

VARIANCE REDUCTION FOR QUANTILE ESTIMATION VIA CORRELATION INDUCTION

Athanasios N. Avramidis

School of Industrial Engineering
Purdue University
West Lafayette, Indiana 47907, U.S.A.

ABSTRACT

We propose correlation induction techniques for reducing the variance of quantile estimators in a finite-horizon simulation experiment. Both single-sample and multisample estimators are considered. If the response is monotone in the random-number inputs that drive the simulation and has a smooth distribution, then the multisample estimators are guaranteed to have asymptotically smaller mean squared error than the direct-simulation estimator. The results of a Monte Carlo study suggest that significant variance reductions can be achieved when estimating quantiles of the completion time of stochastic activity networks.

1 INTRODUCTION

The purpose of many simulation experiments is to estimate the quantile of order p of a target response variable Y , having an unknown cumulative distribution function (CDF) $F(y) \equiv P[Y \leq y]$. That is, we wish to estimate

$$\xi \equiv \inf\{y : F(y) \geq p\} \quad \text{for some } p \in (0, 1).$$

A single replication of the simulation is driven by a random vector $\mathbf{U} \equiv (U_1, \dots, U_d)$, where $\{U_i : i = 1, \dots, d\}$ are independent *random numbers*—that is, random variables uniformly distributed on the interval $(0, 1)$. The response of interest is $Y = f(\mathbf{U})$, where the function $f(\cdot)$ is defined by the simulation code.

In a direct simulation experiment, we perform n independent replications that yield independent and identically distributed (IID) observations $\{Y_i : i = 1, \dots, n\}$ of the target response. In terms of the order statistics

$$Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)},$$

the *direct* estimator of ξ based on n replications is defined as

$$\hat{\xi}_{\text{DI}}(n) \equiv Y_{(\lfloor np \rfloor + 1)}.$$

(See David (1981) and Juritz, Juritz and Stephens (1983) for properties of this estimator.) Since $\hat{\xi}_{\text{DI}}(n)$ is, in general, biased, we define variance reduction in terms of the mean squared error (MSE) criterion; that is, we seek an alternative estimator $\hat{\xi}(n)$ based on n replications such that

$$\text{MSE}[\hat{\xi}(n)] \equiv E[(\hat{\xi}(n) - \xi)^2] < E[(\hat{\xi}_{\text{DI}}(n) - \xi)^2].$$

We will suppress the dependence on the sample size n when it is not essential in the discussion.

The problem of variance reduction for quantile estimators has received relatively little attention in the simulation literature. Lewis and Ressler (1989) consider the method of control variates, extended to allow for nonlinear transformations of the control variable. Having identified a random variable C that is observed in the simulation experiment and has known quantiles, they propose using $C_{(\lfloor np \rfloor + 1)}$, the direct estimator of the p th quantile of C , as a control variable for $\hat{\xi}_{\text{DI}}(n)$. Hsu and Nelson (1990) also use a control variable with known quantiles, even though the estimators they develop are not classical control-variate estimators.

This paper is organized as follows. In Section 2 we discuss correlation induction techniques across multiple samples. Section 3 contains some ideas on correlation induction within a sample. In Section 4 we report the results of a Monte Carlo study that was used to gauge the variance reductions achieved by the proposed techniques.

2 CORRELATION INDUCTION ACROSS SAMPLES

Motivated by the need to obtain an estimate of the variance of the quantile estimator, Schafer (1974) suggests using k independent samples, each consisting of $m = n/k$ independent observations (to simplify the exposition, we assume throughout the paper that n

is an integral multiple of k). Denote the direct estimator of ξ based on the i th sample by $\hat{\xi}_{DI}^{(i)}$. The multisample estimator of ξ based on k samples and a total of n replications is then defined as

$$\hat{\xi}_M(k, n) \equiv k^{-1} \sum_{i=1}^k \hat{\xi}_{DI}^{(i)}.$$

Although the multisample estimator is not aimed at variance-reduction, we introduce it because it simplifies our subsequent discussion. We propose new multisample estimators, where we induce dependence across the samples, while we maintain independence within each sample. The key to the development is the notion of negative quadrant dependence, which was proposed by Lehmann (1966).

