# SIMULATION OPTIMIZATION OF AIRLINE DELAY WITH CONSTRAINTS

David W. Hutchison

Department of Mathematical Sciences
The Johns Hopkins University
3400 N. Charles Street
Baltimore, MD 21218, U.S.A.

Stacy D. Hill

The Johns Hopkins University
Applied Physics Laboratory
11100 Johns Hopkins Road,
Laurel, MD 20723, U.S.A.

## ABSTRACT

Air traffic delay is a growing and expensive problem. We investigated ways to reduce the cost and magnitude of such delays by trading gate delays against more expensive air delays. Air management and planning at this level can be facilitated by simulation, especially for strategies that alter controls on the system. We used the SIMMOD air traffic simulation to model the system. The objective was to determine a set of control measures that achieve the best system performance subject to restrictions on the decision parameters and selected system output measurements. Because observed system performance is "noisy," the problem is a constrained stochastic optimization problem with nonlinear objective function and nonlinear, stochastic constraints, which requires efficient stochastic optimization methods for its solution. Our approach used the simultaneous perturbation stochastic approximation (SPSA) algorithm with a penalty function to handle the difficult constraints.

## 1 INTRODUCTION

Increased air travel has added growing numbers of travelers and flights to an already congested system. The result is an almost inevitable rise in air traffic delay.

The costs and causes of air traffic delay have been documented in previous studies (see, e.g., Odoni 1987), and so have strategies to reduce controllable delays (e.g., Gilbo 1993, 1997b). Such strategies fall into two broad categories—direct and indirect. Direct strategies are those that reduce delay by eliminating its causes, for example reducing congestion at the destination airport by using larger aircraft to reduce the number of flights. Indirect strategies reduce costs by distributing the delay to other parts of the system where it is less expensive, such as holding aircraft at the departure gate to avoid holding on the ground or in the air.

Many of the most effective direct measures are costly and require time to implement and, therefore, offer no immediate resolution. More attractive are incremental strategies that work within the existing system to reduce the effects of congestion or the cost of delays, even though these measures may produce more modest results.

Prominent among indirect strategies are control measures such as gate holding policies, metering aircraft through a control point, and vectoring. These strategies reflect control decisions that are typically made at the individual flight level, often just hours in advance of execution, and based on projected traffic flows (Gilbo 1997a).

Deciding on the control measures to apply to a set of flights for any given day is a difficult nonlinear optimization problem. Problems of this type lend themselves to simulation optimization methods to determine the values of system parameters that yield optimal performance (L'Ecuyer, Giroux, and Glynn 1994). Earlier studies showed that reductions in the cost of delay can be obtained by using a simulation optimization procedure to process delay cost measurements (Kleinman, Hill, and Ilenda 1998). This work was conducted in an unconstrained setting. We extend these results by considering the constrained case.

We present a formulation of the aircraft delay problem and explore the effectiveness of using gate holding delays to reduce costly delays in the air. We accomplish this by constructing a scaled-down simulation model of air traffic control, focusing on the flights between four airports. We add reasonable constraints, and seek an optimal solution. Even though the simulation model is a small-scale system, the results highlight features of a larger-scale problem.

The innovative aspect of this work is the manner in which the constraints are managed. As the constraints of our problem are not simple, we are led to the use of penalty functions. The second goal of this paper is to examine a set of penalty functions for effectiveness (measured by success in finding an acceptable solution) and efficiency (measured by how quickly such solutions are found).

This paper is organized into five sections. In the following section we formulate the problem as a simulation optimization. We use a constrained version of the simultaneous perturbation stochastic approximation method (SPSA) to solve for an optimal set of parameters for some performance measure or loss function of interest. The spe-

cifics of SPSA are developed in section 3. In particular, we discuss the adaptations of the SPSA algorithm to handle complex constraints. Section 4 presents some of the main results of our analysis. We end with conclusions and some observations on directions for future research.

## 2 PRELIMINARIES

The air control system is of such complexity that simulation is often the best method to study its performance. In our case we used the SIMMOD simulation program (ATAC Corporation 1995) to model the flow of 336 flights (departures and arrivals) on a network of four airports. Due to the amount of traffic there is considerable delay in the system. We want to determine the effectiveness of a set of control measures in reducing delay and its costs. SIMMOD is a very flexible discrete event simulation, and there was an array of control measures available. We chose gate hold policies because they are easily implemented and have been studied previously. Gate holds occur when a flight is delayed in departing the gate. The decision parameters in our formulation, therefore, are actual (versus scheduled) aircraft departure times (or, more correctly, the delay between the scheduled time and actual departure). We formulate the problem as a continuous decision parameter optimization problem.

