly lags pGMIA for Styblinski-Tang after a faster start. Since the experimental settings are the same, the difference in performance between these test functions can be attributed to different objective functions. For our target applications of DASSO the initial rapid decrease in the optimality gap is what we desire; we do not expect to achieve anything like convergence when the simulation is high-dimensional and computationally expensive.

5.2. A Higher-Dimensional Test Function

To evaluate the performance of DASSO in problems of significantly higher dimension, we compare it to pGMIA+ (Li and Song 2024), a multi-layer extension of pGMIA for higher-dimensional problems. We test the 100-dimensional Zakharov function with $\mathcal{X} = \{-5, -4, \dots, 5\}^{100}$ which takes values between $[0, 2.54 \times 10^{16}]$. To make the problem stochastic, we add zero-mean normally distributed noise with variance 1.8^2 to objective function values.

In numerical experiments for pGMIA+, Li and Song (2024) set the numbers of dimensions in the first two layers to 3 leaving the 94 dimensions to the last. The sampling decisions in the first two layers are made

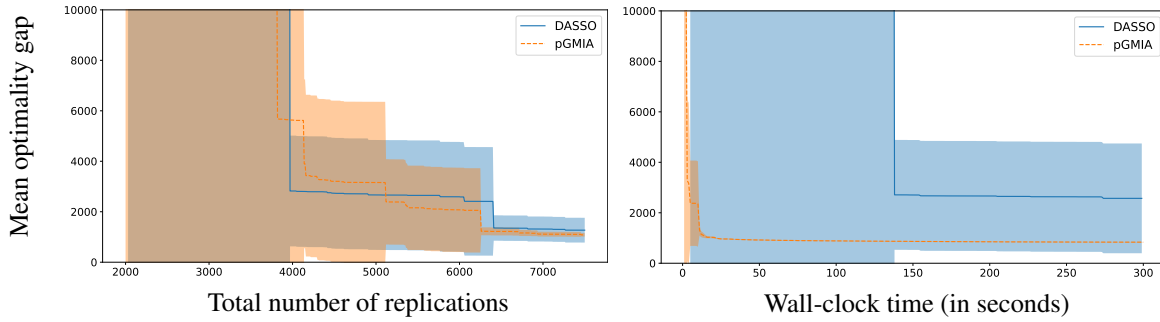


Figure 5 Mean optimality gap of the 100-dimensional Zakharov function vs. total number of replications and wall-clock time across 100 macro-replications for DASSO (blue solid line) with decomposition \mathcal{G}_3 and pGMIA+ (orange dashed line). The shaded area around each curve shows point-wise ± 2 standard error of the average.

by maximizing the CEI while the coordinates in the remaining 94 dimensions are selected randomly. The closest correspondence for DASSO is a decomposition consisting of thirty-two 3-dimensional groups and two 2-dimensional groups. Let \mathcal{G}_3 denote such a decomposition with the sets of component indices being $\mathbf{g}^{(i)} = \{3i - 2, 3i - 1, 3i\}$ for $i = 1, 2, \dots, 32$, $\mathbf{g}^{(33)} = \{97, 98\}$, and $\mathbf{g}^{(34)} = \{99, 100\}$. Furthermore, since the problem size makes the exact CEI maximization computationally infeasible in each dice stage, \mathcal{F} is heuristically constructed as $\left\{ \mathbf{x} \in \mathcal{U}: \mathbf{x}^{(\rho)} \in \hat{\mathcal{F}}^{(\rho)}, \forall \rho \in \mathcal{G}^{(-g)} \right\}$, where $\hat{\mathcal{F}}^{(\rho)}$ is the Pareto frontier $\mathcal{F}^{(\rho)}$ for two randomly chosen ρ s in $\mathcal{G}^{(-g)}$ and a single randomly chosen point from $\mathcal{X}^{(\rho)}$ for the remaining thirty-one groups. This is a heuristic modification with the goal to demonstrate feasibility for DASSO to tackle even higher-dimensional problems.

Figure 5 shows the mean optimality gap of the 100-dimensional Zakharov function as the distance from the optimal value versus the total number of simulation replications (left) and versus the wall-clock time (right), averaged over 100 macro-replications of DASSO and pGMIA+. In terms of progress per replication in this example, DASSO and pGMIA+ outperform each other at different points while in the later iterations, pGMIA+ appears to have smaller optimality gap on average (although statistically indistinguishable from DASSO) with smaller run-to-run variability. Considering the range of the function, 2.54×10^{16} , the optimality gap (approximately 1000) after 7500 replications is negligible for both algorithms. In terms of progress per second, the large number of groups in DASSO causes computational overhead, making the algorithm iterate slower than pGMIA+. If *simulation* overhead was more substantial, rather than negligible as it is for simulating the Zakharov function, the difference would be less.

5.3. Inventory Problem

Consider a multi-product inventory problem where each product follows an (s, S) policy, i.e., once the inventory level of product ρ falls below its reorder point, $s^{(\rho)}$, an order is given to bring the inventory level up to $S^{(\rho)}$. This problem is a variant of the single-product version in Salemi et al. (2019), where the product is subject to periodic demand that follows a Poisson distribution and the optimal policy is $(s, S - s) = (18, 35)$

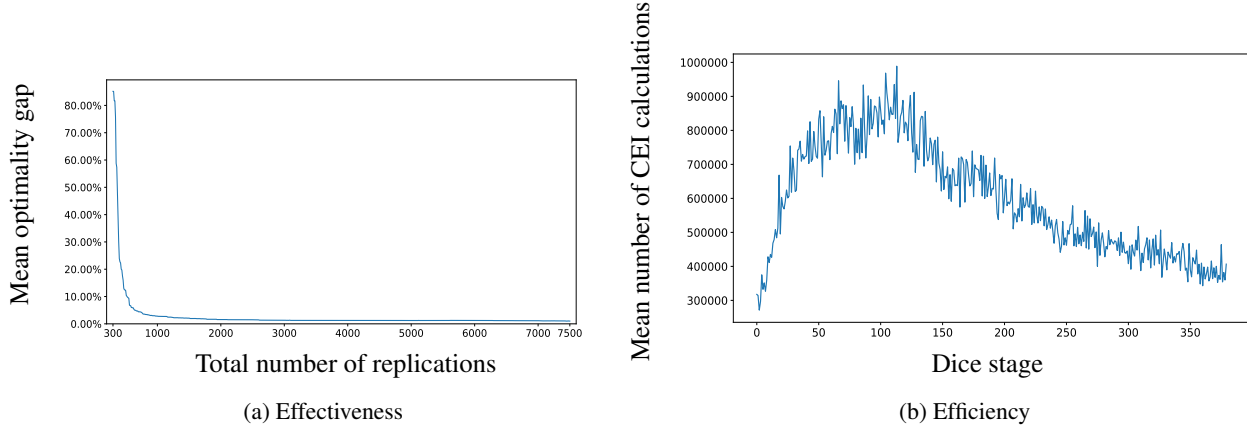


Figure 6 (a) Mean optimality gap vs. total number of replications and (b) mean number of CEI calculations at each dice stage across 30 macro-replications for the natural decomposition in the inventory problem.

under the objective of minimizing the expected average cost per period over a fixed planning horizon. Similarly, in our version, demand for each product follows a common Poisson distribution and the objective is to minimize the expected value of cost defined as $y(\mathbf{x}) = \sum_{\rho \in \mathcal{G}} y^{(\rho)}(\mathbf{x}^{(\rho)}) + y^{(r)}(\mathbf{x})$, where $y^{(\rho)}(\mathbf{x}^{(\rho)})$ represents the expected average cost per period for product ρ and $y^{(r)}(\mathbf{x})$ is a deterministic interaction term that depends on how much feasible solution \mathbf{x} deviates from the optimal solution. In other words, we induce an interaction so that the problem is not separable. Therefore, the optimal policy for each product is $(s^{(\rho)}, S^{(\rho)} - s^{(\rho)}) = (18, 35)$, and the interaction term is $y^{(r)}(\mathbf{x}) = \prod_{\rho \in \mathcal{G}} \sqrt{(s^{(\rho)} - 18)^2 + (S^{(\rho)} - s^{(\rho)} - 35)^2}$. To treat the feasible region as a hyperbox defined on the integer lattice, we let $\mathbf{x}^{(\rho)} = (s^{(\rho)}, S^{(\rho)} - s^{(\rho)})$ for $\rho \in \mathcal{G}$ and $\mathbf{x} = \{(s^{(\rho)}, S^{(\rho)} - s^{(\rho)})\}_{\rho \in \mathcal{G}}$.

The natural decomposition for this multi-product inventory problem is that each group represents a product, and thus is 2-dimensional, one dimension is for $s^{(\rho)}$ and the other is for $S^{(\rho)} - s^{(\rho)}$. We consider a 5-product (i.e., 10-dimensional) problem where the feasible region is defined by constraints $10 \leq s^{(\rho)} \leq 34$ and $20 \leq S^{(\rho)} - s^{(\rho)} \leq 44$, which leads to $n = 25^{10} = 95,367,431,640,625$ feasible solutions. In the experiments, we set the number of macro-replications to 30 and the total number of replications to 7500, i.e., DASSO terminates after obtaining this number of simulation outputs.

Figure 6a depicts the mean optimality gap as a percentage versus total number of replications for 30 macro-replications of the inventory problem with the natural decomposition. Before DASSO starts performing dice and slice stages, the mean optimality gap of the sample-best solution among the first $s = 15$ initial design points is around 85%; recall that these design points are chosen by Latin hypercube sample and simulated for 20 replications. Clearly DASSO obtains rapid improvement: it only takes exploring around 145 (or 38) design points and simulating 2500 (or 650) replications in total to have the mean optimality gap below 1.5% (or 5%). Considering the very large scale of the problem, with more than 95 trillion feasible solutions, the design points explored by the algorithm to obtain such small optimality gaps is an amazingly tiny portion of the feasible solution set.

Recall from Section 4.6 that rather than performing CEI calculation for all 95 trillion feasible solutions at each dice stage, it is sufficient to calculate the CEIs of solutions from the Pareto set $\mathcal{D} \cup \check{\mathcal{F}}$. To illustrate how much smaller $|\mathcal{D} \cup \check{\mathcal{F}}|$ is than the number of feasible solutions $n = |\mathcal{X}|$, Figure 6b shows the mean number of CEI calculations performed, i.e., $|\mathcal{D} \cup \check{\mathcal{F}}|$, at each dice stage averaged over 30 macro-replications of the inventory problem with the natural decomposition. Although the mean number of CEI calculations varies from one dice stage to another, it is no greater than 989 thousand, which is around one hundred-millionth of the number of feasible solutions, i.e., $|\mathcal{D} \cup \check{\mathcal{F}}| \leq 10^{-8}n$. Without such a reduction in the number of CEI calculations it would be computationally impossible to solve a problem of this scale.

The mean total CPU times spent to perform three different tasks—simulation execution, slice stage, and dice stage—averaged over 30 macro-replications of the inventory problem with the natural decomposition are as follows; recall that the algorithm stops after simulating 7500 replications. It takes around 347 CPU seconds (5.8 CPU minutes) on average to run a single macro-replication: 251 seconds for the dice stage, 79 seconds for the slice stage, and 17 seconds for the simulation executions. The dice and slice stage times are independent of the simulation time, which is abnormally small in this constructed example. In other words, even if the simulation took orders of magnitude longer the DASSO algorithm overhead would be the same. These results indicate that DASSO is not only effective but also efficient to solve very large-scale problems.

The results above are for the natural decomposition, where each group represents a product. To explore the importance of decomposition, in addition to the natural one, we consider 10 alternative decompositions with five 2-dimensional groups; notice that the decompositions are the same size. These additional decompositions are different from the natural one in a way that at least one group represents either $(s^{(\rho)}, S^{(\varrho)} - s^{(\varrho)})$, or $(s^{(\rho)}, s^{(\varrho)})$, or $(S^{(\rho)} - s^{(\rho)}, S^{(\varrho)} - s^{(\varrho)})$, for some $\rho \neq \varrho$. Figure 7 illustrates the mean optimality gap as a percentage versus total number of replications for 30 macro-replications of the inventory problem with various decompositions. The results indicate that the performance of DASSO with some alternative decompositions is statistically indistinguishable from that achieved with the natural one, suggesting that an alternative decomposition can be used when there is no knowledge about a natural decomposition. In the next section, we further explore the importance of decomposition.

