

## CONTROL VARIATES IN NONLINEAR REGRESSION

James J. Swain  
School of Industrial and Systems Engineering  
Georgia Institute of Technology  
Atlanta, Georgia 30332

Control variates can be applied to Monte Carlo sampling experiments to improve the precision of the results. This method is especially useful in statistical problems where low order approximators of a particular variate of interest are available and possibly several statistical properties of the variate are to be investigated. In this paper a control variate scheme based on the linear approximator  $\delta$  of the nonlinear parameter estimator  $\hat{\theta}$  is used to improve the precision of the first four moments of  $\hat{\theta}$  and the covariance matrix of the parameter estimates. The control variate method is shown to improve the effectiveness of the Monte Carlo results without substantially increasing the estimation effort, and it is effective over a wide range of nonlinearities. An approximate expression for the effectiveness of the control variate method based on the Beale measure of nonlinearity  $N_{\theta}$  is given.

### 1. INTRODUCTION

In the nonlinear parameter estimation or regression problem, the  $p$ -vector of parameters  $\underline{\theta}$  are estimated from the  $n$  responses  $y_i$  ( $i=1, \dots, n$ ) which are assumed to consist of the true response  $\eta(\underline{x}_i; \underline{\theta}_0)$  plus an additive error  $\epsilon_i$ ,

$$y_i = \eta(\underline{x}_i; \underline{\theta}_0) + \epsilon_i. \quad (1)$$

The  $\epsilon_i$  are assumed to be independently and normally distributed random errors with variance  $\sigma^2$ , and  $\eta(\underline{x}; \underline{\theta})$  is a nonlinear function of the parameters  $\underline{\theta}$ . The unknown parameters  $\underline{\theta}_0$  are most often estimated by the method of least squares. That is, the estimator  $\hat{\underline{\theta}}$  is chosen to minimize the sum of the squared residuals

$$S(\underline{\theta}) = \sum_{i=1}^n (y_i - \eta(\underline{x}_i; \underline{\theta}))^2 = \sum_{i=1}^n e_i^2(\underline{\theta}) \quad (2)$$

so that  $s(\hat{\underline{\theta}}) = \min S(\underline{\theta})$ . Of course, other estimators for  $\underline{\theta}$  can be used, but the least

squares method is the most common and will be treated here.

The nonlinear regression problem is similar to that of linear regression, save that the nonlinearity of the response function complicates the numerical problem of obtaining the estimates and of resolving the statistical properties of the estimators obtained. The solution to the least squares problem is generally not a closed form function of the observed responses  $y_i$  so that iterative procedures are needed to obtain a solution (Bard, 1974; Gallant, 1975). The lack of a closed form solution makes the distribution of the estimator  $\hat{\underline{\theta}}$  intractable, and the sampling distribution of the estimator  $\hat{\underline{\theta}}$  is known exactly only asymptotically. The asymptotic approximation of the sampling distribution can be very misleading in problems with a small sample size. In particular, the finite sample estimator is almost always biased and the shape of the confidence region can differ markedly from the elliptical contours of the asymptotic sample distribution.

The statistical properties of the estimator  $\hat{\theta}$  can be approached using the asymptotic distribution, through a series approximation solution, or via Monte Carlo sampling. Under certain circumstances (see Gallant, 1975, for instance) the distribution of  $\hat{\theta}$  is asymptotically normal with mean  $\theta_0$  and variance matrix  $(F^T(\theta_0)F(\theta_0))^{-1}\sigma^2$ , where  $F(\theta_0)$  is the  $n$  by  $p$  Jacobian matrix of first derivatives of  $\eta(\underline{x};\theta)$  with respect to the parameters  $\theta$  and evaluated at  $\theta_0$ . Note that the variance matrix is the limiting matrix as the sample size  $n$  increases indefinitely. To use the asymptotic approximation in practice one uses the finite sample and the two formulas given above. This solution is essentially the first order approximation to be described.

