# MATHEMATICS FOR SIMULATION

Shane G. Henderson

School of Operations Research and Industrial Engineering
Cornell University
Ithaca, NY 14853, U.S.A.

## ABSTRACT

I survey several mathematical techniques and results that are useful in the context of stochastic simulation. The concepts are introduced through the study of a simple model of ambulance operation to ensure clarity, concreteness and cohesion.

## 1  INTRODUCTION

The incredibly rich set of mathematical tools and techniques that underlie stochastic simulation is the subject of this paper. Of course, the field is far too large to be covered in a single paper, and so I choose to focus the discussion somewhat. For example, there is no discussion in this paper on the vast array of techniques that may be used for input analysis, uniform and nonuniform random variate generation, and sensitivity analysis. For excellent overviews of these and other topics, see Bratley, Fox and Schrage (1987), Law and Kelton (2000), and the tutorials and advanced tutorials in recent proceedings of the Winter Simulation Conference.

Instead, what I attempt to do is to describe a set of mathematical tools and techniques that can be used to explore density estimation of performance measures in both the terminating and steady-state simulation context. It would be very easy to provide a smorgasbord of such results, but such a paper would read like an encyclopedia. Therefore, many of the results are applied to a simple model of ambulance operation that serves to unify the discussion.

This paper is an updated version of Henderson (2000), in which there were 2 main topics. First, in the terminating simulation context, performance measures were rigorously defined through the strong law of large numbers for i.i.d. random variables. The performance of estimators of these performance measures was studied via the central limit theorem. Variants of these results were used to study performance measures that, instead of being expectations of random variables, were *functions* of expectations of random variables. Second, in the steady-state context, performance measures were rigorously defined and analyzed by appealing to asymptotic results for general state-space Markov chains. Lyapunov conditions were used to establish that the asymptotic results held.

All of the performance measures described in Henderson (2000) take the form of an expectation of a random variable, or a differentiable function of a finite number of expectations. Such performance measures are particularly useful when the goal is to compare many different stochastic systems, as they provide a concrete basis for the comparison. However, if the goal is to enhance one's *understanding* of a single stochastic system, then it is often more useful to analyze the *distribution* of certain random variables, perhaps through density estimation techniques.

In this paper we switch the focus to estimating the densities of random variables related to these performance measures. The goal remains the same as in Henderson (2000), namely to demonstrate mathematical tools and techniques that are useful in simulation analysis.

In Section 2 I review some approaches to density estimation in a particularly transparent context, namely that of estimating the density of the completion time of a stochastic activity network. The analysis in this section requires the use of the strong law of large numbers (SLLN) and central limit theorem (CLT). Any treatment of mathematics for simulation would be incomplete without these 2 fundamental results.

Section 2 then sets the stage for the remainder of the paper in which we analyze a simple model of ambulance operation in several ways. The model of ambulance operations is introduced in Section 3. Then, in Section 4 we specialize the model to the terminating simulation context. Even the *definition* of certain performance measures leads to the use of some interesting tools and techniques.

By imposing different assumptions on the model, one obtains a steady-state simulation, where the performance measures are all long-run averages. To rigorously define these performance measures, it is necessary to define an appropriate stochastic process with which to work. A great

deal is known about the class of Markov processes evolving on general (not necessarily countable) state spaces. In Section 5 a general state space Markov chain is defined. To ensure that long-run averages exist, it is necessary to show that this chain is, in a certain sense, positive recurrent.

A very practical approach to establishing that a Markov chain is positive recurrent is to use Lyapunov functions, and this approach is the central mathematical tool illustrated in Section 5. We use Lyapunov theory to show that certain Markov chains are positive recurrent, that our performance measures are well-defined, that certain density estimators are consistent and satisfy central limit theorems, and that confidence intervals obtained through the method of batch means are asymptotically valid. An important consideration in the steady-state context is that of initialization bias. We also use Lyapunov theory to characterize the magnitude of such bias.

The underlying theme of Section 5 is then that Lyapunov functions provide an enormously powerful, and easily applied (at least relative to many other methods) approach to establishing results that underlie steady-state simulation methodology.

Throughout this paper, results are rigorously quoted, and references given for the proofs. To simplify the exposition, it is often the case that results are quoted using stronger hypotheses than are strictly necessary, but tighter hypotheses can be found in the references provided. Notation is reused from section to section, but is consistently applied within each section.

## 2 DENSITY ESTIMATION

Our running example in this section will be that of a stochastic activity network (SAN).

A SAN is a directed graph that represents some set of activities that, taken together, represent some project/undertaking. Each arc in the graph represents a task that needs to be completed, and the (random) length of the arc represents the time required to complete the task. Nodes are used to indicate the precedence relationships between tasks. The time required to complete the project is indicated by the longest path between the designated "source" and "sink" nodes.

**Example 1.** *The simple stochastic activity network in Figure 1 is adapted from Avramidis and Wilson (1996). Nodes 1 and 9 are the source and sink nodes respectively. The labels on the arcs give the mean task durations (the distributions will be specified shortly), and task durations are independent.*

Let $L$ be the (random) network completion time. When does $L$ have a density?

Before answering this question, we first need to understand what we mean when we say that a random variable
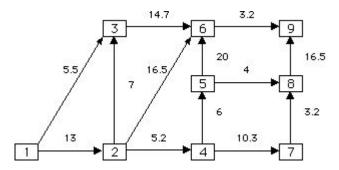


Figure 1: A Stochastic Activity Network with Mean Task Durations.

has a given density. The standard definition is based on the Radon-Nikodym theorem; see Theorem 31.7 on p. 434, and Theorem 32.2 on p. 443 of Billingsley (1986).

**Definition 1.** *We say that a real-valued random variable X has a density if $P(X \in A) = 0$ for all Lebesgue measurable sets A with Lebesgue measure 0. This is equivalent to saying that there exists a nonnegative function f (the density function) with the property that for all $x \in \mathbf{R}$,*

$$F(x) = \int_{-\infty}^{x} f(y)\, dy,$$

*where F is the distribution function of X.*

The first part of this definition may not be familiar to those who have not had a measure-theoretic course in probability. Heuristically speaking, $X$ has a density if the probability that $X$ takes on values in "insignificant" sets is 0. The second part of the definition is perhaps more familiar. We will use the 2 definitions interchangeably in what follows.

**Proposition 1.** *Consider a SAN with a finite number of arcs/tasks, where the individual task durations are independent. Suppose that every path from the source to the sink contains an arc for which the corresponding task duration has a density. Then the time to complete the project has a density.*

**Proof:** Let $P$ be a path from the source to the sink. The time $T$ required to traverse $P$ is a sum of the times required to traverse each arc in the path. One of these times has a density, and since the task durations are independent, it follows that $T$ has a density. Let $b < \infty$ denote the number of such paths from the source to the sink, and let $T_1, \ldots, T_b$ denote the times required to traverse each path. Then $L = \max\{T_1, \ldots, T_b\}$.

