

THE USE OF ANTITHETIC CONTROL VARIATES IN COMPUTER SIMULATIONS

R.C.H. Cheng
Department of Mathematics
UWIST
King Edward VII Avenue
Cardiff CF1 3NU
Britain

ABSTRACT

The method of antithetic variates is a well-known technique for reducing the variability of estimators in computer simulation experiments. However, the usually suggested way of using the method, if incorrectly applied, can lead to an increase in the variance of estimators of certain quantities, such as percentiles. A procedure, based on two non-standard methods of generating samples from the normal distribution, is suggested which does not suffer this weakness. Numerical examples are given showing the ease of implementation and the effectiveness of the procedure.

1. INTRODUCTION

There has been much recent interest in developing computer methods for generating random variates from specified distributions such as the normal, gamma and beta. However, as has been pointed out by Andrews (in the Discussion of Atkinson and Pearce's paper, 1976), the cost of generating variates is often very small relative to the total cost of the simulation experiment. Variance reduction methods which apply to the cost of the simulation as a whole therefore offer much longer savings.

In this paper a variance reduction procedure is suggested, using antithetic variates, which might be applied in a simple and automatic way in a number of computer simulation situations. The procedure is based on two novel algorithms for generating random samples of specified size n from the standard normal distribution: RSNA and RSNB. The procedure and the two algorithms are described in section 3. In section 4, two numerical examples are given making use of the procedure, which illustrate its effectiveness.

This work is an extension of that of Cheng (1981) where similar methods are applied to control variables which are sums of independent and identically distributed random variables.

2. ANTITHETIC CONTROL VARIABLES

Before describing the proposed antithetic procedure, it will be of help to look at its motivation. Its basis is the fairly well-known method of antithetic variates as applied to control variables (see, for example, Tocher 1963), whose rationale is as follows.

Suppose a run of a computer simulation model yields a response y whose mean μ is to be estimated. If two runs are made yielding response values y and y' , then their average

$$\bar{y} = \frac{1}{2}(y + y')$$

estimates μ . If the two runs are independent then the variance of \bar{y} is

$$\text{var}(\bar{y}) = \frac{1}{2} \text{var}(y).$$

The method of antithetic variates aims to reduce this variance by introducing negative correlation between the runs. Then $\text{var}(\bar{y})$ becomes

$$\text{var}(\bar{y}) = \frac{1}{2}(1 + r)\text{var}(y),$$

where $r < 0$ is the correlation between y and y' . This negative correlation can be obtained by first selecting a variable x which is highly correlated with y but whose distribution is known, and whose sampling can be controlled directly in the simulation. Such a variable is usually called a control variable. The idea is then to negatively correlate the control variables x and x' in the two runs, rather than correlate y and y' . Because y and x are highly correlated, the

negative correlation between x and x' will induce a similar correlation between y and y' .

The most frequently suggested way of correlating x and x' is the inverse distribution function transform (IDFT) method. If $F(\cdot)$ is the c.d.f. of x , set

$$x = F^{-1}(u), \quad x' = F^{-1}(1-u) \quad (2.1)$$

where u is a uniform $U(0,1)$ variable. The degree of negative correlation, ρ , obtained between x and x' in this way depends on the shape of the density function of x . It is best (i.e. it is most negative) when x is symmetrical distributed, when $\rho = -1$. If x is not symmetrical then $\rho > -1$. For example Page (1965) shows that, if x is negative exponentially distributed, then $\rho = -0.645$.

In many simulations it is not just μ , the mean of y , that is of interest. Its variance, certain percentiles of its distribution, or even its distribution as a whole may be required. To estimate such quantities a set of runs is needed. It is tempting to apply the IDFT method of negative correlation in blanket fashion by making a set of runs divided into two blocks each of n runs, and to correlate corresponding pairs of runs from the two blocks. Thus if x_j and x'_j are the control variate values in the j th run of the first and second blocks respectively, then the same $U(0,1)$ variate value, u_j , is used in calculating x_j and x'_j as

$$x_j = F^{-1}(u_j), \quad x'_j = F^{-1}(1-u_j), \quad j = 1, 2, \dots, n. \quad (2.2)$$

This is a bad policy however.