5.4. Controlled Test Function

Consider the objective function $y(\mathbf{x}) = (1 - \alpha) \sum_{\rho \in \mathcal{G}} y^{(\rho)}(\mathbf{x}^{(\rho)}) + \alpha \lambda y^{(r)}(\mathbf{x})$ for $\mathbf{x} \in \mathcal{X}$, where $\mathcal{X} = \{\underline{x}, \underline{x} + 1, \dots, \bar{x} - 1, \bar{x}\}^d$ for some integers \underline{x} and \bar{x} and dimension d . The role of $\alpha \in [0, 1]$ is to control how close to separable the objective function is, or equivalently, how much the remainder matters. Let each $y^{(\rho)}(\cdot)$ for $\rho \in \mathcal{G}$ and $y^{(r)}(\cdot)$ be an inverted multivariate normal density function with a shift, generically defined as $f(x_1, x_2, \dots, x_k) = -\gamma_1 \exp\left\{-\gamma_2 \sum_{i=1}^k i x_i^2\right\} + \gamma_1$, where $\gamma_1 = 1000$ and $\gamma_2 = 0.001$; notice that $k = d^{(\rho)}$ for each $\rho \in \mathcal{G}$ and $k = d$ for the remainder term. The function is minimized at $x_i = 0$ for $i = 1, 2, \dots, k$, and thus, $y(\mathbf{x})$ is also minimized at $\mathbf{x} = (0, 0, \dots, 0)^\top$. The shifting term makes the optimal

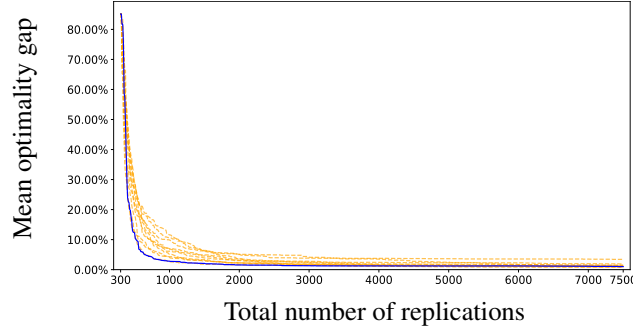


Figure 7 Mean optimality gap vs. total number of replications across 30 macro-replications for 11 different decompositions in the inventory problem. The blue solid line represents the natural decomposition while each orange dashed line represent an alternative decomposition.

objective function value 0 regardless of the value of α . To have the same range of $y(\cdot)$ for any α , we set $\lambda = \sum_{\rho \in \mathcal{G}} y^{(\rho)}(\mathbf{x}^{(\rho)*}) / y^{(r)}(\mathbf{x}^*)$ with solution \mathbf{x}^* that maximizes the objective function. We consider a 12-dimensional problem with $d^{(\rho)} = 2$ for $\rho \in \mathcal{G} = \{1, 2, \dots, 6\}$, where the sets of component indices are $\mathbf{g}^{(1)} = \{1, 2\}$, $\mathbf{g}^{(2)} = \{3, 4\}$, $\mathbf{g}^{(3)} = \{5, 6\}$, $\mathbf{g}^{(4)} = \{7, 8\}$, $\mathbf{g}^{(5)} = \{9, 10\}$, and $\mathbf{g}^{(6)} = \{11, 12\}$; recall that $\mathbf{x}^{(\rho)} = (x_i)_{i \in \mathbf{g}^{(\rho)}}$. Letting $\underline{x} = -2$ and $\bar{x} = 2$, the objective function values range between $[0, 71.57]$, and the total number of feasible solutions is $n = 5^{12} = 244,140,625$. To make the problem stochastic, a zero-mean normally distributed noise with variance 3^2 is added to the objective function values.

Slightly abusing notation, let \mathcal{G}_2 denote the decomposition that consists of $\mathbf{g}^{(1)} = \{1, 2\}$, $\mathbf{g}^{(2)} = \{3, 4\}$, $\mathbf{g}^{(3)} = \{5, 6\}$, $\mathbf{g}^{(4)} = \{7, 8\}$, $\mathbf{g}^{(5)} = \{9, 10\}$, and $\mathbf{g}^{(6)} = \{11, 12\}$; we use the subscript to indicate that each group is 2-dimensional. We treat \mathcal{G}_2 as a natural decomposition since it defines the objective function. We also consider three alternative decompositions: \mathcal{G}_1 consists of $\mathbf{g}^{(\rho)} = \{\rho\}$ for $\rho = 1, 2, \dots, 12$; \mathcal{G}_3 consists of $\mathbf{g}^{(1)} = \{1, 2, 3\}$, $\mathbf{g}^{(2)} = \{4, 5, 6\}$, $\mathbf{g}^{(3)} = \{7, 8, 9\}$, and $\mathbf{g}^{(4)} = \{10, 11, 12\}$; and finally \mathcal{G}_4 consists of $\mathbf{g}^{(1)} = \{1, 2, 3, 4\}$, $\mathbf{g}^{(2)} = \{5, 6, 7, 8\}$, and $\mathbf{g}^{(3)} = \{9, 10, 11, 12\}$. \mathcal{G}_4 is also a natural decomposition because there is no interaction among groups.

Figure 8 depicts the mean optimality gap as a distance from the optimal value versus total number of replications averaged over 50 macro-replications of the stylized problem for $\alpha \in \{0, 0.5, 1\}$ with decompositions \mathcal{G}_1 , \mathcal{G}_2 , \mathcal{G}_3 , and \mathcal{G}_4 . Although \mathcal{G}_4 shows rapid improvement in the early stages of the algorithm, it fails to close the gap as fast as the other decompositions. On the other hand, \mathcal{G}_1 performs the best. Even for $\alpha = 0$, where the objective function value is fully separable, the natural decompositions \mathcal{G}_2 and \mathcal{G}_4 do not show a better performance than \mathcal{G}_1 . For $\alpha = 1$, where the objective function value is not even approximately separable, DASSO still performs well with a good choice of decomposition, such as \mathcal{G}_1 , \mathcal{G}_2 or \mathcal{G}_3 .

We note that the discussion above ignores the computation time overhead of DASSO and focuses on the total number of replications. One might expect that decompositions with smaller-size groups are computationally faster, and thus \mathcal{G}_1 should be the ideal decomposition. However, this is not the case based on Table 1 because \mathcal{G}_1 is much slower than the other decompositions. The table exhibits the mean CPU times spent to

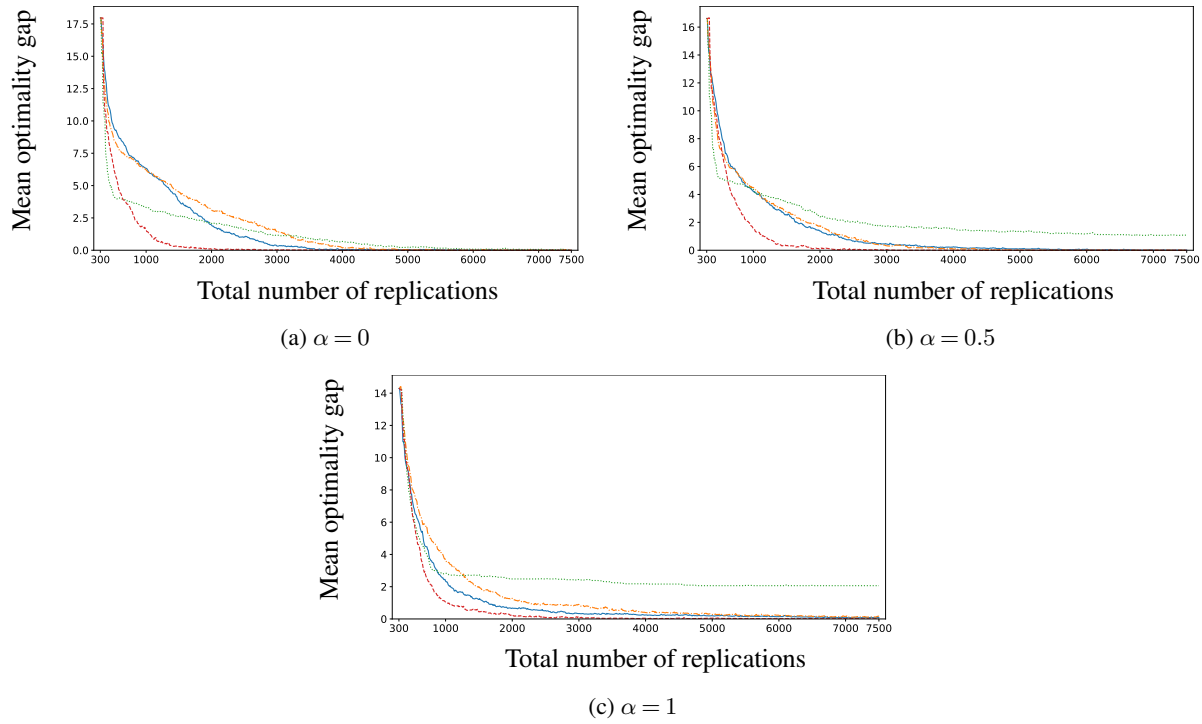


Figure 8 Mean optimality gap vs. total number of replications across 50 macro-replications for 4 different decompositions in the stylized problem with various α values. The red dashed, blue solid, orange dash-dotted, and green dotted lines represent \mathcal{G}_1 , \mathcal{G}_2 , \mathcal{G}_3 , and \mathcal{G}_4 , respectively.

Table 1 Mean CPU times (in seconds) spent for performing three different tasks across 50 macro-replications for 4 different decompositions in the stylized problem with $\alpha = 0$. The percentages in parentheses represent the proportion of each task to total time.

Task	\mathcal{G}_1	\mathcal{G}_2	\mathcal{G}_3	\mathcal{G}_4
Simulation	1.5 (0.2%)	2.3 (4.5%)	2.5 (3.9%)	4.2 (2.5%)
Slice stage	4.6 (0.7%)	10.3 (20.0%)	15.2 (24.0%)	80.8 (47.7%)
Dice stage	698.7 (99.1%)	38.9 (75.5%)	45.7 (72.1%)	84.5 (49.8%)
All	704.8 (100%)	51.5 (100%)	63.4 (100%)	169.5 (100%)

perform three different tasks—simulation execution, slice stage and dice stage—as well as the total time across 50 macro-replications for decompositions \mathcal{G}_1 , \mathcal{G}_2 , \mathcal{G}_3 , and \mathcal{G}_4 in the stylized problem with $\alpha = 0$; recall that the algorithm stops after simulating 7500 replications.

The time spent to perform the slice stages increases as the size of the groups increases, whereas that of the dice stages does not show the same trend. This is mainly because \mathcal{G}_1 requires many more CEI calculations: the mean number of CEI calculations in each dice stage varies between 72 thousand to 5 million for \mathcal{G}_1 , 3.7 thousand to 49 thousand for \mathcal{G}_2 , 1.3 thousand to 5.4 thousand for \mathcal{G}_3 , and 103 to 1.2 thousand for \mathcal{G}_4 . This result also illustrates the importance of the reduction of the number of CEI calculations.

To further investigate the importance of decomposition, we consider two more decompositions with the same size as \mathcal{G}_2 , i.e., they also consists of six 2-dimensional groups. The sets of component indices for the

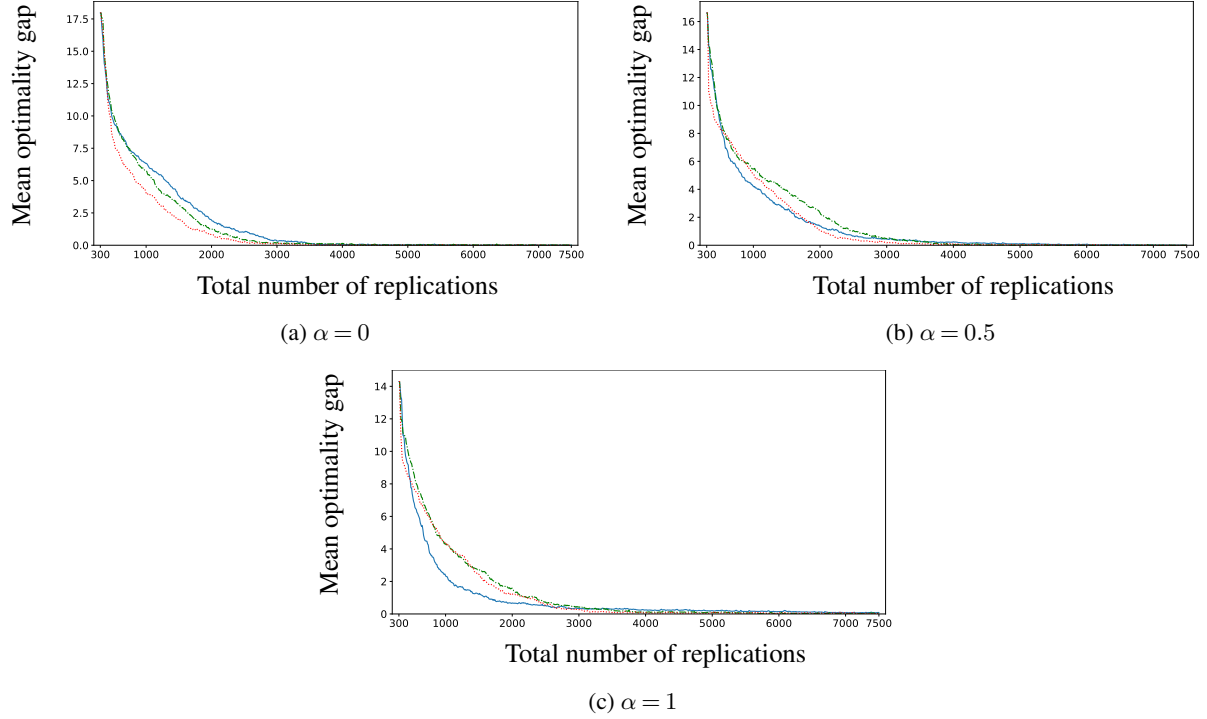


Figure 9 Mean optimality gap vs. total number of replications across 50 macro-replications for 3 different decompositions in the stylized problem with various α values. The blue solid, green dash-dotted, and red dotted lines represent \mathcal{G}_2 , $\hat{\mathcal{G}}_2$, and $\tilde{\mathcal{G}}_2$, respectively.

