

COMPENSATORS AND DERIVATIVE ESTIMATION FOR QUEUEING SYSTEMS

Michael A. Zazanis
Dept. of Industrial Engineering and Management Sciences
Northwestern University
Evanston, IL 60208

1. Abstract

A fairly general method for sensitivity analysis of simulations is proposed involving the use of compensator identities. The case of a GI/GI/1 queue is discussed in some detail. Expressions for the derivatives of state probabilities that provide direct simulation estimates are given for single class closed queueing networks with markovian routing.

2. Introduction

The problem of sensitivity analysis of simulations has in the last few years attracted the attention of a number of authors. One of the methods proposed, infinitesimal perturbation analysis (e.g. see Ho and Cao, 1983 and Suri and Zazanis, 1988), involves the direct differentiation of sample performance measures. It is well known that, depending on the nature of the parameter and the performance measure, the estimates obtained in this way may be biased (Heidelberger et al., 1988). An alternative method using likelihood ratios was proposed by Reiman and Weiss (1989), Glynn (1986), and Rubinstein (1989). The estimators obtained by this method are unbiased for a wide class of systems but in many cases they have large variance. In this paper, a general method for derivative estimation is proposed for a large class of problems that can be cast in a form involving stochastic integrals with respect to a counting process. It uses compensator identities in conjunction with infinitesimal perturbation analysis (IPA) techniques to provide low variance unbiased estimates at the expense of additional computational requirements. For an earlier paper on the same subject we refer the reader to Zazanis (1988).

3. Compensators and sensitivity of simulations

Let $X_i(\theta)$ be a real valued stochastic process depending on a parameter θ , and A_t a counting process (e.g. arrivals to or departures from the system), both defined on a filtered probability space and adapted to $\{F_t\}_{t \geq 0}$. Let λ_t be the F_t -intensity of A_t which may also depend on θ . We will assume throughout the paper that $\lambda_s < K$ for all s a.s. and we will consider performance criteria of the form

$$J(\theta) = E \left[\int_0^t f(X_s(\theta)) dA_s \right] \quad (1.1)$$

and we will derive estimators for the derivative $\frac{d}{d\theta} J(\theta)$. Special emphasis will be given to the case $f(x) = I_B(x)$, where $B = [x, \infty)$ and we will show how to estimate $\frac{d}{d\theta} P_B(\theta)$, where $P_B(\theta)$ is the steady state probability that an arrival finds the system in B .

The key idea in the method we propose here is the use of the compensator identity

$$E \left[\int_0^T f(X_s) dA_s \right] = E \left[\int_0^T f(X_s) \lambda_s ds \right] \quad (1.2)$$

which, assuming X_s to be left continuous, holds for all f for which the expectation exists (e.g. see Bremaud, 1981). (Here and in what follows we suppress the dependence on θ except when we want to draw attention to it). While infinitesimal perturbation analysis (IPA) methods would use $\frac{d}{d\theta} \int_0^T f(X_s) dA_s$,

evaluated along a sample path of the process, to estimate $\frac{d}{d\theta} J(\theta)$, the method we propose here uses $\frac{d}{d\theta} \int_0^T f(X_s) \lambda_s ds$, which in many cases that are important in applications can be evaluated along a sample path. The estimators we propose here have increased computational requirements but are unbiased in many cases where the IPA estimators are biased. The reason for this is that while it may not be permissible to differentiate with respect to θ inside the expectation on the left hand side of (1.2), it is often permissible to do so on the right hand side as we will see in the examples given in the next two sections. For a related but not identical method we refer the reader to Gong and Ho (1987) and Gong and Glasserman (1989).

The method we proposed above is illustrated for the case of a GI/G/1 queue and a single class, single server closed queueing network with renewal service times.

