# ADAPTIVE IMPORTANCE SAMPLING SIMULATION OF QUEUEING NETWORKS

Pieter-Tjerk de Boer

Telematics Systems and Services
Department of Computer Science
University of Twente
P.O. Box 217
7500 AE Enschede
THE NETHERLANDS

Victor F. Nicola

Telematics Systems and Services
Department of Electrical Engineering
University of Twente
P.O. Box 217
7500 AE Enschede
THE NETHERLANDS

Reuven Y. Rubinstein

William Davidson Faculty
of Industrial Engineering
and Management
Technion
Haifa
ISRAEL

## ABSTRACT

In this paper, a method is presented for the efficient estimation of rare-event (overflow) probabilities in Jackson queueing networks using importance sampling. The method differs in two ways from methods discussed in most earlier literature: the change of measure is state-dependent, i.e., it is a function of the content of the buffers, and the change of measure is determined using a cross-entropy-based adaptive procedure. This method yields asymptotically efficient estimation of overflow probabilities of queueing models for which it has been shown that methods using a state-independent change of measure are not asymptotically efficient. Numerical results demonstrating the effectiveness of the method are presented as well.

## 1 INTRODUCTION

During the last decade, there has been much interest in the estimation of rare-event probabilities in queues and networks of queues, with applications to models of telecommunication networks as well as computer and manufacturing systems. Two methods have gained popularity: importance sampling (Heidelberger 1995, Asmussen and Rubinstein 1995), and importance splitting (RESTART) (Villén-Altamirano and Villén-Altamirano 1994), the former of which is used in this paper.

One simple network that received a lot of attention is a set of two or more queues in tandem. Despite its simplicity, a complete analysis of this system is hard due to the behaviour at the state-space boundaries. As a consequence, no importance sampling change of measure that is provably asymptotically efficient is known.

In Parekh and Walrand (1989), an importance sampling procedure was described for estimating the overflow probability of the total population in tandem queues. A simple and static (i.e., state-independent) change of measure was used: exchange the arrival rate with the service rate (of the bottleneck queue, in case of a tandem system). In Sadowsky (1991), the asymptotic efficiency of that method for a single queue was proved. In Frater, Lennon, and Anderson (1991) this heuristic was extended to overflows of the total population in any Jackson network. However, it was shown in Glasserman and Kou (1995) that for two or more queues in tandem, this heuristic does not always give an asymptotically efficient simulation, depending on the values of arrival and service rates. It reasonable is to expect that similar problems will occur with this method in other Jackson networks.

Clearly, by allowing the change of measure to depend on the state of the system (i.e., the content of each of the queues), more efficient importance sampling schemes may be obtained. This approach was recently used in Kroese and Nicola (1999), where the overflow probability of the second queue in a two-node tandem Jackson network is estimated

using a simulation in which the change of measure depends on the content of the first buffer; the functional dependence of the rates on the buffer content is derived from on a Markov additive process representation of the system. Furthermore, in Heegaard (1998) a state-dependent change of measure is used for simulating link overloads in a telecommunications network; again, the functional dependence of the importance sampling rates on the system state is derived using a heuristic. The biggest obstacle to the use of a state-dependent change of measure in general is the problem of determining this dependence: rather specific mathematical models are used in the publications mentioned, making the results very specific to those problems.

As an alternative to avoid the complex mathematical analysis often used to determine a good (state-independent) change of measure, several adaptive methods have been proposed recently; see Devetsikiotis and Townsend (1993b), Devetsikiotis and Townsend (1993a), Al-Qaq, Devetsikiotis, and Townsend (1995), Rubinstein (1997), Rubinstein and Melamed (1998), Rubinstein (1999), Lieber (1999). All of these either try to iteratively minimize the variance of the estimator involved, or a related quantity like the cross-entropy. However, none of these papers consider a state-dependent change of measure for simulation of queueing models.

In this paper, we present an adaptive method for determining a state-dependent change of measure for rare events in queueing problems. This is a rather versatile method:

- due to the adaptiveness, a complex mathematical analysis of the problem is not necessary.
- since the state-dependent change of measure is less restrictive, problems can be solved for which no effective state-independent change of measure exists.

In particular, with this method the probability of overflow of the total network population in Jackson tandem networks can be asymptotically efficiently estimated, even in those cases where Glasserman and Kou (1995) show that the heuristic of exchanging the arrival rate with the bottleneck service rate does not work. In addition, the combination of state-dependence and adaptiveness leads to another useful property: the standard deviation of the estimator can decrease faster than proportional to the square-root of the total simulation effort.

Here we restrict our discussion to the estimation of rare-event probabilities in discrete-time Markov chains (DTMCs). However, we plan to expand the method to more general models in the future.

The rest of this paper is organized as follows. Section 2 provides a summary of the most important aspects of the adaptive method from Rubinstein and Melamed (1998). Section 3 explains the implementation of this algorithm for state-dependent simulation, and discusses some of the problems involved and their solutions. In Section 4, empirical results demonstrate the effectiveness of the method. Concluding remarks and directions for further research are given in Section 5.

## 2 PRINCIPLES OF THE CROSS-ENTROPY METHOD

In this section, we briefly review the cross-entropy method for the adaptive optimization of an importance sampling simulation. Only the aspects that are relevant for the rest of this paper are discussed; for more details, the reader is referred to Rubinstein and Melamed (1998) and Rubinstein (1999).