Now, the maximum of 2 random variables, $X$ and $Y$ say, that have densities, also has a density. To see why, let

$A$ denote an arbitrary (measurable) subset of the real line, and let $Z = \max\{X, Y\}$. Then

$$
\begin{aligned}
P(Z \in A) &\leq P(\{X \in A\} \cup \{Y \in A\}) \\
&\leq P(X \in A) + P(Y \in A)
\end{aligned}
$$

where the first inequality follows since $Z \in A$ implies that at least one of $X$ and $Y$ must be in $A$, and the second is Boole's inequality. Now, we know that $X$ and $Y$ have densities, so if the Lebesgue measure of $A$ is 0, then $P(X \in A) = 0$ and $P(Y \in A) = 0$. Hence $P(Z \in A) = 0$, and so, by Definiton 1, $Z$ has a density.

We can now apply this result inductively to $L = \max\{T_1, \ldots, T_b\}$ to conclude that $L$ has a density. $\square$

Applying this result to Example 1, we see that the network completion time $L$ will have a density if tasks (6, 9) and (8, 9) have densities. But then, how can we estimate this density?

The look-ahead density estimators developed by Henderson and Glynn (2001) are easily analyzed, and have excellent statistical properties. In fact, Example 1 is one of the examples given in that paper. Let us see how to construct a look-ahead density estimator for the completion time $L$ in Example 1.

Let $L_6$ and $L_8$ be the lengths of the longest paths from the source node to nodes 6 and 8 respectively, and let $G$ be their joint distribution function, so that $G(l_6, l_8) = P(L_6 \leq l_6, L_8 \leq l_8)$. Let $T_{69}, T_{89}$ denote the task durations for tasks (6, 9) and (8, 9), let $F_{69}, F_{89}$ be the corresponding distribution functions, and $f_{69}, f_{89}$ be the corresponding densities. Then

$$
\begin{aligned}
&P(L \leq t) \\
&= E\, P(L \leq t | L_6, L_8) \\
&= E\, P(T_{69} \leq t - L_6, T_{89} \leq t - L_8 | L_6, L_8) \\
&= E[P(T_{69} \leq t - L_6 | L_6) P(T_{89} \leq t - L_8 | L_8)] \\
&= E[F_{69}(t - L_6) F_{89}(t - L_8)] \\
&= \iint F_{69}(t - u) F_{89}(t - v)\, dG(u, v) \\
&= \iint \int_{-\infty}^{t} \frac{d}{dx}[F_{69}(x - u) F_{89}(x - v)]\, dx\, dG(u, v) \\
&= \iint \int_{-\infty}^{t} [F_{69}(x - u) f_{89}(x - v) \\
&\qquad + f_{69}(x - u) F_{89}(x - v)]\, dx\, dG(u, v) \\
&= \int_{-\infty}^{t} \iint [F_{69}(x - u) f_{89}(x - v) \\
&\qquad + f_{69}(x - u) F_{89}(x - v)]\, dG(u, v)\, dx \\
&= \int_{-\infty}^{t} E[F_{69}(x - L_6) f_{89}(x - L_8) \\
&\qquad + f_{69}(x - L_6) F_{89}(x - L_8)]\, dx.
\end{aligned}
$$

Thus, we can conclude that $L$ has a density $f$ say, where $f(x)$ is given by

$$
E[F_{69}(x - L_6) f_{89}(x - L_8) + f_{69}(x - L_6) F_{89}(x - L_8)]. \quad (1)
$$

(See Section 4.1 of Avramidis and Wilson (1996) for a related discussion.)

The expression (1) has an intuitive interpretation. The first term in (1) is related to the probability that the longest path from the source to the sink through node 6 has length at most $x$ and at the same time, the longest path from the source to the sink through node 8 is exactly of length $x$. The second term can be interpreted similarly.

The expression (1) immediately suggests a density estimator for $f$. We generate i.i.d. replicates $L_6(i), L_8(i)$ for $i = 1, \ldots, n$, and estimate $f(x)$ by

$$
f_n(x) = \frac{1}{n} \sum_{i=1}^{n} L(i; x),
$$

where

$$
\begin{aligned}
L(i; x) &= F_{69}(x - L_6(i)) f_{89}(x - L_8(i)) \\
&\quad + f_{69}(x - L_6(i)) F_{89}(x - L_8(i)).
\end{aligned}
$$

The consistency of $f_n(x)$ follows from the strong law of large numbers (SLLN).

**Theorem 2 (SLLN).** *If $X_1, X_2, \ldots$ is an i.i.d. sequence of random variables with $E|X_1| < \infty$, then*

$$
\frac{\sum_{i=1}^{n} X_i}{n} \to EX_1 \ a.s.
$$

*as $n \to \infty$.*

For a proof, see p. 290 of Billingsley (1986).

Applying this result to our context, we note that $f_n(x) \to f(x)$ a.s., for all $x$ for which (1) is finite. The set of values $x$ for which this does not hold has Lebesgue measure 0, and so the convergence of $f_n$ to $f$ occurs for almost all $x$. This is the best that can be hoped for, because densities are only defined up to a set of Lebesgue measure 0.

An estimate of the accuracy of this estimator can be obtained through the central limit theorem (CLT).

**Theorem 3 (CLT).** *If $X_1, X_2, \ldots$ is an i.i.d. sequence of random variables with $EX_1^2 < \infty$, then*

$$
\sqrt{n}\left(\frac{1}{n} \sum_{i=1}^{n} X_i - EX_1\right) \Rightarrow \sigma N(0, 1)
$$

as $n \to \infty$, *where* $\sigma^2 = \text{var} \, X_1$, $\Rightarrow$ *denotes weak convergence, and* $N(0, 1)$ *denotes a standard normal random variable.*

For a proof, see p. 367 of Billingsley (1986).

We can apply the CLT to estimate the error in $f_n(x)$ if $E[L(1; x)^2] < \infty$, which is implied by the simpler condition that

$$E[f_{69}^2(x - L_6) + f_{89}^2(x - L_8)] < \infty. \qquad (2)$$

If (2) holds, then the CLT basically establishes that the error in the estimator $f_n(x)$ is asymptotically normally distributed with mean 0 and variance $s^2(x)/n$ where $s^2(x) = \text{var} \, L(1; x)$, and this is the basis for obtaining confidence intervals for $f(x)$. In particular, an approximate 95% confidence interval for $f(x)$ is given by

$$f_n(x) \pm 1.96 \sqrt{\frac{s^2(x)}{n}}. \qquad (3)$$

Of course, $s^2(x)$ must invariably be estimated. The usual estimator is the sample variance

$$s_n^2(x) = \frac{1}{n - 1} \sum_{i=1}^{n} (L(i; x) - f_n(x))^2.$$

The confidence interval that is reported is the same as (3) with $s^2(x)$ replaced with its sample counterpart $s_n^2(x)$. But is the modified confidence interval then valid?

If (2) holds, then the SLLN implies that $s_n^2(x) \to s^2(x)$ a.s. as $n \to \infty$. Hence, by Exercise 29.4 of Billingsley (1986), we have that

$$\begin{pmatrix} n^{1/2}(f_n(x) - f(x)) \\ s_n^2(x) \end{pmatrix} \Rightarrow \begin{pmatrix} s(x)N(0, 1) \\ s^2(x) \end{pmatrix}. \qquad (4)$$

The natural tool to apply at this point is the continuous mapping theorem. For a real-valued function $h$ in $\mathbf{R}^d$, let $D_h$ denote its set of discontinuities (in $\mathbf{R}^d$).