Series approximations can be obtained by approximating the sum of squares function  $S(\theta)$  or its derivative  $\partial S(\theta)/\partial\theta$  (which vanishes at  $\hat{\theta}$ ) at some point such as  $\theta_0$ . The approximator results from minimizing the approximate sum of squares function or solving for the point which solves the equation  $\partial S(\theta)/\partial\theta = 0$ . For instance, the first order approximator  $\hat{\delta}$  for  $\hat{\theta}$  is given by

$$\hat{\delta} = \theta_0 + (F^T(\theta_0) F(\theta_0))^{-1} F^T(\theta_0) \underline{\varepsilon}. \quad (3)$$

This is similar to the asymptotic solution, in that  $\hat{\delta}$  is normally distributed with mean  $\theta_0$  and variance matrix  $(F^T(\theta_0) F(\theta_0))^{-1}\sigma^2$ , although in this case the number of rows in  $F(\theta)$  will be finite. The approximator  $\hat{\delta}$  corresponds to the Gauss approximation method for finding a solution to the nonlinear least squares problem. Higher order approximations can be derived (Box, 1971; Clarke, 1981), but simplifications are introduced in each to make the solution tractable in the multiple parameter case. Note that measures of nonlinearity (Beale, 1960; Bates and Watts, 1980) such as Beale's  $N_\theta$  have been developed to assess the extent to which  $\hat{\delta}$  can be used to approximate the statistical properties of  $\hat{\theta}$ . The measures can readily be adopted to any parameterization of the model, e.g.,  $\psi = \psi(\theta)$ , and Beale denotes the least possible nonlinearity under any parameterization as the intrinsic nonlinearity,  $N_\phi$ .

Monte Carlo methods can in principle be used to overcome the limitations of approximation methods. A direct Monte Carlo algorithm consists of  $N$  repeated independent samples of the error vector  $\underline{\varepsilon}(\underline{\varepsilon}_v, v=1,1,\dots,N)$ , from which (using (1)) the observed vectors  $\underline{y}_v$  are obtained and used to obtain estimates  $\hat{\theta}_v$  by repeated solution of equation (2) for each  $\underline{y}_v$ . The  $N$  random sample points  $\hat{\theta}_v$  can be combined to obtain the sample statistics to estimate the properties of  $\hat{\theta}$ . For instance, the  $k$ th marginal moments of  $\hat{\theta}$ ,  $\mu_k = E(\hat{\theta})^k$  (where the exponentiation is component by component) can be estimated using the statistics  $\hat{\mu}_k$ ,

$$\hat{\mu}_k = N^{-1} \sum_{v=1}^N (\hat{\theta}_v)^k \quad (4)$$

and the variance matrix of the parameter estimators  $\hat{\theta}$ ,  $\text{Var}(\hat{\theta}) = \Sigma_\theta$  can be estimated using the sample statistic  $\hat{\Sigma}_\theta$ ,

$$\hat{\Sigma}_\theta = (N-1)^{-1} \sum_{v=1}^N (\hat{\theta}_v - \hat{\mu}_1)(\hat{\theta}_v - \hat{\mu}_1)^T. \quad (5)$$

The Monte Carlo method is conceptually simple and can be easily extended to nonlinear least estimators other than nonlinear least squares or to errors other than normal, but suffers from the drawback that  $N$  nonlinear estimation problems must be solved and that the standard errors of statistics such as  $\hat{\mu}_1$  decrease only as  $N^{-1/2}$ . Therefore, although arbitrary accuracy can be obtained via Monte Carlo sampling, arbitrarily many estimation problems may have to be solved to obtain that accuracy.

## 2. IMPROVED MONTE CARLO USING CONTROL VARIATES

The efficiency of Monte Carlo sampling can be improved by variance reduction techniques (Hammersley and Handscomb, 1964; Law and Kelton, 1982). These methods can take several forms, among which the best known are: antithetic variates, stratification, conditional expectations, and control variates. Control variates, using the linear approximator  $\hat{\delta}$  (equation (3)) as the control variate is a very natural approach. This method is very simple to

implement, and it leaves the sampling method unchanged so that the sample estimates  $\hat{\theta}_v$  are all independent. The distribution of  $\hat{\delta}$  is known and  $\hat{\delta}$  is easier to compute than  $\hat{\theta}$  so that very little additional work is required to implement the method. In particular, the control  $\hat{\delta}$  can be used as the basis for controls for all the marginal moments and covariances, as well as the quantiles. It is more efficient than direct Monte Carlo at all the levels of nonlinearity commonly encountered in practice.