Firstly, though this policy reduces the variance in the estimate of μ , it tends to increase the variance in the estimate of σ^2 , the variance of y . For example, when x is symmetrically distributed, both the first and second blocks of runs yield exactly the same value for the sample variance of the x 's. In other words no additional information on the variance of x has been obtained from the second block. As y is assumed highly correlated with x , this means that little, if any, additional information has been obtained about the variance of y either. The second block of runs has thus been made without improving the estimate of σ^2 and in effect the antithetic procedure (2.2) has led to variance increase not variance reduction. As percentiles depend on both the mean and variance, the blanket use of the IDFT method can thus lead to variance increase in the estimation of percentiles, particularly for those near the tails of the distribution.

Secondly, even if the estimation of μ is being considered, though use of (2.2) will usually lead to variance reduction, the degree of reduction may not be the best achievable. This is because the overall level of correlation between the two blocks will not exceed that achieved between individual pairs of runs. It turns out that one can do better if runs are not correlated in pairs, but if the two blocks are considered as a whole. The procedure described in the following section does this.

3. THE PROPOSED PROCEDURE.

The antithetic procedure suggested below makes use of two alternative ways of generating a random sample z_1, z_2, \dots, z_n of size n from the standard normal distribution. These will be described first.

The first algorithm: RSNA, has the following property. If z_1, z_2, \dots, z_n is a random sample of standard normal variables, then it is well-known that the sample mean \bar{z} is normal with mean zero and variance equal to n^{-1} , i.e. $\bar{z} \sim N(0, n^{-1})$, and that the sample variance

$$v = (\sum z_i^2 - n\bar{z}^2)/(n-1)$$

is independent of \bar{z} , with $w = (n-1)v$ a χ^2 variable with $n-1$ degrees of freedom i.e. $w \sim \chi^2_{n-1}$.

Algorithm RSNA allows the values of \bar{z} and v to be sampled first, then constructs exact $N(0,1)$ variables z_1, z_2, \dots, z_n in such a way that their sample mean and variance are equal to these sampled values.

Algorithm RSNA. Inputs : $z_0 \sim N(0, n^{-1})$,
 $w_0 \sim \chi^2_{n-1}$

Constants : (These need only be set once.)

$$v = n^{-1}, \quad \alpha_j = [j/(j+1)]^2, \quad j = 0, 1, \dots, v.$$

1. Generate v NID(0,1) variates: v_1, v_2, \dots, v_v (independent of z_0 and w_0)
2. Set $a = (w_0 / \sum_{i=1}^v v_i^2)^{1/2}$; $t_j = av_j$, $j = 1, 2, \dots, v$.
3. Set $z_n = z_0 - \alpha_v t_v$
 $z_j = z_{j+1} + \alpha_j^{-1} t_j - \alpha_{j-1} t_{j-1}$, $j = v, v-1, \dots, 1$.
4. Return with z_1, z_2, \dots, z_n , which are exact NID(0,1) variables with mean $\bar{z} = z_0$ and variance $v = w_0/(n-1)$. \square

A detailed derivation of the Algorithm is given in Cheng (1981). The key is that the t_j 's of step 2 are NID(0,1) variates with the property $\sum t_j^2 = w_0$.

The variates z_1, z_2, \dots, z_n are merely a linear transformation (actually the inverse of Helmert's transformation, see Kendall and Stuart, 1968) of z_0 and the t 's. It is readily verified that z_1, z_2, \dots, z_n are exact NID(0,1) variates, with sample mean and variance as claimed.

The second algorithm: RSNB, has the following property. If a_1, a_2, \dots, a_n is a set of constants such that $\sum a_i^2 = 1$, then clearly $z = \sum a_i z_i$ will be a standard normal variable if z_1, z_2, \dots, z_n are

NID(0,1) variables. RSNB allows the value of z to be sampled first, then constructs z_1, z_2, \dots, z_n in such a way that $\sum a_i z_i$ equals the sampled z value.

Algorithm RSNB. Input : $z \sim N(0,1)$.