**Theorem 4 (Continuous Mapping Theorem).** *Let* $(X_n : n \geq 1)$ *be a sequence of* $\mathbf{R}^d$ *valued random variables with* $X_n \Rightarrow X$ *as* $n \to \infty$ *and let* $h : \mathbf{R}^d \to \mathbf{R}$ *be measurable. If* $P(X \in D_h) = 0$, *then* $h(X_n) \Rightarrow h(X)$ *as* $n \to \infty$.

For a proof, see p. 391 of Billingsley (1986).

Define $h(x, y) = x/y^{1/2}$, and then apply the continuous mapping theorem to (4), to obtain that when $s^2(x) > 0$,

$$\frac{n^{1/2}(f_n(x) - f(x))}{s_n(x)} \Rightarrow N(0, 1)$$

as $n \to \infty$, and so the confidence interval procedure outlined above is indeed valid.

The look-ahead density estimator described above has very appealing and easily derived asymptotic properties. These attractive properties are a result of carefully investigating the model to identify exploitable properties.

One might ask if this somewhat specialized density estimation technique can be complemented by a more general-purpose approach that does not require as much tailoring to specific applications. The field of *nonparametric functional estimation* encompasses several such approaches. Prakasa Rao (1983) is an excellent survey of this field, especially *kernel density estimation*. We will not go into this area in any detail because the mathematical techniques used to analyze kernel density estimators are beyond the scope of this paper.

# 3 A SIMPLE MODEL

We now describe a very simple model that will serve as a vehicle for the concepts to follow. The purpose of the example is therefore simplicity, and certainly not realism, although with a few straightforward extensions, the model could be considered to be quite practical.

Suppose that a single ambulance serves calls in a square region. By translating and rescaling units, we may assume that the square is centred at the origin, with lower left-hand corner at $(-1/2, -1/2)$ and upper right-hand corner at $(1/2, 1/2)$. For simplicity, we assume that the ambulance travels at unit speed within the square. The combined hospital/ambulance base is located at the origin.

Calls arrive (in time) according to a homogeneous Poisson process with rate $\lambda$ calls per hour. The location of a call is independent of the arrival process, and uniformly distributed over the square. To serve a call, the ambulance travels in a Manhattan fashion (i.e., at any given time, movement is restricted to lie only in the $x$ direction or the $y$ direction) from its present location to the location of the call. A random amount of time is then spent at the scene treating the patient, independent of all else. After this scene time is complete, with probability $p$ (independent of all else), the ambulance is required to transport and admit the patient to the hospital, with hospital admission occurring instantaneously once the ambulance reaches the hospital, and with probability $1 - p$ the ambulance is freed for other work.

# 4 TERMINATING SIMULATION

In this section, we assume that the ambulance only receives calls from (say) 7am until 11pm each day. At 11pm, the ambulance completes the call that it is currently serving (if any) and returns to base. We will further assume that if the ambulance is engaged with a call when another call

is received, then some outside agency, such as another emergency service, handles the other call. Finally, we assume that the random variables associated with each day are independent of those for all other days.

We are interested in 3 performance measures.

$u$      Let $U$ be the (random) fraction of a 16 hour day that the ambulance is busy. There is some positive probability that $U = 0$. But $U$ also has a density $u$ on $(0, 1]$.

$\alpha$      The long-run fraction of calls attended by the ambulance.

$r$      The conditional density of the response time to a call given that the ambulance attends the call.

Let $U_i$ denote the fraction of a 16 hour day that the ambulance is busy on day $i$, so that the $U_i$s are i.i.d. replicates of $U$, and $0 \leq U_i \leq 1$. (We do not count any residual time after 11pm needed to complete any call in progress.) Then

$$
\begin{aligned}
P(U_1 = 0) &= P(0 \text{ calls are received}) \\
&= \exp(-16\lambda).
\end{aligned}
$$

But $U_1$ also has a density $u$ on the interval $(0, 1]$, and we can construct a look-ahead density estimator for $u$ by conditioning on all of the random variables that constitute a single day of operation of the ambulance *except* for the time of the arrival of the first call. We will not go into the details here because the mathematics of this approach are exactly as discussed in Section 2 (although some of the implementation details are also interesting), and we will study a look-ahead density estimator for $r$ in detail below.

Let us turn to $\alpha$, the long-run fraction of calls attended by the ambulance.

Let $N_i$ denote the total number of calls received on day $i$, and for $j = 1, \ldots, N_i$, let $A_{ij}$ be 1 if the ambulance is available when the $j$th call arrives on day $i$ and 0 otherwise. Then the number of calls $A_i$ attended by the ambulance on day $i$ is $\sum_{j=1}^{N_i} A_{ij}$. After $n$ days, the fraction of calls attended by the ambulance is given by

$$
\frac{\sum_{i=1}^{n} A_i}{\sum_{i=1}^{n} N_i}. \tag{5}
$$

Dividing both the numerator and denominator of (5) by $n$, and applying the SLLN separately to both the numerator and denominator, we see that

$$
\frac{\sum_{i=1}^{n} A_i}{\sum_{i=1}^{n} N_i} \to \alpha = \frac{E A_1}{E N_1} \text{ a.s.}
$$

as $n \to \infty$. But $E N_1 = 16\lambda$, so we can estimate $\alpha$ by

$$
\alpha_n = \frac{1}{n} \sum_{i=1}^{n} \frac{A_i}{16\lambda}.
$$

The SLLN establishes that $\alpha_n$ is a consistent estimator of $\alpha$, and the CLT allows us to construct confidence intervals for $\alpha$ based on $\alpha_n$.

Finally we develop a look-ahead density estimator for $r$, the density of the response time. Before doing so, let us make sure that we understand what the density $r$ represents. For $x > 0$ and $\delta > 0$ sufficiently small, we want

$$
r(x)\delta \approx \lim_{n \to \infty} \frac{\sum_{i=1}^{n} \sum_{j=1}^{N_i} I(R_{ij} \in (x, x + \delta))}{\sum_{i=1}^{n} A_i} \tag{6}
$$

where $R_{ij}$ is the time required for the ambulance to reach the $j$th call on day $i$ if the ambulance responds to that call, and is (somewhat arbitrarily) -1 otherwise.

The right-hand side of (6) represents the limiting value of the number of calls responded to with a response time in the interval $(x, x + \delta)$ expressed as a fraction of the number of calls that the ambulance responded to. We have written "$\approx$" to reflect the more rigorous notion that if both sides of (6) are divided by $\delta$ and then $\delta \to 0$ from above, then the resulting limits are equal, at least at points $x$ at which $r$ is continuous.

To derive a look-ahead density estimator of $r$ we proceed as follows. For $i \geq 1$ and $j = 1, \ldots, N_i$, let $B_{ij}$ ($C_{ij}$) denote the vector location of the ambulance (new call) at the time at which the $j$th call on day $i$ is received. Let $d(b, c)$ denote the time required for the ambulance to travel from location $b$ to location $c$. For $x > 0$,

$$
P(R_{ij} \in dx \mid A_{ij}, B_{ij}) = A_{ij} P(d(B_{ij}, C_{ij}) \in dx),
$$

i.e., the response time will be $x$ if the ambulance is available when the call arrives and the time required for the ambulance to reach the call from its current location is exactly $x$.