first one, denoted by $\hat{\mathcal{G}}_2$, are $\mathbf{g}^{(1)} = \{1, 10\}$, $\mathbf{g}^{(2)} = \{5, 8\}$, $\mathbf{g}^{(3)} = \{2, 7\}$, $\mathbf{g}^{(4)} = \{4, 11\}$, $\mathbf{g}^{(5)} = \{3, 9\}$, and $\mathbf{g}^{(6)} = \{6, 12\}$. The sets of component indices for the second one, denoted by $\tilde{\mathcal{G}}_2$, are $\mathbf{g}^{(1)} = \{1, 7\}$, $\mathbf{g}^{(2)} = \{2, 6\}$, $\mathbf{g}^{(3)} = \{3, 8\}$, $\mathbf{g}^{(4)} = \{4, 10\}$, $\mathbf{g}^{(5)} = \{5, 12\}$, and $\mathbf{g}^{(6)} = \{9, 11\}$. Figure 9 depicts the mean optimality gap as a distance from the optimal value versus total number of replications for 50 macro-replications of the stylized problem for $\alpha \in \{0, 0.5, 1\}$ with decompositions \mathcal{G}_2 , $\hat{\mathcal{G}}_2$, and $\tilde{\mathcal{G}}_2$.

All decompositions work well. For $\alpha = 0$, where the objective function values are fully separable, the natural decomposition \mathcal{G}_2 is outperformed by $\tilde{\mathcal{G}}_2$, which is not a natural decomposition. On the other hand, for $\alpha = 1$, where the objective function values are not separable at all, the behavior is the opposite, that is, \mathcal{G}_2 outperforms $\tilde{\mathcal{G}}_2$. These results suggest that the natural decomposition is not necessarily the best one, aligned with the results in Section 5.3, but DASSO is robust to the choice.

6. Conclusion

In this paper, we proposed DASSO to tackle high-dimensional DSO problems. DASSO decomposes the prior distribution into an additive form, reducing the problem dimensionality to facilitate efficient posterior updates. This decomposition makes the search much more efficient and computationally possible by avoiding the CEI calculation for each feasible solution: we showed that it is sufficient to calculate the CEIs of only small fraction of solutions at a dice stage. Our numerical results revealed the effectiveness and efficiencies of DASSO on very large-scale problems: it can obtain rapid improvement on a problem with more than

a trillion feasible solutions within a couple of minutes of algorithm overhead. Furthermore, it empirically outperforms the state-of-the-art high-dimensional DSO algorithm pGMIA and BO algorithm Bounce.

Future research includes improving the performance of DASSO via an adaptive stopping condition for the slice stage, inventing a principled way to update the identity of the last group, and exploiting parallel computing: it is easy to parallelize replications for simulating the sample-best and best-CEI solutions and computation of the posterior distributions and Pareto frontiers for each group can also be parallelized. Attaining greater benefit from fast computational linear algebra is also relevant future work. Lastly, although we show that DASSO is capable of tackling a 100-dimensional problem, sampling decisions in the dice stage can be improved from the heuristic we demonstrated here.

Acknowledgments

This research was partially supported by National Science Foundation Grants CMMI-2206973, CMMI-2045400 and CMMI-2417616.

References

- Binois M, Wycoff N (2022) A survey on high-dimensional Gaussian process modeling with application to Bayesian optimization. *ACM Transactions on Evolutionary Learning and Optimization* 2(2):1–26.
- Chen Y, Ryzhov IO (2019) Complete expected improvement converges to an optimal budget allocation. *Advances in Applied Probability* 51(1):209–235.
- Deshwal A, Ament S, Balandat M, Bakshy E, Doppa JR, Eriksson D (2023) Bayesian optimization over high-dimensional combinatorial spaces via dictionary-based embeddings. *International Conference on Artificial Intelligence and Statistics*, 7021–7039 (PMLR).
- Djolonga J, Krause A, Cevher V (2013) High-dimensional Gaussian process bandits. *Advances in Neural Information Processing Systems*, 1025–1033.
- Durrande N, Ginsbourger D, Roustant O (2010) Additive kernels for Gaussian process modeling. *Annales de la Faculté de Sciences de Toulouse* 17.
- Eriksson D, Jankowiak M (2021) High-dimensional Bayesian optimization with sparse axis-aligned subspaces. *Uncertainty in Artificial Intelligence*, 493–503 (PMLR).
- Frazier P, Powell W, Dayanik S (2009) The knowledge-gradient policy for correlated normal beliefs. *INFORMS Journal on Computing* 21(4):599–613.
- Fu MC, et al. (2015) *Handbook of simulation optimization*, volume 216 (Springer).
- Gardner J, Guo C, Weinberger K, Garnett R, Grosse R (2017) Discovering and exploiting additive structure for Bayesian optimization. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 1311–1319 (PMLR).
- Garnett R, Osborne MA, Hennig P (2014) Active learning of linear embeddings for Gaussian processes. *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, 230–239.

- Ginsbourger D, Roustant O, Schuhmacher D, Durrande N, Lenz N (2016) On ANOVA decompositions of kernels and Gaussian random field paths. *Monte Carlo and Quasi-Monte Carlo Methods*, 315–330 (Springer).
- Hernández-Lobato JM, Hoffman MW, Ghahramani Z (2014) Predictive entropy search for efficient global optimization of black-box functions. *Advances in Neural Information Processing Systems* 27.
- Hoang TN, Hoang QM, Ouyang R, Low KH (2018) Decentralized high-dimensional Bayesian optimization with factor graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 3231–3238.
- Hutter F, Hoos HH, Leyton-Brown K (2011) Sequential model-based optimization for general algorithm configuration. *Learning and intelligent optimization: 5th international conference, LION 5, rome, Italy, January 17-21, 2011. selected papers* 5, 507–523 (Springer).
- Hvarfner C, Hutter F, Nardi L (2022) Joint entropy search for maximally-informed Bayesian optimization. *Advances in Neural Information Processing Systems* 35:11494–11506.
- Jian N, Freund D, Wiberg HM, Henderson SG (2016) Simulation optimization for a large-scale bike-sharing system. *Proceedings of the Winter Simulation Conference*, 602–613 (IEEE).
- Jones DR, Schonlau M, Welch WJ (1998) Efficient global optimization of expensive black-box functions. *Journal of Global Optimization* 13(4):455–492.
- Kandasamy K, Schneider J, Póczos B (2015) High dimensional Bayesian optimisation and bandits via additive models. *International Conference on Machine Learning*, 295–304 (PMLR).
- Li X, Song E (2020) Smart linear algebraic operations for efficient Gaussian Markov improvement algorithm. *Proceedings of the Winter Simulation Conference*, 2887–2898 (IEEE).
- Li X, Song E (2024) Projected Gaussian Markov improvement algorithm for high-dimensional discrete optimization via simulation. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* 34(3):1–29.
- Mathesen L, Chandrasekar KK, Li X, Pedrielli G, Candan KS (2019) Subspace communication driven search for high dimensional optimization. *Proceedings of the Winter Simulation Conference*, 3528–3539 (IEEE).
- Mes MR, Powell WB, Frazier PI (2011) Hierarchical knowledge gradient for sequential sampling. *Journal of Machine Learning Research* 12(10).
- Muehlenstaedt T, Roustant O, Carraro L, Kuhnt S (2012) Data-driven Kriging models based on FANOVA-decomposition. *Statistics and Computing* 22:723–738.
- Oh C, Tomczak J, Gavves E, Welling M (2019) Combinatorial Bayesian optimization using the graph cartesian product. *Advances in Neural Information Processing Systems* 32.
- Papenmeier L, Nardi L, Poloczek M (2023) Bounce: Reliable high-dimensional Bayesian optimization for combinatorial and mixed spaces. *Advances in Neural Information Processing Systems* 36:1764–1793.
- Quan N, Yin J, Ng SH, Lee LH (2013) Simulation optimization via kriging: a sequential search using expected improvement with computing budget constraints. *IIE Transactions* 45(7):763–780.

- Rolland P, Scarlett J, Bogunovic I, Cevher V (2018) High-dimensional Bayesian optimization via additive models with overlapping groups. *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, 298–307.
- Rue H, Held L (2005) *Gaussian Markov random fields: theory and applications* (Chapman and Hall/CRC).
- Salemi PL, Song E, Nelson BL, Staum J (2019) Gaussian Markov random fields for discrete optimization via simulation: Framework and algorithms. *Operations Research* 67(1):250–266.
- Saltelli A (2002) Making best use of model evaluations to compute sensitivity indices. *Computer physics communications* 145(2):280–297.
- Semelhago M, Nelson BL, Song E, Wächter A (2021) Rapid discrete optimization via simulation with Gaussian Markov random fields. *INFORMS Journal on Computing* 33(3):915–930.
- Semelhago M, Nelson BL, Wächter A, Song E (2017) Computational methods for optimization via simulation using gaussian Markov random fields. *Proceedings of the Winter Simulation Conference*, 2080–2091 (IEEE).
- Sun L, Hong LJ, Hu Z (2014) Balancing exploitation and exploration in discrete optimization via simulation through a gaussian process-based search. *Operations Research* 62(6):1416–1438.
- Tongarlak MH, Ankenman B, Nelson BL, Borne L, Wolfe K (2010) Using simulation early in the design of a fuel injector production line. *Interfaces* 40(2):105–117.
- Tripathy R, Bilonis I, Gonzalez M (2016) Gaussian processes with built-in dimensionality reduction: Applications to high-dimensional uncertainty propagation. *Journal of Computational Physics* 321:191–223.
- Ulaganathan S, Couckuyt I, Dhaene T, Degroote J, Laermans E (2016) High dimensional Kriging metamodeling utilising gradient information. *Applied Mathematical Modelling* 40(9-10):5256–5270.
- Wan X, Nguyen V, Ha H, Ru B, Lu C, Osborne MA (2021) Think global and act local: Bayesian optimisation over high-dimensional categorical and mixed search spaces. *International Conference on Machine Learning*, 10663–10674 (PMLR).
- Wang Z, Gehring C, Kohli P, Jegelka S (2018) Batched large-scale Bayesian optimization in high-dimensional spaces. *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, 745–754 (PMLR).
- Wang Z, Hutter F, Zoghi M, Matheson D, De Freitas N (2016) Bayesian optimization in a billion dimensions via random embeddings. *Journal of Artificial Intelligence Research* 55:361–387.
- Wang Z, Jegelka S (2017) Max-value entropy search for efficient Bayesian optimization. *International Conference on Machine Learning*, 3627–3635 (PMLR).
- Wang Z, Li C, Jegelka S, Kohli P (2017) Batched high-dimensional Bayesian optimization via structural kernel learning. *Proceedings of the 34th International Conference on Machine Learning*, 3656–3664 (PMLR).
- Xie J, Frazier PI, Chick SE (2016) Bayesian optimization via simulation with pairwise sampling and correlated prior beliefs. *Operations Research* 64(2):542–559.

Ziomek JK, Ammar HB (2023) Are random decompositions all we need in high dimensional Bayesian optimisation?
International Conference on Machine Learning, 43347–43368 (PMLR).

Electronic Companion to Dice and Slice Simulation Optimization

EC.1. Useful Lemmas

LEMMA EC.1 (The Woodbury Matrix Identity).

$$(A + CBC^\top)^{-1} = A^{-1} - A^{-1}C(B^{-1} + C^\top A^{-1}C)^{-1}C^\top A^{-1},$$

where A , B and C are conformable matrices.

LEMMA EC.2 (The Block Matrix Inversion).

$$Q^{-1} = \begin{pmatrix} A & B \\ B^\top & C \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} + A^{-1}B(Q/A)^{-1}B^\top A^{-1} & -A^{-1}B(Q/A)^{-1} \\ -(Q/A)^{-1}B^\top A^{-1} & (Q/A)^{-1} \end{pmatrix},$$

where A and C are invertible block matrices, B is conformable with them for partitioning and $Q/A = C - B^\top A^{-1}B$ is the Schur complement of A in Q .