DEFINITION 1. *The pair (X, Y) is negatively quadrant dependent (NQD) if*

$$P[X \leq x, Y \leq y] \leq P[X \leq x] \cdot P[Y \leq y] \quad \text{for all } x, y.$$

In section 2.1 we induce negative quadrant dependence by using Latin hypercube sampling (LHS), and in section 2.2 by using antithetic sampling (AS).

2.1 Latin Hypercube Sampling

We begin with

DEFINITION 2. *The sample $\{Y_1, \dots, Y_k\}$ is a Latin hypercube sample if it is generated as*

$$Y_i = f\left(\frac{\pi_1(i) - 1 + U_{i1}}{k}, \dots, \frac{\pi_d(i) - 1 + U_{id}}{k}\right)$$

for $i = 1, \dots, k$, where (a) $\{\pi_1(\cdot), \dots, \pi_d(\cdot)\}$ are independent random permutations of $\{1, \dots, k\}$; and (b) $\{U_{ij} : j = 1, \dots, d, i = 1, \dots, k\}$ are independent random numbers that are independent of $\{\pi_1(\cdot), \dots, \pi_d(\cdot)\}$.

Our definition of LHS is more general than the standard definition, introduced by McKay, Beckman, and Conover (1980) and followed by Stein (1987). These authors assume that the nonuniform random variates that drive the simulation are independent, and each of these variates is generated by the method of inversion. We do not make either of these assumptions.

Note that (a) each Y_i has the distribution of Y , since it is generated by arguments that are uniformly distributed on $[0, 1]$ and are independent of each other; and (b) the pairs (Y_i, Y_j) are dependent for all i, j . Now we have a basic result in

PROPOSITION 1. *If $f(\cdot)$ is monotone in each coordinate and $\{Y_1, \dots, Y_k\}$ is a Latin hypercube sample, then (Y_i, Y_j) is NQD for $i \neq j$.*

PROOF. Define

$$(V_{1j}, V_{2j}) = \left(\frac{\pi_j(1) - 1 + U_{1j}}{k}, \frac{\pi_j(2) - 1 + U_{2j}}{k} \right)$$

for $j = 1, \dots, d$. It is straightforward to check that $(\pi_j(1), \pi_j(2))$ is NQD, so it follows from Theorem 1(iii) of Lehmann (1966) that (V_{1j}, V_{2j}) is NQD for each $j = 1, \dots, d$. Moreover, the pairs $\{(V_{1j}, V_{2j}) : j = 1, \dots, d\}$ are independent. Since $f(\cdot)$ is monotone in each coordinate, Theorem 1(ii) of Lehmann (1966) implies that (Y_1, Y_2) is NQD. Finally, we observe that all pairs (Y_i, Y_j) with $i \neq j$ have the same distribution. \square

Now we describe a new quantile estimator based on Latin hypercube sampling. Let $m = n/k$. Obtain k samples $\{Y_j^{(i)} : j = 1, \dots, m\}$ for $i = 1, \dots, k$ such that $\mathcal{Y}_j \equiv \{Y_j^{(i)} : i = 1, \dots, k\}$ is a Latin hypercube sample for each j and $\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_m$ are independent. Observe that each sample consists of m independent observations, but we have induced dependence across the respective observations of samples. Let $\hat{\xi}_{LH}^{(i)}$ denote the direct estimator of ξ based on the i th sample. Define the *Latin hypercube* estimator of ξ based on k subsamples and a total of n replications as

$$\hat{\xi}_{LH}(k, n) \equiv k^{-1} \sum_{i=1}^k \hat{\xi}_{LH}^{(i)}.$$

Note that the direct estimator $\hat{\xi}_{DI}(n)$ is a special case of $\hat{\xi}_{LH}(k, n)$ for $k = 1$. Our main result is

THEOREM 1. *If $f(\cdot)$ is monotone in each coordinate, then $\text{Var}[\hat{\xi}_{LH}(k, n)] \leq \text{Var}[\hat{\xi}_M(k, n)]$ for all k, n .*