Suppose that $C$ represents a random call location, and is uniformly distributed over the square as our model assumes. It immediately follows that for all fixed locations $b$ contained within the square,

$$
P(d(b, C) \in dx) = f(x; b) \, dx \tag{7}
$$

for an easily computed function $f(\cdot; \cdot)$. (The function $f$ gives the total length of all points within the square that are the fixed $L_1$ distance $x$ from the point $v$. It can be visualized by superimposing a diamond representing all points at a fixed $L_1$ distance $x$ from $v$, and measuring the length of the portions of the diamond that fall within the square.)

The above discussion establishes that we can estimate the density $r$ using the estimator

$$r_n(x) = \frac{\sum_{i=1}^n \sum_{j=1}^{N_i} A_{ij} f(x; B_{ij})}{\sum_{i=1}^n A_i}. \qquad (8)$$

Dividing both the numerator and denominator of (8) by $n$, we find that

$$r_n(x) \to \frac{E Y_1(x)}{E A_1}$$

almost surely as $n \to \infty$, where

$$Y_i(x) \stackrel{\triangle}{=} \sum_{j=1}^{N_i} A_{ij} f(x; B_{ij}).$$

So how can we assess the accuracy of the estimator $r_n(x)$? Certainly, the standard central limit theorem cannot be applied, because $r_n(x)$ is a *ratio* of sample means of i.i.d. observations. We first consider a strongly related question, and then return to the problem at hand.

Suppose that $X_1, X_2, \ldots$ is an i.i.d. sequence of random variables with finite mean $\mu = E X_1$. Let $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ denote the sample mean. If the real-valued function $h$ is continuous at $\mu$, it follows that $h(\bar{X}_n) \to h(\mu)$ a.s. as $n \to \infty$. So how does the error $h(\bar{X}_n) - h(\mu)$ behave, for large $n$? Note that for large $n$, $\bar{X}_n$ will be very close to $\mu$, and so the asymptotic behaviour of the error should depend only on the local behaviour of $h$ near $\mu$. Indeed, if $h$ is appropriately differentiable, then Taylor's theorem implies that

$$h(\bar{X}_n) - h(\mu) \approx h'(\mu)(\bar{X}_n - \mu),$$

and so if the $X_i$s have finite variance, then

$$\begin{aligned} n^{1/2}(h(\bar{X}_n) - h(\mu)) & \approx & h'(\mu) n^{1/2}(\bar{X}_n - \mu) \\ & \Rightarrow & \eta N(0, 1) \end{aligned}$$

as $n \to \infty$, where $\eta^2 = h'(\mu)^2 \operatorname{var} X_1$.

This intuitive argument can be formalized, and also generalized to higher dimensions to obtain the following result, sometimes referred to as the delta method.

**Theorem 5.** *Suppose that $(X_n : n \geq 1)$ is an i.i.d. sequence of $\mathbf{R}^d$ valued random variables with $E \|X_1\|_2^2 < \infty$. Let $\mu = E X_1$ denote their common mean, and let $\Lambda$ denote their common covariance matrix. Let $\bar{X}_n$ denote the sample mean of $X_1, \ldots, X_n$. If $h : \mathbf{R}^d \to \mathbf{R}$ is continuously differentiable in a neighbourhood of $\mu$ with nonzero gradient $g = \nabla h(\mu)$*

*at $\mu$, then*

$$n^{1/2}(h(\bar{X}_n) - h(\mu)) \Rightarrow \sigma N(0, 1)$$

*as $n \to \infty$, where $\sigma^2 = g' \Lambda g$.*

For a proof, see p. 122 of Serfling (1980).

To apply this result in our context, let

$$X_i(x) = (Y_i(x), A_i),$$

and define $h(y, a) = y/a$. Theorem 5 then implies that

$$n^{1/2}(r_n(x) - r(x)) \Rightarrow \sigma(x) N(0, 1),$$

where

$$\sigma^2(x) = \frac{E(Y_1(x) - r(x) A_1)^2}{(E A_1)^2}.$$

Using the SLLN, one can easily show that $\sigma^2(x)$ can be consistently estimated by

$$s_n^2(x) = \frac{n^{-1} \sum_{i=1}^n (Y_i(x) - r_n(x) A_i)^2}{(n^{-1} \sum_{i=1}^n A_i)^2},$$

and the same continuous mapping argument discussed earlier establishes that

$$r_n(x) \pm 1.96 s_n(x)/\sqrt{n}$$

is an approximate 95% confidence interval for $r(x)$.

Observe that the estimator $r_n(x)$, being a ratio estimator, is biased. Taylor's theorem can be used to examine this bias. In particular, reverting to our 1-dimensional digression for the moment, Taylor's theorem implies that

$$h(\bar{X}_n) - h(\mu) \approx h'(\mu)(\bar{X}_n - \mu) + \frac{1}{2} h''(\mu)(\bar{X}_n - \mu)^2.$$

Taking expectations, we find that

$$E h(\bar{X}_n) - h(\mu) \approx \frac{1}{2} h''(\mu) \operatorname{var} X_1/n,$$

i.e., we have an explicit expression for the asymptotic bias. As before, this argument can be formalized, and generalized to higher dimensions.

**Theorem 6.** *Suppose that $(X_n : n \geq 1)$ is an i.i.d. sequence of $\mathbf{R}^d$ valued random variables with $E \|X_1\|_2^4 < \infty$. Let $\mu = E X_1$ denote their common mean, and let $\Lambda$ denote their common covariance matrix. Let $\bar{X}_n$ denote the sample mean of $X_1, \ldots, X_n$. If $h : \mathbf{R}^d \to \mathbf{R}$ is such that $h(\bar{X}_n)$ is bounded for all $n$ with probability 1, and twice continuously*

*differentiable in a neighbourhood of $\mu$, then*

$$n(Eh(\bar{X}_n) - h(\mu)) \rightarrow \frac{1}{2} \sum_{i,j=1}^{d} \nabla^2 h(\mu)_{ij} \Lambda_{ij}$$

*as $n \rightarrow \infty$.*

The proof is a slight modification of Theorem 7 in Glynn and Heidelberger (1990).

We would like to apply this result to the estimator $r_n(x)$. Define $h(y, a) = y/a$. The only condition that is not obviously satisfied is that $h(\bar{X}_n)$ is bounded for all $n$ with probability 1. The function $f$ as given in (7) is bounded by $2\sqrt{2}$, and so $h(\bar{X}_n(x)) = r_n(x)$ is also bounded by $2\sqrt{2}$. We have therefore established that the bias in the estimator $r_n(x)$ is of the order $n^{-1}$.

It is reasonable to ask whether this bias is sufficient to noticeably affect the performance of the confidence intervals produced earlier for a given runlength $n$. Recall that the widths of the confidence intervals are of the order $n^{-1/2}$. Thus, the bias decreases at a (much) faster asymptotic rate than the width of the confidence intervals, and so unless runlengths are quite small, it is reasonable to neglect bias.

