

Simulation Optimization: A Review on Theory and Applications

WANG Long-Fei¹ SHI Le-Yuan¹

Abstract Simulation optimization is a very powerful tool in analysis and optimization of complex real systems. In this paper, a tutorial introduction and review of simulation optimization are given. The simulation optimization problems are classified according to the underlying structure of decision variables (discrete or continuous). And some important techniques for simulation optimization are discussed in detail, including their principles, implementation procedures, advantages and disadvantages, and applications. The future research directions are also provided in this paper.

Key words Simulation optimization, gradient-based methods, ranking and selection, optimal computing budget allocation (OCBA), nested partitions method

Citation Wang Long-Fei, Shi Le-Yuan. Simulation optimization: a review on theory and applications. *Acta Automatica Sinica*, 2013, **39**(11): 1957–1968

DOI 10.3724/SP.J.1004.2013.01957

Simulation optimization is a very effective and powerful tool to solve optimization problems arising in a wide variety of complex real systems. The main goal of simulation optimization is to determine optimal parameter values of systems that result in the best performance measures. With the development of the theory and methods of simulation optimization and the computing technology, simulation optimization is receiving considerable attentions and tremendous achievements have been obtained. Reference [1] gives an excellent survey on techniques for simulation optimization. Some other good reviews on the theory, techniques, applications and commercial software developments can be found in [2–18].

In this paper, the simulation optimization problems are classified into two categories: continuous variable optimization problems and discrete variable optimization problems. And some main methods that are used to solve these problems are introduced. The existing studies related to these methods are reviewed and the advantages and disadvantages of these methods are discussed. What is more, some future research directions are summarized. The main purpose of this paper is to give a comprehensive introduction and tutorial rather than an exhaustive literature survey, to the researchers who are interested in this area.

The remainder of the paper is organized as follows. In Section 1, the formal description of simulation optimization problems is given, and the classification according to the underlying structure of decision variables is presented. In Section 2, we review the techniques that are used to solve the continuous variable optimization problems, and discuss methods for solving discrete variable optimization problems in Section 3. This paper ends with some concluding remarks and suggestions for future research directions in Section 4.

1 Problem description and classification

In general, the simulation optimization problems can be stated as

$$\min_{\theta \in \Theta} f(\theta)$$

where θ can be a single variable or a p -dimensional vector of all the decision variables, and Θ is the feasible region.

For simulation optimization problems, we do not have much knowledge on the structure of $f(\theta)$ and the analytical expression of $f(\theta)$ cannot be obtained, or may not even exist. So the objective function must be estimated based on the outputs of simulation runs, such as

$$f(\theta) = E[L(\theta, \omega)]$$

where ω represents the randomness of the simulation systems, and $L(\theta, \omega)$ is the performance value obtained from the outputs of simulation runs. This is an unbiased estimate of the objective function. In general, the goal of simulation optimization is to find the optimal θ that minimizes (or maximizes) $f(\theta)$.

There are many practical problems related to simulation optimization in the real world. For example, θ can represent the number of servers and service rates in a service center, or the number of machines and the size of buffers in a manufacturing system, or the maximum queue length and waiting time in a queueing system, then $f(\theta)$ may indicate the customer satisfaction of the service center, or the production rate of the manufacturing system, or the average waiting time and sojourn time in the queueing system, respectively. The simulation optimization can be seen as a structured and systematic approach to determine the optimal input parameters θ that result in the best performance measures.

Unlike other optimization problems (such as linear programming or mixed integer linear programming), the major difficulties in solving simulation optimization problems include:

1) There does not exist an analytical expression of the objective function $f(\theta)$. What is more, the feasible solution space Θ may not be described explicitly.

2) The existence of randomness causes that it is difficult to estimate the objective function values of solution points. Usually, more than one simulation replication is needed to ensure the estimation accuracy.

3) In many cases, the simulation run is very expensive and time-consuming.

According to the underlying structure of the decision variables, simulation optimization problems can be classified into two kinds: continuous variable optimization problems and discrete variable optimization problems. And these two big groups of problems can be subdivided based

Manuscript received July 1, 2013; accepted August 29, 2013
Supported by National High Technology Research and Development Program of China (863 Program) (2012AA040909)
Invited Articles for the Special Issue for the 50th Anniversary of Acta Automatica Sinica
1. Department of Industrial Engineering and Management, College of Engineering, Peking University, Beijing 100871, China

on the solution methods used. For continuous variable optimization problems, the most important and popular methods include gradient-based methods, stochastic approximation methods, sample path methods and response surface methodology. On the other hand, ranking and selection, multiple comparison procedures, ordinal optimization, optimal computing budget allocation, metaheuristics are often used to solve the discrete variable optimization problems. The classification scheme is showed in Fig. 1.

2 Continuous decision variables

For the problems discussed in this section, variables are continuous and Θ is uncountable and infinite. This kind of problems have attracted a great deal of research attentions. Below we discuss the commonly used methods in the literature.

2.1 Gradient-based approaches

The gradient-based approaches can be seen as the stochastic form of gradient search methods that are used to solve deterministic optimization problems. When using the gradient-based approaches to solve simulation optimization problems, we assume that Θ is continuous and differentiable in θ , and the gradients of the responses of the simulation models to the variables should be estimated first, then the gradient search methods developed for non-linear programming problems (such as Newton's method) are employed to identify the optimum. In the following subsections, four main gradient estimation methods are introduced.

2.1.1 Finite difference estimation

The finite difference estimation (FDE) method is very intuitive. It is based on the derivative or partial derivatives of a continuous function, i.e., for the gradient of $f(\theta)$, it can be estimated by

$$\frac{df(\theta)}{d\theta} = \frac{f(\theta + \Delta\theta) - f(\theta)}{\Delta\theta}$$

For the multiple decision variables problems, the gradient can be estimated by

$$\frac{\partial f(\theta)}{\partial \theta^i} = \frac{f(\theta^1, \dots, \theta^i + \Delta\theta^i, \dots, \theta^p) - f(\theta^1, \dots, \theta^i, \dots, \theta^p)}{\Delta\theta^i}$$

For a p decision variables problem, at least $p + 1$ outputs of simulation runs are needed to estimate the gradient at each point. And because the observations are noisy, generally more than one observation is necessary to make reliable estimations.

2.1.2 Infinitesimal perturbation analysis

Perturbation analysis is first proposed in [19]. In general, there are two kinds of perturbation analysis: finite pertur-

bation analysis (FPA) and infinitesimal perturbation analysis (IPA). FPA is designed to estimate the derivatives of discrete variables, so the discussion in this subsection will focus on IPA.

IPA can be used to estimate all gradients of the objective function from a single simulation run. The idea behind IPA is that if the decision variable is perturbed by an infinitesimal amount, the sensitivity of the response of the objective function can be estimated by tracing related statistics of certain events during a simulation run. A comprehensive and detailed discussion about IPA can be found in [20].

Even though IPA method can be used to estimate all gradients based on a single simulation run, there are some restrictive conditions that have to be satisfied. For example, if the sequence of events that govern the behavior of the system changes, the results of IPA may not be reliable^[21]. Also, when using the IPA method, a complete knowledge of the simulation model is necessary.

There have been extensive researches on IPA, such as the sufficient conditions for unbiased estimates^[22], convergence rates^[23], perturbation analysis via coupling^[24], applications in queueing systems^[25-29], inventory systems^[30-32], manufacturing systems^[33-39], etc.

2.1.3 Frequency domain analysis

In [40], a method called frequency domain analysis (FDA) that estimates the sensitivity and gradients of the performance values or responses of simulation models to the variables is proposed. The main idea of FDA is to oscillate the value of a variable according to a sinusoidal function and approximate the partial derivative of the objective function with respect to the variable using the magnitude of the performance value variation. The vector of variables can be stated as

$$\theta(t) = \theta_0 + \alpha \sin(\omega t)$$

where θ_0 is the variable vector, α is the vector of oscillation amplitudes, and ω is the vector of oscillation frequencies. The selection of an appropriate index for these oscillations is critical for this method. In [41], using global simulation clock as the sinusoidal oscillations index is discussed. And the problem of optimally selecting input frequencies is studied in [42].

Even though the efficiency of FDA is very high for some simulation optimization problems, it also has some disadvantages, such as:

- 1) The indexing problem, such as oscillation indexing problem and sampling indexing problem, is difficult to solve. Reference [43] addresses the problem by providing some guidelines on the selection of the oscillation index and

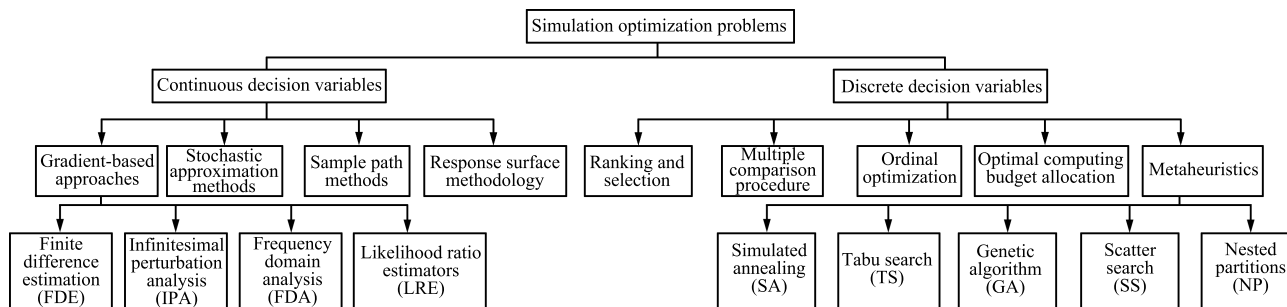


Fig. 1 The classification of simulation optimization problems

compares two different indices: the continuous global simulation clock time and an inherently discrete index.