### 2.1 Basics

Assume that the change of measure (or "tilting") is parameterized by some vector $v$; then the aim of an adaptive importance sampling procedure should be to find the value of $v$ which results in minimal variance for the resulting estimator.

Another approach for choosing $v$ was introduced in Rubinstein (1999). It is well known that always an importance sampling distribution (change of measure) exists which results in a zero-variance estimator, and that this distribution is precisely the original distribution conditioned on the occurrence of the rare event. In practice, this distribution may not be within the family of distributions that can be obtained by the change of measure parameterized by $v$. However, if a simulation distribution is used that is in some sense "close" to the unattainable zero-variance distribution, then a low (but non-zero) variance should be expected. So instead of choosing $v$ such that the variance is minimized explicitly, one could try to devise a procedure that minimizes some distance measure between the distribution under the change of measure given by $v$, and the distribution that would give zero variance. The latter distribution will henceforth be called the "zero-variance distribution".

Before proceeding with details of such a procedure, some more notation needs to be defined. The sample path of one replication of the simulation is denoted by $Z$. The function $I(Z)$ is the indicator function of the occurrence of the rare event in $Z$. We already defined $v$ to denote the tilting vector; consistently with this, $f(Z, v)$ is the probability (or, for continuous systems, the probability density) of the sample path $Z$ under the tilting $v$, with $v = 0$ corresponding to the original (untilted) system. The likelihood ratio associated with the sample path $Z$ and the tilting vector $v$ is denoted by $L(Z, v)$:

$$L(Z, v) = \frac{f(Z, 0)}{f(Z, v)}.$$

Finally, $\mathbb{E}_{\boldsymbol{v}}$ denotes the expectation under the tilting $\boldsymbol{v}$.

A suitable "distance" measure for this procedure is the Kullback-Leibler cross-entropy, which is defined, see Kapur and Kesavan (1992), as

$$CE = \int f(z) \ln \frac{f(z)}{g(z)} dz,$$

where $f(z)$ and $g(z)$ are two density functions whose "distance" is to be calculated. Note that this "distance" measure is not symmetric: in general, exchanging $f$ and $g$ in the above will result in a different value of $CE$.

We want to apply the Kullback-Leibler cross-entropy to measure the distance between the distribution to be used for the simulation (assumed to be of the form $f(z, \boldsymbol{v})$) and the zero-variance distribution. To do this, substitute $g(z) = f(z, \boldsymbol{v})$ (i.e., the distribution to be optimized by changing $\boldsymbol{v}$) and $f(z) = \rho_0 I(z) f(z, 0)$ with normalization factor $\rho_0^{-1} = \int I(z) f(z, 0) dz$ into the above; note that this $f(z)$ is the original distribution conditioned on the rare event (i.e., the zero-variance distribution). Then we need to do the following minimization:

$$\begin{aligned} \boldsymbol{v}^{\dagger} &= \arg\min_{\boldsymbol{v}} \int \rho_0 I(z) f(z, 0) \ln \frac{\rho_0 I(z) f(z, 0)}{f(z, \boldsymbol{v})} dz \\ &= \arg\max_{\boldsymbol{v}} \int I(z) f(z, 0) \ln f(z, \boldsymbol{v}) dz \\ &= \arg\max_{\boldsymbol{v}} \mathbb{E}_0 I(Z) \ln f(Z, \boldsymbol{v}), \end{aligned} \quad (1)$$

where $\boldsymbol{v}^{\dagger}$ denotes the value of $\boldsymbol{v}$ that minimizes the cross-entropy. In the above form the equation is not useful, since we do not know $\mathbb{E}_0 I(Z) \ln f(Z, \boldsymbol{v})$. However, we can rewrite it as follows:

$$\boldsymbol{v}^{\dagger} = \arg\max_{\boldsymbol{v}} \mathbb{E}_{\boldsymbol{v}_j} I(Z) L(Z, \boldsymbol{v}_j) \ln f(Z, \boldsymbol{v}),$$

where $\boldsymbol{v}_j$ is any other tilting vector; we will later interpret it as the tilting vector used during the $j$th iteration of an iterative procedure. The above form can easily be approximated by a sum (stochastic counterpart of the expectation) over $N$ samples from a simulation performed with tilting $\boldsymbol{v}_j$, thus yielding an approximation to $\boldsymbol{v}^{\dagger}$ which we call $\boldsymbol{v}_{j+1}$:

$$\boldsymbol{v}_{j+1} = \arg\max_{\boldsymbol{v}} \sum_{i=1}^{N} I(Z_i) L(Z_i, \boldsymbol{v}_j) \ln f(Z_i, \boldsymbol{v}), \quad (2)$$

where the $Z_i$ are sample paths drawn under the tilting $\boldsymbol{v}_j$.

## 2.2 Algorithm

Now we have all elements for an iterative procedure to approximate the optimal tilting vector $\boldsymbol{v}^{\dagger}$ in the sense that it minimizes the cross-entropy; this may not be equal to

minimizing the variance of the estimator, although in practice the difference turns out to be small:

1. Initialize as follows:
   $j := 1$ (iteration counter)
   $\boldsymbol{v}_1 :=$ initial tilting vector (see below)
2. Simulate $N$ replications with tilting $\boldsymbol{v}_j$, yielding $Z_1 \ldots Z_N$.
3. Find the new tilting vector $\boldsymbol{v}_{j+1}$ from the maximization (2).
4. Increment $j$ and repeat steps 2–4, until the tilting vector has converged (i.e., $\boldsymbol{v}_{j+1} \approx \boldsymbol{v}_j$).