PROOF. It follows from Proposition 1 that the pair $(Y_j^{(1)}, Y_j^{(2)})$ is NQD for each $j = 1, \dots, m$. Let $l = \lfloor mp \rfloor + 1$. Define

$$g(x_1, \dots, x_m) \equiv \text{the } l\text{th smallest in } \{x_1, \dots, x_m\}$$

and observe that $g(\cdot)$ is increasing in each coordinate. Since $\hat{\xi}_{LH}^{(i)} = g(Y_1^{(i)}, \dots, Y_m^{(i)})$ for $i = 1, 2$ and the pairs $\{(Y_1^{(1)}, Y_1^{(2)}), \dots, (Y_m^{(1)}, Y_m^{(2)})\}$ are independent, Theorem 1(ii) of Lehmann implies that the pair $(\hat{\xi}_{LH}^{(1)}, \hat{\xi}_{LH}^{(2)})$ is NQD. By Lemma 3 of Lehmann (1966), $\text{Cov}(\hat{\xi}_{LH}^{(1)}, \hat{\xi}_{LH}^{(2)}) \leq 0$. Thus

$$\begin{aligned} \text{Var}[\hat{\xi}_{LH}(k, n)] &= k^{-1} \left[\text{Var}(\hat{\xi}_{LH}^{(1)}) + (k-1) \text{Cov}(\hat{\xi}_{LH}^{(1)}, \hat{\xi}_{LH}^{(2)}) \right] \\ &\leq k^{-1} \text{Var}(\hat{\xi}_{LH}^{(1)}) \\ &= \text{Var}[\hat{\xi}_M(k, n)]. \quad \square \end{aligned}$$

Clearly

$$E[\hat{\xi}_{\text{LH}}(k, n)] = E[\hat{\xi}_{\text{DI}}(n/k)] = E[\hat{\xi}_{\text{M}}(k, n)].$$

Since the mean squared error of an estimator $\hat{\xi}$ is $\text{MSE}(\hat{\xi}) = (E[\hat{\xi}] - \xi)^2 + \text{Var}(\hat{\xi})$, it follows that

$$\text{MSE}[\hat{\xi}_{\text{LH}}(k, n)] \leq \text{MSE}[\hat{\xi}_{\text{M}}(k, n)] \quad \text{for all } k, n. \quad (1)$$

Next we wish to compare $\hat{\xi}_{\text{LH}}(k, n)$ with the direct estimator $\hat{\xi}_{\text{DI}}(n)$. A comparison of the bias and variance of these two estimators for finite n seems impossible in general. However, as n becomes large, the asymptotic behavior of the bias and variance of $\hat{\xi}_{\text{DI}}(n)$ (and thus of $\hat{\xi}_{\text{M}}(k, n)$ for fixed k) is well-known. Let $Q(x) \equiv \inf\{y : F(y) \geq x\}$ for $x \in (0, 1)$ be the quantile function of Y , and assume that $Q(\cdot)$ is continuous at p and has three continuous derivatives at p . Denote the first derivative of $Q(\cdot)$ by $Q'(\cdot)$. Zelterman (1987) proves that “for p not near 0 or 1”,

$$\begin{aligned} \lim_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{DI}}(n)] &= \lim_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{M}}(k, n)] \\ &= p(1-p)Q'^2(p) \end{aligned} \quad (2)$$

for any fixed k . (Under conditions given in van Zwet (1964), (2) holds for all $p \in (0, 1)$; it follows easily from Lemmata 3.2.2 and 3.2.3.) Combining (1) and (2) we get

$$\limsup_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{LH}}(k, n)] \leq \lim_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{DI}}(n)]$$

for any fixed k . That is, for any fixed k , $\hat{\xi}_{\text{LH}}(k, n)$ has asymptotically (as n becomes large) no larger MSE than $\hat{\xi}_{\text{DI}}(n)$. An unresolved problem is how to choose k in order to minimize $\text{MSE}[\hat{\xi}_{\text{LH}}(k, n)]$; some rough guidelines are given in section 4.