2) The incorporation of FDA with independently built simulation models may be difficult.

3) It may be impossible to induce sinusoidal oscillations to some variables.

2.1.4 Likelihood ratio estimators

Likelihood ratio estimators (LRE, also called the score function method) assumes that the performance measure function is $L(\mathbf{Y})$, where \mathbf{Y} is a random vector with joint cumulative distribution function $F(\theta, \cdot)$ and density function $f'(\theta, \cdot)$, and dependence on θ enters only through the random vector \mathbf{Y} , thus

$$f(\theta) = E[L(\mathbf{Y})] = \int L(\mathbf{y})dF(\theta, \mathbf{y})$$

By differentiating the above equation with respect to θ , we can estimate the derivative of the performance measure. And the derivative of $f(\theta)$ can be written as

$$\frac{\partial f(\theta)}{\partial \theta} = E_{\theta} \left[L(\mathbf{Y}) \frac{\partial \ln f'(\theta, \mathbf{Y})}{\partial \theta} \right]$$

In [44], an overview of LRE is given. And its applications in the discrete-time and continuous stochastic systems are discussed in [45]. Reference [46] shows that both the sensitivities (derivatives, gradients, Hessians, etc.) and the performance measure can be estimated simultaneously from the same simulation when using LRE method. But it also has some disadvantages, such as it can not be applied to structural parameters, and the variance of the estimator increases with the increase in the run length^[47]. Reference [48] gives detailed discussion about the LRE method. More research about the applications of the LRE method to different situations, its advantages and disadvantages, and the comparisons with other methods can be found in [5, 49–53].

In this section, the main gradient estimation methods, including FDE, IPA, FDA and LRE, are introduced and discussed. For more reviews and summaries, please refer to [54–55].

2.2 Stochastic approximation methods

Stochastic approximation (SA) methods are very popular and can be employed to solve many kinds of simulation optimization problems. The first SA method is proposed in [56–57].

The basic idea of SA is using noisy observations to build the regression function of a stochastic response surface and finding the optimal solution through some recursive procedures. The original recursive formula for a single variable can be stated as

$$\theta_{n+1} = \Pi_{\Theta}(\theta_n - a_n \hat{\nabla} f_n)$$

where a_n is a series of real numbers that satisfy $\sum a_n < +\infty$, $\sum a_n^2 < +\infty$. θ_n is the estimated value at the n th iteration, $\hat{\nabla} f_n$ is the estimate of the gradient ∇f_n at the n th iteration, and Π_{Θ} is a projection onto Θ . It has been proved that as n approaches infinity, θ_n approaches the value such that the regression function of the stochastic response is minimized (for a minimization problem) or maximized (for a maximization problem). But the biggest disadvantage of SA method is that the convergence rate may be very low, i.e., a large number of iterations are necessary to identify the optimum.

As one of the most important methods to solve simulation optimization problems, SA has attracted a great deal of research attentions. In [58], two SA algorithms are proposed. And an alternative proof for the convergence of stochastic approximation algorithms is provided in [59–60]. To deal with the problem of low convergence rate and divergence, [61] studies how to allocate the total available computational budget to the successive SA iterations, and [62] develops a variant of stochastic approximation defined over a growing sequence of compact sets. There are also some excellent researches on the variance reduction in stochastic approximation estimates^[63], the incorporation of different gradient estimators into SA method^[64–71], the applications in queueing systems^[72–74] and manufacturing systems^[75], and the methods that are used to accelerate the convergence rates^[76–78]. Reference [79] gives a comprehensive and detailed discussion of the principles, properties, algorithms and applications of stochastic approximation methods.

2.3 Sample path methods

Sample path method is also called stochastic counterpart method, or sample average approximation method^[48]. When using this method, some simulation replications should be performed first and the expected value of the objective function is estimated by the average of the observations. Then the deterministic optimization techniques are used to solve the problem. This method is very effective to deal with the difficulties faced by stochastic approximation method, such as low convergence rates, absence of robust stopping rules and complicated constraints.

The conditions under which the sample path method converges are given in [80]. In [81], the sample path method is used to set release times for jobs with due dates in a stochastic production flow line. And in [47], a sample path algorithm for optimizing simulation models with rare events is proposed. The convergence rates, stopping rules, and computational complexity of the sample path method are discussed in [82]. There are also some researches on the applications of sample path methods to stochastic root-finding problems^[83] and the incorporation of statistical inferences^[84].

2.4 Response surface methodology

Response surface methodology (RSM) is first proposed in [85–86] for the exploration and exploitation of stochastic response functions. The idea behind the response surface methodology is building an approximate functional relationship between the input variables and the output objective, i.e., fitting a series of regression models to the responses of the simulation model by evaluating it at several points, and then optimizing the regression function. Because the first-order regression models generally provide good fit locally, when using this method, first-order models are usually fitted before reaching the vicinity of the optimum, and then other metamodels (such as higher order regression models) are used to provide good global fit.

The biggest advantage of RSM is that many statistical methods, especially statistical design of experiments, can be incorporated into RSM. Early applications of RSM to simulation optimization are given in [87–90]. And [91] gives a survey of the RSM researches from 1966 to 1988. To get better fits when using RSM, some methods are integrated into it, such as quasi-Newton method^[92], gradient deflection and second-order strategies^[93]. The combinations of RSM and regression analysis, statistical designs

and the steepest descent (ascent) method are discussed in [94].

RSM has also been applied to decision support systems^[95–96] and computer-integrated manufacturing systems^[97]. For more discussions and researches about RSM methods, please refer to [98–99].

3 Discrete decision variables

For the problems discussed in this section, the feasible region Θ is finite or countably infinite. So the methods discussed in the previous section may be inapplicable. According to the structure of the feasible region (finite and small, finite and large, countably infinite), different solution methods are introduced in the following section.

3.1 Ranking and selection

For many simulation optimization problems, there is a given set of alternatives, and the goal is to select the best one. Because the set of alternatives is fixed, there is no need to search the solution space. What should do is to calculate or estimate the performance of alternatives and compare them. If the problem is deterministic, the goal can be achieved by enumeration. But due to the existence of randomness, the number of simulation replications for each alternative is usually more than one to get an accurate estimation. So the main problem is how to efficiently allocate the finite computation efforts to each alternative to get a reliable conclusion.

In general, the problem mentioned above can be classified into two categories. And to give a clear description, the notations should be introduced: let $\Theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ denote the set of solutions, let θ^* denote the best solution, and $f(\theta_i)$ the objective function, P the probability of correct selection (the probability of selecting the best solution). So the two kinds of problems can be described as follows:

- 1) Minimize the number of simulation replications subject to P exceeding a given level;
- 2) Maximize P subject to a given simulation budget constraint.

One of the most important methods to deal with these problems is indifferent zone ranking and selection procedure. First define the difference in $f(\theta)$ that is less than $\delta > 0$ to be insignificant, then allocation the simulation replications to each solutions carefully to assure that the probability of correct selection is greater than or equal to a prespecified value P^0 , i.e.,

$$Prob\{|f(\theta) - f(\theta^*)| < \delta\} \geq P^0$$

To achieve the desired probability guarantee, [100–101] use a two-stage procedure to determine how many simulation replications are allocated to each solution. The main idea is to estimate the mean and variance of each solution in the first stage and then use the variances to determine how many more simulation replications are needed in the second stage. But there is a restriction that the simulation runs should be conducted independently so that the outputs from each run are independent.

There is also another popular method, called subset selection or screening. In subset selection, the objective is not to select the best solution, but to reduce the feasible region to a small subset of solutions, i.e., identify a subset that contains the best solution. Some early studies on this problem are given in [102–103], in which two basic assumptions should be satisfied, i.e., the simulation output is normal with common variance and the same number of

simulation observations is used for each solution. Reference [104] develops two procedures for screening a set of populations with unknown moments. And [105] also presents procedures the assumption of common variance is not required. In [106], the combinations of statistical subset selection and indifference-zone ranking procedures are investigated. In this kind of approaches, the subset-selection procedures are used to screen out the obviously inferior systems, and the best system is distinguished from the less obviously inferior systems through the indifference-zone procedures. For more researches on ranking and selection, see [107].

3.2 Multiple comparison procedures

Similar to ranking and selection, multiple comparison procedures (MCPs) are also efficient methods to find the optimal alternative over a finite set $\Theta = \{\theta_1, \theta_2, \dots, \theta_k\}$, $k < +\infty$. The MCPs can be classified into three kinds: all pairwise multiple comparisons (MCA), multiple comparisons with the best (MCB)^[108–109], and multiple comparisons with a control (MCC)^[110–112].

