# EFFICIENT SUBOPTIMAL RARE-EVENT SIMULATION

Xiaowei Zhang

Department of Management Science and Engineering
Stanford University
Stanford, C.A. 94305, U.S.A.

Jose Blanchet

Department of Statistics
Harvard University
Cambridge, M.A. 02138, U.S.A

Peter W. Glynn

Department of Management Science and Engineering
Stanford University
Stanford, C.A. 94305, U.S.A.

## ABSTRACT

Much of the rare-event simulation literature is concerned with the development of asymptotically optimal algorithms. Because of the difficulties associated with applying these ideas to complex models, this paper focuses on sub-optimal procedures that can be shown to be much more efficient than conventional crude Monte Carlo. We provide two such examples, one based on "repeated acceptance/rejection" as a mean of computing tail probabilities for hitting time random variables and the other based on filtered conditional Monte Carlo.

## 1 INTRODUCTION

The rare-event simulation problem is concerned with using simulation to compute $\alpha = P(A)$, where $A$ is a "rare-event" (and hence has probability $P(A)$ close to zero). It is well known that if $\alpha$ is computed via crude Monte Carlo (i.e., by estimating $\alpha$ via the proportion of independent simulations on which the event $A$ occurs), the number of trials $n$ required to estimate $\alpha$ to a given relative precision scales is in proportion to $1/P(A)$. As a consequence, crude Monte Carlo (CMC) is a highly inefficient algorithm for computing $\alpha$ when $P(A)$ is small. Much of the rare-event simulation literature is concerned with the development of asymptotically optimal algorithms. Because of the difficulties associated with applying these ideas to complex models, this paper focuses on sub-optimal procedures that can be shown to be much more efficient than conventional crude Monte Carlo. We provide two such examples, one based on "repeated acceptance/rejection" as a mean of computing tail probabil-

ities for hitting time random variables and the other based on filtered conditional Monte Carlo.

It follows that for very rare events (say, of order $10^{-4}$ or smaller), there is great interest (both practically speaking and academically) in using modified simulation algorithms capable of computing $\alpha$ more efficiently (i.e., using a so-called "efficiency improvement technique"). Accordingly, there is now a substantial literature on such efficient rare-event simulation algorithms; see, for example, Bucklew (2004) and Juneja and Shahabuddin (2007).

Given a family of problem instances $(P(A_\theta) : \theta \in \Lambda)$, let $Z(\theta)$ be a random variable having mean equal to $P(A_\theta)$ for each $\theta \in \Lambda$. The family of r.v.'s $(Z(\theta) : \theta \in \Lambda)$ is said to have *bounded relative variance* (BRV) if

$$\sup_{\theta \in \Lambda} \frac{\mathrm{var}Z(\theta)}{P(A_\theta)^2} < \infty.$$

Bounded relative variance implies, via Chebyshev's inequality, that the number of observations $n$ required to compute $P(A_\theta)$ to a given relative precision is bounded in $\theta \in \Lambda$. In particular, this implies that if $P(A_\theta) \downarrow 0$ as $\theta \to \theta_0$, then $EZ^2(\theta) = O(P(A_\theta)^2)$ as $\theta \to \theta_0$. Because the Cauchy-Schwarz inequality implies that $EZ^2(\theta) \geq P(A_\theta)^2$, it follows that $EZ^2(\theta)$ is within a constant multiple, uniformly in $\theta \in \Lambda$, of the smallest possible second moment. Hence, the family $(EZ^2(\theta) : \theta \in \Lambda)$ is within a constant multiple of optimality.

An extensive theory of asymptotic optimality based on the concept of BRV (and a weaker notion known as "logarithmic optimality") has developed in response to rare-event simulation problems in several applications domains. In order to guarantee asymptotic optimality, one clearly

needs an upper bound on $EZ^2(\theta)$. Many different methods are available for computing such bounds, varying from ad hoc ideas to the systematic use of Lyapunov bounds; see Blanchet and Glynn (2007a). But it is equally evident that one needs lower bounds on $P(A_\theta)^2$; we provide a new Lyapunov based lower bound in Section 4 of this paper. Note, however, that the square of such a lower bound must be within a constant multiple of the upper bound, in order that the family $(Z(\theta) : \theta \in \Lambda)$ possesses the BRV property. In fact, both bounds, must be with a constant multiple of the actual $P(A_\theta)^2$. Hence, asymptotic optimality can typically be verified only in settings in which the actual probability $P(A_\theta)$ can be calculated (without simulation) to within a constant multiple.