2.2 Antithetic Sampling

Obtain two *antithetic* samples $\mathcal{Y}_1 \equiv \{Y_i^{(1)} : i = 1, \dots, n/2\}$ and $\mathcal{Y}_2 \equiv \{Y_i^{(2)} : i = 1, \dots, n/2\}$ as follows:

$$Y_i^{(1)} = f(U_{i1}, \dots, U_{id}),$$

$$Y_i^{(2)} = f(1 - U_{i1}, \dots, 1 - U_{id})$$

for $i = 1, \dots, n/2$, where $\{U_{ij} : j = 1, \dots, d, i = 1, \dots, n/2\}$ are independent random numbers. Let $\hat{\xi}_{\text{AV}}^{(1)}$, $\hat{\xi}_{\text{AV}}^{(2)}$ denote the direct estimators of ξ based on samples \mathcal{Y}_1 and \mathcal{Y}_2 respectively. Define the *antithetic variate* estimator of ξ based on a total of n replications as

$$\hat{\xi}_{\text{AV}}(n) \equiv \frac{1}{2} \left(\hat{\xi}_{\text{AV}}^{(1)} + \hat{\xi}_{\text{AV}}^{(2)} \right).$$

In analogy with Theorem 1, we have

THEOREM 2. *If $f(\cdot)$ is monotone in each coordinate, then $\text{Var}[\hat{\xi}_{\text{AV}}(n)] \leq \text{Var}[\hat{\xi}_{\text{M}}(2, n)]$ for all n .*

PROOF. It is straightforward to check that if U is a random number, then $(U, 1 - U)$ is NQD. Since $f(\cdot)$ is monotone in each coordinate, Theorem 1(ii) of Lehmann (1966) shows that $(Y_i^{(1)}, Y_i^{(2)})$ is NQD for each $i = 1, \dots, n/2$. By reasoning as in Theorem 1, we see that $\text{Cov}(\hat{\xi}_{\text{AV}}^{(1)}, \hat{\xi}_{\text{AV}}^{(2)}) \leq 0$, which in turn implies $\text{Var}[\hat{\xi}_{\text{AV}}(n)] \leq \text{Var}[\hat{\xi}_{\text{M}}(2, n)]$. \square

As in section 2.1, we get

$$\text{MSE}[\hat{\xi}_{\text{AV}}(n)] \leq \text{MSE}[\hat{\xi}_{\text{M}}(2, n)] \quad \text{for all } n,$$

and, under the conditions mentioned in section 2.1,

$$\limsup_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{AV}}(n)] \leq \lim_{n \rightarrow \infty} n \text{MSE}[\hat{\xi}_{\text{DI}}(n)].$$

That is, $\hat{\xi}_{\text{AV}}(n)$ has asymptotically no larger MSE than $\hat{\xi}_{\text{DI}}(n)$.

3 CORRELATION INDUCTION WITHIN A SAMPLE

Multisample estimators are usually more biased than single-sample estimators; and bias can be the dominant factor of MSE when the size of each subsample is small, especially when we estimate extreme quantiles (Juritz, Juritz and Stephens 1983). In this case a single-sample estimator may be desirable.

Let $\mathcal{Y} \equiv \{Y_1, \dots, Y_n\}$ be a Latin hypercube sample. Define $\hat{\xi}_{\text{S,LH}}(n)$, the *single-sample Latin hypercube* estimator of ξ based on n replications as the direct estimator of ξ based on the sample \mathcal{Y} —that is, $\hat{\xi}_{\text{S,LH}}(n)$ is the $([np] + 1)$ st smallest observation in \mathcal{Y} . This estimator is fundamentally different from those discussed so far: it is an order statistic from a sample of *dependent* observations. In contrast, so far we have considered order statistics from samples of independent observations ($\hat{\xi}_{\text{DI}}$), and averages of independent ($\hat{\xi}_{\text{M}}$) or dependent ($\hat{\xi}_{\text{LH}}, \hat{\xi}_{\text{AV}}$) observations of such order statistics. The idea is that a Latin hypercube sample is more representative of the underlying distribution of the response Y than a random sample, because the marginal distribution of each input random number is sampled more thoroughly. Thus the appropriate order statistic from a Latin hypercube sample should better approximate the quantile of interest.