EC.2. Proofs

Proof of Proposition 1 Suppose that $P\{\mathbb{X} = \mathbf{x}\} = \prod_{i=1}^d P\{\mathbb{X}_i = x_i\}$, $\forall \mathbf{x} = [x_1, x_2, \dots, x_d]^\top \in \mathcal{X}$, and $P\{\mathbb{X}_i = x\} = 1/n_i$, $\forall x \in \mathcal{X}_i$, for each component dimension \mathbb{X}_i of \mathbb{X} . Thus, $P\{\mathbb{X} = \mathbf{x}\} = 1/n$, $\forall \mathbf{x} \in \mathcal{X}$.

$$(a) E[y(\mathbb{X})] = \sum_{\mathbf{x} \in \mathcal{X}} P\{\mathbb{X} = \mathbf{x}\} y(\mathbf{x}) = \sum_{\mathbf{x} \in \mathcal{X}} y(\mathbf{x})/n = \bar{y}.$$

(b) The proof is by induction on u . Let $u = \{i\}$ for some $i \in d$. Then, replacing $\{i\}$ with i for notational convenience,

$$E[y_i(\mathbb{X})] = E\left[\frac{1}{n_{-i}} \sum_{\mathbf{x}_{-i} \in \mathcal{X}_{-i}} (y(\mathbb{X}) - \bar{y})\right] = \frac{1}{n_{-i}} \sum_{\mathbf{x}_{-i} \in \mathcal{X}_{-i}} (E[y(\mathbb{X})] - \bar{y}) = 0$$

since $E[y(\mathbb{X})] = \bar{y}$ from (a). Now, suppose that $E[y_v(\mathbb{X})] = 0$ for all $v \subset u$ with $v \neq \emptyset$. Then,

$$\begin{aligned} E[y_u(\mathbb{X})] &= E\left[\frac{1}{n_{-u}} \sum_{\mathbf{x}_{-u} \in \mathcal{X}_{-u}} \left(y(\mathbb{X}) - \sum_{v \subset u: v \neq \emptyset} y_v(\mathbb{X}) - y_\emptyset(\mathbb{X})\right)\right] \\ &= \frac{1}{n_{-u}} \sum_{\mathbf{x}_{-u} \in \mathcal{X}_{-u}} \left(E[y(\mathbb{X})] - \sum_{v \subset u: v \neq \emptyset} E[y_v(\mathbb{X})] - \bar{y}\right) \\ &= \frac{1}{n_{-u}} \sum_{\mathbf{x}_{-u} \in \mathcal{X}_{-u}} (\bar{y} - \bar{y}) = 0, \end{aligned}$$

completing the induction.

(c) Consider subsets $u, v \subseteq d$ with $u \neq v$. Without loss of generality, suppose $v \subset u$. Using the law of total expectation, $E[y_u(\mathbb{X})y_v(\mathbb{X})] = E[E[y_u(\mathbb{X})y_v(\mathbb{X}) | \mathbb{X}_v]] = E[E[y_u(\mathbb{X}) | \mathbb{X}_v]y_v(\mathbb{X})]$ since $y_v(\mathbb{X})$ depends on \mathbb{X} only through \mathbb{X}_v . Since $E[y_u(\mathbb{X}) | \mathbb{X}_v] = 0$ from (b), we have $E[y_u(\mathbb{X})y_v(\mathbb{X})] = 0$.

(d) $\text{Var}[y(\mathbb{X})] = \text{Var}\left[\sum_{u \subseteq d} y_u(\mathbb{X})\right] = \sum_{u \subseteq d} \text{Var}[y_u(\mathbb{X})]$ since $y_u(\mathbb{X})$ and $y_v(\mathbb{X})$ are uncorrelated for $u \neq v$ from (c). \square

Proof of Theorem 1 To establish the conditional distribution of $(\mathbb{Y}_{\mathcal{U}}^{(\rho)}, \mathbb{Y}_{\mathcal{D}}^{(\rho)})$ given observed $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$, for each $\rho \in \mathcal{G}$, we follow similar steps as in Salemi et al. (2019). We first derive the joint distribution of $(\mathbb{Y}_{\mathcal{U}}^{(\rho)}, \mathbb{Y}_{\mathcal{D}}^{(\rho)}, \mathbb{Y}_{\mathcal{D}}^{\epsilon})$ and then apply Lemma 2.1 of Rue and Held (2005).

Since $\mathbb{Y}^{(\rho)}$ is a GMRF and independent of the intrinsic noise, $\mathbb{Y}_{\mathcal{U}}^{(\rho)}$ and $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ are conditionally independent, given $\mathbb{Y}_{\mathcal{D}}^{(\rho)}$. From Theorem 2.5 in Rue and Held (2005), the conditional distribution of $\mathbb{Y}_{\mathcal{U}}^{(\rho)}$ given $\mathbb{Y}_{\mathcal{D}}^{(\rho)} = \mathbf{y}_{\mathcal{D}}^{(\rho)}$ is

$$\mathcal{N}\left(-[\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1}\mathbf{Q}_{\mathcal{UD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}, [\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1}\right).$$

From the definition of $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ and the independence assumption, the conditional distribution of $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ given $\mathbb{Y}_{\mathcal{D}}^{(\rho)} = \mathbf{y}_{\mathcal{D}}^{(\rho)}$ is

$$\mathcal{N}\left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{T}_{\mathcal{DD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}, [\mathbf{E}^{(\rho)}]^{-1}\right).$$

Recall $\mathbf{E}^{(\rho)} = \left[\boldsymbol{\Sigma}_{\mathcal{DD}} - \mathbf{T}_{\mathcal{DD}}^{(\rho)}\boldsymbol{\Sigma}_{\mathcal{DD}}^{(\rho)}[\mathbf{T}_{\mathcal{DD}}^{(\rho)}]^{\top} + \boldsymbol{\Sigma}^{\epsilon}\right]^{-1}$. Moreover, using standard results for the inverse of a partitioned matrix,

$$\mathbb{Y}_{\mathcal{D}}^{(\rho)} \sim \mathcal{N}\left(\vec{\mathbf{0}}_{|\mathcal{D}^{(\rho)}|}, \boldsymbol{\Sigma}_{\mathcal{DD}}^{(\rho)}\right).$$

Then, the joint distribution $f(\mathbf{y}_{\mathcal{U}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{\epsilon})$ satisfies

$$\begin{aligned} f(\mathbf{y}_{\mathcal{U}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{\epsilon}) &\propto \exp\left\{-\frac{1}{2}\left(\mathbf{y}_{\mathcal{U}}^{(\rho)} + [\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1}\mathbf{Q}_{\mathcal{UD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}\right)^{\top}\mathbf{Q}_{\mathcal{UU}}^{(\rho)}\left(\mathbf{y}_{\mathcal{U}}^{(\rho)} + [\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1}\mathbf{Q}_{\mathcal{UD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}\right)\right\} \\ &\quad \times \exp\left\{-\frac{1}{2}[\mathbf{y}_{\mathcal{D}}^{(\rho)}]^{\top}\left(\mathbf{Q}_{\mathcal{DD}}^{(\rho)} - \mathbf{Q}_{\mathcal{DU}}^{(\rho)}[\mathbf{Q}_{\mathcal{UU}}^{(\rho)}]^{-1}\mathbf{Q}_{\mathcal{UD}}^{(\rho)}\right)\mathbf{y}_{\mathcal{D}}^{(\rho)}\right\} \\ &\quad \times \exp\left\{-\frac{1}{2}\left[\mathbf{y}_{\mathcal{D}}^{\epsilon} - \left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{T}_{\mathcal{DD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}\right)\right]^{\top}\mathbf{E}^{(\rho)}\left[\mathbf{y}_{\mathcal{D}}^{\epsilon} - \left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{T}_{\mathcal{DD}}^{(\rho)}\mathbf{y}_{\mathcal{D}}^{(\rho)}\right)\right]\right\}. \end{aligned}$$

A significant quantity of matrix algebra leads to

$$\begin{aligned} f(\mathbf{y}_{\mathcal{U}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{(\rho)}, \mathbf{y}_{\mathcal{D}}^{\epsilon}) &\propto \exp\left\{-\frac{1}{2}\begin{pmatrix}\mathbf{y}_{\mathcal{U}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}^{\top} \begin{pmatrix} \mathbf{Q}_{\mathcal{UU}}^{(\rho)} & \mathbf{Q}_{\mathcal{UD}}^{(\rho)} & \mathbf{0}_{|\mathcal{U}^{(\rho)}| \times |\mathcal{D}|} \\ \mathbf{Q}_{\mathcal{DU}}^{(\rho)} & \mathbf{Q}_{\mathcal{DD}}^{(\rho)} + [\mathbf{T}_{\mathcal{DD}}^{(\rho)}]^{\top}\mathbf{E}^{(\rho)}\mathbf{T}_{\mathcal{DD}}^{(\rho)} & -[\mathbf{T}_{\mathcal{DD}}^{(\rho)}]^{\top}\mathbf{E}^{(\rho)} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}^{(\rho)}|} & -\mathbf{E}^{(\rho)}\mathbf{T}_{\mathcal{DD}}^{(\rho)} & \mathbf{E}^{(\rho)} \end{pmatrix} \begin{pmatrix}\mathbf{y}_{\mathcal{U}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}\right. \\ &\quad \left. + \begin{pmatrix}\mathbf{y}_{\mathcal{U}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{(\rho)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}^{\top} \begin{pmatrix} \vec{\mathbf{0}}_{|\mathcal{U}^{(\rho)}|} \\ -[\mathbf{T}_{\mathcal{DD}}^{(\rho)}]^{\top}\mathbf{E}^{(\rho)}\boldsymbol{\mu}_{\mathcal{D}} \\ \mathbf{E}^{(\rho)}\boldsymbol{\mu}_{\mathcal{D}} \end{pmatrix}\right\}. \end{aligned}$$

Finally, following Lemma 2.1 in Rue and Held (2005), the conditional distribution of $\mathbb{Y}^{(\rho)}$ given $\mathbb{Y}_{\mathcal{D}}^{\epsilon} = \mathbf{y}_{\mathcal{D}}^{\epsilon}$ is as in the theorem's statement. \square

Proof of Theorem 2 To establish the conditional distribution of $(\mathbb{W}_{\mathcal{U}}^{(g)}, \mathbb{W}_{\mathcal{D}}^{(g)})$ given the observed $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$, we follow similar steps as in Theorem 1. We first derive the joint distribution of $(\mathbb{W}_{\mathcal{U}}^{(g)}, \mathbb{W}_{\mathcal{D}}^{(g)}, \mathbb{Y}_{\mathcal{D}}^{\epsilon})$ and then apply Lemma 2.1 of Rue and Held (2005).

Since $\mathbb{W}^{(g)}$ is independent of the intrinsic noise, $\mathbb{W}_{\mathcal{U}}^{(g)}$ and $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ are conditionally independent, given $\mathbb{W}_{\mathcal{D}}^{(g)}$. The conditional distribution of $\mathbb{W}_{\mathcal{U}}^{(g)}$ given $\mathbb{W}_{\mathcal{D}}^{(g)} = \mathbf{w}_{\mathcal{D}}^{(g)}$ is

$$\mathcal{N}\left(\vec{\mathbf{0}}_{|\mathcal{U}|}, \sigma_g^2 \mathbf{I}_{|\mathcal{U}|}\right).$$

From the definition of $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ and the independence assumption, the conditional distribution of $\mathbb{Y}_{\mathcal{D}}^{\epsilon}$ given $\mathbb{W}_{\mathcal{D}}^{(g)} = \mathbf{w}_{\mathcal{D}}^{(g)}$ is

$$\mathcal{N}\left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{w}_{\mathcal{D}}^{(g)}, [\mathbf{E}^{(g)}]^{-1}\right).$$

Recall that $\mathbf{E}^{(g)} = [\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} - \sigma_g^2 \mathbf{I}_{|\mathcal{D}|} + \boldsymbol{\Sigma}^{\epsilon}]^{-1}$. Moreover,

$$\mathbb{W}_{\mathcal{D}}^{(g)} \sim \mathcal{N}\left(\vec{\mathbf{0}}_{|\mathcal{D}|}, \sigma_g^2 \mathbf{I}_{|\mathcal{D}|}\right).$$

Then, the joint distribution $f(\mathbf{w}_{\mathcal{U}}^{(g)}, \mathbf{w}_{\mathcal{D}}^{(g)}, \mathbf{y}_{\mathcal{D}}^{\epsilon})$ satisfies

$$\begin{aligned} f(\mathbf{w}_{\mathcal{U}}^{(g)}, \mathbf{w}_{\mathcal{D}}^{(g)}, \mathbf{y}_{\mathcal{D}}^{\epsilon}) &\propto \exp\left\{-\frac{1}{2} \frac{1}{\sigma_g^2} [\mathbf{w}_{\mathcal{U}}^{(g)}]^{\top} \mathbf{I}_{|\mathcal{U}|} \mathbf{w}_{\mathcal{U}}^{(g)}\right\} \times \exp\left\{-\frac{1}{2} \frac{1}{\sigma_g^2} [\mathbf{w}_{\mathcal{D}}^{(g)}]^{\top} \mathbf{I}_{|\mathcal{D}|} \mathbf{w}_{\mathcal{D}}^{(g)}\right\} \\ &\times \exp\left\{-\frac{1}{2} \left[\mathbf{y}_{\mathcal{D}}^{\epsilon} - \left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{w}_{\mathcal{D}}^{(g)}\right)\right]^{\top} \mathbf{E}^{(g)} \left[\mathbf{y}_{\mathcal{D}}^{\epsilon} - \left(\boldsymbol{\mu}_{\mathcal{D}} + \mathbf{w}_{\mathcal{D}}^{(g)}\right)\right]\right\}. \end{aligned}$$