## 5 STEADY-STATE SIMULATION

We now turn to useful mathematical techniques and results for steady-state simulation analysis. For this purpose, we will modify the assumptions of the previous section on the dynamics of the ambulance model. In particular, in addition to the assumptions given in Section 3, we assume that the ambulance operates 24 hours a day, 7 days a week. Furthermore, calls that arrive while the ambulance is busy are queued, and answered in first-in first-out order. Once the current call is complete, the ambulance then attends to the next call. Recall that a call is completed either at the scene (with probability $1 - p$), or when the ambulance drops the patient off at the hospital (with probability $p$).

For this model, 2 of the previous 3 performance measures are still relevant, but because the ambulance is now handling *all* calls, the fraction of calls answered by the ambulance ($\alpha$) is no longer of interest. For convenience, and also to refine the statement of the performance measures to our new setting, we restate the performance measures.

$\beta$     The long-run utilization of the ambulance, i.e., the percentage of time that the ambulance is occupied with a call.

$r$     The long-run density of the response time to a call.

Note that $\beta$ is a deterministic constant, while $r$ is a density function.

Both performance measures involve the term "long-run". In order that such long-run measures exist, it is first

necessary that the ambulance model be stable. In order to be able to make statements about the stability, or lack thereof, of the model, it is first necessary to define an appropriate stochastic process from which our performance measures can be derived. Statements about the stability of the model really relate to the stability of the stochastic process.

There are typically a host of stochastic processes that may be defined from the elements of a simulation. The choice of stochastic process depends partly on the performance measures in question. Given that two of our measures are related to response time, it is natural to consider a stochastic process that yields information on response times. Furthermore, for mathematical convenience, it is often helpful to ensure that one's stochastic process is Markov.

For $n \geq 1$, let $T_n$ denote the time at which the $n$th call is received, and define $T_0 = 0$. For $n \geq 1$, let $W_n$ be the *residual workload* of the ambulance at time $T_n+$, i.e., just after the $n$th call is received. By residual workload at some time $t$, we mean the amount of time required for the ambulance to complete any current call, along with calls that might also be queued at time $t$. We assume that the ambulance is idle at the hospital at time $T_0 = 0$, so that $W_0 = 0$.

Unfortunately, $(W_n : n \geq 0)$ is not a Markov process because the response time for a future call, and hence the workload, depends on the location of the ambulance when the ambulance clears the previous workload. So if we also keep track of the location coordinates of the ambulance $B_n = (B_n(1), B_n(2))$ at the instant at which the workload $W_n$ is first cleared, then the resulting process $Z = (Z_n : n \geq 0)$ is Markov, where $Z_n = (W_n, B_n)$.

The process $Z$ is a general state space Markov chain, and evolves on the state space

$$S = [0, \infty) \times \left[ -\frac{1}{2}, \frac{1}{2} \right]^2 .$$

The first step in ensuring that our "long-run" performance measures are defined is to establish that $Z$ exhibits some form of positive recurrence. One way to achieve this is to verify that the chain $Z$ satisfies the following condition, which certainly deserves some explanation!

To avoid a potential confusion between general results and those for our particular model, we will state general results in terms of a Markov chain $\Phi = (\Phi_n : n \geq 0)$ evolving on a state space $\mathcal{S}$.

**The Lyapunov Condition**     There exists a $B \subseteq \mathcal{S}$, positive scalars $a < 1, c$, and $\delta$, an integer $m \geq 1$, a probability distribution $\varphi$ on $\mathcal{S}$, and a function $V : \mathcal{S} \rightarrow [1, \infty)$ such that

1. $P(\Phi_m \in \cdot | \Phi_0 = z) \geq \delta \varphi(\cdot)$ for all $z \in B$, and
2. $E(V(\Phi_1) | \Phi_0 = z) \leq aV(z) + cI(z \in B)$ for all $z \in \mathcal{S}$.

**89**

The Lyapunov condition (sometimes called a Foster-Lyapunov condition) is a stronger condition than we really require, but it simplifies the presentation considerably. The function $V$ is called a Lyapunov (think of energy) function. The second requirement basically states that when the chain $\Phi$ lies outside of the set $B$, the energy in the system tends to decrease, and when the chain lies inside $B$, the energy in the system cannot become too big on the next step. This condition implies that the set $B$ gets hit infinitely often. Of course, if one takes $B = \mathcal{S}$, the entire state space, then this requirement is trivially satisfied. The first condition is needed to ensure that the set $B$ is not too "big".

In any case, the point is that if a chain $\Phi$ satisfies the Lyapunov condition, then $\Phi$ is appropriately positive recurrent. The precise statement is as follows.

**Theorem 7.** *If a discrete time Markov chain $\Phi$ is aperiodic and satisfies the Lyapunov condition, then it is $V$-uniformly ergodic. In particular, $\Phi$ has a unique stationary probability distribution.*

For a proof, see Theorem 16.0.1 of Meyn and Tweedie (1993).

So the question then is, does our chain $Z$ satisfy the Lyapunov condition? The answer is yes, and it is instructive to go through a proof. However, on a first reading one may skip the following development up to the statement of Proposition 8 without loss of continuity.

For many systems, the function $V$ may be taken to be $e^{\gamma v}$, where $v$ is some measure of the work in the system. In fact, as we now show, one may take $V(w, b) = e^{\gamma w}$ for some yet to be determined constant $\gamma > 0$.

Consider what happens on a single transition of the chain $Z$, starting from the point $Z_n = (w, b)$, where $n \geq 0$. The workload decreases at unit rate, at least until it hits 0, until the arrival of the next call over an interval of length $\tau_{n+1} = T_{n+1} - T_n$. At time $T_{n+1}$ a new call arrives at location $C_{n+1}$ say, and adds some work to the workload. In particular, there will be some travel time $\eta_{n+1}$ to the scene of the call, some time $U_{n+1}$ spent at the scene, and then potentially some travel time $\xi_{n+1}$ to transport the patient to the hospital. If the patient requires transport to the hospital then $B_{n+1} = (0, 0)$, which is the location of the hospital. If not, then $B_{n+1} = C_{n+1}$, which is the location of the new call, and $\xi_{n+1} = 0$.

So for $n \geq 0$, the new workload $W_{n+1}$ is given by $[W_n - \tau_{n+1}]^+ + Q_{n+1}$, where $[x]^+ = \max\{x, 0\}$, and

$$Q_n = \eta_n + U_n + \xi_n. \tag{9}$$

We assume that the scene time sequence $(U_n : n \geq 1)$ is i.i.d. and independent of all else, and that the same is true of the call location sequence $(C_n : n \geq 1)$.