Choosing the initial tilting vector $\boldsymbol{v}_1$ in step 1 is not trivial. The most obvious choice is to set $\boldsymbol{v}_1 = 0$, i.e., use the original transition probabilities. However, with that choice the rare event of interest will typically not be observed, making (2) unusable. In Rubinstein (1999), this is solved by introducing an additional step in the algorithm, in which the rare event is temporarily modified into a less rare event. In the present paper, a different approach is used: we choose $\boldsymbol{v}_1$ on the basis of heuristics like exchanging the arrival rate with the bottleneck service rate. Although such a heuristic by itself does not produce an asymptotically efficient simulation in the cases considered here, it does provide a convenient starting point for the iterative process.

## 3 STATE-DEPENDENT TILTING

One application of the adaptive procedure described above is finding a "static" change of measure for queueing problems, i.e., finding the optimal arrival and service rates for simulation of a buffer overflow. In that case, the vector $\boldsymbol{v}$ just contains one component for every rate that is allowed to change. Indeed, for many problems this turns out to work well; see Rubinstein and Melamed (1998) and de Boer (2000) for examples. However, for many other problems no static change of measure seems to exist that gives an efficient simulation; for those systems, a less restrictive change of measure should be used, which can be obtained by allowing the arrival and service rates in the simulation to depend on the state. In this section we will do precisely that for DTMC models of queueing networks.

### 3.1 Principles

A DTMC model is completely described by its initial probability distribution and its set of transition probabilities: the probabilities of going from one state to another. Since many DTMC models (e.g., for queueing systems) are derived from continuous time models with exponential arrival and service time distributions (CTMCs), the transition probabilities are typically calculated from transition rates: the probability of going from state $i$ to state $j$ is given by $\lambda_{ij} / \sum_k \lambda_{ik}$,

where $\lambda_{ij}$ is the transition rate from state $i$ to state $j$, and $k$ in $\sum_k$ runs over all states. In fact, for the cross-entropy calculations done in this section, it is more convenient to work with rates; the transition probabilities can trivially be calculated from the rates by normalizing their sum to 1. Collectively, all rates $\lambda_{ij}$ will be referred to as a vector $\boldsymbol{\lambda}$.

In DTMC models, only one type of tilting is possible: changing the transition probabilities. Equivalently, one can change the transition rates of the corresponding CTMC model and calculate the transition probabilities for the DTMC from those, as shown above. It turns out that the latter approach is slightly simpler. So the aim is to find a set of transition rates $\boldsymbol{\lambda}$ which minimizes the cross-entropy.

Before deriving the actual cross-entropy minimization formula, let us first build a mathematical description of one replication $Z$ of a DTMC simulation. Define the sequence $z_i$, $i = 1, 2, 3, \ldots$, which denotes the state of the system just before the $i$th transition in this replication $Z$. Denote by $\lambda_{lm}$ the rate (or probability) of going from state $l$ to state $m$. Then obviously the probability of the $i$th step is

$$\frac{\lambda_{z_i z_{i+1}}}{\sum_k \lambda_{z_i k}},$$

where $k$ runs over all states (or only those states that can be reached in one step from state $z_i$, since all other $\lambda_{z_i k}$ are 0). The total probability of the sample path $Z$ is

$$\Pr(Z) = \prod_i \frac{\lambda_{z_i z_{i+1}}}{\sum_k \lambda_{z_i k}},$$

where $i$ runs over all steps in the sample path.

Substitute the above expression for the probability of a sample path into equation (1); then we get the following expression for the optimal transition rate vector $\boldsymbol{\lambda}^\dagger$:

$$\boldsymbol{\lambda}^\dagger = \arg\max_{\boldsymbol{\lambda}} \mathbb{E}_0 I(Z) \ln \prod_i \frac{\lambda_{z_i z_{i+1}}}{\sum_k \lambda_{z_i k}}$$
$$= \arg\max_{\boldsymbol{\lambda}} \mathbb{E}_0 I(Z) \sum_i \left( \ln \lambda_{z_i z_{i+1}} - \ln \sum_k \lambda_{z_i k} \right).$$

To find the maximum in the right-hand side, set the derivative with respect to $\lambda_{lm}$ to 0, for any two states $l$ and $m$:

$$0 = \mathbb{E}_0 I(Z) \sum_{i:z_i=l} \left( \frac{1_{(z_{i+1}=m)}}{\lambda_{lm}} - \frac{1}{\sum_k \lambda_{lk}} \right),$$

or, equivalently:

$$\frac{1}{\lambda_{lm}^\dagger} \mathbb{E}_0 I(Z) \sum_{i:z_i=l} 1_{(z_{i+1}=m)} = \frac{1}{\sum_k \lambda_{lk}^\dagger} \mathbb{E}_0 I(Z) \sum_{i:z_i=l} 1.$$

Thus, we find the following expression for the optimal transition probability $q_{lm}$ from state $l$ to state $m$:

$$q_{lm} = \frac{\lambda_{lm}^\dagger}{\sum_k \lambda_{lk}^\dagger}$$
$$= \frac{\mathbb{E}_0 I(Z) \sum_{i:z_i=l} 1_{(z_{i+1}=m)}}{\mathbb{E}_0 I(Z) \sum_{i:z_i=l} 1}. \tag{3}$$