A significant quantity of matrix algebra leads to

$$\begin{aligned} f(\mathbf{w}_{\mathcal{U}}^{(g)}, \mathbf{w}_{\mathcal{D}}^{(g)}, \mathbf{y}_{\mathcal{D}}^{\epsilon}) &\propto \exp\left\{-\frac{1}{2} \begin{pmatrix} \mathbf{w}_{\mathcal{U}}^{(g)} \\ \mathbf{w}_{\mathcal{D}}^{(g)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}^{\top} \begin{pmatrix} \frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & \frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{D}|} + \mathbf{E}^{(g)} & -\mathbf{E}^{(g)} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & -\mathbf{E}^{(g)} & \mathbf{E}^{(g)} \end{pmatrix} \begin{pmatrix} \mathbf{w}_{\mathcal{U}}^{(g)} \\ \mathbf{w}_{\mathcal{D}}^{(g)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}\right. \\ &\quad \left. + \begin{pmatrix} \mathbf{w}_{\mathcal{U}}^{(g)} \\ \mathbf{w}_{\mathcal{D}}^{(g)} \\ \mathbf{y}_{\mathcal{D}}^{\epsilon} \end{pmatrix}^{\top} \begin{pmatrix} \vec{\mathbf{0}}_{|\mathcal{U}|} \\ -\mathbf{E}^{(g)} \boldsymbol{\mu}_{\mathcal{D}} \\ \mathbf{E}^{(g)} \boldsymbol{\mu}_{\mathcal{D}} \end{pmatrix}\right\}. \end{aligned}$$

Finally, following Lemma 2.1 in Rue and Held (2005), the conditional distribution of $\mathbb{W}^{(g)}$ given $\mathbb{Y}_{\mathcal{D}}^{\epsilon} = \mathbf{y}_{\mathcal{D}}^{\epsilon}$ is as in the theorem's statement. \square

Proof of Proposition 2 For $\mathbf{x}, \check{\mathbf{x}} \in \mathcal{U}$, notice that $m^{(g)}(\mathbf{x}) = m^{(g)}(\check{\mathbf{x}}) = 0$ as well as $v^{(g)}(\mathbf{x}) = v^{(g)}(\check{\mathbf{x}}) = \sigma_g^2$ and $c^{(g)}(\check{\mathbf{x}}, \mathbf{x}) = c^{(g)}(\tilde{\mathbf{x}}, \check{\mathbf{x}}) = 0$ from the block-diagonal structure of $\bar{\boldsymbol{\Sigma}}^{(g)}$.

(R1) Suppose that $\mathbf{x}^{(\rho)} = \check{\mathbf{x}}^{(\rho)}$ for all $\rho \in \mathcal{G}^{(-g)}$. Then, $m(\mathbf{x}) = m(\check{\mathbf{x}})$ and $v(\tilde{\mathbf{x}}, \mathbf{x}) = v(\tilde{\mathbf{x}}, \check{\mathbf{x}})$, and thus $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) = \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$.

(R2) Suppose for some group $\rho \in \mathcal{G}^{(-g)}$ that the following conditions hold:

$$(C2.1) \quad \mathbf{x}^{(\varrho)} = \check{\mathbf{x}}^{(\varrho)} \text{ for all } \varrho \in \mathcal{G} \setminus \{\rho, g\},$$

$$(C2.2) \quad m^{(\rho)}(\mathbf{x}^{(\rho)}) \geq m^{(\rho)}(\check{\mathbf{x}}^{(\rho)}), \text{ and}$$

$$(C2.3) \quad v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \mathbf{x}^{(\rho)}) \leq v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \check{\mathbf{x}}^{(\rho)}), \text{ where } v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \cdot) = v^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}) + v^{(\rho)}(\cdot) - 2c^{(\rho)}(\tilde{\mathbf{x}}^{(\rho)}, \cdot).$$

Conditions (C2.1) and (C2.2) together imply that $m(\mathbf{x}) \geq m(\check{\mathbf{x}})$. Conditions (C2.1) and (C2.3) together imply that $v(\tilde{\mathbf{x}}, \mathbf{x}) \leq v(\tilde{\mathbf{x}}, \check{\mathbf{x}})$. Then, $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$ from the fact that $\text{CEI}(\tilde{\mathbf{x}}, \cdot)$ is decreasing in $m(\cdot)$ and increasing in $v(\tilde{\mathbf{x}}, \cdot)$ as

$$\frac{\partial \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})}{\partial m(\mathbf{x})} = -\Phi\left(\frac{m(\tilde{\mathbf{x}}) - m(\mathbf{x})}{\sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})}}\right) < 0$$

and

$$\frac{\partial \text{CEI}(\tilde{\mathbf{x}}, \mathbf{x})}{\partial v(\tilde{\mathbf{x}}, \mathbf{x})} = \frac{1}{2\sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})}} \phi\left(\frac{m(\tilde{\mathbf{x}}) - m(\mathbf{x})}{\sqrt{v(\tilde{\mathbf{x}}, \mathbf{x})}}\right) > 0.$$

(R3) For $\mathbf{x}, \check{\mathbf{x}} \in \mathcal{D}$, suppose that $\mathbf{x}^{(\rho)} = \check{\mathbf{x}}^{(\rho)}$ for all $\rho \in \mathcal{G}^{(-g)}$, $m^{(g)}(\mathbf{x}) \geq m^{(g)}(\check{\mathbf{x}})$, and $v^{(g)}(\tilde{\mathbf{x}}, \mathbf{x}) \leq v^{(g)}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$. These imply that $m(\mathbf{x}) \geq m(\check{\mathbf{x}})$ and $v(\tilde{\mathbf{x}}, \mathbf{x}) \leq v(\tilde{\mathbf{x}}, \check{\mathbf{x}})$, leading to $\text{CEI}(\tilde{\mathbf{x}}, \mathbf{x}) \leq \text{CEI}(\tilde{\mathbf{x}}, \check{\mathbf{x}})$. \square

EC.3. Global Convergence of DASSO

We establish the global convergence of DASSO in Theorem EC.1 under the very mild conditions in Assumption EC.1 and a small tweak to the algorithm in Assumption EC.2.

ASSUMPTION EC.1. $y(\mathbf{x}) > -\infty$ and $0 < \text{Var}[Y(\mathbf{x})] < \infty$, for all $\mathbf{x} \in \mathcal{X}$.

ASSUMPTION EC.2. *Each slice stage is performed for a finite number of iterations. Further, the prior precision matrix of $\mathbb{Y}^{(\rho)}$ is adjusted to be $\alpha_t \mathbf{Q}^{(\rho)}$ for $\rho \in \mathcal{G}^{(-g)}$, where $\alpha_t \rightarrow \infty$ as $t \rightarrow \infty$, and t represents a counter that advances with the number of iterations completed in the algorithm.*

The first part of Assumption EC.2 is to avoid a failure to explore solutions with different values of the first $g - 1$ components by getting stuck in a slice stage; this is easily enforced by the slice-stage stopping criterion. The second part is to eventually eliminate the effect of the first $g - 1$ components as more simulation outputs are obtained; this can be achieved through multiplying the precision matrix by α_t , which diverges to infinity as the algorithm progresses (although this is not the only possible adjustment). Without this assumption, the uncertainty caused by the model structure (as we observe the overall objective function values, rather than individual values for each group) does not disappear, making the algorithm simulate only a subset of solutions eventually without exploring the others. Theorem EC.1 guarantees the asymptotic convergence of DASSO to the global optimum with probability one regardless of the choice of decomposition and how the identity of the last group is updated.

THEOREM EC.1. *Under Assumptions EC.1 and EC.2, the DASSO algorithm without a stopping condition converges to the global optimum with probability one as the number of iterations goes to infinity.*

We confess that Theorem EC.1 is strictly of academic interest, since we do not expect to approach anything like convergence in the class of problems we consider.

EC.3.1. Proof

Recall that $y(\mathbf{x})$ denote the objective function value at solution $\mathbf{x} \in \mathcal{X}$. Letting $\mathcal{X}_{\min} = \{\mathbf{x} \in \mathcal{X} : y(\mathbf{x}) = y_{\min}\}$ denote the set of optimal solutions, where $y_{\min} = \min_{\mathbf{x} \in \mathcal{X}} y(\mathbf{x})$ is the optimal objective function value, the aim is to find an optimal solution $\mathbf{x} \in \mathcal{X}_{\min}$. We allow $|\mathcal{X}_{\min}| \geq 1$. To prove the theorem, we will show that each solution will be simulated infinitely often with probability one as the number of iterations goes to infinity. Then, the main result will follow by the strong law of large numbers.

We fix a sample path but suppress it in the notation. Let t denote the dice-stage iteration of the algorithm. Without loss of generality, we set the number of slice-stage iterations for each dice stage to 1 and consider t as the iteration of the algorithm as well. Moreover, without loss of generality, we assume that the identity of the last group g is fixed. Or equivalently, since at least one group must be chosen as the last group infinitely often, we consider a subsequence of iterations on which the identity of the last group g is one

such a group chosen infinitely often. The statistics used to compute the posterior distribution and the CEI values at iteration t are conditional on the simulation outputs obtained up to, but not including, iteration t . To ease notation, we write these statistics with subscript t ; for example, \tilde{x}_t denotes the sample-best solution at iteration t .

Let $\mathcal{A} \subset \mathcal{D}$ denote the set of the feasible solutions that will be simulated infinitely often in this sample path, i.e., $\mathcal{A} = \{x \in \mathcal{X} : r_t(x) \rightarrow \infty \text{ as } t \rightarrow \infty\}$. Since at least one solution must be simulated infinitely often as $t \rightarrow \infty$, \mathcal{A} is non-empty. Assume that $\mathcal{A}^c = \mathcal{X} \setminus \mathcal{A}$ is also non-empty; notice that $\mathcal{U} \subset \mathcal{A}^c$ since $\mathcal{A} \subset \mathcal{D}$. Therefore, T , the last iteration at which a solution in \mathcal{A}^c is simulated, is finite. Since the sample-best solution \tilde{x}_t is simulated at each dice stage, $\tilde{x}_t \in \mathcal{A}$ for all $t > T$. Although \mathcal{D}_t and \mathcal{U}_t change as t increases, they remain the same for $t > T$. Therefore, we suppress the dependency of \mathcal{D} and \mathcal{U} on t for notational simplicity.

Each slice stage employs the GMIA algorithm of Salemi et al. (2019). Since GMIA without a stopping condition is proven to simulate each solution infinitely often with probability one, if the first $g - 1$ components are fixed to z in the slice stage infinitely often, then all solutions in $\mathcal{X}_z = \{x \in \mathcal{X} : x^{(-g)} = z\}$ are simulated infinitely often. That is, if $x \in \mathcal{A} \setminus \{\tilde{x}_t\}$ for $t > T$, then $\check{x} \in \mathcal{A}$ for all $\check{x} \in \mathcal{X}$ such that $\check{x}^{(-g)} = x^{(-g)}$ and therefore simulated infinitely often. Alternatively, if $x \in \mathcal{A}^c$, then $\check{x} \in \mathcal{A}^c$ for all $\check{x} \in \mathcal{X} \setminus \{\tilde{x}_t\}$ such that $\check{x}^{(-g)} = x^{(-g)}$ for $t > T$. Therefore, none of the solutions in \mathcal{A}^c is chosen in a dice stage after iteration T , i.e., $\arg \max_{x \in \mathcal{X} \setminus \{\tilde{x}_t\}} \text{CEI}_t(\tilde{x}_t, x) \notin \mathcal{A}^c$, for $t > T$.

We need the following the asymptotic results for the conditional mean and conditional variance; the proof is at the end of this section.