When using the multiple comparison procedures, the difference between the estimates of the performance measures of a pair of solutions is computed and the confidence interval for a given confidence level is constructed. For example, for a pair of solutions θ_i and θ_j , if the $(1 - \alpha)100\%$ confidence interval for $f(\theta_i) - f(\theta_j)$ is strictly negative, then we can say that θ_i is a better solution than θ_j (for a minimization problem).

The existing researches about MCP mainly focus on the incorporation of variance reduction techniques^[113–114], the applications to steady-state simulation^[115–118], and the connection between indifference zone procedures and MCB^[119–121].

3.3 Ordinal optimization

Ordinal optimization is a very effective and efficient method used to solve problems with a large number of decision variables. It is first proposed in [122].

The basic idea of ordinal optimization is to concentrate on finding good, better, or best designs rather than estimate accurately the performance values of these designs. This idea is called “goal soften”, and it is intuitive that finding the ordering among alternative designs is much easier than estimating the performance value of every design and ranking order. Reference [123–124] give more discussions and explanations about the goal softening in ordinal optimization. And a comprehensive review is given in [125].

One of the significant advantages of ordinal optimization is the exponential convergence rate. In [126–128], the convergence properties of ordinal optimization are investigated. And [127] shows that for a regenerative system, the probability of obtaining a desired solution by using ordinal optimization converges with an exponential rate while the variance of the performance measures converges with the rate $O(1/t^2)$, where t is the simulation time. There are also some researches on the impact of correlation on ordinal optimization^[129], subset selection rules^[130], the lower bound on the probability that the selected subset contains at least one good design and the probability that the best of the selected subset is very close to the true best design^[131], iterative ordinal optimization procedure^[132] and the methods that can be used to enhance the efficiency of ordinal optimization^[133–135].

For a comprehensive and detailed introduction and discussion about ordinal optimization, see [136].

3.4 Optimal computing budget allocation

As discussed in Subsection 3.1, if a set of alternatives is given, and the number of the alternatives is small, then the key issue is how to efficiently allocate the computing efforts to each alternative. One simple and intuitive method is to gradually increase the number of simulation replications for each alternative until the variance decreases to a sufficiently small value, which can be used to give a satisfactory conclusion. There is also a simple allocation method, i.e., allocating the computing budget to each alternative equally. But these methods are not efficient. Consider one situation that one or several alternatives have very small variances, then there is no need to take many replications to estimate their performance. And if the means of some alternatives are very large, i.e., these alternatives are obviously inferior to others (for a minimization problem), then these alternatives may be screened out and there is also no need to allocate many simulation replications to them, even though the variances of them may be large. So it can be seen that the number of simulation replications allocated to each alternative should be dependent on the means and variances of them, and it decreases with the increase in means and the decrease in variance (for a minimization problem).

The idea discussed above is included in the optimal computing budget allocation (OCBA) approach. The OCBA is introduced in [133–134, 137–138]. In [134], the simulation computing budget allocation problem is formulated as an optimization problem. To describe the formulation clearly, some notations should be defined first. Let N_i denote the number of simulation replications that are allocated to the i th alternative, T the total computing budget. Let $P\{CS\}$ denote the probability of correct selection, i.e., the probability of selecting the best alternative. Then the problem can be formulated as

$$\begin{aligned} \max \quad & P\{CS\} \\ \text{s.t.} \quad & N_1 + N_2 + \dots + N_k = T \\ & N_i \in N, \quad i = 1, \dots, k \end{aligned}$$

where N is the set of non-negative integers.

To solve this problem, $P\{CS\}$ must be estimated. Based on the estimation technique that approximates $P\{CS\}$ for ordinal comparison developed in [131], a cost-effective sequential approach is proposed. In this approach, the computing budget allocation is based on the combination of the means and variances of alternatives. The details of the sequential algorithm are showed below.

Step 1. Input k, T, Δ, n_0 .

Step 2. Initialize, and let l be 0, perform n_0 simulation replications for each alternative, i.e., let

$$N_1^l = N_2^l = \dots = N_k^l = n_0$$

Step 3. If $\sum_{i=1}^k N_i^l < T$, go to Step 4, otherwise, go to Step 7.

Step 4. Based on the outputs of the simulation runs, calculate sample means \bar{f}_i , sample standard deviation s_i , $i = 1, 2, \dots, k$ for each alternative, and $b = \arg \min_i \bar{f}_i$.

Step 5. Increase the computing budget by Δ and calculate the new budget allocation, $N_1^{l+1}, N_2^{l+1}, \dots, N_k^{l+1}$, according to

$$\frac{N_i^{l+1}}{N_j^{l+1}} = \frac{\frac{s_i^2}{(\bar{f}_b - \bar{f}_i)^2}}{\frac{s_j^2}{(\bar{f}_b - \bar{f}_j)^2}}, \quad i \neq j \neq b$$

$$N_b^{l+1} = s_b \sqrt{\sum_{i=1, i \neq b}^k \left(\frac{N_i^{l+1}}{s_i} \right)^2}$$

Step 6. Perform $\max(N_i^{l+1} - N_i^l, 0)$ additional simulations for alternative i , $i = 1, \dots, k$; and let $l = l + 1$; go to Step 3.

Step 7. Based on the outputs of the simulation runs, identify the optimal alternative $b = \arg \min_i \bar{f}_i$.

In [139], the dual problem of the original problem discussed above is formulated. In the dual problem, the objective is to minimize the computing efforts subject to a specified $P\{CS\}$, such as 95%, the formulation is showed below:

$$\begin{aligned} \min \quad & N_1 + N_2 + \dots + N_k \\ \text{s.t.} \quad & P\{CS\} \geq P^* \\ & N_i \in N, \quad i = 1, \dots, k \end{aligned}$$

This problem coincides with the first kind of problems that ranking and selection is used to solve. Because the approximation solutions of the problem are the same as the original problem, so the OCBA approach can also be applied to this problem.

In [140], the OCBA approach is compared with other popular methods through extensive experimentations, and the results indicate that the efficiency of OCBA is very high.

When using the OCBA approach, we assume that the sample performances $L(\theta, \omega)$ are independently and normally distributed. In [141–142], the OCBA problem in which the simulated designs or system performances are correlated is considered. And in [143], the non-normal distribution of the underlying random variables is discussed. There are also some researches about other kinds of objective functions, such as minimizing the expected opportunity cost^[144–145], minimizing variance^[146], and selecting an optimal subset of top- m solutions^[147]. In [148], the incorporation of OCBA and the selection of Gaussian process models are used on different noisy mathematical test functions. For more researches on the OCBA approach, please refer to [149].

3.5 Metaheuristics

Most of the current commercial simulation software packages contain the optimization modules. Rather than making statistical estimation, these optimization modules incorporate some search methods to find the optimal values of input parameters. It should be noted that metaheuristics are the most commonly used methods embedded in simulation software.

When combining the metaheuristics with simulation models, the latter can be seen as a black box, i.e., some input parameters are given to the black box, then the simulation models will give some feedbacks or responses, which can be used to guide the search process in metaheuristics. Fig. 2, which is proposed in [16], gives a good description of the relationship between the simulation models and metaheuristics. In this section, five metaheuristics are introduced and discussed, but an emphasis is placed on the nested partitions method, which is a newly-developed and powerful method for global optimization.

3.5.1 Simulated annealing

Simulated annealing is a kind of global optimization method based on the simulation of the physical annealing process, which is first proposed in [150] to solve combinatorial optimization problems. When implementing the sim-

ulated annealing method, the search process moves from one solution to the next until the terminating condition is satisfied. To avoid to be trapped in a local optimum, the inferior solutions may be accepted with a probability, i.e., for each iteration, the probability that a solution will be accepted is

$$\text{Prob}\{\text{Accept}\} = \begin{cases} 1, & \text{if } f(\theta^c) < f(\theta^k) \\ e^{-\frac{f(\theta^c) - f(\theta^k)}{T_k}}, & \text{otherwise} \end{cases}$$

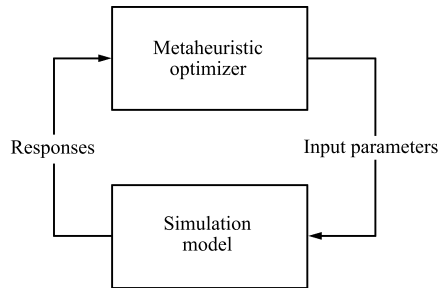


Fig. 2 The relationship between simulation models and metaheuristics

In other words, the candidate solution will be accepted if it is better than the current one. And if it is inferior, it may be accepted or rejected. The probability of being accepted is inversely proportional to the difference between the performance values of the two solutions, and proportional to the temperature T_k .

Even though simulated annealing is commonly used to solve the deterministic optimization problems, there are numerous studies about its applications to the area of simulation optimization. Reference [151] extends the basic convergence results for the simulated annealing algorithm to a stochastic optimization problem where the objective function is stochastic and can be evaluated only through Monte Carlo simulation. And [152] proposes a modified simulated annealing algorithm designed for solving discrete stochastic optimization problems. Rather than the decreasing temperature used in the original simulated annealing algorithm, the constant temperature is used in the new algorithm. And two approaches for estimating the optimal solution are considered and it is shown that both are guaranteed to converge almost surely to the global optimal solutions. To improve the performance of simulated annealing for discrete variable simulation optimization, [153] proposes a method in which portions of the search procedure are based on inferred statistical knowledge of the system.