Of course, the expectations in the right-hand side are generally not known, but we can approximate them as follows:

$$q_{lm} = \frac{\mathbb{E}_{\boldsymbol{\lambda}_j} I(Z) L(Z, \boldsymbol{\lambda}_j) \sum_{i:z_i=l} 1_{(z_{i+1}=m)}}{\mathbb{E}_{\boldsymbol{\lambda}_j} I(Z) L(Z, \boldsymbol{\lambda}_j) \sum_{i:z_i=l} 1}$$
$$\approx \frac{\sum_{Z=Z_1}^{Z_N} I(Z) L(Z, \boldsymbol{\lambda}_j) \sum_{i:z_i=l} 1_{(z_{i+1}=m)}}{\sum_{Z=Z_1}^{Z_N} I(Z) L(Z, \boldsymbol{\lambda}_j) \sum_{i:z_i=l} 1}, \tag{4}$$

where $\sum_{Z=Z_1}^{Z_N}$ is a sum over the sample paths from $N$ replications, simulated with transition rates $\boldsymbol{\lambda}_j$ (i.e., from the $j$th iteration). Note that the factor $\sum_{i:z_i=l} 1$ in the denominator is just the number of visits to state $l$ during replication $Z$, and that $\sum_{i:z_i=l} 1_{(z_{i+1}=m)}$ in the numerator is the number of those visits in which the transition to state $m$ was chosen next. Consequently, the right-hand side of (4) can be interpreted as the *observed* conditional (on the occurrence of the rare event) probability of the transition from state $l$ to state $m$; this is not surprising, since it is known that using the true conditional distributions for importance sampling yields a zero-variance estimator, as discussed before.

## 3.2 Practical Problems

Using the adaptive importance sampling algorithm from Section 2.2 with state-dependent parameters chosen according to (4) seems very simple. There are, however, practical difficulties. The cause of these is the enormous number of states that a typical queueing network can have. For example, a network with three queues and an overflow level of 50 for the total network population has 23461 states. This is the total number of ways to distribute among three distinct queues a total of 1 customer (3 ways), 2 indistinguishable customers (6 ways), 3 indistinguishable customers (10 ways), up to 50 indistinguishable customers. Doubling the overflow level to 100 multiplies this number of states by almost 8. If the rare event of interest is the overflow of one particular queue, other queues in the network can have an infinite size, thus making the number of states infinite.

One of the consequences of the enormous state space is that a lot of data needs to be stored: this takes a lot of memory capacity; but with present-day computers and the size of the queueing networks studied here, this is typically not a problem (except if the state space is infinite,

of course). However, manipulating such a lot of data (e.g., in the smoothing techniques that will be discussed later) can be prohibitively time-consuming.

The accuracy of the estimations in the right-hand side of (4) is more problematic. The only sample paths that give a contribution to the sums in the numerator and denominator are those that reach the rare event (because of the $I(Z)$ factor) and pass through the state $l$ (because of the summation over $i$ for which $z_i = l$). The factor $I(Z)$ will typically not be a problem: the tilting used in the $j$th iteration is usually such that the event of interest is no longer rare. However, the tilting will not favor visits to states that are away from some optimal path to the rare event of interest. If the state space is multi-dimensional, this means that many states will not be visited often or at all, even under a tilting that makes the event of interest non-rare. States that are not visited at all during the $N$ replications of a simulation yield 0/0 (undefined) in the right-hand side of (4). And states that are visited only a few times make the quotient of sums a bad approximation of the quotient of expectations.

There is in fact a rather fundamental risk here: suppose the transition from some state $l$ to another state $m$ happens in only 10 % of all visits to state $l$, and state $l$ is visited only 5 times during the $N$ replications of a simulation. Then it is quite likely that in none of those 5 visits to state $l$, a transition to state $m$ will be made. Consequently, using (4) to choose the simulation parameters for the next iteration would set the rate (probability) of this transition to 0, thus making the transition impossible. Then in the next simulation, surely no transitions from state $l$ to state $m$ will be observed, so this rate will again be set to 0 for the next iteration: it will remain at 0 forever, even though that is wrong if the transition has a non-zero probability in the untilted system, thus possibly resulting in a biased estimator.

The only case in which the above does not give a biased estimator is when the rare event of interest can no longer be reached after that particular transition has been made. As a matter of fact, all paths $Z$ which contain such a transition necessarily have $I(Z) = 0$; as a consequence, (4) will automatically set the rate of such a transition to zero for the next iteration. Therefore, after the first iteration, *all* sample paths will reach the rare event.

### 3.3 Dealing with a Large Number of States

In this section, three techniques will be outlined to deal with the problems caused by the large state space. For details, the reader is referred to de Boer (2000).

The basic idea of these techniques is the assumption that the optimal transition probabilities for a particular state are typically close to those of other "similar" states in its neighbourhood. If this is the case, the estimates of the transition probabilities for a given state may be improved by also including observations from sample paths passing through an appropriate set of such "similar" states. Of course, this introduces an error, since the optimal probabilities are probably not exactly equal. On the other hand, since more samples are used, the accuracy of the estimation increases. Furthermore, treating several states as if they were one state saves memory for storing the transition probabilities. This is necessary for systems with an infinite number of states.