The details of the control variate method are given in Swain (1982) and Swain and Schmesier (1983). The control estimators for the marginal moments  $\underline{\mu}_k$ , are given by

$$\hat{\underline{\mu}}_k(B_k) = \hat{\underline{\mu}}_k - B_k(\hat{\delta}_k - E\hat{\delta}_k) \quad (6)$$

where

$$\hat{\delta}_k = N^{-1} \sum_{v=1}^N (\delta_v)^k$$

is the sample  $k$ th moment for  $\hat{\delta}$  and the  $p$  by  $p$  matrix  $B_k$  is the control weight matrix. The method depends upon the covariation between  $\hat{\theta}_v$  and  $\delta_v$  (they both depend upon the same vector of errors,  $\epsilon_v$ ), so that when  $\hat{\delta}_k$  is greater than its expectation, it can be assumed that  $\hat{\underline{\mu}}_k$  exceeds its expectation as well, and a correction proportional to  $-(\hat{\delta}_k - E\hat{\delta}_k)$  should be made to  $\hat{\underline{\mu}}_k$ . It can be shown that the variance of  $\hat{\underline{\mu}}_k(B_k)$  is minimized for the choice

$$B_k^* = \text{Cov}(\hat{\theta}^k, (\delta)^k) \text{Var}^{-1}((\delta)^k),$$

and that the decrease in the control variate variance is largest when the correlation between  $\hat{\delta}$  and  $\hat{\theta}$  is greatest. The covariance term in equation (7) for  $B_k^*$  is generally unknown, so  $B_k^*$  must be specified in some other way. Experience with the method has shown that the choice  $B_k = I$  ( $I$  is the  $p$  by  $p$  identity matrix) is nearly optimal and leads to a simple, unbiased estimator for  $\underline{\mu}_k$ . This is also consistent with the asymptotic case, since  $B_k^*$  (Eqn(7)) tends to the

identity matrix as the sample size increases indefinitely.

A similar control variate scheme can also be given for  $\Sigma_\theta$ . Let  $\underline{S}_\theta$  be a vector of length  $m=p(p+1)/2$  containing the lower diagonal elements of  $\hat{\Sigma}_\theta$  stored by row (e.g.,  $\hat{\sigma}_{11}, \hat{\sigma}_{12}, \dots, \hat{\sigma}_{pp}$ ). Let  $\underline{S}_\delta$  be a similar  $m$ -vector for  $\Sigma_\delta$ . Then

$$\underline{S}_\theta(C) = \underline{S}_\theta - C(\underline{S}_\delta - E(\underline{S}_\delta)) \quad (8)$$

is a control variate estimator for  $\Sigma_\theta$  and again by experience the choice of the  $m$  by  $m$  identity matrix for  $C$  is nearly optimal.

### 3. RESULTS

The efficiency of two Monte Carlo procedures can be computed as the ratio of their precisions, with precision measured as the inverse of the variance. For the multiple parameter case the scalar measure adopted is the determinant of the variance matrices, which is also known as the generalized variance. Then the efficiency of a control variate estimator for the  $k$ th moment compared to the crude estimator of the  $k$ th moment is given by

$$E_k = (1 / |\text{Var} \hat{\underline{\mu}}_k(I)|) / (1 / |\text{Var} \hat{\underline{\mu}}_k(0)|) \\ = |\text{Var} \hat{\underline{\mu}}_k(0)| / |\text{Var} \hat{\underline{\mu}}_k(I)|$$

since direct Monte Carlo corresponds to the use of a control variates with a 0 weighting matrix. Note that efficiencies in excess of 1 indicate that the control variance is less than that of the direct Monte Carlo estimator.

The sample efficiencies ( $k=1$ ) are plotted in figure 1 versus the nonlinearity measure  $N_\theta$ , showing that the control variate estimator is more efficient than direct Monte Carlo for a wide range of nonlinearities, and the efficiency becomes infinite as the nonlinear estimator approaches the behavior of a linear estimator (i.e., asymptotically). The range of  $N_\theta$  depicted is typical of the values encountered in practice.