3.5.2 Tabu search

Tabu search can be seen as a special search procedure that is constrained by the tabu list at each step. The main idea behind this method is that an adaptive memory is used to forbid backtracking moves. The basic framework of tabu search is developed in [154–155].

In tabu search method, a fixed-length list of explored moves is maintained, and the solutions are tabu if they require the moves in the list. The restriction can help the search process to escape from the local optima. For example, if the current solution is θ , and θ' is the best solution in the neighborhood of θ and θ' is not tabu, then the process will move from θ to θ' . The opposite move $\theta' \rightarrow \theta$ will be added into the tabu list and the oldest move in the list will be deleted.

Tabu search procedures have been applied to simulation

optimization in a number of papers. Reference [156] uses the tabu search method with random moves to solve optimal engineering design problems. The reliability and efficiency of the algorithm are investigated by using some standard test functions and the results show that it outperforms the random search and a composite genetic algorithm on the example problems. The combination of simulation of stochastic inventory system simulator and the tabu search is considered in [157].

In addition to these researches, tabu search has also been applied to other different kinds of simulation optimization problems, including determining buffer location and size^[158], determining the number of kanbans and lotsizes^[159], identifying the optimal number of kanbans in a just-in-time system^[160], flow-shop scheduling problems^[161], unmanned aerial vehicles routing problems^[162], and so on.

3.5.3 Genetic algorithm

Genetic algorithm (GA) is a global optimization method based on the natural selection principles. The main difference between the GA and other random search methods is that a population of solutions in GA is manipulated rather than a single solution. In each iteration, a population of solutions is used to generate the next population based on the operators such as selection, crossover and mutation. When selecting solutions, the main principle is that each solution should have a chance to be selected and the probability of being selected as the better solution should be higher. The crossover operation is to take two solutions (parents) to generate new solutions (children). And to escape the local optima, the chromosome of some individuals (solutions) will change slightly in the mutation stage. The genetic algorithm hopefully converges to the global optimal solution.

Since being introduced in [163], a great deal of research has been done and there are many discussions about its application to simulation optimization. Readers interested in this topic, please refer to [164–168].

3.5.4 Scatter search

Similar to GA, scatter search is also designed to operate on a set of points, called reference points. It is first introduced in [169] as a heuristic for integer programming. In each iteration of the scatter search algorithm, new set of points are generated based on the weighted linear combination of the previous points, and some generalized rounding mechanisms should be used to assure the points satisfy integer feasibility conditions.

Reference [170] presents a template for scatter search and its generalized form called path relinking method. The basic design to implement scatter search consists of five methods: a diversification generation method, an improvement method, a reference set update method, a subset generation method, and a solution combination method^[171].

There have been many applications of scatter search in recent years, and a comprehensive discussion and overview can be found in [172].

3.5.5 Nested partitions

The nested partitions (NP) method is a newly-developed and novel optimization framework for solving complex system optimization problems. The advantages of the NP method include flexibility, convergence to a global optimum, high compatibility with parallel computer structures and so on. So it is suitable to solve the simulation optimization problems that arise in the field of manufacturing system, service system, product design, healthcare, etc.

When implementing the NP method, there is a most

promising region in each iteration. The most promising region is partitioned into M sub-regions, and the entire surrounding region is aggregated into one. So, within each iteration, there are $M + 1$ disjoint regions that need to be explored. Some feasible solutions should be sampled from each region by some random sampling schemes. Based on the objective function values or estimated performance values of these solutions, the promising index of each region is calculated. The promising index can be used to determine which region will be the most promising region in the next iteration. If the best promising index corresponds to one of the M sub-regions, the sub-region will be the next most promising region. If the surrounding region is found to be the best, then the method will backtrack to a larger region that contains the previous most promising region. The larger region will be the most promising region in the next iteration. The process will continue until the terminating condition is satisfied.

Mathematically, the NP method will be described as follows. Firstly, some notations and terminologies should be defined. If region η is a sub-region of region σ , then σ is a super-region of η . And let $\sigma(k)$ denote the most promising region of the k th iteration, $d(k)$ the depth of $\sigma(k)$. The depth of the whole solution space Θ is 0. If the feasible solution space Θ contains finite solutions, then there exist some sub-regions, each of which contains only one solution. The singleton regions are called the regions of maximum depth. And if the feasible solutions in Θ is infinite, the maximum depth is the depth of the smallest desired subsets. Let d^* denote the maximum depth. The “generic nested partitions algorithm” or “pure nested partitions” is given below (similar algorithm can also be found in [173]).

1) Partition the most promising region $\sigma(k)$ into M sub-regions $\sigma_1(k), \sigma_2(k), \dots, \sigma_M(k)$, and aggregate the complimentary region $\Theta \setminus \sigma(k)$ into one region $\sigma_{M+1}(k)$.

If we have some knowledge about the structure of the feasible region and the objective function, problem-specific partitioning schemes can be designed and implemented to make the good solutions clustered together in the same sub-regions. Thus the good sub-regions and solutions will be identified quickly and the convergence rate will be improved. Reference [174] gives more discussion about partitioning methods when the NP method is used to solve the traveling salesman problem.

2) Random sampling. Randomly generate N_j sample solutions from each of the region $\sigma_j(k)$, $j = 1, 2, \dots, M + 1$:

$$\theta_1^j, \theta_2^j, \dots, \theta_{N_j}^j, \quad j = 1, 2, \dots, M + 1$$

Calculate the corresponding performance values:

$$f(\theta_1^j), f(\theta_2^j), \dots, f(\theta_{N_j}^j), \quad j = 1, 2, \dots, M + 1$$

Sampling is a very important factor in determining the efficiency of the NP method. To make a correct move in each iteration, or to increase the probability of correct move, good solutions should be sampled in each region. In general, three ways can be used to increase the probability of making a correct move when sampling: a) biasing the sampling distribution so that good solutions are more likely to be selected; b) using some heuristics or metaheuristics (SA, GA, TS, etc.) to search for good solutions; c) using a sufficiently large sample. But whichever sampling method is used must ensure that each solution in the sampling region should be selected with a positive probability.

3) Calculate promising index. For each region σ_j , $j = 1, 2, \dots, M + 1$, calculate the promising index as the best

performance value within the region:

$$I(\sigma_j) = \min_{i=1,2,\dots,N_j} f(\theta_i^j), \quad j = 1, 2, \dots, M + 1$$

There are also some other definitions of the promising index, such as the sample mean^[174].

When using the NP method, suitable promising index should be defined carefully so that more information will be contained and the quality of each region will be better reflected. Because the performance of each sampled point is estimated based on the outputs of simulation runs, it is very important to efficiently allocate the computing budget to each solution and OCBA can be employed in this step.

4) Move. Identify the index of the region with the best performance value:

$$j_k^* = \arg \min_{j=1,2,\dots,M+1} I(\sigma_j), \quad j = 1, 2, \dots, M + 1$$

If more than one regions are equally promising, the tie can be broken arbitrarily. If this index corresponds to a region that is a sub-region of $\sigma(k)$, i.e., $j_k^* \leq M$, then let this sub-region be the most promising region in the next iteration:

$$\sigma(k + 1) = \sigma_{j_k^*}(k)$$

Otherwise, if the index corresponds to the complementary region, i.e., $j_k^* = M + 1$, backtrack to the super-region of the current most promising region (the previous most promising region):

$$\sigma(k + 1) = \sigma(k - 1)$$

or backtrack to the entire solution space:

$$\sigma(k + 1) = \sigma$$

There have been extensive researches on the NP methods, such as the study about its convergence rate to a global optimum and the stopping criteria^[175–177], the combination of NP method and OCBA^[178], and the applications of NP method to production design^[179], supply chain network optimization^[180], beam angle and dose optimization^[181], and scheduling^[182–183], etc. A comprehensive introduction and review on the nested partitions method are given in [173].

Because metaheuristics are usually designed for the combinatorial optimization problems, it is very important to investigate the effect of the simulation noise on these algorithms. The convergence rates of these methods are also worthy of deeper studies. In [184], these issues are investigated and some metaheuristics are discussed in detail.

4 Conclusion

Simulation optimization is a very active field of research. In recent years, extensive researches on simulation optimization has been conducted and a great deal of excellent achievements have been obtained. In this paper, we have given a review on the theory and techniques of simulation optimization. According to the underlying structure of decision variables, the simulation optimization problems are classified into two categories: continuous variable optimization problems and discrete variable optimization problems. For each kind of problems, several important solution methods or techniques are introduced. The principles, implementation procedures, applications, advantages and disadvantages of these techniques are discussed, and existing researches on these topics are reviewed. It is worth mentioning that there are also some other useful methods that

are not covered in this paper, such as Nelder-Mead simple/complex search methods^[185–187], Hooke-Jeeves pattern search methods^[188–189], gradient surface methods^[190], ant colony algorithms^[191], etc.