Constants : (These need only be set once; $a_j, j=1, 2, \dots, n$ are assumed given with

$$\Sigma a_j^2 = 1.) \quad b_j = (\sum_{i=j}^n a_i^2)^{1/2}, r_j = a_j/b_j, s_j = b_j/b_{j+1},$$

$$j = 1, 2, \dots, v \equiv n-1.$$

1. Set $t_1 = z$
2. Set $z_j = r_j t_j + s_j v_j, t_{j+1} = (t_j - r_j z_j)/s_j,$
 $j = 1, 2, \dots, v$
 where v_1, v_2, \dots, v_v are v NID(0,1) variables
 (independent of z).
3. Set $z_n = t_n$
4. Return with z_1, z_2, \dots, z_n , which are exact
 NID(0,1) variables satisfying the condition $\Sigma a_i z_i = z.$ □

Clearly z_1, z_2, \dots, z_n are linear combinations of z and v_1, v_2, \dots, v_v , and so are normal. Direct calculation shows that they are independent standard normal and moreover that $\Sigma a_i z_i = z.$

In what follows, the simulation is assumed to be set out in two blocks of runs each consisting of n independent but identical runs. Assume first that the control variable is a standard normal variable which will be denoted by $z.$ This assumption will be considerably relaxed later.

The value of z is recorded in each run of the first block of runs as: $z_1, z_2, \dots, z_n.$ This allows their sample mean and variance

$$z_0 = \bar{z}, v_0 = (\sum_{i=1}^n z_i^2 - n\bar{z}^2)/(n-1)$$

to be obtained at the end of this block of runs. It is well-known that z is normally distributed $N(0, n^{-1})$ variable, whilst $(n-1)v_0 (=w_0, \text{ say})$ is a χ^2 variable, with $(n-1)$ degrees of freedom, that is independent of $z_0.$

The basic idea is to apply the antithetic technique to z_0 and w_0 only. Antithetic versions of z_0 and w_0 can easily be obtained using the IDFT method of equation (2.1). Because z is symmetrically distributed about zero, this yields

$$z_0' = -z_0 \tag{3.1}$$

whilst

$$w_0' = G^{-1}[1 - G(w_0)]. \tag{3.2}$$

where G is the c.d.f of a χ^2 variable with $(n-1)$ degrees of freedom. (Here and in what follows, a prime will indicate a quantity associated with the second block of runs). The functions G and G^{-1} are awkward to obtain explicitly. An approximation for (3.2) is

$$w_0' = v[2[1-2/(9v)] - (w_0/v)^{1/3}]^3 \tag{3.3}$$

where $v = n-1.$ This formula is explained in Cheng (1981). It makes use of the Wilson-Hilferty approximation twice in such a way that errors tend to cancel so that it is accurate down to small values of $n.$ For example, when $n=4$ the probability for which w_0' is a quantile is never in error from the exact probability by more than 0.004.

Once z_0' and w_0' are obtained, the second block of runs must be made with control variate values

z_1', z_2', \dots, z_n' generated in the n runs in such a way that their sample mean $\bar{z}' = z_0'$ and variance $v_0' = w_0'/(n-1).$ The point of algorithm RSNA is now evident, as it allows precisely such a set of antithetic z' variables to be obtained, using z_0' and w_0' as inputs.

This method of correlation has the effect of correlating both the means and the variances of the responses in the two blocks. This leads to an improvement not only in the estimate of the mean of y but in estimates of quantities which depend on the variance of $y,$ such as percentiles.

The above procedure is only a prototype, as it is restricted by requiring the control variable to be normally distributed. A more general assumption is that a control variable x can be found which is itself dependent on p input variables $x_1, x_2, \dots, x_p:$

$$x = f(x_1, x_2, \dots, x_p) \tag{3.4}$$

The x_i 's will not be required to be normal or even identically distributed; but they will need to be independent.

The prototype procedure cannot be applied to x directly. However it can still be applied if a normal variable z is first constructed from $x.$ This variable will be called the "normal version of x ". A simple way to do this is to find an approximation to f that is linear in x_1, x_2, \dots, x_p (using a Taylor series expansion for example) and then to replace each x_i by a normal variable with the same mean μ_i and variance σ_i^2 as $x_i.$ Thus

$$x \approx f(\mu_1, \mu_2, \dots, \mu_p) + \Sigma (\partial f / \partial x_i) x_i \tag{3.5}$$

and replacing x_i by $(\mu_i + \sigma_i z_i),$ where z_i is $N(0,1)$

$$\hat{z} = f(\underline{\mu}) + \Sigma [\partial f(\underline{\mu}) / \partial x_i] (\mu_i + \sigma_i z_i).$$