4. Derivative estimators for steady state probabilities in a GI/G/1 queue

Consider a GI/G/1 queue with input process (S_i, X_i) , $i=1, 2, \dots$, where S_i is the interarrival time between the i^{th} and $(i+1)^{th}$ customers, and X_i is the load brought by the i^{th} customer to the system. We will assume the input to be i.i.d. with distribution $F(x, \theta)$ depending on a parameter θ . We also assume that the distribution G of the interarrival times is absolutely continuous with density g , and that the corresponding hazard rate, $\frac{g(x)}{1-G(x)}$, is bounded by some $K < \infty$ for all $x \geq 0$.

To be specific, suppose that the first customer arrives to an empty server at time $t = 0$. We will consider the number of customers in the system process $N_t(\theta)$, as θ varies in an interval

$[a, b]$, such that $\sup_{\theta \in [a, b]} \int_0^\infty x dF(x, \theta) < EA_1$ (to ensure stability)

and $\sup_{\theta \in [a, b]} \int_0^\infty x^2 dF(x, \theta) < \infty$. To construct a family of sample paths parametrized by θ on the same probability space let $X(\theta + \delta) = F^{-1}(\theta + \delta, F(\theta, X))$. We will assume that $\frac{dX_i}{d\theta}$ exists

a.s. and that $\sup_{\theta \in [a, b]} E \left(\frac{dX_i}{d\theta} \right)^2 < \infty$. Notice that $\{N_t(\theta); t \geq 0\}$ determines $\{N_t(\eta); t \geq 0\}$ for all $\eta \in [a, b]$. $\{N_t(\theta); t \geq 0\}$ is defined to be *left continuous*. We will denote by $\{A_s; s \geq 0\}$ and $\{D_s; s \geq 0\}$ the counting processes associated with arrivals and departures from the system. (These are

both defined to be right continuous). Let us define for convenience three additional (left continuous) processes related to N_s , namely $\{Z_s; s \geq 0\}$, the time since the last arrival process, $\{\lambda_s; s \geq 0\}$, the stochastic intensity process, and finally $\{Y_s; s \geq 0\}$, where

$$Y_t = \sum_{i=0}^{D_t} \frac{dX_i}{d\theta} - \int_0^t Y_s 1_{(N_s=1)} dD_s. \quad (2.1)$$

Intuitively, if δ is a small perturbation to θ , $Y_t \delta$ is how much ahead (or behind) schedule would the server be as a result of this.

Theorem 1: Let

$$J(\theta, t) = \frac{1}{t} E \left[\int_0^t 1_{(N_s \geq k)} dA_s \right]. \quad (2.2)$$

Then $J(\theta, t)$ is differentiable on $[a, b]$ and

$$\frac{d}{d\theta} J(\theta, t) = \frac{1}{t} E \left[\int_0^t \lambda_s Y_s 1_{(N_s=k)} dD_s \right]. \quad (2.3)$$

Sketch of Proof: Let $\{u_i\}$ and $\{v_i\}$ be the sequences of upcrossings from level $k-1$ to k and downcrossings from level k to $k-1$ respectively. Also, let $\Lambda = \int_0^t \lambda_s ds$ be the compensator associated with the arrival process. Using the identity (1.2) we have

$$\begin{aligned} J(\theta, t) &= \frac{1}{t} E \left[\int_0^t 1_{(N_s \geq k)} d\Lambda_s \right] \\ &= \frac{1}{t} \left[\sum_{\{i: 0 \leq v_i < t\}} \Lambda_{v_i} - \sum_{\{i: 0 \leq u_i < t\}} \Lambda_{u_i} \right]. \end{aligned}$$

Differentiate with respect to θ taking into account that only the v_i 's in the above expression depend on θ . From an easy dominated convergence argument and the fact that $\frac{d}{d\theta} \Lambda_{v_i} = \lambda_{v_i} \frac{d}{d\theta} v_i$ we obtain:

$$\frac{d}{d\theta} J(\theta, t) = \frac{1}{t} \left[\sum_{\{i: 0 \leq v_i < t\}} \lambda_{v_i} \frac{d}{d\theta} v_i \right]$$

This, together with the fact that $Y_{v_i} = \frac{d}{d\theta} v_i$ establishes (2.3).