Equipped with this Lindley-type recursion for the workload, we can now attempt to identify conditions under which the Lyapunov condition will hold. If $z = (w, b)$, then $E[V(Z_1)|Z_0 = z]$ is given by

$$
\begin{aligned}
& E e^{\gamma([w-\tau_1]^+ + Q_1)} \\
= \ & E e^{\gamma[w-\tau_1]^+} E e^{\gamma Q_1} \\
\leq \ & [E e^{\gamma(w-\tau_1)} + P(w - \tau_1 < 0)] E e^{\gamma Q_1} \\
\leq \ & e^{\gamma w}[E e^{-\gamma \tau_1} + e^{-(\lambda+\gamma)w}] E e^{\gamma(3+U_1)} \quad (10) \\
= \ & e^{\gamma w}[1 + \frac{\lambda + \gamma}{\lambda} e^{-(\lambda+\gamma)w}] E e^{\gamma(3+U_1-\tau_1)} \\
= \ & V(z)[1 + \frac{\lambda + \gamma}{\lambda} e^{-(\lambda+\gamma)w}] \phi(\gamma), \quad (11)
\end{aligned}
$$

where $\phi$ is the moment generating function of $3 + U_1 - \tau_1$. Equation (10) follows since the ambulance travels at unit rate, and the distances it can travel are such that $\eta_1 \leq 2$, and $\xi_1 \leq 1$. (Recall that the ambulance travels distances as measured by the Manhattan metric.)

Assuming that $E e^{tU_1}$ is finite in a neighbourhood of 0, i.e., $U_1$ has a moment generating function defined near 0, then we have that $\phi(0) = 1$, and

$$\phi'(0) = E(U_1 + 3 - \tau_1).$$

So if $EU_1 + 3 < E\tau_1$, then $\phi'(0) < 0$, and so $\phi(t) < 1$ for $t > 0$ in some neighbourhood of 0. So choose $\gamma > 0$ so that $\phi(\gamma) < 1$.

Now, there is some $K > 0$ such that if $w > K$, then

$$[1 + \frac{\lambda + \gamma}{\lambda} e^{-(\lambda+\gamma)w}] \phi(\gamma) < 1. \tag{12}$$

Furthermore, for $w \leq K$ we have that

$$E[V(Z_1)|Z_0 = z] \leq E e^{\gamma(K+3+U_1)} < \infty. \tag{13}$$

Thus, if we take $B = [0, K] \times [-\frac{1}{2}, \frac{1}{2}]^2$, then it follows from (11), (12) and (13) that the second requirement in the Lyapunov condition is met.

It remains to check the first requirement. Observe that if the time till the next call is large enough and the ambulance transports the current patient to the hospital, then the ambulance will be at its base when the next call arrives. So if $\tau_1 > K + U_1 + 3$ and the current patient requires transport to the hospital, then independent of $Z_0 = z \in B$, the next call will be served immediately by the ambulance from the base. In fact, the chain regenerates at such times. Let

$$\delta = p P(\tau_1 > K + U_1 + 3)$$

and $\varphi$ denote the distribution of $Z_1 = (W_1, B_1)$ assuming that at time $T_1-$ the ambulance is free and located at the

origin. Then we have that for all $z \in B$,

$$P(z, \cdot) \geq \delta\varphi(\cdot),$$

and the first requirement in the Lyapunov condition is satisfied.

In summary then, we have established that $Z$ satisfies the Lyapunov condition. It is straight-forward to show that $Z$ is aperiodic, and so we arrive at the following result.

**Proposition 8.** *If $U_1$ possesses a moment generating function in a neighbourhood of 0, and $EU_1 + 3 < E\tau_1$, then the chain $Z$ is $V$-uniformly ergodic, where $V(w, b) = e^{\gamma w}$, for some $\gamma > 0$.*

The stability condition

$$EU_1 + 3 < E\tau_1$$

has a very nice interpretation in terms of the model. The left-hand side of the inequality gives an upper bound on the expected amount of work (time at the scene + travel time to the scene + travel time from the scene to the hospital) brought in by an arriving call, whereas the right-hand side gives the expected amount of time that the ambulance has available between calls to deal with this work. This condition can certainly be weakened by being more careful about defining how much work each call brings to the system, but this is not something that we will pursue further.

The main point is that Proposition 8 gives *easily verifiable* conditions under which the system is stable. While it may have appeared somewhat difficult to verify the Lyapunov condition, the argument used is actually quite straightforward, and we will see that the payoff is easily worth the effort. Based on this result, we can now define our performance measures rigorously, and also construct estimators that are consistent and satisfy central limit theorems.

As in Section 4, the rigorous definition of our performance measures is based on the strong law of large numbers. For simplicity, we state this theorem under stronger hypotheses than are really necessary.

**Theorem 9 (MCSLLN).** *Let $\Phi$ be a $V$-uniformly ergodic Markov chain on state space $\mathcal{S}$ with stationary probability distribution $\pi$. Let $h : \mathcal{S} \to \mathbf{R}$ be a real-valued function on $\mathcal{S}$. If $\pi|h| = \int_{\mathcal{S}} |h(x)|\pi(dx) < \infty$, then*

$$\frac{1}{n} \sum_{i=0}^{n-1} h(\Phi_i) \to \pi h \quad a.s.$$

*as $n \to \infty$.*

For a proof, see Theorem 17.0.1 of Meyn and Tweedie (1993).

We turn now to the performance measure $\beta$, the long-run utilization of the ambulance. The actual utilization of the ambulance over the time interval $[0, T_n)$, i.e., up until the time of the $n$th arrival is

$$\frac{n^{-1} \sum_{i=0}^{n-1} \min\{W_i, \tau_{i+1}\}}{n^{-1} \sum_{i=1}^{n} \tau_i}. \tag{14}$$

Now, the SLLN for i.i.d. random variables implies that the denominator converges to $\lambda^{-1}$. We would like to apply the MCSLLN to the numerator, but it is not yet in an appropriate form. However, using a simple device we can fix this difficulty. In essence, we are going to apply filtering; see Glasserman (1993). We have that

$$
\begin{aligned}
E \min\{w, \tau_1\} &= wP(\tau_1 > w) + E\tau_1 I(\tau_1 \leq w) \\
&= \lambda^{-1}(1 - e^{-\lambda w}),
\end{aligned}
$$

and so we replace (14) by

$$\beta_n = \frac{1}{n} \sum_{i=0}^{n-1} (1 - e^{-\lambda W_i}). \tag{15}$$

Notice that $\beta_n$ is in exactly the form that we need to apply the MCSLLN, with $h(w, b) = 1 - e^{-\lambda w}$, which is bounded, and so we find that

$$\beta_n \to \beta \text{ a.s.}$$

as $n \to \infty$. This then is a rigorous definition of $\beta$, and also a proof that the estimator $\beta_n$ is (strongly) consistent.

Let us turn now to $r$, the steady-state density of the response time to a call. For $n \geq 0$, the response time $R_{n+1}$ to the call arriving at time $T_{n+1}$ is the sum of the workload $[W_n - \tau_{n+1}]^+$ just before the arrival of the call and the time $\eta_{n+1}$ for the ambulance to travel to the location of the new call. So the response time $R_{n+1}$ depends not only on $Z_n = (W_n, B_n)$, but also on the time interval $\tau_{n+1}$.

This dependence of the response time on the time between calls causes some difficulties in our analysis. There are several ways to solve these difficulties, but perhaps the "cleanest" is to expand the state space of our Markov chain so that it is "sufficiently rich" to supply all of the needed information.