Given the difficulties associated with analytically calculating $P(A_\theta)$ to such a precision in many applied domains, this suggests focusing on computational strategies that, while not necessarily optimal, are at least demonstrably better than crude Monte Carlo. Our focus in this paper is on describing a couple of such methods that, while more efficient than CMC, are rarely optimal procedures (in the sense of possessing BRV). In addition to describing these methods, it is our hope that this paper will stimulate additional research on identifying rare-event simulation algorithms that are provably efficient (even if suboptimal). The first such method is "repeated acceptance/rejection", in which step-by-step conditioning is used to compute the tail probability for a "first hitting time" random variable. This method has been previously introduced and studied by Glasserman and Staum (2001) in the setting of barrier options. Our discussion, in Section 2, notes that this method is an example of an efficient suboptimal rare-event simulation algorithm. We also provide a new (simpler) proof of this method's efficiency, and provide a new analytical characterization of its efficiency in the above tail probability computational setting by appealing to Perron-Frobenius theory.

## 2  REPEATED ACCEPTANCE / REJECTION

Let $X = (X_m : m \geq 0)$ be a finite state irreducible discrete-time Markov chain living on state space $S$. We will describe a variance reduction technique for computing rare-event probabilities of the form $\alpha_n = P_x(T > n)$ when $n$ is large. Here, $T = \inf\{m \geq 0 : X_m \in C^c\}$. is the *hitting time* of some non-empty subset $C^c \subseteq S$. We assume that for each $y \in C$, $r(y) = P_y(X_1 \in C)$ can be computed.

Algorithm 1
1.  $m \leftarrow 0, X_m \leftarrow x, L \leftarrow 1$.
2.  $L \leftarrow Lr(X_m)$.
3.  Generate $Z$ from the probability distribution $P(X_{m+1} \in \cdot | X_m)$.
4.  If $Z \in C^c$, go to 3.
    Else, $m \leftarrow m+1$, $X_m \leftarrow Z$.

5.  If $m < n$, go to 2.
    Else, return $L$.

Note that Steps 2 and 3 can be viewed as an acceptance/rejection algorithm in which the target distribution is $P(X_{m+1} \in \cdot | X_m, X_{m+1} \in C)$. Thus, if one can efficiently generate transitions from the transition matrix $Q = (Q(y,z) : y, z \in C)$ with entries

$$Q(y,z) = \frac{P(y,z)}{r(y)}$$

using an alternative (more efficient) method, this alternative variate generation scheme can (and should) be used instead.

In the form specified by Algorithm 1, paths of length $n$ are generated via an acceptance/rejection step at each of the $n$ transitions, leading to our terminology "repeated acceptance/rejection". To analyze this algorithm, let $\widetilde{P}_y(\cdot)$ and $\widetilde{E}_y(\cdot)$ denote, respectively, the probability and expectation operators, conditional on $X_0 = y$, associated with generating paths of $X$ under the Markovian dynamics corresponding to the transition matrix $Q$. The probability $\alpha_n$ can then be represented as

$$\alpha_n = \widetilde{E}_x L_n$$

where

$$L_n = \prod_{i=0}^{n-1} r(X_i).$$

We now proceed to analyze the variance of the random variable $L_n$ under $\widetilde{P}_x$. Note that if $\widetilde{\text{var}}_y(\cdot)$ is the variance under $\widetilde{P}_y$, then

$$\begin{aligned}
\widetilde{\text{var}}_x L_n &= \widetilde{E}_x L_n^2 - \alpha_n^2 \\
&= E_x L_n I(\tau > n) - \alpha_n^2.
\end{aligned}$$

Since $L_n \leq 1$, it follows that

$$\begin{aligned}
\widetilde{\text{var}}_x L_n &\leq E_x I(\tau > n) - \alpha_n^2 \\
&= \alpha_n(1 - \alpha_n) \\
&= \text{var}_x I(\tau > n).
\end{aligned}$$

So, repeated acceptance/rejection clearly beats crude Monte Carlo.