Finally, we make the following suggestions for future research:

1) In the literature, many techniques and algorithms have been proposed. But there is not enough research on the comparisons between them. So it is very important to build and maintain a testbed for the simulation optimization problems^[192], including benchmark problems and standard test functions to evaluate the performances of different techniques, which can help researchers to compare their algorithms.

2) As pointed in [13], there seems to be a gap between academic theory and commercial applications. Most of the existing optimization modules in simulation software search for good solutions based on metaheuristics (SA, GA, TS, etc.), other methods and techniques, such as gradient search methods, statistical inference techniques, have not been integrated into these modules yet.

3) The efficiency of the simulation optimization software needs to be improved. Because of the complexity and huge scale of real problems, it is very difficult and time-consuming to solve them using the existing simulation software packages. Some advanced techniques, such as parallel computing and factor screening techniques, should be integrated into those software packages. And more research efforts are necessary to improve the efficiency of the simulation optimization systems.

References

- 1 Tekin E, Sabuncuoglu I. Simulation optimization: a comprehensive review on theory and applications. *IIE Transactions*, 2004, **36**(11): 1067–1081
- 2 Meketon M S. Optimization in simulation: a survey of recent results. In: Proceedings of the 1987 Winter Simulation Conference. Piscataway, NJ: IEEE, 1987. 58–67
- 3 Jacobson S H, Schruben L W. Techniques for simulation response optimization. *Operations Research Letters*, 1989, **8**(1): 1–9
- 4 Azadivar F. A tutorial on simulation optimization. In: Proceedings of the 1992 Winter Simulation Conference. Piscataway, NJ: IEEE, 1992. 198–204
- 5 Fu M C. Optimization via simulation: a review. *Annals of Operations Research*, 1994, **53**(1): 199–247
- 6 Carson Y, Maria A. Simulation optimization: methods and applications. In: Proceedings of the 1997 Winter Simulation Conference. Atlanta, Georgia: IEEE, 1997. 118–126
- 7 Andradóttir S. A review of simulation optimization techniques. In: Proceedings of the 1998 Winter Simulation Conference. Piscataway, NJ: IEEE, 1998. 151–158
- 8 Bowden R O, Hall J D. Simulation optimization research and development. In: Proceedings of the 1998 Winter Simulation Conference. Washington, DC: IEEE, 1998. 1693–1698
- 9 Azadivar F. Simulation optimization methodologies. In: Proceedings of the 1999 Winter Simulation Conference. Phoenix, AZ: IEEE, 1999. 93–100
- 10 Swisher J R, Hyden P D, Jacobson S H, Schruben L W. A survey of simulation optimization techniques and procedures. In: Proceedings of the 2000 Winter Simulation Conference. Orlando, FL: IEEE, 2000. 119–128
- 11 Law A M, McComas M G. Simulation-based optimization. In: Proceedings of the 2000 Winter Simulation Conference. San Diego, CA, USA: IEEE, 2000. 46–49
- 12 Fu M C. Simulation optimization. In: Proceedings of the 2001 Winter Simulation Conference. Arlington, VA, USA: IEEE, 2001. 53–61
- 13 Fu M C. Optimization for simulation: theory vs. practice. *INFORMS Journal on Computing*, 2002, **14**(3): 192–215
- 14 Ólafsson S, Kim J. Simulation optimization. In: Proceedings of the 2002 Winter Simulation Conference. Piscataway, NJ: IEEE, 2002. 79–84
- 15 Law A M, McComas M G. Simulation optimization: simulation-based optimization. In: Proceedings of the 2002 Winter Simulation Conference. San Diego, CA, USA: IEEE, 2002. 41–44
- 16 April J, Glover F, Kelly J P, Laguna M. Practical introduction to simulation optimization. In: Proceedings of the 2003 Winter Simulation Conference. New Orleans, LA, USA: IEEE, 2003. 71–78
- 17 Fu M C, Glover F W, April J. Simulation optimization: a review, new developments, and applications. In: Proceedings of the 2005 Winter Simulation Conference. Orlando, FL: IEEE, 2005. 83–95
- 18 Fu M C, Chen C H, Shi L. Some topics for simulation optimization. In: Proceedings of the 2008 Winter Simulation Conference. Austin, TX: IEEE, 2008. 27–38
- 19 Ho Y C, Eyster M A, Chien T T. A gradient technique for general buffer storage design in a production line. *International Journal of Production Research*, 1979, **17**(6): 557–580
- 20 Ho Y C, Cao X R. *Perturbation Analysis of Discrete Event Dynamic Systems*. Norwell, MA, USA: Kluwer Academic Pub, 1991
- 21 Heidelberger P. Limitations of Infinitesimal Perturbation Analysis. International Business Machine Research Report RC 11891, IBM Watson Research Labs, Yorktown Heights, NY, 1986
- 22 Glasserman P. Structural conditions for perturbation analysis derivative estimation: finite-time performance indices. *Operations Research*, 1991, **39**(5): 724–738
- 23 L'Ecuyer P, Perron G. On the convergence rates of IPA and FDC derivative estimators. *Operations Research*, 1994, **42**(4): 643–656
- 24 Dai L L. Perturbation analysis via coupling. *IEEE Transactions on Automatic Control*, 2000, **45**(4): 614–628
- 25 Ho Y C, Suri R, Cao X R, Diehl G W, Dille J W, Zazanis M. Optimization of large multiclass (non-product-form) queueing networks using perturbation analysis. *Large Scale Systems*, 1984, **7**: 165–180
- 26 Ho Y C. On the perturbation analysis of discrete-event dynamic systems. *Journal of Optimization Theory and Applications*, 1985, **46**(4): 535–545
- 27 Ho Y C, Hu J Q. An infinitesimal perturbation analysis algorithm for a multiclass G/G/1 queue. *Operations Research Letters*, 1990, **9**(1): 35–44
- 28 Chong E K P, Ramadge P J. Optimization of queues using an infinitesimal perturbation analysis-based stochastic algorithm with general update times. *SIAM Journal on Control and Optimization*, 1993, **31**(3): 698–732
- 29 Fu M C, Hu J Q. Smoothed perturbation analysis derivative estimation for Markov chains. *Operations Research Letters*, 1994, **15**(5): 241–251
- 30 Fu M C. Sample path derivatives for (s, S) inventory systems. *Operations Research*, 1994, **42**(2): 351–364
- 31 Bashyam S, Fu M C. Application of perturbation analysis to a class of periodic review (s, S) inventory systems. *Naval Research Logistics (NRL)*, 1994, **41**(1): 47–80
- 32 Bashyam S, Fu M C. Optimization of (s, S) inventory systems with random lead times and a service level constraint. *Management Science*, 1998, **44**(12-Part-2): S243–S256
- 33 Donohue K L, Spearman M L. Improving the design of stochastic production lines: an approach using perturbation analysis. *International Journal of Production Research*, 1993, **31**(12): 2789–2806