This is of the form $z = c_0 + \Sigma c_i z_i;$ for the purposes of antithetic correlation it is easiest to work with the standardised normal variable.

$$z = (\hat{z} - c_0) / (\Sigma c_i^2)^{1/2} = \Sigma a_i z_i, \text{ say (where } \Sigma a_i^2 = 1) \tag{3.6}$$

The prototype procedure can now be applied with z replacing x as the control variable. When $x_{ij},$ the i th input variable in the j th run, is generated, a normal version z_{ij} of x_{ij} must also be produced that is strongly correlated with $x_{ij}.$ At the end of the j th run the control variate z_j is obtained from the z_{ij} using (3.6). The prototype procedure can then be applied to produce an antithetic set of z_j' for the second block of runs. In this second block the process of going from x_{ij} to z_j is reversed. In the j th run, $z = z_j'$ is used as input to Algorithm RSNB to produce a set of normal variates x_{ij}' satisfying $\Sigma a_i z_{ij}' = z_j'. The z_{ij}' are then used as normal versions of the input variates x_{ij}' actually used in the runs. Strong correlation between z_{ij}' and x_{ij}' is obtained by using z_{ij}' to produce a uniform variate u_{ij}' = \Phi(z_{ij}') (where \Phi is the standard normal c.d.f) which is then used to generate x_{ij}' by the IDFT method. The whole procedure is as follows:$

Procedure A

1. Make the first block of runs. Generate input variates and normal versions by the IDFT method.

$$x_{ij} = F_i^{-1}(u_{ij}), z_{ij} = \Phi^{-1}(u_{ij}), i=1,2,\dots,p; \\ j=1,2,\dots,n.$$

Set $z_j = \sum a_i z_{ij}, j=1,2,\dots,n.$

2. Calculate the sample mean and variance of z_j and find their antithetic versions by (3.1) and (3.3). Then use Algorithm RSNA to produce a set of antithetic $z_j', j=1,2,\dots,n.$

3. Make the second block of runs. In the j th run use z_j' as input to Algorithm RSNB to construct $z_{1j}', z_{2j}', \dots, z_{pj}'$. The input variates of x_{ij}' are then obtained as

$$x_{ij}' = F_i^{-1}[\Phi(z_{ij}')] \quad i=1,2,\dots,p$$

where F_i is the c.d.f. of x_i and Φ is the standard normal c.d.f. □

4. EXAMPLES

The two examples described below have been kept deliberately simple as they are intended to be illustrative only. There is no difficulty in using the methods for more complicated and realistic examples of the same sort.

Example 1. Bury (1975, §15.24) describes the simulation of the current gain of a certain transistor amplifier defined as

$$G = \beta R_0 (R_0 + R_C)^{-1},$$

where β and R_0 are characteristics of the transistor and R_0 is an external resistance. The quantities β, R_0 and R_C are assumed to be random variables with 0 distributions as given in Table 1. It is desired to find the c.d.f. of the current gain G .

The simulation was organised as follows. A run consisted of sampling a set of β, R_0, R_C values and calculating G for this amplifier "specimen". A block of n such determinations constituted the first block of runs. A second block of n was then obtained. However, for comparison two different versions of the second block were carried out.

Version A of the second block was obtained using Procedure A. The control variable was a normal version of G itself. Expanding G in a Taylor series about the mean values of β, R_0 and R_C gives

$$G \sim (\partial G / \partial \beta) \beta + (\partial G / \partial R_0) R_0 + (\partial G / \partial R_C) R_C \\ = m_0 (m_0 + m_C)^{-1} \beta + m_\beta (m_0 + m_C)^{-2} [m_C R_0 - m_0 R_C].$$

Table 1.