Note that the "error" discussed above does not imply that the resulting estimate of the rare-event probability will be biased; in principle that estimate will be unbiased as long as the correct likelihood ratios are used. Rather, it means that the used transition probabilities deviate from the optimal transition probabilities, so the estimate has a larger variance than without this error. In fact, such errors and the associated non-optimal variance are always present, even if no grouping of states is used, due to the fact that the transition probabilities are estimated by simulation and thus subject to statistical errors.

#### 3.3.1 Local Average

The local average technique tries to automatically choose the optimal amount of grouping, separately for every state. It does this as follows:

First, just the observations obtained at the state itself are used. If this gives good enough (see below) estimates of the transition probabilities out of this state, then the estimates are accepted. If not, the observations from a set of neighbouring states are combined with those from this state; if the transition probability estimates are now good enough, these are accepted. If not, this is repeated with ever larger sets of neighbouring states, until the resulting transition probability estimates are good enough.

The test for deciding whether the transition probability estimates are good enough comprises several aspects. First of all, the number of visits to the state (including the states with which it is being grouped): if the state or group of states has been visited too few times, its transition probability estimates cannot be trusted. Secondly, no probabilities that theoretically should be non-zero, are set to zero: if this happens, again the results cannot be trusted. Thirdly, one can construct an estimator for the variance of the transition probability estimates, and compare its value to some threshold to decide whether or not the transition probability estimate is acceptable.

#### 3.3.2 Boundary Layers

The boundary layer technique is based on the observation that when a queue's content is sufficiently large, the optimal transition probabilities tend to become nearly independent of that queue's content. Thus, all states in which a queue

contains $B$ or more customers are grouped together. When drawing this in a picture of the state space, layers are seen along the boundaries; hence the name. See Figure 1 for $B = 3$ in a two-dimensional state-space. Choosing the optimal number of boundary layers $B$ seems to be done best by trial and error: using too few gives a less efficient simulation, since the resulting change of measure is less dependent on the state.
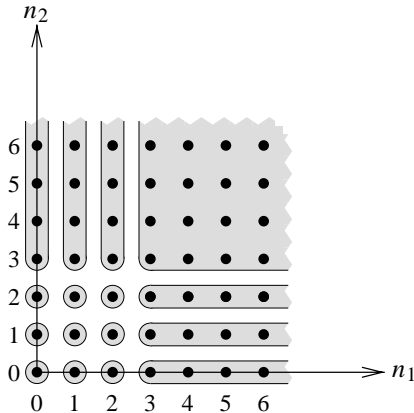


Figure 1: Grouping of States Using Three Boundary Layers in the State Space of a Two-queue System; $n_i$ = Level of $i$th Queue

### 3.3.3 Smoothing Using Splines

After applying the above two methods, the transition probabilities can still be rather "noisy" functions of the state; they are simulation results, after all. It might be beneficial to replace the noisy data by a smooth function fitted through it. The form of the optimal transition probability functions is not known in general, so fitting a flexible generic function to the data is the best one can do. We have succesfully used cubic splines for this smoothing; basically, this means that the state space is divided into pieces, and on every piece a third-order polynomial of the coordinates (i.e., the contents of the queues) is fitted to the data. Choosing the size of the pieces is a compromise between noise reduction and accurace of the fit.

### 3.3.4 Combination

In practice, two or all three of the above methods are combined. Before the simulations are started, the number of boundary layers is chosen; this is very effective at reducing the amount of data to be stored and processed. Next, the simulation is performed. Following this, the transition rates are calculated on the basis of the simulation results, using equation (4); the local average technique is used to group neighbouring states where necessary to obtain reliable estimates of the transition probabilities. Finally, the spline smoothing can be applied, if needed or desired. If the results

after the local average step are already relatively good, spline smoothing may worsen the accuracy by imposing an unsuitable form on the data; on the other hand, if the data is rather noisy, the spline smoothing usually improves its accuracy. We will see examples of both in Section 4.

### 3.4 The Variance of the Estimator

It can be shown (see de Boer (2000)) that if the state-dependent transition probabilities given by (3) are used, the resulting estimator for the rare event probability has zero variance. In practice, the variance will not be zero, due to the fact that one cannot obtain the exact transition probabilities that satisfy (3): instead, simulation results are used in equation (4) to approximate the optimal transition probabilities, thus causing them to have a statistical error. Furthermore, the techniques for dealing with the large state space limit the accuracy of the transition probabilities.

Now consider what happens if the number of replications per iteration is, say, quadrupled. If this had no influence on the transition probability estimates, the relative error of the rare-event probability estimator would obviously improve by a factor of $\sqrt{4} = 2$. However, errors in the estimates of the transition probabilities would also improve, by up to a factor of 2 if the statistical error in them is dominant. This means that they become closer to the optimal (zero-variance) transition probabilities, causing the estimate of the rare-event probability to improve; under some assumptions it can be argued that this improvement is linear in the reduction of the statistical error in the transition probabilities. Therefore, the error in the rare-event probability estimate decreases by up to a total factor of 4, i.e., up to linear in the number of replications used per iteration. We will demonstrate this experimentally in the next section.

## 4 EXPERIMENTAL RESULTS

In this section, overflows in a simple Jackson network will be considered. The network consists of four queues in tandem, with arrival and service rates chosen in the region where the standard state-independent change of measure (exchanging the arrival rate with the bottleneck service rate) does not work well according to Glasserman and Kou (1995): the arrival rate is 0.09, the service rates of the first through fourth queue are 0.23, 0.227, 0.227 and 0.226, respectively. The rare event of interest is the total network population reaching a high level, starting from 0 and before returning to 0 again.