Accordingly, define $\tilde{Z} = (\tilde{Z}_n : n \geq 0)$, where, for $n \geq 0$, $\tilde{Z}_n = (W_n, B_n, \tau_{n+1})$. Using techniques that are very similar to those used for the chain $Z$, we can show that $\tilde{Z}$ is a $\tilde{V}$-uniformly ergodic chain on the state space

$$\tilde{S} = [0, \infty) \times \left[-\frac{1}{2}, \frac{1}{2}\right]^2 \times [0, \infty),$$

where

$$\tilde{V}(w, b, t) = e^{\theta[w-t]^+}$$

for some $\theta > 0$.

We are now equipped to both define the density $r$ and derive a look-ahead density estimator for it. Recall that for $n \geq 0$,

$$R_{n+1} = [W_n - \tau_{n+1}]^+ + \eta_{n+1}.$$

So

$$
\begin{aligned}
&P(R_1 \in dx | \tilde{Z}_0 = z = (w, b, t)) \\
&= P(\eta_1 \in (x - [w-t]^+, x - [w-t]^+ + dx) | \tilde{Z}_0 = z) \\
&= P(d(b, C_1) \in (x - [w-t]^+, x - [w-t]^+ + dx)) \\
&= f(x - [w-t]^+; b) \, dx, \qquad (16)
\end{aligned}
$$

where the function $f(\cdot; \cdot)$ was introduced in (7).

If we define $g(x, z) = g(x, w, b, t) = f(x - [w - t]^+; b)$, then we can write (16) as

$$P(R_1 \in dx | \tilde{Z}_0 = z) = g(x, z) \, dx.$$

It follows that we can define $r(x) = \pi g(x, \cdot)$ where $\pi$ is the stationary distribution of $\tilde{Z}$, i.e., that $r$ defined in this way is a density of the steady-state response time. This follows immediately from results in Henderson and Glynn (2001), and is very similar to the calculation in Section 2 showing that the network completion time $L$ has a density. We omit the details.

So we define $r(x) = \pi g(x, \cdot)$, and this expression suggests that we might estimate $r(x)$ by

$$r_n(x) = \frac{1}{n} \sum_{i=0}^{n-1} g(x, Z_n).$$

Recall that $f$ is bounded (by $2\sqrt{2}$), and therefore so is $g$. So the MCSLLN establishes that under the conditions of Proposition 8, $r_n(x) \to r(x)$ almost surely as $n \to \infty$, for all $x > 0$. Hence, we have rigorously defined the density $r$, and established that it can be consistently estimated by $r_n$.

We summarize the above discussion with the following proposition.

**Proposition 10.** *Under the conditions of Proposition 8, the utilization $\beta$ and the density $r$ are well-defined, and the estimators $\beta_n$ and $r_n$ are strongly consistent (as $n \to \infty$).*

So we now turn to the error in the estimators. As before, this can be assessed through confidence intervals that derive from a central limit theorem. Again, in order for simplicity, we state the Markov chain central limit theorem under stronger conditions than are strictly necessary.

For a function $h : S \to \mathbf{R}$ with $\pi |h| < \infty$, let $\bar{h}(\cdot) = h(\cdot) - \pi h$ denote the function centred by its steady-state mean. Also, let $E_\nu$ denote the expectation operator over the path space of a Markov chain under initial distribution $\nu$.

**Theorem 11 (MCCLT).** *Suppose that the chain $\Phi$ satisfies the Lyapunov condition and is aperiodic. Then, for any function $h : \mathcal{S} \to \mathbf{R}$ with $h(z)^2 \leq V(z)$ for all $z$,*

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=0}^{n-1} h(\Phi_i) - \pi h \right) \Rightarrow \sigma N(0, 1),$$

*where $\pi$ is the stationary probability distribution of $\Phi$, and*

$$\sigma^2 = E_\pi[\bar{h}(\Phi_0)^2] + 2 \sum_{k=1}^{\infty} E_\pi[\bar{h}(\Phi_0)\bar{h}(\Phi_k)]. \qquad (17)$$

For a proof, see Theorem 17.0.1 of Meyn and Tweedie (1993).

We immediately obtain the following result.

**Proposition 12.** *Under the conditions of Proposition 8,*

$$\sqrt{n}(\beta_n - \beta) \Rightarrow \sigma_\beta N(0, 1)$$

*as $n \to \infty$, for an appropriately defined constant $\sigma_\beta^2$. In addition, for all $x > 0$,*

$$\sqrt{n}(r_n(x) - r(x)) \Rightarrow \sigma(x) N(0, 1)$$

*as $n \to \infty$, for an appropriately defined constant $\sigma(x)$.*

Thus, just as in the terminating simulation case, the error in the estimators $\beta_n$ and $r_n(x)$ is approximately normally distributed with mean 0 and variance on the order of $n^{-1}$.

This result serves as a foundation for constructing confidence intervals for our performance measures. One approach is to estimate the variance constants directly using the regenerative method, which is certainly easily applied to our example. But the method of batch means is, at least currently, more widely applicable, and the preferred method in commercial simulation software, and so we instead consider this approach.

Suppose that we have a sample path $\Phi_0, \Phi_1, \ldots, \Phi_{n-1}$. Divide this sample path into $m$ batches of size $b$, where for convenience we assume that $n = mb$, so that the $k$th batch consists of observations $\Phi_{(k-1)b}, \ldots, \Phi_{kb-1}$. Now, for $k = 1, \ldots, m$, let $M_k$ be the sample mean over the $k$th

batch, i.e.,

$$M_k = \frac{1}{b} \sum_{i=(k-1)b}^{kb-1} h(\Phi_i),$$

and let $\bar{M}_m$ denote the sample mean of the $m$ batch means $M_1, \ldots, M_m$. Finally, let

$$s_m^2 = \frac{1}{m-1} \sum_{i=1}^{m} (M_k - \bar{M}_m)^2$$

denote the sample variance of the $M_k$'s. The method of batch means provides a confidence interval for $\pi h$ of the form $\bar{M}_m \pm t s_m / \sqrt{m}$, for some constant $t$, and relies on the assumption that for large $n$, $(\bar{M}_m - \pi h)/(s_m/\sqrt{m})$ is approximately $t$-distributed, with $m-1$ degrees of freedom.

The MCCLT above implies that as $n \to \infty$ with $m$, the number of batches, held fixed, all of the batch means are asymptotically normally distributed with mean $\pi h$, and variance $m\sigma^2/n$. If each of the batch means are also asymptotically independent, then a standard result (see p. 173 of Rice 1988 for example) shows that the above confidence interval methodology is valid.

But how can we be sure that this asymptotic independence of the batch means will hold? A sufficient condition that supplies both the asymptotic independence, together with asymptotic normality, is that the chain $\Phi$ satisfy a functional central limit theorem; see Glynn and Iglehart (1990), from which much of the following discussion is adapted.