But given the fact that $L_n$ is a product of $n$ terms, each of which is less than or equal to 1, the expectation $\widetilde{E}_x L_n^2$ is actually smaller than $E_x I(\tau > n)$ by a factor that is geometric in the time horizon $n$.

The degree of variance reduction can be quantified more precisely by appealing to Perron-Frobenius theory. Let $B = (B(y,z) : y, z \in C)$ be the restriction of $P$ to $C$ (so that $B(y,z) = P(y,z)$). If $B$ is irreducible, the Perron-

Frobenius theorem guarantees existence of a positive column eigenvector $v = (v(z) : z \in C)$ and a positive eigenvalue $\lambda$ such that $Bv = \lambda v$. Note that

$$B^n(y,z) = P_y(\Gamma > n, X_n = z).$$

Because $B^n v = \lambda^n v$ with $v$ strictly positive, we may conclude that $E_x I(\Gamma > n) = \theta(\lambda^n)$ as $n \to \infty$ (where $a_n = \theta(b_n)$ means that there exist positive finite constants $c_1$ and $c_2$ such that $c_1 \le a_n/b_n \le c_2$ for $n \ge 1$.)

Similarly, let $G = (G(y,z) : y, z \in C)$ be the matrix in which $G(y,z) = r(y)P(y,z)$. Again, if $G$ is irreducible, the Perron-Frobenius theory guarantees that $\widetilde{E}_x L_n^2 = \theta(\gamma^n)$ as $n \to \infty$, where $\gamma$ is the Perron-Frobenius eigenvalue of $G$. Because $\widetilde{E}_x L_n^2 \ge P_x(\Gamma > n)^2$, we conclude that $\gamma \ge \lambda^2$, while the inequality $\gamma \le \lambda$ follows from the fact that CMC has higher variance than repeated acceptance/rejection.

We expect repeated acceptance/rejection to be most efficient in settings in which the $r(X_i)$'s tend to be small over most of the time horizon $[0,n]$. For example, in simulating a stable queue, the method is likely to provide large variance reductions when computing (for example) the probability that the duration of the busy period exceeds $n$. On the other hand, the degree of variance reduction is likely to be smaller when computing the probability that no buffer overflow has occurred by time $n$. In the latter case, the stability of the queue tends to move the system rapidly away from those regions in which $r(\cdot)$ tends to be small, whereas for busy period computations, the tendency of the natural dynamics is to push the system towards regions in which $r(\cdot)$ is small.

It is worth noting that variants of this algorithm can also be applied when $r(\cdot)$ is unknown. In particular, $r(x)$ can be estimated by drawing a Bernoulli $(P_x(X_1 \in C))$ indicator random variable $I(x)$. A simple inductive argument establishes that $\prod_{j=0}^{n-1} I_j(X_j)$ has expectation $\alpha_n$, where $I_j(X_j)$ has distribution, conditional on $X_0, X_1, \ldots, X_j$, equal to $I(X_j)$. But $\prod_{j=0}^{n-1} I_j(X_j) \stackrel{\mathscr{D}}{=} I(T > n)$ (where $\stackrel{\mathscr{D}}{=}$ denotes "equality in distribution"), so that CMC is recovered and no variance reduction is achieved. However, we can also choose to estimate $r(x)$ via a binomial $(k(x), P_x(X_1 \in C))$ random variable. This leads to an estimator of $\alpha_n$ based on a product of such binomial random variables. Another estimator of $\alpha_n$ based on products of negative binomial random variables can also be derived; see Glasserman and Staum (2001) for a discussion and related numerical results.

We close this section with a brief discussion of computational efficiency for repeated acceptance/rejection. It should be noted that while $\text{var}\, L_n \le \text{var}\, I(T > n)$, the computational effort per observation for repeated acceptance/rejection is greater. Any comprehensive discussion of computational efficiency should reflect this additional work. It is well known (see, for example, Glynn and Whitt 1992) that the appropriate figure of merit for comparing computational efficiency of two competing Monte Carlo algorithms is

(average computer time per observation)

$\times$ (variance per observation);

the most efficient algorithm is the one with the smaller figure of merit. We have already computed the variances. A rough measure of the average computer time per observation $I(T > n)$ is $E_x \min(T,n)$.

On the other hand, the total number of acceptance/rejection steps required to generate $L_n$ is, on average, $E_x \sum_{j=0}^{n-1} 1/r(X_j)$. Because the ratio of variances decays geometrically in $n$, whereas, the ratio of average time per observation will often increase additively, it follows that repeated acceptance/rejection will generally be more computational efficient than CMC when the time horizon $n$ is large.