LEMMA EC.3. *For $x \in \mathcal{A}$, we have $\lim_{t \rightarrow \infty} v_t(\tilde{x}_t, x) = 0$ and $\lim_{t \rightarrow \infty} m_t(x) = y(x)$. Assuming $\mathcal{A}^c \neq \emptyset$, for $x \in \mathcal{A}^c$, we have $\liminf_{t \rightarrow \infty} v_t(\tilde{x}_t, x) > 0$.*

Let $y_{\min}^{\mathcal{A}} = \min_{x \in \mathcal{A}} y(x)$, and $\mathcal{A}_{\min} = \{x \in \mathcal{A} : y(x) = y_{\min}^{\mathcal{A}}\}$, the set of optimal solutions in \mathcal{A} . Using Lemma EC.3, we show that $\lim_{t \rightarrow \infty} m_t(\tilde{x}_t) = y_{\min}^{\mathcal{A}}$: If $y(x) = y_{\min}^{\mathcal{A}}$ for all $x \in \mathcal{A}$, then it immediately follows from Lemma EC.3 that $m_t(\tilde{x}_t) \rightarrow y_{\min}^{\mathcal{A}}$ as $\tilde{x}_t \in \mathcal{A}$ for all $t > T$. If $y(x) \neq y_{\min}^{\mathcal{A}}$ for some $x \in \mathcal{A}$, otherwise, let ε be a constant such that $0 < \varepsilon < \min_{\mathcal{A} \setminus \mathcal{A}_{\min}} y(x) - y_{\min}^{\mathcal{A}}$. We can find such a constant as there are only a finite number of feasible solutions in \mathcal{A} . Since $r_t(x) \rightarrow \infty$ as $t \rightarrow \infty$ for all $x \in \mathcal{A}$, i.e., every solution in \mathcal{A} is simulated infinitely often, there exists some $T^* > T$ such that for all $x \in \mathcal{A}$, we have $|\bar{Y}_{\mathcal{D},t}(x) - y(x)| < \varepsilon/2$ for $t \geq T^*$ by the strong law of large numbers. Then, by the definition of $\tilde{x}_t = \arg \min_{x \in \mathcal{A}} \bar{Y}_{\mathcal{D},t}(x)$ for all $t > T$ and the choice of ε , we have $y(\tilde{x}_t) = y_{\min}^{\mathcal{A}}$ for all $t \geq T^*$. Therefore, together with Lemma EC.3, we have $m_t(\tilde{x}_t) \rightarrow y_{\min}^{\mathcal{A}}$.

Recall that $\text{CEI}_t(\tilde{x}_t, x)$ in the dice stage for $x \in \mathcal{X} \setminus \{\tilde{x}_t\}$ is

$$\text{CEI}_t(\tilde{x}_t, x) = (m_t(\tilde{x}_t) - m_t(x)) \Phi \left(\frac{m_t(\tilde{x}_t) - m_t(x)}{\sqrt{v_t(\tilde{x}_t, x)}} \right) + \sqrt{v_t(\tilde{x}_t, x)} \phi \left(\frac{m_t(\tilde{x}_t) - m_t(x)}{\sqrt{v_t(\tilde{x}_t, x)}} \right).$$

Because $\text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x})$ is a non-negative increasing function in $m_t(\tilde{\mathbf{x}}_t) - m_t(\mathbf{x})$, and $\text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}) \rightarrow 0$ as $m_t(\tilde{\mathbf{x}}_t) - m_t(\mathbf{x}) \rightarrow -\infty$ provided $\lim_{t \rightarrow \infty} v_t(\tilde{\mathbf{x}}_t, \mathbf{x}) > 0$, Lemma EC.3 guarantees $\liminf_{t \rightarrow \infty} \text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}) > 0$ for $\mathbf{x} \in \mathcal{A}^c$. On the other hand, for $\mathbf{x} \in \mathcal{A}$, we have $\text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}) \rightarrow 0$ as $t \rightarrow \infty$ by Lemma EC.3 since $\lim_{t \rightarrow \infty} v_t(\tilde{\mathbf{x}}_t, \mathbf{x}) = 0$ and $\lim_{t \rightarrow \infty} (m_t(\tilde{\mathbf{x}}_t) - m_t(\mathbf{x})) \leq 0$, which follows from the fact that $\lim_{t \rightarrow \infty} m_t(\tilde{\mathbf{x}}_t) = y_{\min}^{\mathcal{A}} = \min_{\mathbf{x} \in \mathcal{A}} y(\mathbf{x})$. Hence, we will eventually have

$$\min_{\mathbf{x} \in \mathcal{A}^c} \text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}) > \max_{\mathbf{x} \in \mathcal{A}} \text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}),$$

and therefore, $\arg \max_{\mathbf{x} \in \mathcal{X} \setminus \{\tilde{\mathbf{x}}_t\}} \text{CEI}_t(\tilde{\mathbf{x}}_t, \mathbf{x}) \in \mathcal{A}^c$ and DASSO would choose a solution in \mathcal{A} .

However, this contradicts the existence of T . Hence, \mathcal{A}^c is empty, and thus $\mathcal{X} = \mathcal{A}$ and $y_{\min}^{\mathcal{A}} = y_{\min}$. Therefore, together with the fact that $m_t(\tilde{\mathbf{x}}_t) \rightarrow y_{\min}^{\mathcal{A}}$, we conclude that $m_t(\tilde{\mathbf{x}}_t) \rightarrow y_{\min}$, as $t \rightarrow \infty$. Since this will occur on almost all sample paths, the convergence is with probability 1.

Proof of Lemma EC.3 Under Assumption EC.2, $\mathbf{Q}_t^{(\rho)} = \alpha_t \mathbf{Q}^{(\rho)}$ for each group $\rho \in \mathcal{G}^{(-g)}$. Then, since $\mathbf{E}_t^{(\rho)}$ is non-negative and $\alpha_t \rightarrow \infty$ as $t \rightarrow \infty$,

$$\bar{\Sigma}_t^{(\rho)} = [\bar{\mathbf{Q}}_t^{(\rho)}]^{-1} \rightarrow \mathbf{0}_{n^{(\rho)} \times n^{(\rho)}}, \text{ and thus } \mathbf{m}_t^{(\rho)} \rightarrow \vec{\mathbf{0}}_{n^{(\rho)}}. \quad (\text{EC.1})$$

Moreover, $\Sigma_{\mathcal{D}\mathcal{D},t}^{(\rho)} = \frac{1}{\alpha_t} \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)} \rightarrow \mathbf{0}_{|\mathcal{D}| \times |\mathcal{D}|}$, and thus $\Sigma_{\mathcal{D}\mathcal{D},t} \rightarrow \sigma_g^2 \mathbf{I}_{|\mathcal{D}|}$, as $t \rightarrow \infty$. Equation (EC.1) implies that $\lim_{t \rightarrow \infty} m_t(\mathbf{x}) = \lim_{t \rightarrow \infty} m_t^{(g)}(\mathbf{x})$ and $\liminf_{t \rightarrow \infty} v_t(\tilde{\mathbf{x}}_t, \mathbf{x}) = \liminf_{t \rightarrow \infty} [v_t^{(g)}(\tilde{\mathbf{x}}_t) + v_t^{(g)}(\mathbf{x}) - 2c_t^{(g)}(\tilde{\mathbf{x}}_t, \mathbf{x})]$, for $\mathbf{x} \in \mathcal{X}$.

Let $\mathcal{B} = \mathcal{D} \setminus \mathcal{A}$; notice that $\mathcal{X} = \mathcal{U} \cup \mathcal{A} \cup \mathcal{B}$. For $t > T$, reordering the elements of $\mathbf{E}_t^{(g)}$ and Σ_t^ϵ , we can partition them as

$$\mathbf{E}_t^{(g)} = \begin{pmatrix} \mathbf{E}_{\mathcal{B}\mathcal{B},t}^{(g)} & \mathbf{E}_{\mathcal{B}\mathcal{A},t}^{(g)} \\ \mathbf{E}_{\mathcal{A}\mathcal{B},t}^{(g)} & \mathbf{E}_{\mathcal{A}\mathcal{A},t}^{(g)} \end{pmatrix} \text{ and } \Sigma_t^\epsilon = \begin{pmatrix} \Sigma_{\mathcal{B}\mathcal{B},t}^\epsilon & \mathbf{0}_{|\mathcal{B}| \times |\mathcal{A}|} \\ \mathbf{0}_{|\mathcal{A}| \times |\mathcal{B}|} & \Sigma_{\mathcal{A}\mathcal{A},t}^\epsilon \end{pmatrix},$$

respectively, where $\Sigma_{\mathcal{B}\mathcal{B},t}^\epsilon$ and $\Sigma_{\mathcal{A}\mathcal{A},t}^\epsilon$ are diagonal matrices whose diagonal element corresponding to solution \mathbf{x} is $S_t^2(\mathbf{x})/r_t(\mathbf{x})$. Notice that the diagonal element corresponding to \mathbf{x} of $\Sigma_{\mathcal{B}\mathcal{B},t}^\epsilon$ is $S_T^2(\mathbf{x})/r_T(\mathbf{x})$ for $t > T$ whereas that of $\Sigma_{\mathcal{A}\mathcal{A},t}^\epsilon$ converges to 0 as $t \rightarrow \infty$. Recall that $\mathbf{E}_t^{(g)} = [\Sigma_{\mathcal{D}\mathcal{D},t} - \sigma_g^2 \mathbf{I}_{|\mathcal{D}|} + \Sigma_t^\epsilon]^{-1}$, where $\Sigma_{\mathcal{D}\mathcal{D},t} \rightarrow \sigma_g^2 \mathbf{I}_{|\mathcal{D}|}$ as $t \rightarrow \infty$. Therefore, the diagonal elements of $\mathbf{E}_{\mathcal{A}\mathcal{A},t}^{(g)}$ diverge to ∞ as $t \rightarrow \infty$.

Letting

$$\mathbf{X}_t = \begin{pmatrix} \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c, t} & \mathbf{X}_{\mathcal{A}^c \mathcal{A}, t} \\ \mathbf{X}_{\mathcal{A} \mathcal{A}^c, t}^\top & \mathbf{E}_{\mathcal{A}\mathcal{A}, t}^{(g)} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{B}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{A}|} \\ \mathbf{0}_{|\mathcal{B}| \times |\mathcal{U}|} & \mathbf{E}_{\mathcal{B}\mathcal{B}, t}^{(g)} & \mathbf{E}_{\mathcal{B}\mathcal{A}, t}^{(g)} \\ \mathbf{0}_{|\mathcal{A}| \times |\mathcal{U}|} & \mathbf{E}_{\mathcal{A}\mathcal{B}, t}^{(g)} & \mathbf{E}_{\mathcal{A}\mathcal{A}, t}^{(g)} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & \mathbf{E}_t^{(g)} \end{pmatrix}$$

with

$$\mathbf{X}_{\mathcal{A}^c \mathcal{A}^c, t} = \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{B}|} \\ \mathbf{0}_{|\mathcal{B}| \times |\mathcal{U}|} & \mathbf{E}_{\mathcal{B}\mathcal{B}, t}^{(g)} \end{pmatrix}, \text{ and } \mathbf{X}_{\mathcal{A}^c \mathcal{A}, t} = \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{A}|} \\ \mathbf{E}_{\mathcal{B}\mathcal{A}, t}^{(g)} \end{pmatrix},$$

we have $\bar{\mathbf{Q}}_t^{(g)} = (\sigma_g^2)^{-1} \mathbf{I}_n + \mathbf{X}_t$. Let $\mathbf{A}_t = (\sigma_g^2)^{-1} \mathbf{I}_{|\mathcal{A}^c|} + \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c, t}$, $\mathbf{B}_t = \mathbf{X}_{\mathcal{A}^c \mathcal{A}, t}$, and $\mathbf{C}_t = (\sigma_g^2)^{-1} \mathbf{I}_{|\mathcal{A}|} + \mathbf{E}_{\mathcal{A} \mathcal{A}, t}^{(g)}$ denote the corresponding components of $\bar{\mathbf{Q}}_t^{(g)}$. Using block matrix inversion, see Lemma EC.2 in Section EC.1,

$$\bar{\Sigma}_t^{(g)} = [\bar{\mathbf{Q}}_t^{(g)}]^{-1} = \begin{pmatrix} \mathbf{A}_t & \mathbf{B}_t \\ \mathbf{B}_t^\top & \mathbf{C}_t \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}_t^{-1} + \mathbf{A}_t^{-1} \mathbf{B}_t (\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \mathbf{B}_t^\top \mathbf{A}_t^{-1} & -\mathbf{A}_t^{-1} \mathbf{B}_t (\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \\ -(\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \mathbf{B}_t^\top \mathbf{A}_t^{-1} & (\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \end{pmatrix},$$

where $\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t = \mathbf{C}_t - \mathbf{B}_t^\top \mathbf{A}_t^{-1} \mathbf{B}_t$ is the Schur complement of \mathbf{A}_t in $\bar{\mathbf{Q}}_t^{(g)}$. Since the diagonal elements of $\mathbf{E}_{\mathcal{A} \mathcal{A}, t}^{(g)}$ diverges to ∞ as $t \rightarrow \infty$, so do those of \mathbf{C}_t . Therefore, we have $(\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \rightarrow \mathbf{0}_{|\mathcal{A}| \times |\mathcal{A}|}$, and thus

$$\bar{\Sigma}_t^{(g)} \rightarrow \begin{pmatrix} \left(\left(\frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{A}^c|} + \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} \right)^{-1} & \mathbf{0}_{|\mathcal{A}^c| \times |\mathcal{A}|} \\ \mathbf{0}_{|\mathcal{A}| \times |\mathcal{A}^c|} & \mathbf{0}_{|\mathcal{A}| \times |\mathcal{A}|} \end{pmatrix}, \quad (\text{EC.2})$$

where $\mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} = \lim_{t \rightarrow \infty} \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c, t}$. Moreover,

$$(\bar{\mathbf{Q}}_t^{(g)} / \mathbf{A}_t)^{-1} \mathbf{E}_{\mathcal{A} \mathcal{A}, t}^{(g)} = \left[\mathbf{E}_{\mathcal{A} \mathcal{A}, t}^{(g)} + \frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{A}|} - \mathbf{X}_{\mathcal{A}^c \mathcal{A}, t}^\top \left(\frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{A}^c|} + \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c, t} \right)^{-1} \mathbf{X}_{\mathcal{A}^c \mathcal{A}, t} \right]^{-1} \mathbf{E}_{\mathcal{A} \mathcal{A}, t}^{(g)} \rightarrow \mathbf{I}_{|\mathcal{A}|},$$

and thus

$$[\bar{\mathbf{Q}}_t^{(g)}]^{-1} \begin{pmatrix} \mathbf{0}_{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{0}_{|\mathcal{U}| \times |\mathcal{D}|} \\ \mathbf{0}_{|\mathcal{D}| \times |\mathcal{U}|} & \mathbf{E}_t^{(g)} \end{pmatrix} = [\bar{\mathbf{Q}}_t^{(g)}]^{-1} \mathbf{X}_t \rightarrow \begin{pmatrix} \left(\frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{A}^c|} + \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} \right)^{-1} \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} & \mathbf{0}_{|\mathcal{A}^c| \times |\mathcal{A}|} \\ \mathbf{0}_{|\mathcal{A}| \times |\mathcal{A}^c|} & \mathbf{I}_{|\mathcal{A}|} \end{pmatrix}. \quad (\text{EC.3})$$