**Definition 2.** *Let $\Phi$ be a Markov chain on state space $\mathcal{S}$, and let $h : \mathcal{S} \to \mathbf{R}$. Define the continuous time process $Y = (Y(t) : t \geq 0)$ by $Y(t) = \Phi_{\lfloor t \rfloor}$. For $0 \leq t \leq 1$, let*

$$\bar{Y}_n(t) = n^{-1} \int_0^{nt} h(Y(s)) \, ds$$

*and set*

$$\zeta_n(t) = n^{1/2}(\bar{Y}_n(t) - \kappa t),$$

*for some constant $\kappa$. We say that $\Phi$ satisfies a functional central limit theorem (FCLT) if there exists an $\eta > 0$ such that $\zeta_n \Rightarrow \eta B$ as $n \to \infty$, where $B$ denotes a standard Brownian motion.*

Observe that if $\Phi$ satisfies a FCLT, then the $j$th batch mean $M_j$ can be expressed as

$$\begin{aligned} M_j &= m[\bar{Y}_n(j/m) - \bar{Y}_n((j-1)/m)] \\ &= \kappa + n^{-1/2} m(\zeta_n(j/m) - \zeta_n((j-1)/m)). \end{aligned}$$

Since the increments of Brownian motion are normally distributed, the FCLT then implies that the $M_j$'s are asymptotically normally distributed with mean $\kappa$ and variance $m\eta^2/n$, which is a conclusion that we had already reached. But the increments of Brownian motion are also independent, which implies that the $M_j$'s are asymptotically independent, and this is the final result needed to ensure that the batch means confidence methodology outlined above is asymptotically valid.

So when can we be sure that $\Phi$ satisfies a FCLT? One sufficient condition is the following result.

**Theorem 13.** *Suppose that $\Phi$ satisfies the Lyapunov condition, and $h$ is such that $h(z)^2 \leq V(z)$ for all $z$. If the constant $\sigma^2$ defined in (17) above is positive, then $\Phi$ satisfies a functional central limit theorem with $\kappa = \pi h$ and $\eta^2 = \sigma^2$.*

For a proof, see Theorems 17.4.4 and 17.5.3 of Meyn and Tweedie (1993).

Notice that we have already established that the conditions of Theorem 13 hold for our estimators. Thus, we immediately arrive at the conclusion that the method of batch means will yield asymptotically valid confidence intervals.

As in the terminating simulation case, the performance of these confidence interval procedures for finite $n$ may be negatively impacted by bias. Of course, the bias depends on the initial distribution $\mu$ say of the chain. The bias in the estimator $\beta_n$ is $E_\mu \beta_n - \beta$, with a similar expression for $r_n(x)$ for each $x > 0$.

We give the appropriate calculations for $\beta$, as those for $r$ are the same. Let $h(w, b) = 1 - e^{-\lambda w}$. Borrowing a technique from Glynn (1995), we see that the bias in $\beta_n$ under initial distribution $\mu$ is

$$\begin{aligned} E_\mu &\frac{1}{n} \sum_{i=0}^{n-1} [h(Z_i) - \pi h] \\ &= \frac{1}{n} E_\mu \sum_{i=0}^{\infty} [h(Z_i) - \pi h] - \frac{1}{n} E_\mu \sum_{i=n}^{\infty} [h(Z_i) - \pi h] \\ &= \frac{c}{n} + o(n^{-1}) \end{aligned}$$

provided that

$$c = E_\mu \sum_{i=0}^{\infty} [h(Z_i) - \pi h] < \infty. \tag{18}$$

So the bias in the estimator $\beta_n$ will be of the order $n^{-1}$ if (18) holds. This result holds in great generality. We in fact have the following result.

**Theorem 14.** *Suppose that $\Phi$ satisfies the Lyapunov condition and is aperiodic. Let $\pi$ be the stationary probability distribution of $\Phi$. If $h(z)^2 \leq V(z)$ for all $z$, and $\mu V < \infty$, then*

$$c = E_\mu \sum_{i=0}^{\infty} [h(\Phi_i) - \pi h] < \infty,$$

*and*

$$E_\mu \frac{1}{n} \sum_{i=0}^{n-1} h(\Phi_i) - \pi h = \frac{c}{n} + O(q^n),$$

*as $n \to \infty$, where $q < 1$.*

The proof of this result is a straightforward extension of Theorem 16.0.1 of Meyn and Tweedie (1993).

We can conclude from this result that if the initial conditions are chosen appropriately (e.g., if $Z_0$ and $\tilde{Z}_0$ are chosen to be deterministic), then the bias of our estimators is of the order $n^{-1}$.

Since the width of the batch mean confidence intervals is of the order $n^{-1/2}$, and the bias in the estimators is of the order $n^{-1}$, it follows that bias will typically only be an important factor for small sample sizes.

## ACKNOWLEDGMENTS

## REFERENCES

Avramidis, A. N., J. R. Wilson. 1996. Integrated variance reduction strategies for simulation. *Operations Research* 44: 327–346.

Billingsley, P. 1986. *Probability and Measure, 2nd ed.* Wiley, New York.

Bratley, P., B. L. Fox, and L. E. Schrage. 1987. *A Guide to Simulation*. 2d ed. New York: Springer-Verlag.

Glasserman, P. 1993. Filtered Monte Carlo. *Mathematics of Operations Research* 18:610–634.

Glynn, P. W. and P. Heidelberger. 1990. Bias properties of budget constrained simulations. *Operations Research* 38: 801–814.

Glynn, P. W. and D. L. Iglehart. 1990. Simulation output analysis using standardized time series. *Mathematics of Operations Research* 15:1–16.

Glynn, P. W. 1995. Some new results on the initial transient problem. Proceedings of the 1995 Winter Simulation Conference. C. Alexopoulos, K. Kang, W. R. Lilegdon, D. Goldsman, eds. IEEE, Piscataway NJ. 165–170.

Henderson, S. G. 2000 Mathematics for simulation. Proceedings of the 2000 Winter Simulation Conference. J. A. Joines, R. R. Barton, K. Kang, and P. A. Fishwick, eds. IEEE, Piscataway NJ. 137-146.

Henderson, S. G. and P. W. Glynn. 2001. Computing densities for Markov chains via simulation. *Mathematics of Operations Research.* To appear.

Law, A. M. and W. D. Kelton. 2000. *Simulation Modeling and Analysis, 3rd ed.* McGraw-Hill, New York.

Meyn, S. P. and R. L. Tweedie. 1993. *Markov Chains and Stochastic Stability.* Springer-Verlag, New York.

Prakasa Rao, B. L. S. 1983. *Nonparametric Functional Estimation.* Academic Press.

Rice, J. A. 1988. *Mathematical Statistics and Data Analysis.* Wadsworth and Brooks/Cole, Pacific Grove, California.

Serfling, R. J. 1980 *Approximation Theorems of Mathematical Statistics.*

## AUTHOR BIOGRAPHY

**SHANE G. HENDERSON** is an assistant professor in the School of Operations Research and Industrial Engineering at Cornell University. He has previously held positions in the Department of Industrial and Operations Engineering at the University of Michigan and the Department of Engineering Science at the University of Auckland. He is an associate editor for the ACM Transactions on Modeling and Computer Simulation, and the assistant newsletter editor for the INFORMS College on Simulation. His research interests include discrete-event simulation, queueing theory and scheduling problems. His e-mail address is <shane@orie.cornell.edu>, and his web page is <www.orie.cornell.edu/~shane>.