The apparent relation between  $E_1$  and  $N_\theta$  is not surprising, since the Beale measure  $N_\theta$  explicitly measures the appropriateness of  $\hat{\delta}$  as

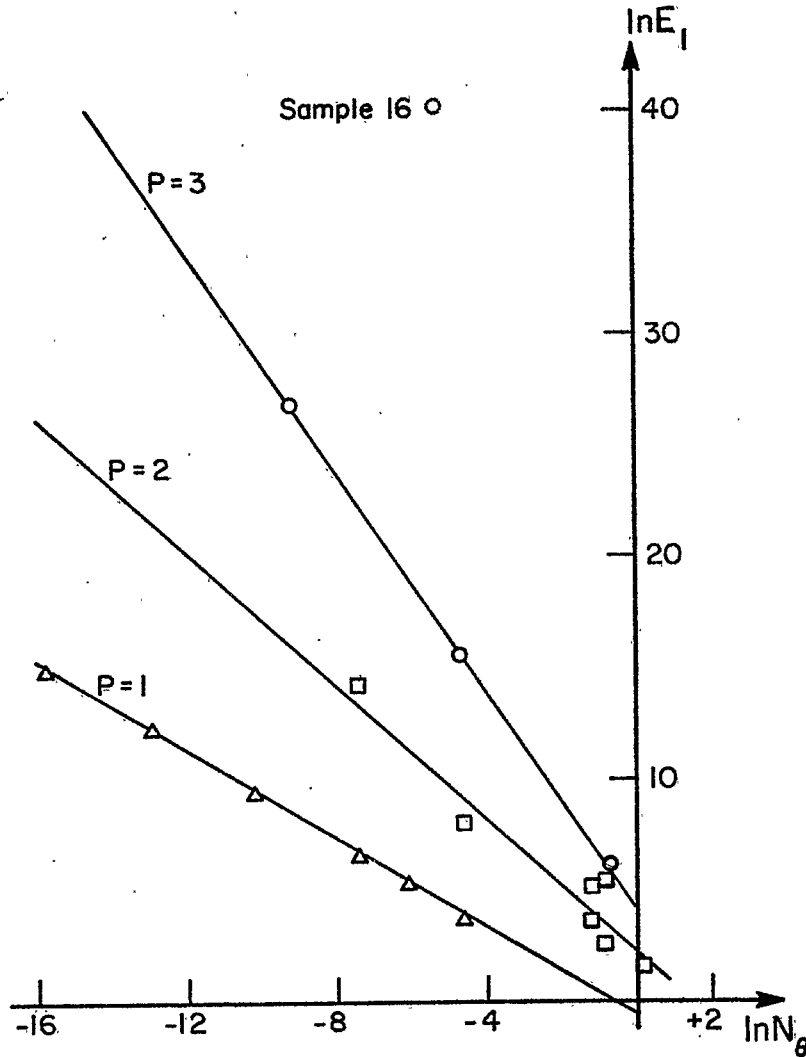


Figure 1. Summary of sampling efficiencies for the control estimator  $\hat{\mu}_1(I)$  for 25 samples taken from Swain (1982). Some points are multiple samples. Sample 16 has one parameter which behaves like a linear estimator and therefore has higher than expected efficiency.

an approximator for  $\hat{\theta}$  and  $\delta$  is the control variate used here. Moreover, an approximate relation between the two (for  $k=1$ ) is given (Swain, 1982) by

$$E_1 = (2N_\theta)^{-P} \quad (9)$$

and the degree of this fit is given in figure 2.

The relation between efficiency and  $N_\theta$ , equation (9), allows a prediction in advance of sampling of how efficient the control variate scheme is likely to be. However, since the

control variable strategy is efficient for all but the most extreme nonlinearities ( $N_\theta$  in excess of  $1/F_{p,n-p;\alpha}$  can be considered extreme, according to Beale) and because the method requires only relatively simple linear computations, it will always be advantageous to use the method. In addition, equation (9) suggests that the efficiency of the method can be further improved by a suitable choice of a parameter transformation. An upper bound on the possible efficiency using control variables can be given by substituting the intrinsic