- 34 Yan H, Yin G, Lou S X C. Using stochastic optimization to determine threshold values for the control of unreliable manufacturing systems. *Journal of Optimization Theory and Applications*, 1994, **83**(3): 511–539
- 35 Liberopoulos G, Caramanis M. Infinitesimal perturbation analysis for second derivative estimation and design of manufacturing flow controllers. *Journal of Optimization Theory and Applications*, 1994, **81**(2): 297–327
- 36 Cheng D W. On the design of a tandem queue with blocking: modeling, analysis, and gradient estimation. *Naval Research Logistics (NRL)*, 1994, **41**(6): 759–770
- 37 Heidergott B. Sensitivity analysis of a manufacturing workstation using perturbation analysis techniques. *International Journal of Production Research*, 1995, **33**(3): 611–622
- 38 Brooks C A, Varaiya P. Using augmented infinitesimal perturbation analysis for capacity planning inintree ATM networks. *Discrete Event Dynamic Systems*, 1997, **7**(4): 377–390
- 39 Heidergott B. Optimisation of a single-component maintenance system: a smoothed perturbation analysis approach. *European Journal of Operational Research*, 1999, **119**(1): 181–190
- 40 Schruben L W, Coglianò V J. Simulation sensitivity analysis: a frequency domain approach. In: Proceedings of the 1981 Winter Simulation Conference. Piscataway, NJ: IEEE, 1981. 455–459
- 41 Jacobson S H, Morrice D, Schruben L W. The global simulation clock as the frequency domain experiment index. In: Proceedings of the 1988 Winter Simulation Conference, San Diego, CA, USA: IEEE, 1988. 558–563
- 42 Jacobson S H, Buss A H, Schruben L W. Driving frequency selection for frequency domain simulation experiments. *Operations Research*, 1991, **39**(6): 917–924
- 43 Hazra M M, Morrice D J, Park S K. A simulation clock-based solution to the frequency domain experiment indexing problem. *IEEE Transactions*, 1997, **29**(9): 769–782
- 44 Glynn P W. Likelihood ratio gradient estimation: an overview. In: Proceedings of the 1987 Winter Simulation Conference. Piscataway, NJ: IEEE, 1987. 366–375
- 45 Glynn P W. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 1990, **33**(10): 75–84
- 46 Rubinstein R Y. Sensitivity analysis and performance extrapolation for computer simulation models. *Operations Research*, 1989, **37**(1): 72–81
- 47 Rubinstein R Y. Optimization of computer simulation models with rare events. *European Journal of Operational Research*, 1997, **99**(1): 89–112
- 48 Rubinstein R Y, Shapiro A. *Discrete Event Systems: Sensitivity Analysis and Stochastic Optimization by the Score Function Method*. New York: Wiley, 1993
- 49 Nakayama M K, Goyal A, Glynn P W. Likelihood ratio sensitivity analysis for Markovian models of highly dependable systems. *Operations Research*, 1994, **42**(1): 137–157
- 50 Nakayama M K, Shahabuddin P. Likelihood ratio derivative estimation for finite-time performance measures in generalized semi-Markov processes. *Management Science*, 1998, **44**(10): 1426–1441
- 51 Andradóttir, S. Optimization of the transient and steady-state behavior of discrete event systems. *Management Science*, 1996, **42**(5): 717–737
- 52 Fu M C, Hu J Q, Nagi R. Comparison of gradient estimation techniques for queues with non-identical servers. *Computers and Operations Research*, 1995, **22**(7): 715–729
- 53 Fu M C, Hu J Q. Efficient design and sensitivity analysis of control charts using Monte Carlo simulation. *Management Science*, 1999, **45**(3): 395–413
- 54 Fu M C. 2006. Gradient estimation. *Handbooks in Operations Research and Management Science: Simulation, Chapter 19* (Henderson S G, Nelson B L (editor)). Amsterdam: Elsevier. 575–616
- 55 Fu M C. What you should know about simulation and derivatives. *Naval Research Logistics (NRL)*, 2008, **55**(8): 723–736
- 56 Robbins H, Monro S. A stochastic approximation method. *The Annals of Mathematical Statistics*, 1951, **22**(3): 400–407
- 57 Kiefer J, Wolfowitz J. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, 1952, **23**(3): 462–466
- 58 Glynn P W. Stochastic approximation for Monte Carlo optimization. In: Proceedings of the 1986 Winter Simulation Conference. Piscataway, NJ: IEEE, 1986. 356–365
- 59 Kouritzin M A. On the convergence of linear stochastic approximation procedures. *IEEE Transactions on Information Theory*, 1996, **42**(4): 1305–1309
- 60 Kulkarni S R, Horn C S. An alternative proof for convergence of stochastic approximation algorithms. *IEEE Transactions on Automatic Control*, 1996, **41**(3): 419–424
- 61 L'Ecuyer P, Yin G. Budget-dependent convergence rate of stochastic approximation. *SIAM Journal on Optimization*, 1998, **8**(1): 217–247
- 62 Andradóttir S. A stochastic approximation algorithm with varying bounds. *Operations Research*, 1995, **43**(6): 1037–1048
- 63 Kleinman N L, Spall J C, Naiman D Q. Simulation-based optimization with stochastic approximation using common random numbers. *Management Science*, 1999, **45**(11): 1570–1578
- 64 Fu M C, Ho Y C. Using perturbation analysis for gradient estimation, averaging and updating in a stochastic approximation algorithm. In: Proceedings of the 1988 Winter Simulation Conference. San Diego, CA, USA, 1988. 509–517
- 65 Fu M C. Convergence of a stochastic approximation algorithm for the GI/G/1 queue using infinitesimal perturbation analysis. *Journal of Optimization Theory and Applications*, 1990, **65**(1): 149–160
- 66 L'Ecuyer P, Glynn P W. Stochastic optimization by simulation: convergence proofs for the GI/G/1 queue in steady-state. *Management Science*, 1994, **40**(11): 1562–1578
- 67 Hill S D, Fu M C. Optimizing discrete event systems with the simultaneous perturbation stochastic approximation algorithm. In: Proceedings of the 33rd IEEE Conference on Decision and Control. Lake Buena Vista, FL: IEEE, 1994. 2631–2632
- 68 Fu M C, Hill S D. Optimization of discrete event systems via simultaneous perturbation stochastic approximation. *IEEE Transactions*, 1997, **29**(3): 233–243
- 69 Fu M, Hu J Q. *Conditional Monte Carlo: Gradient Estimation and Optimization Applications*. Boston: Kluwer Academic Publishers, 1997
- 70 Chong E K P, Ramadge P J. Optimal load sharing in soft real-time systems using likelihood ratios. *Journal of Optimization Theory and Applications*, 1994, **82**(1): 23–48
- 71 Andradóttir S. A scaled stochastic approximation algorithm. *Management Science*, 1996, **42**(4): 475–498
- 72 Tang Q Y, Chen H F, Han Z J. Convergence rates of Perturbation-Analysis-Robbins-Monro-Single-Run algorithms for single server queues. *IEEE Transactions on Automatic Control*, 1997, **42**(10): 1442–1447
- 73 Tang Q Y, L'Ecuyer P, Chen H F. Asymptotic efficiency of perturbation-analysis-based stochastic approximation with averaging. *SIAM Journal on Control and Optimization*, 1999, **37**(6): 1822–1847
- 74 Tang Q Y, Chen H F. Central limit theorems for stochastic optimization algorithms using infinitesimal perturbation analysis. *Discrete Event Dynamic Systems*, 2000, **10**(1–2): 5–32

- 75 Hasan C N, Spearman M L. Optimal material release times in stochastic production environments. *International Journal of Production Research*, 1999, **37**(6): 1201–1216
- 76 Andradóttir S. A new algorithm for stochastic optimization. In: Proceedings of the 1990 Winter Simulation Conference. New Orleans, LA: IEEE, 1986. 364–366
- 77 Andradóttir S. A projected stochastic approximation algorithm. In: Proceedings of the 1991 Winter Simulation Conference. Phoenix, AZ: IEEE, 1991. 954–957
- 78 Yakowitz S. A globally convergent stochastic approximation. *SIAM Journal on Control and Optimization*, 1993, **31**(1): 30–40
- 79 Kushner H J, Yin G. *Stochastic Approximation and Recursive Algorithms and Applications (2nd edition)*. New York: Springer-Verlag, 2003
- 80 Gurkan G, Yonca-Ozge A, Robinson T M. Sample-path optimization in simulation. In: Proceedings of the 1994 Winter Simulation Conference. Lake Buena Vista, FL, USA: IEEE, 1994. 247–254
- 81 Homem-de-Mello T, Shapiro A, Spearman M L. Finding optimal material release times using simulation-based optimization. *Management Science*, 1999, **45**(1): 86–102
- 82 Kleywegt A J, Shapiro A, Homem-de-Mello T. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on Optimization*, 2002, **12**(2): 479–502
- 83 Chen H F, Schmeiser B W. Retrospective approximation algorithms for stochastic root finding. In: Proceedings of the 1994 Winter Simulation Conference. Lake Buena Vista, FL, USA: IEEE, 1994. 255–261
- 84 Shapiro A. Simulation based optimization. In: Proceedings of the 1996 Winter Simulation Conference. Washington, DC: IEEE, 1996. 332–336
- 85 Box G E P, Wilson K B. On the experimental attainment of optimum conditions. *Journal of the Royal Statistical Society. Series B (Methodological)*, 1951, **13**(1): 1–45
- 86 Box G E P. The exploration and exploitation of response surfaces: some general considerations and examples. *Biometrics*, 1954, **10**(1): 16–60
- 87 Biles W E. A gradient-regression search procedure for simulation experimentation. In: Proceedings of the 1974 Winter Simulation Conference. Piscataway, NJ: IEEE, 1974. 491–497
- 88 Smith D E. Automatic optimum-seeking program for digital simulation. *Simulation*, 1976, **27**(1): 27–31
- 89 Daugherty A F, Turnquist M A. Simulation optimization using response surfaces based on spline approximations. In: Proceedings of the 1978 Winter Simulation Conference. Piscataway, NJ: IEEE, 1978. 183–193
- 90 Wilson J R. Future directions in response surface methodology for simulation. In: Proceedings of the 1987 Winter Simulation Conference. Piscataway, NJ: IEEE, 1987. 378–381
- 91 Myers R H, Khuri A I, Carter W H Jr. Response surface methodology: 1966-1988. *Technometrics*, 1989, **31**(2): 137–157
- 92 Safizadeh M H, Signorile R. Optimization of simulation via quasi-Newton methods. *ORSA Journal on Computing*, 1994, **6**(4): 398–408
- 93 Joshi S, Sherali H D, Tew J D. An enhanced response surface methodology (RSM) algorithm using gradient deflection and second-order search strategies. *Computers and Operations Research*, 1998, **25**(7–8): 531–541
- 94 Kleijnen J P C. 1998. Experimental design for sensitivity analysis, optimization, and validation of simulation models. *Handbook of Simulation, Chapter 6 (Banks J editor)*. New York: Wiley. 173–223
- 95 Kleijnen J P C. Simulation and optimization in production planning: a case study. *Decision Support Systems*, 1993, **9**(3): 269–280
- 96 Kleijnen J P C. Sensitivity analysis and optimization of system dynamics models: regression analysis and statistical design of experiments. *System Dynamics Review*, 1995, **11**(4): 275–288
- 97 Shang J S, Tadikamalla P R. Output maximization of a CIM system: simulation and statistical approach. *International Journal of Production Research*, 1993, **31**(1): 19–41
- 98 Barton R R, Meckesheimer M. 2006. Metamodelbased simulation optimization. *Handbooks in Operations Research and Management Science: Simulation, Chapter 18 (Henderson S G, Nelson B L editor)*. Amsterdam: Elsevier. 535–574
- 99 Kleijnen J P C. *Design and Analysis of Simulation Experiments*. New York: Springer, 2007
- 100 Dudewicz E J, Dalal S R. Allocation of observations in ranking and selection with unequal variances. *Sankhyā: The Indian Journal of Statistics, Series B*, 1975, **37**(1): 28–78
- 101 Rinott Y. On two-stage selection procedures and related probability-inequalities. *Communications in Statistics-Theory and Methods*, 1978, **7**(8): 799–811
- 102 Gupta S S, Sobel M. On a statistic which arises in selection and ranking problems. *The Annals of Mathematical Statistics*, 1957, **28**(4): 957–967
- 103 Gupta S S. On some multiple decision (selection and ranking) rules. *Technometrics*, 1965, **7**(2): 225–245
- 104 Sullivan D W, Wilson J R. Restricted subset selection procedures for simulation. *Operations Research*, 1989, **37**(1): 52–71
- 105 Nelson B L, Swann J, Goldsman D, Song W. Simple procedures for selecting the best simulated system when the number of alternatives is large. *Operations Research*, 2001, **49**(6): 950–963
- 106 Boesel J, Nelson B L, Kim S H. Using ranking and selection to “clean up” after simulation optimization. *Operations Research*, 2003, **51**(5): 814–825
- 107 Kim S H, Nelson B L. 2006. Selecting the best system. *Handbooks in Operations Research and Management Science: Simulation, Chapter 17 (Henderson S G, Nelson B L editor)*. Amsterdam: Elsevier. 501–534
- 108 Hsu J C. Constrained simultaneous confidence intervals for multiple comparisons with the best. *The Annals of Statistics*, 1984, **12**(3): 1136–1144
- 109 Hsu J C. Sample size computation for designing multiple comparison experiments. *Computational Statistics and Data Analysis*, 1988, **7**(1): 79–91
- 110 Dunnett C W. A multiple comparison procedure for comparing several treatments with a control. *Journal of the American Statistical Association*, 1955, **50**(272): 1096–1121
- 111 Bofinger E, Lewis G J. Two-stage procedures for multiple comparisons with a control. *American Journal of Mathematical and Management Sciences*, 1992, **12**: 253–275
- 112 Damerjji H, Nakayama M K. Two-stage multiple-comparison procedures for steady-state simulations. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 1999, **9**(1): 1–30
- 113 Yang W N, Nelson B L. Optimization using common random numbers, control variates and multiple comparisons with the best. In: Proceedings of the 1989 Winter Simulation Conference. Washington, DC, USA: IEEE, 1989. 444–449
- 114 Yang W N, Nelson B L. Using common random numbers and control variates in multiple-comparison procedures. *Operations Research*, 1991, **39**(4): 583–591
- 115 Goldsman L, Nelson B L. Batch-size effects on simulation optimization using multiple comparisons with the best. In: Proceedings of the 1990 Winter Simulation Conference. New Orleans, LA, 1990. 288–293
- 116 Yuan M, Nelson B L. Multiple comparisons with the best for steady-state simulation. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 1993, **3**(1): 66–79