## 3 FILTERED CONDITIONAL MONTE CARLO

We discuss here, in the rare-event simulation setting, a variance reduction technique known variously as extended conditional Monte Carlo (see Section 2.6 of Bratley, Fox, and Schrage 1987) or filtered (conditional) Monte Carlo (see Glasserman 1993). To illustrate the idea, we consider its application in the setting of the (stable) single-server GI/G/1 queue under a first-in-first-out (FIFO) queue discipline. In particular, let $\alpha_w$ be the probability that some customer experiences a delay greater than $w$ within a busy cycle. To be more precise, let $W_n$ be the waiting time (exclusive of service) of the $i$'th customer, and suppose $W_0 = 0$. Note that $W = (W_n : n \ge 0)$ is a Markov chain living on state space $[0, \infty)$; see, for example, Asmussen (2003). Let $\tau = \inf\{n \ge 1 : W_n = 0\}$ be the index of the customer that initiates the second busy cycle. We are interested in computing

$$\alpha_w = P_0(T_w < \tau),$$

where $T_w = \inf\{n \ge 1 : W_n > w\}$ is the first hitting time for level $w$. To applying filtered conditional Monte Carlo here, note that

$$
\begin{aligned}
\alpha_w &= \sum_{n=1}^{\infty} P_0(T_w = n, T_w < \tau) \\
&= \sum_{n=1}^{\infty} E_0 P_0(T_w = n, T_w < \tau | W_0, \ldots, W_{n-1}) \\
&= \sum_{n=1}^{\infty} E_0 I(\min(T_w, \tau) > n - 1) P_0(W_n > w | W_n) \\
&= \sum_{n=0}^{\infty} E_0 I(\min(T_w, \tau) > n) \bar{F}(w - W_n)
\end{aligned}
$$

$$= \quad E_0 \sum_{n=0}^{\min(T_w,\tau)-1} \bar{F}(w - W_n),$$

where $\bar{F}(x) = P(V - \chi > x)$ and $V, \chi$ are independent realization of a typical service time and interarrival time for the queue.

Set $X = V - \chi$. If $X$ is a random variable bounded above by some constant $c$, then

$$\sum_{n=0}^{\min(T_w,\tau)-1} \bar{F}(w - W_n) = \sum_{n=0}^{\min(T_w,\tau)-1} \bar{F}(w - W_n) I(T_{w-c} < \tau), \tag{1}$$

since the $\bar{F}(w - W_n)$'s vanish unless $w - W - n \leq c$. Hence, the above estimator inherits many of the same inefficiencies as associated with computing $P_0(T_{w-c} < \tau)(= \alpha_{w-c})$ via CMC. Similarly, if the random variable $X$ has a light right tail (i.e., $P(X > x)$ can be bounded above by some multiple of a decaying exponential in $x$), we should not expect a variance reduction (over CMC) that gets progressing larger when $w$ increases.

On the other hand, if the random variable $X$ has a heavy tail, the $\bar{F}(w - W_n)$'s are (over a typical busy cycle) of the same order of magnitude. To be specific, assume that $X$ has a Pareto-like right tail, go that $P(X > x) \sim cx^{-\alpha}$ as $x \to \infty$ for some $c \in (0, \infty)$ and $\alpha > 1$. (Here, $f(x) \sim g(x)$ as $x \to \infty$ means that $f(x)/g(x) \to 1$ as $x \to \infty$.) In this case,

$$\frac{\sum_{n=0}^{\min(T_w,\tau)-1} \bar{F}(w - W_n)}{\bar{F}(w)} \Rightarrow \tau$$

as $w \to \infty$. Furthermore, under our tail assumptions, $\alpha_w = O(\bar{F}(w))$ as $w \to \infty$ (see, for instance, Blanchet, Glynn, and Liu 2007). So, the estimator (1) is, with high probability, of the same order of magnitude as $\alpha_w$, suggesting that the coefficient of variance of the estimator, as a function of $w$, should be better bounded than that of the crude Monte Carlo estimator. This intuition is verified in Blanchet and Glynn (2007b), where it is shown that

$$\frac{\text{var}\left(\sum_{n=0}^{\min(T_w,\tau)-1} \bar{F}(w - W_n)\right)}{\text{var } I(T_w < \tau)} = O\left(\frac{1}{w}\right),$$

as $w \to \infty$.