For all experiments, the boundary layer technique was used to reduce the enormous state space; 10 boundary layers turned out to work well, but possibly fewer would have been sufficient. Furthermore, the local average technique was used. The spline-smoothing was only used in some cases, as indicated below.

## 4.1 Results for Overflow Level 50

The results for an overflow level of 50 are presented in Figure 2, both without and with spline-based smoothing. Along the horizontal axis, the iteration number is indicated. Vertically, the estimate of the overflow probability and its relative error (standard deviation from the simulation, divided by the estimate itself) are shown as two lines in the graph. At the first iteration, a static tilting according to the well-known heuristic of exchanging the arrival rate with the bottleneck service rate was used, to get things started. In the experiments without spline-based smoothing (upper graph), $10^4$ replications were used per iteration up to the 23rd iteration; the 23rd iteration was performed twice, once with $10^4$ and once with $10^5$ replications, and all later iterations used $10^5$ replications. With spline-based smoothing (lower graph), the switch from $10^4$ to $10^5$ replications per iteration was made at the 9th instead of the 23rd iteration.

Obviously, the spline smoothing is quite beneficial to the convergence in this case: without splines, the convergence is rather slow and irregular, with a major excursion
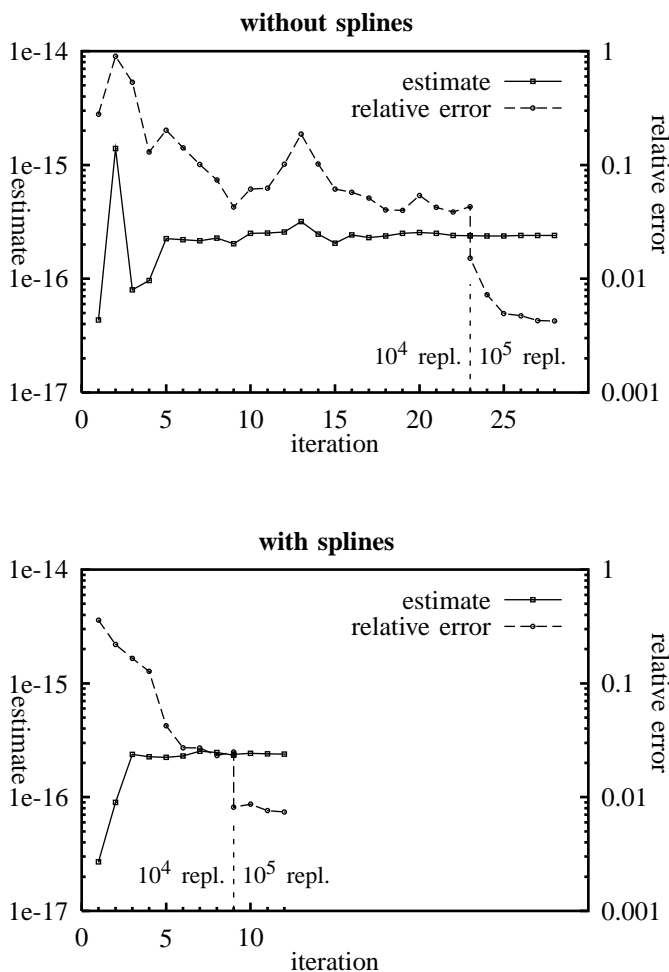


Figure 2: Results for the Four-node Tandem Queue, Overflow Level = 50

around the 13th iteration, whereas with spline smoothing the convergence is quick and monotonic, and the resulting relative error at $10^4$ replications is smaller by almost a factor of 2.

Next, note what happens when the number of replications is increased: at the 23rd (without splines) and 9th (with splines) iteration, the same simulation was done with $10^4$ and $10^5$ replications; the relative error of the latter clearly is about a factor of $\sqrt{10}$ smaller, as it should. However, without splines the relative error continues to decrease in the next iteration: this is a consequence of the fact that these later iterations have better transition probabilities because those have been obtained with $10^5$ instead of $10^4$ replications, as discussed in Section 3.4. In the end, the relative error has decreased by a factor of 10 in total. With splines, this does not happen: the relative error does not significantly decrease further, and in fact is higher than without splines; apparently, the spline form does not fit the optimal state-dependence well enough.