Theorem 2: Let P_k be the customer stationary probability of k or more customers in the system. Then

$$\frac{d}{d\theta} P_k = \frac{1}{EQ_1} E \left[\int_0^{T_1} \lambda_s Y_s 1_{(N_s=k)} dD_s \right]. \quad (2.4)$$

where Q_1 and T_1 is the number of customers in the first busy period and its the length respectively.

Sketch of proof: It is enough to show that $\frac{d}{d\theta} \lim_{t \rightarrow \infty} J(\theta, t) = \lim_{t \rightarrow \infty} \frac{d}{d\theta} J(\theta, t)$, the rest following from a standard regenerative argument. This is guaranteed from a standard theorem (e.g. see Bartle, 1978, p.204) provided that $\lim_{t \rightarrow \infty} J(\theta_0, t)$ exists for some $\theta_0 \in [a, b]$ and that $\frac{d}{d\theta} J(\theta, t)$ converges uniformly in $[a, b]$ to some limit. The first condition is obviously satisfied because of the regenerative nature of the system and the second using a modification of Lorden's inequality for renewal-reward processes (Zazanis, 1990).

If $d_i = \sum_{j=1}^i X_j$ is the i 'th departure time, (2.4) can be written in a form suitable for regenerative simulation:

$$\frac{d}{d\theta} P_k = \frac{1}{EQ_1} E \left[\sum_{i=1}^{Q_1} 1_{(N_{d_i}=k)} \frac{g(Z_{d_i})}{1-G(Z_{d_i})} \left(\sum_{j=1}^i \frac{dX_j}{d\theta} \right) \right] \quad (2.5)$$

Let us note here that the application of the Dominated Convergence Theorem in (2.2) to get (2.3) would not have been possible had we not used the compensator identity (1.2). Thus the IPA estimate would be biased in that case.

5. A Single Class Closed Queueing Network

Consider a single server, single class CQN with M stations, N customers, and renewal service times. To simplify the notation we restrict ourselves to the case of tandem networks (The extension to general markovian routing is straightforward). Let N_t^m be the number of customers at station m at time t (again left continuous). Let $A_t^m (D_t^m)$ be the arrival (departure) counting process at station m . Clearly $A_t^m = D_t^{m-1}$ for $m=2, 3, \dots, M$, and $A_t^1 = D_t^M$. Let F_m be the service time distribution of the m 'th server. We will assume that the corresponding hazard rates exist and are bounded by a finite constant K as in §2 and we will denote the stochastic intensity associated with $\{A_s^m\}_{s \geq 0}$ by $\{\lambda_s^m\}_{s \geq 0}$. Suppose that the service time distribution of server i , $F_i(\theta, x)$ depends on a parameter θ . Define

$$\begin{aligned} Y_t^m &= \int_0^t Y_s^{m-1} 1_{(N_s^m=0)} dA_s^m \\ &\quad - \int_0^t Y_s^m 1_{(N_s^m=1)} dD_s^m \quad m \neq i, \end{aligned} \quad (3.1)$$

and

$$\begin{aligned} Y_t^i &= \int_0^t Y_s^{i-1} 1_{(N_s^i=0)} dA_s^i \\ &\quad - \int_0^t Y_s^i 1_{(N_s^i=1)} dD_s^i + \sum_{k=0}^{D_t^i} \frac{dX_k^i}{d\theta}, \end{aligned} \quad (3.2)$$

where X_k^i is the k 'th service time of the i 'th service station.

Theorem 3: Let $J(\theta, t) = \frac{1}{t} E \left[\int_0^t 1_{(N_s \geq k)} dA_s^m \right]$. Then

$$\begin{aligned} \frac{d}{d\theta} J(\theta, t) &= \frac{1}{t} E \left[\int_0^t 1_{(N_s^m=k-1)} Y_s^{m-1} \lambda_s^{m-1} dA_s^m \right] \\ &\quad - \frac{1}{t} E \left[\int_0^t 1_{(N_s^m=k)} Y_s^m \lambda_s^{m-1} dD_s^m \right]. \end{aligned} \quad (3.2)$$