Variable	Distribution & Method of Generation	Parameters	Mean (m) Standard Deviation (s)
β	Weibull : $\beta = \sigma_\beta [-\log(1-u)]^{1/\lambda_\beta}$ where $u \sim U(0,1)$	$\sigma_\beta = 110$ $\lambda_\beta = 4.3$	$m_\beta = 100$ $s_\beta = 26.3$
R_0	Lognormal : $R_0 = \exp[\mu_0 + \sigma_0 z]$ where $z \sim N(0,1)$	$\mu_0 = 9.2$ $\sigma_0 = 0.12$	$m_0 = 10,000\Omega$ $s_0 = 1,200\Omega$
R_C	Normal : $R_C = \mu_C + \sigma_C z$ where $z \sim N(0,1)$	$\mu_C = 10,000$ $\sigma_C = 600$	$m_C = 10,000\Omega$ $s_C = 600\Omega$

The normal version of this is

$$\hat{z} = \text{constant} + m_0 (m_0 + m_C)^{-1} s_\beta z_\beta + m_\beta (m_0 + m_C)^{-2} \\ \times [m_C s_0 z_0 - m_0 s_C z_C].$$

Substituting in the values of the parameters and standardising the normal variate gives the control variable used in Procedure A as :

$$z = .969z_\beta + .221z_0 - .111z_C.$$

Version B of the second block consisted of an independent block identical to the first block.

The means and variances of the G 's of each block will be written as

$$\bar{G}, \bar{G}_A, \bar{G}_B, V = \Sigma G^2 / n - \bar{G}^2, V_A = \Sigma G_A^2 / n - \bar{G}_A^2,$$

$$V_B = \Sigma G_B^2 / n - \bar{G}_B^2.$$

Combined estimates like $(\bar{G} + \bar{G}_A) / 2$ and $(V + V_A) / 2$ can then be obtained by pooling the first block with each version of the second block. Also the c.d.f. of G can be estimated by pooling the values of G obtained from the first block with those from version A or version B of the second block. To see the effect that Procedure A has of reducing the variance of pooled estimators it is necessary to replicate the entire above experiment a number of times. Table 2 shows the means and variance of pooled estimators obtained from 1000 replicates. As will be seen the variance of the estimates of the mean of G is reduced from 1.9 to 0.01, and of the variance of G is reduced from 637 to 28.

From each replicate two empirical c.d.f.'s can be obtained by pooling the first block with each version of the second block. For each replicate the values of the two c.d.f.'s were calculated at a selection of G values. Table 2 summarises these results by giving the mean values of these c.d.f.'s (averaged over the 1000 replicates) at different G values and also the sample variances of the c.d.f.'s about these means. As is seen, use of the antithetic procedure A leads to a variance reduction over most of the range of G values. Thus for example, the proportion of gains under 40 is estimated, with 95% confidence limits, as

$$P_A(40) = .1410 \pm 1.96 \times (.000435/1000)^{1/2} = .1410 \pm .0013$$

if procedure A is used, compared with

$$P_B(40) = .1411 \pm 1.96 \times (.00119/1000)^{1/2} = .1411 \pm .0022$$

if independent second blocks are used.

Table 2. Results of 1000 replicates of a simulation experiment estimating the gain of a transistor amplifier. Block size n=50.

	$\frac{1}{2}(G+G_A)$	$\frac{1}{2}(G+G_B)$	$\frac{1}{2}(V+V_A)$	$\frac{1}{2}(V+V_B)$
Mean of 1000 values.	49.847	49.827	180.1	180.2
Variance of 1000 values	.011	1.88	28.	637.
	Mean of 1000 Empirical C.D.F. Values		Variance of 1000 Empirical C.D.F. Values	
With Second Block :	Antithetic	Independent	Antithetic	Independent
G			x 100	x 100
10	.0000	.0001	.000020	.000079
20	.0042	.0041	.0038	.0042
30	.0358	.0355	.020	.036
40	.1410	.1411	.043	.119
50	.3581	.3591	.075	.226
60	.6414	.6425	.083	.234
70	.8661	.8672	.043	.123
80	.9710	.9709	.017	.030
90	.9967	.9961	.0033	.0039
100	.9998	.9998	.00020	.00027

Example 2. Battersby(1970)described an oil refinery maintenance problem involving critical path analysis. The 18 jobs and their durations are

Job	A	B	C	D	E	F	G	H	K	L	M	N	P	Q
Duration	16	16	8	6	16	40	24	16	16	24	8	4	36	12
Job	R	S	T	U										
Duration	8	16	8	24.										