Since $\tilde{\mathbf{x}}_t \in \mathcal{A}$ for all $t > T$, Equation (EC.2) implies that $\lim_{t \rightarrow \infty} c_t^{(g)}(\tilde{\mathbf{x}}_t, \mathbf{x}) = 0$ for $\mathbf{x} \in \mathcal{X}$. It also implies that $\lim_{t \rightarrow \infty} v_t^{(g)}(\mathbf{x}) = 0$ for $\mathbf{x} \in \mathcal{A}$ while $\lim_{t \rightarrow \infty} v_t^{(g)}(\mathbf{x}) > 0$ for $\mathbf{x} \in \mathcal{A}^c$. Therefore, $\lim_{t \rightarrow \infty} v_t(\tilde{\mathbf{x}}_t, \mathbf{x}) = 0$ for $\mathbf{x} \in \mathcal{A}$ while $\liminf_{t \rightarrow \infty} v_t(\tilde{\mathbf{x}}_t, \mathbf{x}) > 0$ for $\mathbf{x} \in \mathcal{A}^c$.

Let \mathbf{k}_t be a vector whose element corresponding to solution \mathbf{x} is 0 if $\mathbf{x} \in \mathcal{U}$ and $\bar{Y}_{\mathcal{D}, t}(\mathbf{x}) - \beta_0$ otherwise, i.e., if $\mathbf{x} \in \mathcal{D}$. Reordering the elements of \mathbf{k}_t , we can partition it as

$$\mathbf{k}_t = \begin{pmatrix} \mathbf{k}_{\mathcal{A}^c, t} \\ \mathbf{k}_{\mathcal{A}, t} \end{pmatrix}.$$

As $t \rightarrow \infty$, $\mathbf{k}_{\mathcal{A}, t} \rightarrow \mathbf{y}_{\mathcal{A}} - \beta_0 \vec{\mathbf{1}}_{|\mathcal{A}|}$ by the strong law of large numbers whereas $\mathbf{k}_{\mathcal{A}^c, t} = \mathbf{k}_{\mathcal{A}^c}$ for all $t > T$, where the element of $\mathbf{k}_{\mathcal{A}^c}$ corresponding to solution \mathbf{x} is 0 if $\mathbf{x} \in \mathcal{U}$ and $\bar{Y}_{\mathcal{D}, T}(\mathbf{x}) - \beta_0$ otherwise, i.e., if $\mathbf{x} \in \mathcal{B}$. Then, Equation (EC.3) implies that

$$\mathbf{m}_t^{(g)} = [\bar{\mathbf{Q}}_t^{(g)}]^{-1} \mathbf{X}_t \mathbf{k}_t \rightarrow \begin{pmatrix} \left(\left(\frac{1}{\sigma_g^2} \mathbf{I}_{|\mathcal{A}^c|} + \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} \right)^{-1} \mathbf{X}_{\mathcal{A}^c \mathcal{A}^c} \mathbf{k}_{\mathcal{A}^c} \right) \\ \mathbf{y}_{\mathcal{A}} - \beta_0 \vec{\mathbf{1}}_{|\mathcal{A}|} \end{pmatrix}.$$

Recall from Equation (EC.1) that $m_t^{(\rho)}(\mathbf{x}) \rightarrow 0$ as $t \rightarrow \infty$ for all $\mathbf{x} \in \mathcal{X}$. Therefore, $\lim_{t \rightarrow \infty} m_t(\mathbf{x}) = \lim_{t \rightarrow \infty} \left(\beta_0 + \sum_{\rho \in \mathcal{G}(-g)} m_t^{(\rho)}(\mathbf{x}) + m_t^{(g)}(\mathbf{x}) \right) = y(\mathbf{x})$ for $\mathbf{x} \in \mathcal{A}$. \square

EC.4. A Simple Example to Illustrate the Linear Dependence on the Rows of \mathbb{Y} without $\mathbb{Y}^{(r)}$

Consider a 2-dimensional problem with four solutions

$$\mathcal{X} = \left\{ \begin{pmatrix} 1 \\ 3 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 1 \\ 4 \end{pmatrix}, \begin{pmatrix} 2 \\ 4 \end{pmatrix} \right\},$$

and decompose the objective function value into two groups with $\mathcal{X}^{(1)} = \{1, 2\}$ and $\mathcal{X}^{(2)} = \{3, 4\}$. The transformation matrices can be constructed as

$$\mathbf{T}^{(1)} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \text{ and } \mathbf{T}^{(2)} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}.$$

Then, without $\mathbb{Y}^{(r)}$,

$$\mathbb{Y} = \beta_0 \mathbf{I}_n + \sum_{\rho \in \mathcal{G}} \mathbf{T}^{(\rho)} \mathbb{Y}^{(\rho)} = \begin{pmatrix} \beta_0 + \mathbb{Y}^{(1)}(1) + \mathbb{Y}^{(2)}(3) \\ \beta_0 + \mathbb{Y}^{(1)}(2) + \mathbb{Y}^{(2)}(3) \\ \beta_0 + \mathbb{Y}^{(1)}(1) + \mathbb{Y}^{(2)}(4) \\ \beta_0 + \mathbb{Y}^{(1)}(2) + \mathbb{Y}^{(2)}(4) \end{pmatrix}.$$

Notice that the rows of \mathbb{Y} are linearly dependent; for example, the first row can be obtained by summing the last two rows and subtracting the second row.

EC.5. Parameters

Recall that the additive GMRF model \mathbb{Y} in (4) consists of location parameter β_0 , GMRF $\mathbb{Y}^{(\rho)}$ for $\rho \in \mathcal{G}^{(-g)}$ and random vector $\mathbb{W}^{(g)}$. GMRF $\mathbb{Y}^{(\rho)}$ is characterized by its precision matrix, $\mathbf{Q}^{(\rho)}$. We use a vector $\boldsymbol{\theta}^{(\rho)}$ of parameters to construct $\mathbf{Q}^{(\rho)}$; see Section EC.5.1. On the other hand, random vector $\mathbb{W}^{(g)}$ is characterized by only its variance σ_g^2 as it is a zero-mean random vector with covariance matrix $\sigma_g^2 \mathbf{I}_n$. Benefiting from the additive structure of objective function values, we propose a strategy to choose the initial design points in Section EC.5.2, and use the simulation outputs of the initial design points to estimate $\boldsymbol{\theta}^{(\rho)}$'s and σ_g^2 via maximum likelihood in Section EC.5.3. We discuss an alternative method that can estimate σ_g^2 in Section EC.5.4.

In addition to the parameters mentioned above, the conditional distribution of \mathbb{Y} given $\mathbb{Y}_D^\epsilon = \bar{\mathbf{Y}}_D$ depends on the covariance matrix, $\boldsymbol{\Sigma}^\epsilon$, of the stochastic noise. We simulate all solutions independently (i.e., no common random numbers) so that $\boldsymbol{\Sigma}^\epsilon$ is a diagonal matrix whose diagonal element corresponding to solution \mathbf{x} is $\sigma^2(\mathbf{x})/r(\mathbf{x})$, where $r(\mathbf{x})$ is the number of replications obtained at \mathbf{x} . Since $\sigma^2(\mathbf{x})$ is unknown, the corresponding diagonal element is estimated by $S^2(\mathbf{x})/r(\mathbf{x})$, where $S^2(\mathbf{x}) = \sum_{j=1}^{r(\mathbf{x})} [Y_j(\mathbf{x}) - \bar{Y}_D(\mathbf{x})]^2 / (r(\mathbf{x}) - 1)$ is the sample variance estimate of $\sigma^2(\mathbf{x})$.

EC.5.1. Precision Matrices of GMRFs

For each $\rho \in \mathcal{G}$, we define the set of neighbors of $\mathbf{x}_i^{(\rho)} \in \mathcal{X}^{(\rho)}$ as $\mathcal{N}^{(\rho)}(\mathbf{x}_i^{(\rho)}) = \{\mathbf{x}_j^{(\rho)} \in \mathcal{X}^{(\rho)} : \|\mathbf{x}_i^{(\rho)} - \mathbf{x}_j^{(\rho)}\|_2 = 1\}$ as in Salemi et al. (2019), making $\mathbf{Q}^{(\rho)}$ very sparse as the fraction of nonzero elements in $\mathbf{Q}^{(\rho)}$ is no more than $(2d^{(\rho)} + 1)/n^{(\rho)}$. The nonzero elements are specified by a vector of parameters $\boldsymbol{\theta}^{(\rho)} = [\theta_0^{(\rho)}, \theta_1^{(\rho)}, \theta_2^{(\rho)}, \dots, \theta_{d^{(\rho)}}^{(\rho)}]^\top$; we suppress the dependency of $\mathbf{Q}^{(\rho)}$ on $\boldsymbol{\theta}^{(\rho)}$ to simplify the notation. The (i, j) th element of $\mathbf{Q}^{(\rho)}$ is

$$[\mathbf{Q}^{(\rho)}]_{ij} = \begin{cases} \theta_0^{(\rho)}, & \text{if } i = j, \\ -\theta_0^{(\rho)} \theta_l^{(\rho)}, & \text{if } |\mathbf{x}_i^{(\rho)} - \mathbf{x}_j^{(\rho)}| = \mathbf{e}_l, \\ 0, & \text{otherwise,} \end{cases}$$

where \mathbf{e}_l is the l th standard basis vector and $|\cdot|$ is the element-wise absolute value operator. Notice that $\theta_0^{(\rho)}$ is the conditional precision of each point in the group, and thus it must be positive, i.e., $\theta_0^{(\rho)} > 0$. Also, notice that $\theta_l^{(\rho)}$ is the conditional correlation between points that neighbor in the l th dimension, in the group. To have nonnegative conditional correlations, we restrict the values of $\theta_1^{(\rho)}, \theta_2^{(\rho)}, \dots, \theta_{d^{(\rho)}}^{(\rho)}$ to be nonnegative, i.e., $\theta_l^{(\rho)} \geq 0$ for all $1 \leq l \leq d^{(\rho)}$. Finally, to guarantee that $\mathbf{Q}^{(\rho)}$ positive-definite, we force it to be diagonally dominant, i.e., $\sum_{l=1}^{d^{(\rho)}} \theta_l^{(\rho)} < 0.5$, which is sufficient but not necessary for positive definiteness.

EC.5.2. Construction of Initial Design Points

For the purpose of estimating the $\boldsymbol{\theta}^{(\rho)}$'s via maximum likelihood in the next section, we approximate the objective function values as being purely additive, i.e., $y^{(r)}(\mathbf{x}) = 0$ for all $\mathbf{x} \in \mathcal{X}$. This approximation allows us to obtain values that depend only on a single group by taking the *difference* between objective function values at certain feasible solutions. Notice that obtaining parameters in this way has no impact on the convergence of DASSO; DASSO converges to the optimal solution for any set of legitimate GMRF parameters. And, of course, the GMRF parameters are artifacts that facilitate BO; there are no “true” values. Parameters obtained under this approximation are sufficient to capture local group-by-group behavior of the objective function.