Figure 3 serves to give an idea of how the transition probabilities depend on the state in this particular problem. Of course, since we have up to five transition probabilities and a four-dimensional state space, it is hardly feasible to give a complete picture. Therefore, only the probability of the transition corresponding to a service completion at the first queue is shown, as a function of the contents $n_1$ and $n_2$ of the first and second queues, respectively, while the third and fourth queues are empty. Clearly, the splines
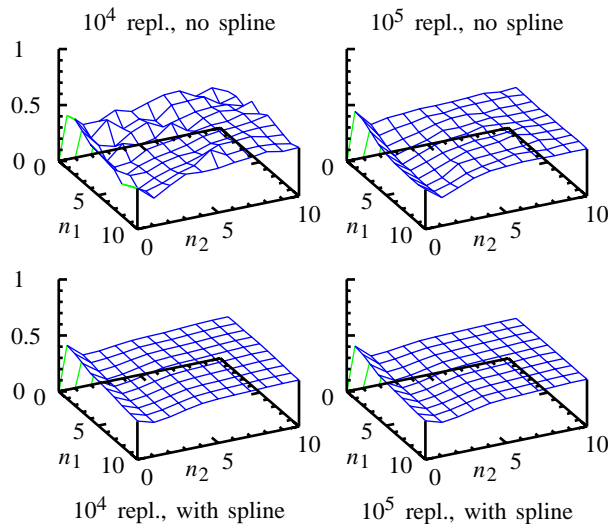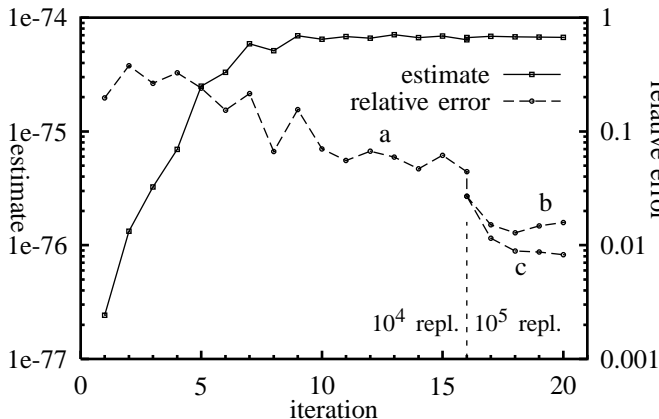


Figure 3: State-dependent Transition Probabilities

perform a very effective smoothing: most of the noise disappears. On the other hand, the splines used here are apparently not able to completely follow the true functions: the "dip" at $n_2 = 1$ is much deeper without splines (only sufficiently visible in the $10^5$-replications plot) than with splines. This agrees with the experimental observation that at $10^5$ replications, the final estimate is more accurate when

the transition probabilities are not restricted by applying splines.

## 4.2  Results for Overflow Level 200

For the case of an overflow level of 200, Figure 4 shows the simulation results. For this problem, all three techniques (local average, 10 boundary layers, and splines) were used initially (up to iteration 16). After convergence had been achieved, the number of replications was increased, resulting in branches b (with splines) and c (without splines) in the graph.



| Branch | Iterations | Description |
|--------|-----------|-------------|
| a | 1–16 | $10^4$ replications, splines |
| b | 16–20 | $10^5$ replications, splines |
| c | 16–20 | $10^5$ replications, no splines |

Figure 4: Results for the Four-node Tandem Queue, Overflow Level = 200

It seems as if the convergence process can be divided into two phases. During the first phase, the estimate is quite inaccurate (typically too low), but it approaches the correct value; in the present example, this phase comprises iterations 1 through 7. During the second phase, the estimate stays correct, and the relative error decreases to its final value; in the present example this happens during iterations 7 through 11. These phases can also be recognized in the results with overflow level 50 in Figure 2.

Note, like before, the strong decrease of the relative error after increasing the number of replications by a factor of 10, and the fact that switching off spline smoothing at that point is beneficial.

## 4.3  Asymptotic Efficiency

Results from the above experiments, and from repetitions of those experiments at overflow levels 25 and 100, are shown in Table 1. All of these experiments used the same number of replications per iteration ($10^5$) and no splines in the final iterations. It is clear from the table that the

relative error grows with the overflow level, but clearly less than exponentially fast, while the probability of interest does decrease exponentially fast. This demonstrates the asymptotic efficiency of the method for this problem.

Table 1: Test of Asymptotic Efficiency

| level | exact | estimate | rel.error |
|-------|-------|----------|-----------|
| 25 | $3.5283 \cdot 10^{-07}$ | $3.504 \cdot 10^{-07}$ | 0.0026 |
| 50 | – | $2.396 \cdot 10^{-16}$ | 0.0042 |
| 100 | – | $1.422 \cdot 10^{-35}$ | 0.0044 |
| 200 | – | $6.722 \cdot 10^{-75}$ | 0.0082 |

The table also shows an exact (numerical) calculation of the overflow probability for an overflow level of 25. Comparing this with the simulation estimate shows a good agreement. No exact numbers could be calculated for higher overflow levels due to the large state space involved.

## 5  CONCLUSIONS AND FURTHER RESEARCH

In this paper, we have proposed an importance sampling simulation method with two important features: the change of measure is completely state-dependent, and a cross-entropy-based adaptive method is used to approximate the optimal change of measure. To show the method's performance, we have applied it to estimate the overflow probability of the total population of a Jackson network consisting of four queues in tandem. This simulation has been shown to be asymptotically efficient, at a parameter setting at which asymptotically efficient simulation is not obtained with state-independent tilting. Furthermore, the method's interesting property that the relative error can decrease faster than proportional to the square root of the total simulation effort has been demonstrated.

The method has also been applied successfully to other rare-event problems in Jackson networks, like overflows in networks with random routing and feedback, bounded queues, and overflows of non-bottleneck queues; see de Boer (2000).

However, all of the systems considered so far are modelled by DTMCs, and the number of queues is not too large to avoid state space explosion. This indicates two obvious directions for future work: extension to non-DTMC systems, and developing more efficient methods for handling large state spaces. Furthermore, it may be possible to improve the method's convergence by combining observations from several iterations.

In the present paper, the good performance of the method has only been demonstrated experimentally. Another direction for further research would therefore be providing more solid mathematical foundations, such as a proof of the convergence of the tilting vector.