Theorem 4: Let P_k^m be the steady state probability that an arriving customer to a station m finds k or more customers. Then

$$\frac{d}{d\theta} P_k^m = \lim_{t \rightarrow \infty} \frac{d}{d\theta} J(\theta, t). \quad (3.3)$$

6. Applications to derivative estimation for simulations

In Zazanis and Suri (1986), and Fox and Glynn (1987) it has been shown that the classical finite difference estimators have very poor performance compared with direct estimators such as the one described in §2. Hence, when available and easy to implement these estimators are preferable. In particular, notice that (2.5) can be used to estimate $\frac{d}{d\theta}P_k$ while observing a single sample path of the system by simply keeping track of two quantities, namely the age of the busy period, Y_i , and the age of the arrival process, Z_i . When θ is a location (scale) parameter of $F(\theta, x)$, Y_i becomes the discrete (continuous) age of the busy period at the i 'th departure epoch, making the implementation of (2.5) and (3.1) very simple.

Here of course we take advantage of the fact that the stochastic intensity is simply the hazard rate, and the only part of the history of the process necessary to determine λ_t is the age of the arrival process at time t . For most models used in practice, one would be able to compute the stochastic intensity easily. This would be the case for instance for superpositions of renewal processes (in which case one would of course need to know the ages of all the arrival processes involved), for Markov renewal processes (in which case one would need to know the state of the underlying Markov chain and the time since the last arrival), for interrupted renewal processes (such as the output from an upstream server) etc.

REFERENCES

- Bartle, R.G. (1976). *The Elements of Real Analysis*, Wiley.
- Bremaud, P. (1981). *Point Processes and Queues*, Springer.
- Cao, X.R. (1985). "Convergence of Parameter Sensitivity Estimates in a Stochastic Experiment", *IEEE Trans. Aut. Control*, Vol.30, No.9, 845-853.
- Chen, H.P. and M.A. Zazanis (1990). "Lorden's Inequality in Regenerative Simulation", Technical Report, IE/MS Dept, Northwestern Univ. Evanston, IL 60208.
- Fox, B.L., and P.W. Glynn (1989). "Replication Schemes for Limiting Expectations", *Probability in Engineering and the Information Sciences*,
- Glasserman, P. and W.B. Gong (1989). "Derivative Estimates from Discontinuous Realizations: Smoothing Techniques", AT&T Technical Report.
- Glynn, P.W. and J.L. Sanders, (1986). "Optimization of Stochastic Systems", *Proceedings of the Winter Simulation Conference*.
- Gong, W.B. and Y.C. Ho (1987). "Smoothed perturbation analysis for discrete event systems", *IEEE Trans. Aut. Control* AC-32,10,pp. 858-866.
- Heidelberger, P., X.R. Cao, M.A. Zazanis, and R. Suri (1988). "Convergence Properties of Infinitesimal Perturbation Analysis", *Management Science*, 34, 11, 1281-1301.
- Ho, Y.C., and X. Cao (1983). "Perturbation Analysis and Optimization of Queueing Networks", *J. Optim. Theory Applic.* 40, 4, 559-582.
- Reiman M.I. and A. Weiss, (1989). "Sensitivity Analysis of Simulations via Likelihood Ratios", *Opns. Res.*, 37, 5, 830-844.
- Rubinstein, R.Y., (1987). "The score function approach for sensitivity analysis of computer simulation models", *Opns. Res.*, 37, 1, 72-81.
- Suri, R. and Zazanis, M.A. (1988). "Perturbation Analysis gives strongly consistent sensitivity estimates for the M/G/1 queue", *Management Science* 34, 1, 39-64.
- Zazanis, M.A., and R. Suri (1986). "Comparison of Perturbation Analysis with Conventional Sensitivity Estimates for Stochastic Systems", Working Paper # 86-123, Dept. of Ind. Engr., Univ. of Wisconsin-Madison.
- Zazanis, M.A. (1988). "Compensators and sensitivity analysis of queueing systems", *Proceedings of the 26th Allerton Conference on Communication, Control, and Computing*, pp. 549-554.