The duration of the project can be written as

$$Y = \max(J+P, E+Q, F+R, G+W+S)+T+U$$

where $J = \max(A+\max(B,C) + D,E)$, $W = \max(N, H+\max(K,L)+ M)$ and where for simplicity the duration of jobs have been denoted by their letters. Suppose that job durations are random. As the example is for illustration only, suppose for simplicity that they are all lognormally distributed with means as shown, but all with the same skewness. The i th job duration can then be sampled as $d_i = m_i \exp[-\sigma^2/2 + \sigma z_i]$, where m_i is the mean of the i th job duration, σ is the same for all jobs and z_i is a standard normal variable. In the simulation $\sigma = 0.3$.

As in the first example simulation runs were made

in blocks of three : the first block and two versions of the second block. A simple control variate to use is the sum of the job durations on the critical path obtained when job durations are set equal to their means. With the parameter values as given this yields the control variate as

$$X = G+H + L+ M+ S+ T+ U.$$

A normal version of X is obtained by replacing each job duration in X by a normal variate with the same mean and variance. As all the normal versions of the d_i 's used in forming X have variances proportional to their means(with same constant of proportionality)it is easily seen that the standardized normal version of X is

$$z = \frac{\sum m_i (\sum m_i^2)^{-\frac{1}{2}} z_i}{\sum m_i (\sum m_i^2)^{-\frac{1}{2}} z_i} = .493(z_G+z_L+z_U) + .329(z_H+z_S) + .164(z_M+z_T).$$

Table 3 gives results analogous to those of Example 1 appearing in Table 2. Again it is seen that procedure A is very effective in reducing estimates of the mean, variance and c.d.f. of Y, the overall project duration.

Table 3. Results of 1000 replicates of a simulation experiment estimating Y, the duration of a project. Block size n = 50.

	$\frac{1}{2}(Y+Y_A)$	$\frac{1}{2}(Y+Y_B)$	$\frac{1}{2}(V+V_A)$	$\frac{1}{2}(V+V_B)$
Mean of 1000 values	122.99	123.02	199.	199.
Variance of 1000 values	.42	2.04	425.	1004.
	Mean of 1000 Empirical C.D.F. Values		Variance of 1000 Empirical C.D.F Values	
With Second Block	Antithetic	Independent	Antithetic	Independent
Y			x 100	x 100
90	.0005	.0004	.00047	.00042
110	.0899	.0895	.052	.076
120	.3036	.3024	.112	.201
125	.4456	.4443	.120	.242
130	.5858	.5876	.113	.244
135	.7108	.7109	.093	.200
140	.8096	.8095	.078	.160
150	.9306	.9298	.038	.068
160	.9781	.9776	.016	.023
180	.9982	.9985	.0019	.0016
190	.9996	.9996	.00038	.00041

5. CONCLUSIONS.

Procedure A offers a simple method of general applicability for reducing the variance of estimates of both the mean and variance of a response variable. This leads to improved estimation of the c.d.f. of the response variable as a whole. The procedure does depend on being able to generate certain input variates by the IDFT method. Moreover variates in the antithetic block will take longer to be generated using Procedure A than if an independent block identical to the first is taken and this must be set against the reduction in variance of estimators. Two features tend to reduce this effect however. Firstly, input variates often require normal variates to be generated as an intermediate step so that generation of normal versions of certain variates in Procedure A is not always disadvantageous. Secondly, as remarked in the introduction, often the generation of variables is a relatively small part of the simulation experiment as a whole, so that the increase in variate generation time is far outweighed by the improvement in overall cost.

REFERENCES

- Atkinson, A.C. and M.C.Pearce (1976), The computer generation of beta, gamma and normal random variables (with Discussion), *Journal of the R.S.S.*, Vol.139, pp.431-461.
- Battersby, A.(1970), *Network Analysis for Planning and Scheduling*. (3rd edn), Macmillan, London, p.21.
- Bury, K.V.(1975), *Statistical Models in Applied Science*, Wiley, New York, p.546.
- Cheng, R.C.H. (1981), *The Use of Antithetic Variates in Computer Simulations*, MATH REPORT 81-1, Department of Mathematics, UWIST, Cardiff, Britain.
- Kendall, M.G. and A.Stuart (1968), *The Advanced Theory of Statistics*, Vol.1, Griffin, London.
- Page, E.S.(1965), *On Monte Carlo Methods in Congestion Problems II - Simulation of Queueing Systems*, *Operations Research*, Vol.13, pp 300-305.
- Tocher, K.D.(1963), *The Art of Simulation*, English Universities Press, London.