- 117 Nakayama M K. Selecting the best system in steady-state simulations using batch means. In: Proceedings of the 1995 Winter Simulation Conference. Piscataway, NJ: IEEE, 1995. 362–366
- 118 Nakayama M K. Multiple-comparison procedures for steady-state simulations. *The Annals of Statistics*, 1997, **25**(6): 2433–2450
- 119 Matejcek F J, Nelson B L. Simultaneous ranking, selection and multiple comparisons for simulation. In: Proceedings of the 1993 Winter Simulation Conference. Piscataway, NJ: IEEE, 1993: 386–392
- 120 Goldsman D, Nelson B L. Ranking, selection and multiple comparisons in computer simulation. In: Proceedings of the 1994 Winter Simulation Conference. Lake Buena Vista, FL, USA: IEEE, 1994: 192–199
- 121 Matejcek F J, Nelson B L. Two-stage multiple comparisons with the best for computer simulation. *Operations Research*, 1995, **43**(4): 633–640
- 122 Ho Y C, Sreenivas R S, Vakili P. Ordinal optimization of DEDS. *Discrete Event Dynamic Systems*, 1992, **2**(1): 61–88
- 123 Ho Y C, Deng M. The problem of large search space in stochastic optimization. In: Proceedings of the 33rd IEEE Conference on Decision and Control. Lake Buena Vista, FL: IEEE, 1994. 1470–1475
- 124 Lee L H, Lau T W E, Ho Y C. Explanation of goal softening in ordinal optimization. *IEEE Transactions on Automatic Control*, 1999, **44**(1): 94–99
- 125 Ho Y C. Overview of ordinal optimization. In: Proceedings of the 33rd IEEE Conference on Decision and Control. Lake Buena Vista, FL: IEEE, 1994. 1975–1977
- 126 Dai L. Convergence properties of ordinal comparison in the simulation of discrete event dynamic systems. *Journal of Optimization Theory and Applications*, 1996, **91**(2): 363–388
- 127 Xie X. Dynamics and convergence rate of ordinal comparison of stochastic discrete-event systems. *IEEE Transactions on Automatic Control*, 1997, **42**(4): 586–590
- 128 Dai L, Chen C H. Rates of convergence of ordinal comparison for dependent discrete event dynamic systems. *Journal of Optimization Theory and Applications*, 1997, **94**(1): 29–54
- 129 Deng M, Ho Y C, Hu J Q. Effect of correlated estimation errors in ordinal optimization. In: Proceedings of the 1992 Winter Simulation Conference. Piscataway, NJ: IEEE, 1992: 466–474
- 130 Lau T W E, Ho Y C. Universal alignment probabilities and subset selection for ordinal optimization. *Journal of Optimization Theory and Applications*, 1997, **93**(3): 455–489
- 131 Chen C H. A lower bound for the correct subset-selection probability and its application to discrete-event system simulations. *IEEE Transactions on Automatic Control*, 1996, **41**(8): 1227–1231
- 132 Deng M, Ho Y C. Iterative ordinal optimization and its applications. In: Proceedings of the 36th IEEE Conference on Decision and Control. San Diego, CA: IEEE, 1997. 3562–3567
- 133 Chen H C, Chen C H, Yücesan E. Computing efforts allocation for ordinal optimization and discrete event simulation. *IEEE Transactions on Automatic Control*, 2000, **45**(5): 960–964
- 134 Chen C H, Lin J, Yücesan E, Chick S E. Simulation budget allocation for further enhancing the efficiency of ordinal optimization. *Discrete Event Dynamic Systems*, 2000, **10**(3): 251–270
- 135 Shi L, Chen C H. A new algorithm for stochastic discrete resource allocation optimization. *Discrete Event Dynamic Systems*, 2000, **10**(3): 271–294
- 136 Ho Y C, Zhao Q C, Jia Q S. *Ordinal Optimization: Soft Optimization for Hard Problems*. New York: Springer, 2007
- 137 Chen H C, Dai L Y, Chen C H, Yücesan E. New development of optimal computing budget allocation for discrete event simulation. In: Proceedings of the 1997 Winter Simulation Conference. Washington, DC: IEEE, 1997. 334–341
- 138 Chen H C, Chen C H, Lin J, Yücesan E. An asymptotic allocation for simultaneous simulation experiments. In: Proceedings of the 1999 Winter Simulation Conference. Phoenix, AZ: IEEE, 1997. 359–366
- 139 Chen C H, Yücesan E. An alternative simulation budget allocation scheme for efficient simulation. *International Journal of Simulation and Process Modelling*, 2005, **1**(1): 49–57
- 140 Branke J, Chick S E, Schmidt C. Selecting a selection procedure. *Management Science*, 2007, **53**(12): 1916–1932
- 141 Fu M C, Hu J Q, Chen C H, Xiong X. Optimal computing budget allocation under correlated sampling. In: Proceedings of the 2004 Winter Simulation Conference. Washington, DC, USA: IEEE, 2004. 595–603
- 142 Fu M C, Hu J Q, Chen C H, Xiong X. Simulation allocation for determining the best design in the presence of correlated sampling. *INFORMS Journal on Computing*, 2007, **19**(1): 101–111
- 143 Glynn P, Juneja S. A large deviations perspective on ordinal optimization. In: Proceedings of the 2004 Winter Simulation Conference. Washington, DC, USA: IEEE, 2004. 577–585
- 144 Chick S E, Wu Y Z. Selection procedures with frequentist expected opportunity cost bounds. *Operations Research*, 2005, **53**(5): 867–878
- 145 He D H, Chick S E, Chen C H. Opportunity cost and OCBA selection procedures in ordinal optimization for a fixed number of alternative systems. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, 2007, **37**(5): 951–961
- 146 Trailovic L, Pao L Y. Computing budget allocation for efficient ranking and selection of variances with application to target tracking algorithms. *IEEE Transactions on Automatic Control*, 2004, **49**(1): 58–67
- 147 Chen C H, He D, Fu M, Lee L H. Efficient simulation budget allocation for selecting an optimal subset. *INFORMS Journal on Computing*, 2008, **20**(4): 579–595
- 148 Bartz-Beielstein T, Friese M, Zaefferer M, Naujoks B, Flasch O, Konen W, Koch P. Noisy optimization with sequential parameter optimization and optimal computational budget allocation. In: Proceedings of the 13th Annual Conference Companion on Genetic and Evolutionary Computation. New York, USA: ACM, 2011. 119–120
- 149 Chen C H, Lee L H. *Stochastic Simulation Optimization: An Optimal Computing Budget Allocation*. Singapore: World Scientific Publishing Company, 2011
- 150 Kirkpatrick S, Gelatt C D Jr, Vecchi M P. Optimization by simulated annealing. *Science*, 1983, **220**(4598): 671–680
- 151 Alkhamis T M, Ahmed M A, Tuan V K. Simulated annealing for discrete optimization with estimation. *European Journal of Operational Research*, 1999, **116**(3): 530–544
- 152 Alrefaei M H, Andradóttir S. A simulated annealing algorithm with constant temperature for discrete stochastic optimization. *Management Science*, 1999, **45**(5): 748–764
- 153 Rosen S L, Harmonosky C M. An improved simulated annealing simulation optimization method for discrete parameter stochastic systems. *Computers and Operations Research*, 2005, **32**(2): 343–358
- 154 Glover F. Tabu search — part I. *INFORSA Journal on Computing*, 1989, **1**(3): 190–206
- 155 Glover F. Tabu search — part II. *INFORSA Journal on Computing*, 1990, **2**(1): 4–32
- 156 Hu N F. Tabu search method with random moves for globally optimal design. *International Journal for Numerical Methods in Engineering*, 1992, **35**(5): 1055–1070
- 157 Lopez-Garcia L, Posada-Bolivar A. A simulator that uses Tabu search to approach the optimal solution to stochastic inventory models. *Computers and Industrial Engineering*, 1999, **37**(1): 215–218