# REFERENCES

Al-Qaq, W. A., M. Devetsikiotis, and J. K. Townsend. 1995. Stochastic gradient optimization of importance sampling for the efficient simulation of digital communication systems. *IEEE Transactions on Communications* 43:2975–2985.

Asmussen, S. and R. Rubinstein. 1995. Complexity properties of steady-state rare-events simulation in queueing models. In *Advances in Queueing: Theory, Methods and Open Problems*, ed. J. Dshalalow, 429–462. CRC Press.

de Boer, P. T. 2000. Analysis and efficient simulation of queueing models of telecommunications systems. Ph. D. thesis, University of Twente. In preparation.

Devetsikiotis, M. and J. K. Townsend. 1993a. An algorithmic approach to the optimization of importance sampling parameters in digital communication system simulation. *IEEE Transactions on Communications* 41:1464–1473.

Devetsikiotis, M. and J. K. Townsend. 1993b. Statistical optimization of dynamic importance sampling parameters for efficient simulation of communication networks. *IEEE/ACM Transactions on Networking* 1:293–305.

Frater, M. R., T. M. Lennon, and B. D. O. Anderson. 1991. Optimally efficient estimation of the statistics of rare events in queueing networks. *IEEE Transactions on Automatic Control* 36:1395–1405.

Glasserman, P. and S.-G. Kou. 1995, January. Analysis of an importance sampling estimator for tandem queues. *ACM Transactions on Modeling and Computer Simulation* 5(1):22–42.

Heegaard, P. E. 1998. A scheme for adaptive biasing in importance sampling. *AEÜ International Journal of Electronics and Communications* 52:172–182.

Heidelberger, P. 1995. Fast simulation of rare events in queueing and reliability models. *ACM Transactions on Modeling and Computer Simulation* 5:43–85.

Kapur, J. N. and H. K. Kesavan. 1992. *Entropy Optimization Principles with Applications*. Academic Press.

Kroese, D. P. and V. F. Nicola. 1999. Efficient simulation of a tandem jackson network. In *Proceedings of the 1999 Winter Simulation Conference*, ed. P. A. Farrington, H. B. Nembhard, D. T. Sturrock, and G. W. Evans, 411–419.

Lieber, D. 1999. The cross-entropy method for estimating probabilities of rare events. Ph. D. thesis, William Davidson Faculty of Industrial Engineering and Management, Technion, Israel.

Parekh, S. and J. Walrand. 1989. A quick simulation method for excessive backlogs in networks of queues. *IEEE Transactions on Automatic Control* 34:54–66.

Rubinstein, R. Y. 1997. Optimization of computer simulation models with rare events. *European Journal of Operations Research* 99:89–112.

Rubinstein, R. Y. 1999. Rare event simulation via cross-entropy and importance sampling. In *Second International Workshop on Rare Event Simulation, RESIM'99*, 1–17.

Rubinstein, R. Y. and B. Melamed. 1998. *Modern Simulation and Modeling*. Wiley.

Sadowsky, J. S. 1991. Large deviations theory and efficient simulation of excessive backlogs in a $GI/GI/m$ queue. *IEEE Transaction on Automatic Control* 36:1383–1394.

Villén-Altamirano, M. and J. Villén-Altamirano. 1994. RESTART: A straightforward method of fast simulation of rare event. In *Proceedings of the 1994 Winter Simulation Conference*, 282–289.

# AUTHOR BIOGRAPHIES

**PIETER-TJERK DE BOER** received the M.S. degree in applied physics (specializing in theoretical physics) in 1996 from the University of Twente, The Netherlands. Since then, he has been working towards a Ph.D. degree at the Department of computer science, University of Twente. His research interests include rare event simulation, importance sampling, queueing theory, and large deviations theory, with applications to performance analysis of telecommunication networks.

**VICTOR F. NICOLA** holds the Ph.D. degree in computer science from Duke University, North Carolina, the B.S. and the M.S. degrees in electrical engineering from Cairo University, Egypt, and Eindhoven University of Technology, The Netherlands, respectively. From 1979, he held faculty and research staff positions at Eindhoven University and at Duke University. In 1987, he joined IBM Thomas J. Watson Research Center, Yorktown Heights, New York, as a Research Staff Member. Since 1993, he has been an Associate Professor at the Department of Electrical Engineering, University of Twente, The Netherlands. His research interests include performance and reliability modeling, fault-tolerance, queueing theory, analysis and simulation methodologies, with applications to computer systems and telecommunication networks.

**REUVEN Y. RUBINSTEIN** Prof. Reuven Rubinstein is with the Faculty of Industrial Engineering and Management of the Technion since 1973. His fields of interest are stochasic models, stochastic optimization and simulation. He published over 80 papers and 4 books on simulation and stochastic optimization, all with Wiley. He was the head of operations research division at the Technion for 4 years. He has visited many universities and research centers around the world, among them University of Illinois, Urbana

(1978–79 academic year), Harvard University (1985–86 academic year), George Washington University (1986–87 academic year), IBM Research Center (Summer 1980), Bell Laboratories, Holmdel, NJ (Summers 1989 and 1990), NEC (February 1992). Motorola US (1997, 6 months), The Institute of Statistical Mathematics (1997-98, 4 months, Tokyo). He is a Technion Management Chair Professor since 1998.