In general, the application of filtered conditional Monte Carlo in the rare-event setting leads to estimators in which the one-step transition probabilities of entering the rare set appear. For models in which the dynamics of the system, conditional on the rare event, are largely determined by a single unusual transition (rather than by a long series of slightly unusual transitions), filtered conditional Monte Carlo can lead to significant variance reduction relative crude Monte Carlo. In particular, the one-step conditional probability integrates out the randomness that is present in the single unusual transition, leading to substantial possible variance reduction.

## 4 A LOWER BOUND FOR $P(A)$

We describe here a Lyapunov-based method for deriving lower bounds on a certain class of rare-event probabilities. In particular, we are concerned here with exit probabilities of the form $P_x(X_T \in A, T < \infty)$, where $T = \inf\{n \geq 0 : X_n \in A \cup B\}$ is the first exit time form $A^c \cap B^c$. Here, $X = (X_n : n \geq 0)$ is an $S$-valued Markov chain, where the state space $S$ can be either discrete or continuous. Many rare-event simulation problems can be re-formulated as special cases of exit problems; see, for example, Blanchet and Glynn (2007a).

**Theorem 1** *Suppose there exist two non-negative finite-valued functions $g_1$ and $g_2$ satisfying the pair of inequalities*

$$E_x g_1(X_1) I(X_1 \in A^c \cap B^c) \geq g_1(X_1) - P_x(X_1 \in A), \quad x \in A^c \cap B^c$$

*and*

$$E_x g_2(X_1) I(X_1 \in A^c \cap B^c) \leq g_2(x) - g_1(x), \quad x \in A^c \cap B^c.$$

*Then, $P_x(X_T \in A, T < \infty) \geq g_1(x)$ for $x \in A^c \cap B^c$.*

<u>Proof.</u> Set $u^*(x) = P_x(X_T \in A, T < \infty)$. Then,

$$
\begin{aligned}
u^*(x) &= \sum_{n=0}^{\infty} \int_{A^c \cap B^c} P_x(X_n \in dy, T > n) P_y(X_1 \in A) \\
&= \sum_{n=0}^{\infty} \int_{A^c \cap B^c} K^n(x, dy) b(y)
\end{aligned}
$$

where $K = (K(x, dy) : x, y \in A^c \cap B^c)$ is the non-negative operator defined by

$$K(x, dy) = P_x(X_1 \in dy)$$

for $x, y \in A^c \cap B^c$ and $b = (b(y) : y \in A^c \cap B^c)$ is given by

$$b(y) = P_y(X_1 \in A)$$

for $y \in A^c \cap B^c$. The operator $K^n$ is then given (inductively) by

$$K^n(x, dy) = \int_{A^c \cap B^c} P_x(X_1 \in dy) K^{n-1}(y, dz)$$

for $n \geq 1$, with $K^0(x, dy) = \delta_x(dy)$ (i.e., a unit point mass measure at $x$).

We therefore need a lower bound on

$$\sum_{n=0}^{\infty}(K^n b)(x).$$

Note that the second Lyapunov inequality asserts that

$$(Kg_2)(x) \leq g_2(x) - g_1(x) \qquad (2)$$

for $x \in A^c \cap B^c$. Applying the non-negative operator $K^{n-1}$ to both sides of (2) yields the inequality

$$(K^n g_2)(x) \leq (K^{n-1} g_2)(x)$$

for $n \geq 1$, and hence $(K^n g_2)(\cdot)$ is finite-valued for $n \geq 0$. The inequality (2) also implies that

$$(K^n g_1)(x) \leq (K^n g_2)(x) - (K^{n+1} g_2)(x), \qquad (3)$$

where the right-hand side of (3) is now guaranteed to be the difference of two finite-valued quantities. We may now sum the inequalities (3) for $n = 0, 1, \ldots, m$, yielding

$$\sum_{n=0}^{m}(K^n g_1)(x) \leq g_2(x) - (K^{m+1} g_2)(x)$$
$$\leq g_2(x).$$

Sending $m \to \infty$, we obtain the upper bound

$$\sum_{n=0}^{\infty}(K^n g_1)(x) \leq g_2(x). \qquad (4)$$

It follows that $(K^n g_1)(x) \to 0$ as $n \to \infty$. On the other hand, the first Lyapunov inequality asserts that