4 MONTE CARLO RESULTS

We performed a Monte Carlo experiment to estimate the MSE reductions (with respect to the direct-simulation estimator $\hat{\xi}_{\text{DI}}$) achieved by the estimators

$\hat{\xi}_{\text{LH}}$, $\hat{\xi}_{\text{AV}}$, and $\hat{\xi}_{\text{S,LH}}$ in the context of stochastic activity network simulation. Stochastic activity networks (SANs) are often used to model projects whose activities have precedence constraints. The graph-theoretic structure of a stochastic activity network is described by the pair $(\mathcal{N}, \mathcal{A})$, where $\mathcal{N} = \{1, \dots, \nu\}$ is the set of nodes (vertices) in the network and $\mathcal{A} = \{(a_j, b_j) : \text{activity } j \text{ has start node } a_j \in \mathcal{N} \text{ and end node } b_j \in \mathcal{N}, j = 1, \dots, p\}$. The network is assumed to be acyclic, with source node $r \in \mathcal{N}$ and sink node $s \in \mathcal{N}$. Each activity j has a random duration V_j , so the input random variates are $\{V_j : j = 1, \dots, p\}$, and the probabilistic structure of the network is described by the joint distribution of the random vector (V_1, \dots, V_p) . Let α denote the number of directed r -to- s paths, and let $A(\ell)$ denote the index set of activities on the ℓ th path, so $A(\ell) = \{j : \text{activity } j \text{ is on the } \ell\text{th directed } r\text{-to-}s \text{ path}\}$ for $\ell = 1, \dots, \alpha$. The duration of the ℓ th path is the random variable $P_\ell = \sum_{j \in A(\ell)} X_j$, and the *network completion time* is $T \equiv \max\{P_1, \dots, P_\alpha\}$. We consider the problem of estimating selected quantiles of the network completion time.

The SAN we used as an example represents the construction of a rock-fill dam and was taken from Antill and Woodhead (1982), Figure 8.5(b), page 189. For each activity duration V_i , the associated distribution was taken to be either (a) a normal distribution with a specified mean μ_i and standard deviation $\sigma_i = \mu_i/4$ whose tail was truncated below the value 0; or (b) an exponential distribution with a specified mean μ_i . We chose the exponential distribution as the nonnormal alternative for reasons elaborated in Avramidis, Bauer, and Wilson (1991). The set of activities with durations as in (a) was taken to be $\{(1,3), (2,6), (2,4), (8,11), (10,13), (12,18), (16,17), (17,21), (17,23), (17,19), (18,19), (23,24)\}$. In addition, we assumed that the activity durations are independent; and we generated all these durations by inversion. As a result, the response function $f(\cdot)$ is monotone in each coordinate, and thus the results of the previous sections apply here.

Table 1 shows the resulting MSE ratios when estimating the quantile of order p of the network completion time. As expected from the results of the previous sections, $\hat{\xi}_{\text{LH}}$ and $\hat{\xi}_{\text{AV}}$ achieve MSE reductions with respect to $\hat{\xi}_{\text{DI}}$. Note, however, that $\hat{\xi}_{\text{LH}}(k, n)$ performs significantly better than $\hat{\xi}_{\text{AV}}(n)$ for $k > 2$. This behavior was observed in several experiments (not reported here), so we would recommend $\hat{\xi}_{\text{LH}}$ with $k > 2$ over $\hat{\xi}_{\text{AV}}$.

To use $\hat{\xi}_{\text{LH}}(k, n)$, a practitioner would probably have to choose k given the total number of replications n . For fixed n , variance is typically decreasing in k (due to the more complete stratification), while

Table 1: Estimated $\text{MSE}[\hat{\xi}_{\text{DI}}(n)]/\text{MSE}[\hat{\xi}(n)]$ for Various Estimators $\hat{\xi}(n)$ and $n = 2048$