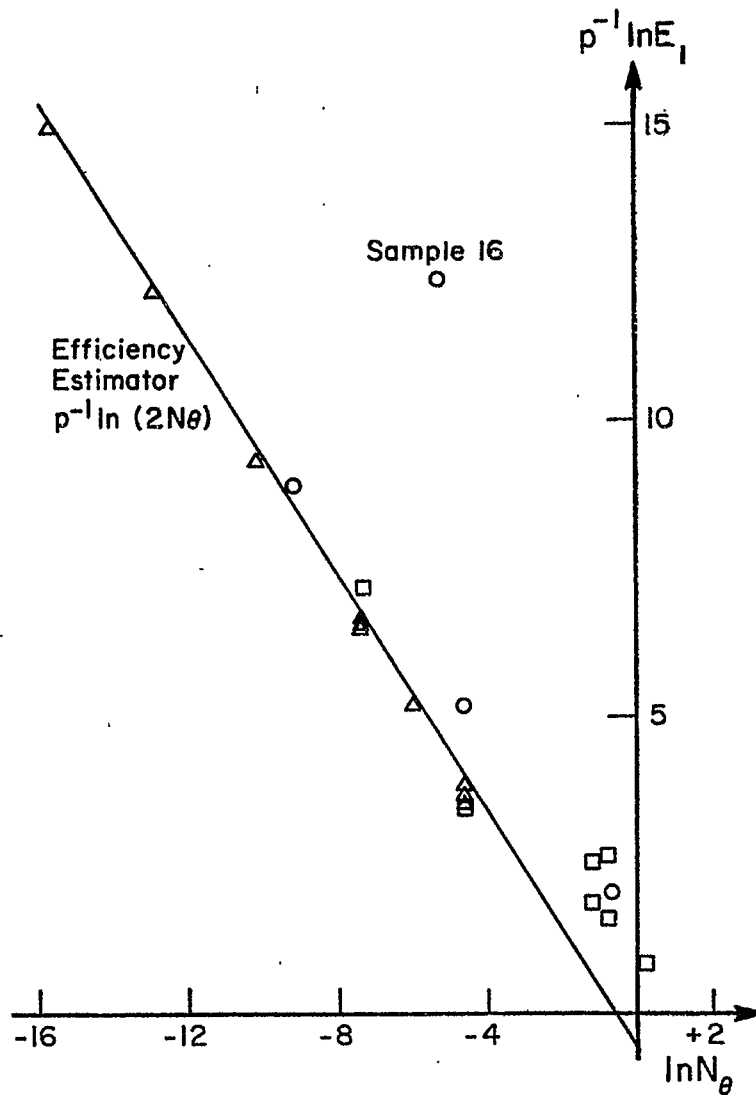


Figure 2. Approximate efficiency formula for the control estimator  $\hat{\mu}_1(I)$ ,  $p^{-1} \ln E_1 = -\ln 2N_\theta$ . The 25 samples are taken from Swain (1982).

nonlinearity  $N_\phi$  in place of  $N_\theta$  in equation (9). For instance, in certain cases a linear approximation based upon the transformed parameters  $\underline{\psi} = (\theta)^k$  leads to additional efficiency in the control variate estimator for the marginal moments (Swain, 1982).

#### REFERENCES

- Bard, Y. (1974), Nonlinear Parameter Estimation, Academic, New York.
- Bates, D.M., Watts, D.G. (1980), Relative Curvature Measures of Nonlinearity, J. Royal Statistical Society, Series B, Vol. 42, No. 1, pp. 1-25.
- Beale, E.M.L. (1960), Confidence Regions in Nonlinear Estimation, J. Royal Statistical Society, Series B, Vol 22, No. 1, pp. 41-88.
- Box, M.J. (1971), Bias in Nonlinear Estimation, J. Royal Statistical Society, Series B, Vol. 33, No. 1, pp. 171-201.

- Clarke, G.P.Y. (1980), Moments of the Least Squares Estimators in a Non-linear Model, J. Royal Statistical Society, Series B, Vol. 42, No. 2, pp. 227-237.
- Gallant, A.R. (1975), Nonlinear Regression, The American Statistician, Vol. 29, No. 2, pp. 73-81.
- Hammersley, J.M., Handscomb, D.C. (1964), Monte Carlo Methods, Chapman and Hall, London.
- Law, A.M., Keltyon, W.D. (1982), Simulation Modeling and Analysis, McGraw-Hill, New York.
- Swain, J.J. (1982), Monte Carlo Estimation of the Sampling Distribution of Nonlinear Parameter Estimators, Unpublished Ph.D. Dissertation, Purdue University, West Lafayette, Indiana.
- Swain, J.J., Schmeiser, B.W. (1983), Monte Carlo Estimation of the Sampling Distribution of Nonlinear Model Parameter Estimators, Technical Report C-83-1, School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, Georgia 17pp. (Also Purdue University Industrial Engineering Technical Memorandum 83-2).