$$(Kg_1)(x) \geq g_1(x) - b(x),$$

so that

$$b(x) \geq g_1(x) - (Kg_1)(x).$$

Again, applying $K^n$ to both sides of the above inequality, we conclude that

$$(K^n b)(x) \geq (K^n g_1)(x) - (K^{n+1} g_1)(x). \qquad (5)$$

Note that (4) guarantees that $(K^n g_1)(x)$ is finite-valued for each $n \geq 0$, so that the right-hand side of (5) is the difference of two finite-valued quantities. Summing the inequalities (5) for $n = 0, \ldots, m$, we find that

$$\sum_{n=0}^{m}(K^n b)(x) \geq g_1(x) - (K^{m+1} g_1)(x). \qquad (6)$$

Since $(K^n g_1)(x) \to 0$ as $n \to \infty$, we conclude, from (6), that by sending $m \to \infty$,

$$u^*(x) = \sum_{n=0}^{\infty}(K^n b)(x) \geq g_1(x),$$

proving the lower bound. $\square$

## REFERENCES

Asmussen, S. 2003. *Applied probability and queues*. 2nd ed. New York: Springer-Verlag.

Blanchet, J., and P. W. Glynn. 2007a. Efficient rare-event simulation for the maximum of the heavy-tailed random walk. Submitted for publication.

Blanchet, J., and P. W. Glynn. 2007b. Filtered monte carlo via taboo representations for heavy-tailed multiserver queues. In preparation.

Blanchet, J., P. W. Glynn, and J. C. Liu. 2007. Lyapunov bounds, fluid heuristics and state-dependent importance sampling for a heavy-tailed G/G/1 queue. Submitted for publication.

Bratley, P., B. L. Fox, and L. Schrage. 1987. *A guide to simulation*. 2nd ed. New York: Springer-Verlag.

Bucklew, J. 2004. *Introduction to rare-event simulation*. New York: Springer-Verlag.

Glasserman, P. 1993. Filtered monte carlo. *Mathematics of Operations Research* 18 (3): 610–634.

Glasserman, P., and J. Staum. 2001. Conditioning on one-step survival in barrier option simulations. *Operations Research* 49:923–937.

Glynn, P. W., and W. Whitt. 1992. The asymptotic efficiency of simulation estimators. *Operations Research* 40.

Juneja, S., and P. Shahabuddin. 2007. Rare event simulation techniques: An introduction and recent advances. In *Handbook on Simulation*, ed. S. Henderson and B. Nelson, 291–350. Elsevier.

## AUTHOR BIOGRAPHIES

**XIAOWEI ZHANG** is currently a second year Ph.D. student in the Department of Management Science and Engineering at Stanford University. Xiaowei Zhang graduated with B.Sc. in Mathematics from Nankai University, China. His research interests include applied probability, stochastic modeling and simulation.

**JOSE BLANCHET** is Assistant Professor of Statistics at Harvard University. Jose holds a M.Sc. in Engineering-Economics System and Operations Research and a Ph.D. in Management Science and Engineering, both from Stanford University. He also holds two B.Sc. degrees: one in Actuarial Science and another one in Applied Mathematics from ITAM (Mexico). Jose worked for two

years as an analyst in Protego Financial Advisors, a leading investment bank in Mexico. He has research interests in applied probability, computational finance, performance engineering, queue theory, risk management, rare-event analysis, statistical inference, stochastic modeling, and simulation.

**PETER W. GLYNN** received his Ph.D. in Operations Research from Stanford University in 1982. He then joined the faculty of the University of Wisconsin at Madison, where he held a joint appointment between the Industrial Engineering Department and Mathematics Research Center, and courtesy appointments in Computer Science and Mathematics. In 1987, he returned to Stanford, where he is now the Thomas Ford Professor of Engineering in the Department of Management Science and Engineering. He also has a courtesy appointment in the Department of Electrical Engineering and serves as Director of the Stanford Institute of Computational and Mathematical Engineering. He is a member of INFORMS and a fellow of the Institute of Mathematical Statistics and his research interests in computational probability, simulation, queueing theory, statistical inference for stochastic processes, and stochastic modeling.