	Order p				
	0.05	0.25	0.50	0.75	0.95
$\hat{\xi}_{\text{AV}}(n)$	1.05	1.20	1.24	1.02	1.06
$\hat{\xi}_{\text{LH}}(2, n)$	1.07	1.27	1.22	1.06	1.05
$\hat{\xi}_{\text{LH}}(4, n)$	1.25	1.43	1.43	1.16	1.15
$\hat{\xi}_{\text{LH}}(8, n)$	1.21	1.60	1.82	1.47	1.26
$\hat{\xi}_{\text{LH}}(16, n)$	1.25	1.66	2.03	1.99	1.41
$\hat{\xi}_{\text{LH}}(32, n)$	1.21	1.66	1.96	2.32	1.91
$\hat{\xi}_{\text{LH}}(64, n)$	0.77	1.34	1.91	2.33	1.84
$\hat{\xi}_{\text{S,LH}}(n)$	1.36	1.80	2.48	2.41	2.34

bias is typically increasing in k (due to the smaller subsample size $m = n/k$). The net effect is that $\text{MSE}[\hat{\xi}_{\text{LH}}(k, n)]$ is typically decreasing for $k \leq k_0$ and increasing for $k \geq k_0$, with the critical value k_0 being an increasing function of n . For example, from Table 1 we see that $k_0 \approx 16$ for $p = 0.50$. As a first-order heuristic, we recommend using $k = O(n^{1/2})$, which we have experimentally found to be a fairly robust strategy for a wide range of values of n and p .

Although no theoretical guarantee for $\hat{\xi}_{\text{S,LH}}$ is currently available, the results shown here and further Monte Carlo experiments suggest that $\hat{\xi}_{\text{S,LH}}$ not only yields significant MSE reductions, but also dominates both $\hat{\xi}_{\text{AV}}$ and $\hat{\xi}_{\text{LH}}$. The development of the asymptotic properties of $\hat{\xi}_{\text{S,LH}}$ is the subject of ongoing research.

ACKNOWLEDGMENTS

The author thanks Professors James Wilson of North Carolina State University and Bruce Schmeiser of Purdue University for their comments. This work was supported by the Purdue Research Foundation under Grant No. 690-1287-1571.

REFERENCES

- Antill, J. M. and R. W. Woodhead. 1982. *Critical Path Methods in Construction Practice*. New York: Wiley.
- Avramidis, A. N., K. W. Bauer, and J. R. Wilson. 1991. Simulation of stochastic activity networks using path control variates. *Naval Research Logistics* 38:183-201.
- David, H. A. *Order Statistics*. 1981. 2d ed. New York: Wiley.

- Hsu, J. C. and B. L. Nelson. 1990. Control Variates for Quantile Estimation. *Management Science* 36:835-851.
- Juritz, J. M., J. W. F. Juritz, and M. A. Stephens. 1983. On the accuracy of simulated percentage points. *Journal of the American Statistical Association* 78:441-444.
- Lehmann, E. L. 1966. Some Concepts of Dependence. *Annals of Mathematical Statistics* 37:1137-1153.
- Lewis, P. A. and R. L. Ressler. 1989. Variance Reduction of Quantile Estimates via Nonlinear Controls. In *Proceedings of the 1989 Winter Simulation Conference*, ed. E. MacNair, K. Musselman, and P. Heidelberger, 450-454. Institute of Electrical and Electronics Engineers, Piscataway, New Jersey.
- McKay, M. D., R. J. Beckman, and W. J. Conover. 1979. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* 21:239-245.
- Schafer, R. E. 1974. On assessing the Precision of Simulation. *Journal of Statistical Computation and Simulation* 3:67-69.
- Stein, M. 1987. Large sample properties of simulations using Latin hypercube sampling. *Technometrics* 29:143-151.
- van Zwet, W. R. 1964. Convex Transformations of Random Variables. Mathematical Center Tracts 7, Mathematisch Centrum, Amsterdam.
- Zelterman, D. 1987. Estimating Percentage Points by Simulation. *Journal of Statistical Computation and Simulation* 27:107-125.

AUTHOR BIOGRAPHY

ATHANASSIOS N. AVRAMIDIS is a Ph.D. candidate in the School of Industrial Engineering at Purdue University. He received a diploma in mechanical engineering from the University of Thessaloniki (Greece) in 1987, and he received an M.S. in industrial engineering from Purdue University in 1989. His research interests are in probabilistic and statistical issues in the design and analysis of simulation experiments, particularly in input modeling, output analysis, and variance reduction techniques. He is also interested in applied probability and operations research applications. He is a student member of ORSA, TIMS, and IIE.