### 2.1 Problem Formulation

Suppose that $\Theta \subset R^p$ and that $\theta \in \Theta$ is a vector with components representing system parameters under control. In our case, for example, the components of $\theta$ are the departure times for each flight. Let $L(\theta)$ be the performance measure or loss function of interest (for example, the expected cost of total flight delays during a specified period of operation). The exact values of loss are unavailable and are estimated by simulation. Our objective is to optimize system performance, i.e. find

$$\min_{\theta \in \Theta} L(\theta). \tag{1}$$

subject to relevant constraints on $\theta$, using only the output measurements of the simulation.

The question of relevant constraints is pertinent. Unconstrained problems are scarce in practice. Departure times are constrained in that flights may be delayed, but for obvious reasons cannot depart early. The restriction on early departures is a *hard* constraint, in that it represents a physical limitation of the system, and parameter values outside these constraints are invalid. Moreover, flights should not be delayed too long, leading to *soft* constraints as upper bounds. We allow parameters to take on values violating soft constraints during optimization, though the final solution must satisfy them. In practice, soft con-

straints are often treated as hard constraints for convenience. In addition, policy or performance *goals* can be modeled as constraints. Examples include limiting airway capacity, implementation of metering strategies, or targeted savings in costs or delays. These goal-constraints are often based on measures of simulation outcomes. The constraints may be highly nonlinear, and generally are implicitly defined. In our problem we consider as one such constraint the goal to effect a 20% reduction in air delay measured as a system output.

Constraints such as the latter type are best handled with a penalty function, and it is our intent to examine a range of penalty functions to investigate comparative effectiveness and efficiency.

### 2.2 Stochastic Approximation

The decision variables are continuous and the solution space may be assumed closed and convex, so the problem lends itself to solution with a gradient-based optimization method. We consider the standard minimization problem, which, it is assumed, is equivalent to finding the root θ* of the equation

$$g(\theta) = \frac{\partial L(\theta)}{\partial \theta} = 0 .$$

The form of $L(\theta)$ and $g(\theta)$ are unknown, and only measurements of $L(\theta)$ are available. In our case SIMMOD is a stochastic simulation, so these measurements are noisy (and the actual process is stochastic).

Our approach to this problem is to use stochastic approximation (SA). SA is a class of algorithms used to minimize (maximize) a function when there is randomness in the optimization process (Andradottir 1998). First introduced by Robbins and Monro (1951) and Kiefer and Wolfowitz (1952), the method has been the subject of considerable research, expanding its applicability and power greatly (see, e.g., Fu 1994, L'Ecuyer, Giroux, and Glynn 1994, Shapiro 1996).

The updating algorithm for stochastic approximation has the form

$$\hat{\theta}_{k+1} = \pi_\Theta \left( \hat{\theta}_k - a_k \hat{g}_k \left( \hat{\theta}_k \right) \right) \tag{2}$$

where $\hat{g}_k(\theta)$ is an estimate of the gradient $g(\theta)$ at iteration $k$ and $\pi_\Theta$ is the projection operator that maps points in $R^p$ to their nearest neighbors in $\Theta$. If the problem is unconstrained, $\Theta = R^p$. When the constrains are known (i.e., not random) and linear, the projection has a simple form. The step-size sequence $a_k$ is nonnegative, decreasing and converges to zero. The generic iterative form of (2) is analogous to the familiar steepest descent algorithm for deterministic problems. The estimate $\hat{\theta}_k$ converges to the

optimizer $\theta^*$, under suitable conditions on the loss function and its gradient (see, e.g., Kushner and Yin 1997).

Difficulty arises when the constraints are nonlinear, stochastic, or both, in which case $\pi_\Theta$ may be complex. Under these conditions, a penalty function may be used to transform the problem into an equivalent unconstrained problem (or one with much simpler constraints) (Bazaraa, Sherali, and Shetty 1993, Fiacco and McCormick 1990).

Since only noisy measurements of the loss function are available, we use stochastic approximation in the Kiefer-Wolfowitz setting. Thus, the gradient estimate is obtained from measurements of $L(\theta)$ (rather than from measurements of $g(\theta)$, as in the Robbins-Monro setting). A common—and computationally inefficient—method of estimating the gradient from loss function measurements is the method of finite-differences. One-sided finite differences, for example, requires $p+1$ function evaluations (symmetric finite-differences requires $2p$ function evaluations). If $p$ is large and the function evaluations difficult or time-consuming, the computational effort could be substantial, since the estimate must be computed at each iteration in (2). To address this difficulty, Spall (1987, 1992) developed the simultaneous perturbation algorithm, which requires only two function evaluations to estimate the gradient at each iteration, regardless of the dimension of $\theta$. The major advantage of SPSA is the reduction in computations required to achieve an optimal solution by reducing the number of required simulation experiments.