To construct such solutions, first consider a subset of feasible solutions, denoted by \mathcal{S} . The size of \mathcal{S} , i.e., $s = |\mathcal{S}|$, is the sample size used for estimation. For each solution $\mathbf{x} \in \mathcal{S}$, pick g solutions $(\mathbf{x}_{-1}, \mathbf{x}_{-2}, \dots, \mathbf{x}_{-g})$ such that $\mathbf{x}^{(\varrho)} \neq \mathbf{x}_{-\rho}^{(\varrho)}$ if $\varrho = \rho$, and $\mathbf{x}^{(\varrho)} = \mathbf{x}_{-\rho}^{(\varrho)}$ otherwise, for $\rho, \varrho \in \mathcal{G}$. In other words, \mathbf{x} and $\mathbf{x}_{-\rho}$ differ only in their lower dimensional components corresponding to group ρ . Therefore, $y(\mathbf{x}) - y(\mathbf{x}_{-\rho}) = y^{(\rho)}(\mathbf{x}^{(\rho)}) - y^{(\rho)}(\mathbf{x}_{-\rho}^{(\rho)})$, which depends only on group ρ (under the assumption stated above). Thus, while estimating $\boldsymbol{\theta}^{(\rho)}$ for group ρ , we can focus on solution pairs $\{(\mathbf{x}, \mathbf{x}_{-\rho})\}_{\mathbf{x} \in \mathcal{S}}$, that is, the difference between their objective function values, and disregard the other $s(g-1)$ solutions; notice that $|\mathcal{D}| = s(g+1)$ initially and only $2s$ solutions are used to estimate $\boldsymbol{\theta}^{(\rho)}$. Of course, we need the other (disregarded) solutions, but to create such pairs for the groups other than ρ .

For $\mathbf{x} \in \mathcal{S}$ and $\rho \in \mathcal{G}$, let $\mathbf{b}^{(\rho)}$ be a row vector of size $1 \times (g+1)$ whose elements are associated with solutions $(\mathbf{x}, \mathbf{x}_{-1}, \mathbf{x}_{-2}, \dots, \mathbf{x}_{-g})$, where the elements corresponding to \mathbf{x} and $\mathbf{x}_{-\rho}$ are 1 and -1 , respectively, and the rest is 0. Notice that

$$\mathbf{b}^{(\rho)} [y(\mathbf{x}), y(\mathbf{x}_{-1}), \dots, y(\mathbf{x}_{-g})]^\top = y(\mathbf{x}) - y(\mathbf{x}_{-\rho}) = y^{(\rho)}(\mathbf{x}^{(\rho)}) - y^{(\rho)}(\mathbf{x}_{-\rho}^{(\rho)})$$

since \mathbf{x} and $\mathbf{x}_{-\rho}$ differ only in their lower dimensional components corresponding to group ρ . Similar to row vector $\mathbf{b}^{(\rho)}$, we construct a matrix $\mathbf{B}^{(\rho)}$ of size $s \times s(g+1)$ whose rows and columns are associated with the solutions in \mathcal{S} and in \mathcal{D} , respectively, where each row has exactly one 1 and one -1 while the rest are 0. In particular, in each row, the elements with 1 and -1 are for \mathbf{x} and $\mathbf{x}_{-\rho}$ of the corresponding solution pair. Notice that $\mathbf{B}^{(\rho)} \vec{\mathbf{1}}_{|\mathcal{D}|} = \vec{\mathbf{0}}_s$. Further, for $\varrho \neq \rho$, notice that $\mathbf{B}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\varrho)} = \mathbf{0}_{s \times |\mathcal{D}^{(\varrho)}|}$ since $\mathbf{x}^{(\varrho)} = \mathbf{x}_{-\rho}^{(\varrho)}$.

EC.5.3. MLEs for GMRFs

As mention in the previous section, we approximate the objective function values as purely additive. In other words, $\mathbb{Y}^{(r)}$ is excluded from \mathbb{Y} in (3) as we focus on the objective function value differences. Then, since $\mathbf{B}^{(\rho)} \vec{\mathbf{1}}_{|\mathcal{D}|} = \vec{\mathbf{0}}_s$ and $\mathbf{B}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\varrho)} = \mathbf{0}_{s \times |\mathcal{D}^{(\varrho)}|}$ for $\varrho \neq \rho$, notice that $\mathbf{B}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^{(\rho)} + \mathbf{B}^{(\rho)} \epsilon$ depends only on group ρ and the stochastic noise. Moreover, $\mathbf{B}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^\epsilon \sim \mathcal{N}(\vec{\mathbf{0}}_s, \Sigma_B^{(\rho)}(\theta^{(\rho)}))$, where $\Sigma_B^{(\rho)}(\theta^{(\rho)}) = \mathbf{B}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{B}^{(\rho)} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top + \mathbf{B}^{(\rho)} \Sigma_\epsilon [\mathbf{B}^{(\rho)}]^\top$; recall that $\Sigma_{\mathcal{D}\mathcal{D}}^{(\rho)}$ is a function of $\mathbf{Q}^{(\rho)}$, which depends on $\theta^{(\rho)}$. Then, the log-likelihood function of $\theta^{(\rho)}$ given $\mathbf{B}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(\rho)} \bar{\mathbf{Y}}_{\mathcal{D}}$ is

$$\mathcal{L}(\theta^{(\rho)} \mid \mathbf{B}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(\rho)} \bar{\mathbf{Y}}_{\mathcal{D}}) \propto \frac{1}{2} \log |[\Sigma_B^{(\rho)}(\theta^{(\rho)})]^{-1}| - \frac{1}{2} (\mathbf{B}^{(\rho)} \bar{\mathbf{Y}}_{\mathcal{D}})^\top [\Sigma_B^{(\rho)}(\theta^{(\rho)})]^{-1} \mathbf{B}^{(\rho)} \bar{\mathbf{Y}}_{\mathcal{D}}.$$

Thus, the MLE $\hat{\theta}^{(\rho)}$ of $\theta^{(\rho)}$ can be obtained by solving

$$\hat{\theta}^{(\rho)} = \arg \max_{\theta^{(\rho)} \in \Theta^{(\rho)}} \mathcal{L}(\theta^{(\rho)} \mid \mathbf{B}^{(\rho)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(\rho)} \bar{\mathbf{Y}}_{\mathcal{D}}),$$

where $\Theta^{(\rho)} = \{\theta^{(\rho)} : \theta_0^{(\rho)} > 0, \theta_l^{(\rho)} \geq 0 \text{ for } 1 \leq l \leq d^{(\rho)}, \text{ and } \sum_{l=1}^{d^{(\rho)}} \theta_l^{(\rho)} < 0.5\}$ is a set of values of $\theta^{(\rho)}$ that makes $\mathbf{Q}^{(\rho)}$ positive-definite. Of course, this approach causes loss of information, but makes it easier to estimate the parameters for each group individually.

Similarly, notice that $\mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^{(g)} + \mathbf{B}^{(g)} \epsilon$ depends only on the last group g and the stochastic noise. Moreover, $\mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^\epsilon \sim \mathcal{N}(\vec{\mathbf{0}}_s, \Sigma_B^{(g)}(\sigma_g^2))$, where $\Sigma_B^{(g)}(\sigma_g^2) = \sigma_g^2 \mathbf{B}^{(g)} [\mathbf{B}^{(g)}]^\top + \mathbf{B}^{(g)} \Sigma_\epsilon [\mathbf{B}^{(g)}]^\top$. Then, the log-likelihood function of σ_g^2 given $\mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(g)} \bar{\mathbf{Y}}_{\mathcal{D}}$ is

$$\mathcal{L}(\sigma_g^2 \mid \mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(g)} \bar{\mathbf{Y}}_{\mathcal{D}}) \propto \frac{1}{2} \log |[\Sigma_B^{(g)}(\sigma_g^2)]^{-1}| - \frac{1}{2} (\mathbf{B}^{(g)} \bar{\mathbf{Y}}_{\mathcal{D}})^\top [\Sigma_B^{(g)}(\sigma_g^2)]^{-1} \mathbf{B}^{(g)} \bar{\mathbf{Y}}_{\mathcal{D}}.$$

Thus, the MLE $\hat{\sigma}_g^2$ of σ_g^2 can be obtained by solving

$$\hat{\sigma}_g^2 = \arg \max_{\sigma_g^2 > 0} \mathcal{L}(\sigma_g^2 \mid \mathbf{B}^{(g)} \mathbb{Y}_{\mathcal{D}}^\epsilon = \mathbf{B}^{(g)} \bar{\mathbf{Y}}_{\mathcal{D}}).$$

EC.5.4. An Alternative Method for Estimation of Random-Effect Vectors

Notice that $\text{Var}[\mathbb{Y}(\mathbf{x})] = \sum_{\rho \in \mathcal{G}(-g)} \text{Var}[\mathbb{Y}^{(\rho)}(\mathbf{x}^{(\rho)})] + \text{Var}[\mathbb{W}^{(g)}(\mathbf{x})]$ for each $\mathbf{x} \in \mathcal{X}$ since $\mathbb{Y}^{(\rho)}$'s and $\mathbb{W}^{(g)}$ are independent. To estimate variance $\sigma_g^2 = \text{Var}[\mathbb{W}^{(g)}(\mathbf{x})]$, it suffices to estimate the total variance, i.e., $\text{Var}[\mathbb{Y}(\mathbf{x})]$, and the first-order effects, i.e., $\text{Var}[\mathbb{Y}^{(\rho)}(\mathbf{x}^{(\rho)})]$ for each ρ . The former can be done with any subset of solutions such as \mathcal{S} , e.g., an estimate of $\text{Var}[\mathbb{Y}(\mathbf{x})]$ is

$$\widehat{\text{Var}}(\mathbb{Y}(\cdot)) = \frac{1}{s} \sum_{\mathbf{x} \in \mathcal{S}} (\bar{Y}_{\mathcal{D}}(\mathbf{x}))^2 - \left(\frac{1}{s} \sum_{\mathbf{x} \in \mathcal{S}} \bar{Y}_{\mathcal{D}}(\mathbf{x}) \right)^2.$$

On the other hand, the latter requires a special construction of initial design points for better estimates (Saltelli 2002). For such construction, for $\mathbf{x} \in \mathcal{S}$, consider g solutions $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_g)$ such that $\mathbf{x}^{(\varrho)} = \mathbf{x}_\rho^{(\varrho)}$ if $\varrho = \rho$, and $\mathbf{x}^{(\varrho)} \neq \mathbf{x}_\rho^{(\varrho)}$ otherwise, for $\rho, \varrho \in \mathcal{G}$. In other words, \mathbf{x} and \mathbf{x}_ρ differ in their lower dimensional components except the ones corresponding to group ρ . Then, using $s = |\mathcal{S}|$ such solution pairs, $\text{Var}[\mathbb{Y}^{(\rho)}(\mathbf{x}^{(\rho)})]$ can be estimated as

$$\widehat{\text{Var}}(\mathbb{Y}^{(\rho)}(\cdot)) = \frac{1}{s-1} \sum_{\mathbf{x} \in \mathcal{S}} \bar{Y}_{\mathcal{D}}(\mathbf{x}) \bar{Y}_{\mathcal{D}}(\mathbf{x}_\rho) - \left(\frac{1}{s} \sum_{\mathbf{x} \in \mathcal{S}} \bar{Y}_{\mathcal{D}}(\mathbf{x}) \right)^2.$$

Notice that the feasible solutions used for this estimation are different than the ones used for the estimation of $\theta^{(\rho)}$'s. Therefore, considering both estimations, $|\mathcal{D}| = s(2g+1)$ initially. Since this method requires additional initial design points, we use MLE to estimate variance σ_g^2 in our numerical experiments, as explained in the previous section.

EC.5.5. MLEs for Location Parameters

For fixed $\theta^{(\rho)}$'s, the MLE for β_0 is

$$\hat{\beta}_0 = \left(\bar{\mathbf{I}}_{|\mathcal{D}|} [\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} + \boldsymbol{\Sigma}^\epsilon]^{-1} \bar{\mathbf{I}}_{|\mathcal{D}|}^\top \right)^{-1} \bar{\mathbf{I}}_{|\mathcal{D}|} [\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} + \boldsymbol{\Sigma}^\epsilon]^{-1} \bar{\mathbf{Y}}_{\mathcal{D}}.$$

Notice that $\hat{\beta}_0$ is a function of $\theta^{(\rho)}$'s, but we suppress the dependency. Also, notice that $\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}$ is a function of σ_g^2 , leading to different estimates depending on which group is considered as the last group (group g). Therefore, while initially estimating β_0 , we compute $\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}$ by excluding group g from the additive model, i.e., $\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}} = \sum_{\rho \in \mathcal{G}} \mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)} \boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}^{(\rho)} [\mathbf{T}_{\mathcal{D}\mathcal{D}}^{(\rho)}]^\top$. However, if β_0 is re-estimated after obtaining additional simulation output, $\boldsymbol{\Sigma}_{\mathcal{D}\mathcal{D}}$ can be computed by taking into account the current last group. Similarly, for fixed $\theta^{(\rho)}$'s, we use MLE to estimate β_z with the simulation outputs of the solutions in \mathcal{D}_z .