- 158 Lutz C M, Roscoe Davis K, Sun M. Determining buffer location and size in production lines using Tabu search. *European Journal of Operational Research*, 1998, **106**(2): 301–316
- 159 Martin A D, Chang T M, Yih Y, Kincaid R K. Using Tabu search to determine the number of kanbans and lotsizes in a generic kanban system. *Annals of Operations Research*, 1998, **78**: 201–217
- 160 Dengiz B, Alabas C. Simulation optimization using Tabu search. In: Proceedings of the 2000 Winter Simulation Conference. Orlando, FL: IEEE, 2000. 805–810
- 161 Yang T, Kuo Y, Chang I. Tabu-search simulation optimization approach for flow-shop scheduling with multiple processors — a case study. *International Journal of Production Research*, 2004, **42**(19): 4015–4030
- 162 Ryan J L, Bailey T G, Moore J T, Carlton W B. Reactive tabu search in unmanned aerial reconnaissance simulations. In: Proceedings of the 1998 Winter Simulation Conference. Washington, DC: IEEE, 1998. 873–880
- 163 Goldberg D E. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Boston: Addison-Wesley, 1989
- 164 Azadivar F, Tompkins G. Simulation optimization with qualitative variables and structural model changes: a genetic algorithm approach. *European Journal of Operational Research*, 1999, **113**(1): 169–182
- 165 Azadivar F, Wang J. Facility layout optimization using simulation and genetic algorithms. *International Journal of Production Research*, 2000, **38**(17): 4369–4383
- 166 Wang L. A hybrid genetic algorithm — neural network strategy for simulation optimization. *Applied Mathematics and Computation*, 2005, **170**(2): 1329–1343
- 167 Pasandideh S H R, Niaki S T A. Multi-response simulation optimization using genetic algorithm within desirability function framework. *Applied Mathematics and Computation*, 2006, **175**(1): 366–382
- 168 Daniel J, Rajendran C. A simulation-based genetic algorithm for inventory optimization in a serial supply chain. *International Transactions in Operational Research*, 2005, **12**(1): 101–127
- 169 Glover F. Heuristics for integer programming using surrogate constraints. *Decision Sciences*, 1977, **8**(1): 156–166
- 170 Glover F. A template for scatter search and path relinking. In: Proceedings of the 1998 Selected Papers from the Third European Conference on Artificial Evolution. Berlin Heidelberg: Springer, 1998. 3–54
- 171 Glover F, Laguna M, Martí R. Fundamentals of scatter search and path relinking. *Control and Cybernetics*, 2000, **39**(3): 653–684
- 172 Laguna M, Martí R, Martí R C. *Scatter Search: Methodology and Implementation in C*. New York: Springer, 2003
- 173 Shi L, Ólafsson S. *Nested Partitions Method, Theory and Applications*. New York: Springer, 2009
- 174 Shi L, Ólafsson S, Sun N. New parallel randomized algorithms for the traveling salesman problem. *Computers and Operations Research*, 1999, **26**(4): 371–394
- 175 Shi L, Ólafsson S. Nested partitions method for global optimization. *Operations Research*, 2000, **48**(3): 390–407
- 176 Shi L Y, Ólafsson S. Convergence rate of the nested partitions method for stochastic optimization. *Methodology and Computing in Applied Probability*, 2000, **2**(1): 37–58
- 177 Shi L Y, Ólafsson S. Nested partitions method for stochastic optimization. *Methodology and Computing in Applied Probability*, 2000, **2**(3): 271–291
- 178 Shi L Y, Chen C H. A new algorithm for stochastic discrete resource allocation optimization. *Discrete Event Dynamic Systems*, 2000, **10**(3): 271–294
- 179 Shi L Y, Ólafsson S, Chen Q. An optimization framework for product design. *Management Science*, 2001, **47**(12): 1681–1692
- 180 Shi L Y, Meyer R R, Bozday M, Andrew J. A nested partitions framework for solving large-scale multicommodity facility location problems. *Journal of Systems Science and Systems Engineering*, 2004, **13**(2): 158–179
- 181 D'Souza W D, Zhang H H, Nazareth D P, Shi L Y, Meyer R R. A nested partitions framework for beam angle optimization in intensity-modulated radiation therapy. *Physics in Medicine and Biology*, 2008, **53**(12): 3293–3307
- 182 Yau H, Shi L Y. Nested partitions for the large-scale extended job shop scheduling problem. *Annals of Operations Research*, 2009, **168**(1): 23–39
- 183 Wu W, Wei J H, Guan X H, Shi L Y. A hybrid nested partitions algorithm for scheduling flexible resource in flow shop problem. *International Journal of Production Research*, 2012, **50**(10): 2555–2569
- 184 Ólafsson S. 2006. Metaheuristics. *Handbooks in Operations Research and Management Science: Simulation, Chapter 21* (Henderson S G, Nelson B L (editor)). Amsterdam: Elsevier, 2006. 633–654
- 185 Nelder J A, Mead R. A simplex method for function minimization. *The Computer Journal*, 1965, **7**(4): 308–313
- 186 Azadivar F, Lee Y H. Optimization of discrete variable stochastic systems by computer simulation. *Mathematics and Computers in Simulation*, 1988, **30**(4): 331–345
- 187 Barton R R, Ivey J S. Nelder-Mead simplex modifications for simulation optimization. *Management Science*, 1996, **42**(7): 954–973
- 188 Hooke R, Jeeves T A. A direct search solution of numerical and statistical problems. *Journal of the ACM (JACM)*, 1961, **8**(2): 212–229
- 189 Chen M C, Tsai D M. Simulation optimization through direct search for multi-objective manufacturing systems. *Production Planning and Control*, 1996, **7**(6): 554–565
- 190 Ho Y C, Shi L, Dai L, Gong W B. Optimizing discrete event dynamic systems via the gradient surface method. *Discrete Event Dynamic Systems*, 1992, **2**(2): 99–120
- 191 Dorigo M, Stützle T. *Ant Colony Optimization*. Cambridge, MA: MIT Press, 2004
- 192 Pasupathy R, Henderson S G. A testbed of simulation-optimization problems. In: Proceedings of the 2006 Winter Simulation Conference. Monterey, CA: IEEE, 2006. 255–263



WANG Long-Fei Ph.D. candidate in the Department of Industrial Engineering and Management, Peking University. He received his bachelor degree from Peking University in 2007. His research interest covers optimization theory and methodology, operation management and system optimization.
E-mail: wanglongfei@pku.edu.cn



SHI Le-Yuan Professor in the Department of Industrial Engineering and Management, Peking University. She received her Ph.D. degree in applied mathematics from Harvard University in 1992, her master degree in engineering from Harvard University in 1990, her master degree in applied mathematics from Tsinghua University in 1985, and her bachelor degree in mathematics from Nanjing Normal University in 1982. Her research interest covers theory and development of large-scale optimization algorithms, discrete event simulation methodology, and modeling and analysis of discrete dynamic systems, with applications to complex systems such as supply chain networks, manufacturing systems, communication networks, and financial engineering. She is a member of INFORMS, and a fellow of IEEE. Corresponding author of this paper. E-mail: leyuan@coe.pku.edu.cn