# 3 SIMULTANEOUS PERTURBATION

The theoretical basis for SPSA was developed by Spall (1987, 1992) and expanded in subsequent work (see Fu and Hill 1997, Sadegh and Spall 1998, Spall 2000). The method relies on a computationally efficient estimate of the gradient of $L(\theta)$.

## 3.1 Gradient Estimate

Let $\hat{g}_k(\theta)$ denote the simultaneous perturbation estimate of $g(\theta)$ at iteration $k$ and $\hat{\theta}_k$ the estimate of $\theta^*$. Let $\Delta_k \in R^p$ be a random *perturbation* vector at the $k$-th iteration. The components of $\Delta_k$ are usually taken to be independent $\pm 1$ Bernoulli variables. (More generally, the components of $\Delta_k$ are independent and independently distributed random variables that are bounded, symmetrically distributed about 0, and whose inverses have finite absolute first moments.) We take measurements of $L(\theta)$ at the two values:

$$y\left(\hat{\theta}_k \pm c_k \Delta_k\right) = L\left(\hat{\theta}_k \pm c_k \Delta_k\right) + \varepsilon_k^{(\pm)}$$

where $\varepsilon_k^{(+)}, \varepsilon_k^{(-)}$ are measurement error terms.

The simultaneous perturbation estimate $\hat{g}_k(\theta)$ of the gradient has $j$-th component, $1 \le j \le p$,

$$\hat{g}_k^{(j)}\left(\hat{\theta}_k\right) = \frac{y\left(\hat{\theta}_k + c_k \Delta_k\right) - y\left(\hat{\theta}_k + c_k \Delta_k\right)}{2 c_k \Delta_{kj}}. \tag{5}$$

Note that, in contrast to finite-difference gradient estimates, only the denominator varies in (5). Under suitable conditions (see Spall 1992 for details), the gradient estimate satisfies

$$E[\hat{g}_k(\hat{\theta}_k) \mid \hat{\theta}_k] = g(\hat{\theta}_k) + O(c_k^2).$$

Furthermore, $\hat{\theta}_k \rightarrow \theta^*$ almost surely as $k \rightarrow \infty$.

## 3.2 Penalty Functions

An approach to constrained stochastic optimization is discussed in Wang and Spall (1999), Sadegh (1997). They showed SPSA converges almost surely for explicit constraints using projection. More difficult constraints can be handled by penalty function, transforming the problem into one that is unconstrained or mildly constrained. Pflug (1981) showed that stochastic approximation using penalty functions converges almost surely, and Wang and Spall (1999) have extended this result to SPSA.

We handle the constraints for this problem in several ways. The non-negativity constraint is managed by recognizing that departure times are translation invariant, that is, scaling the departure times by an additive constant does not affect the outcomes. Consequently we remove the non-negativity constraints, solve the problem, and then translate the departure times to ensure no flight departs earlier than its originally scheduled time.

Removing the non-negativity constraints demanded we address the soft upper bounds on gate hold delays somewhat differently. We treated these as hard constraints, but with respect to the minimum and maximum hold times, $\left|\max(\theta_k) - \min(\theta_k)\right| \le D$. We handled these inequalities with projection.

We placed goal-defining constraints (e.g., achieve a 20% reduction in air delay time) in a penalty function. Following the method of Wang and Spall (1999), we use a class of penalty functions defined by

$$P(x) = \frac{1}{\beta} \max(x, 0)^\beta$$

where $1 \le \beta \le 2$. At $\beta = 1$, this is the absolute value function and quadratic at $\beta = 2$. We investigated a collection of

penalty functions by varying beta in the range $1 \le \beta \le 2$. Previous work on problems with better-behaved constraints showed that penalty functions with $\beta > 1$ quickly diverged, these penalty functions being sensitive to noise. The case for $\beta = 2$ is particularly unstable while $\beta = 1$ is stable and always converges (Wang and Spall 1999). We want to investigate values for $\beta > 1$, below which the method is generally stable.

According to standard procedure, we modify the constrained optimization in (1) by adding a penalty term to obtain the more mildly constrained problem below

$$\min_{\theta \in \Theta} L(\theta) + r_k P(\theta).$$

The penalty weights $r_k$ are nonnegative and $r_k \to \infty$. The SPSA algorithm takes on the form

$$\hat{\theta}_{k+1} = \pi_\Theta \left( \hat{\theta}_k - a_k \hat{g}_k(\hat{\theta}_k) - a_k r_k \hat{h}_k(\hat{\theta}_k) \right).$$

Here $\hat{h}_k(\hat{\theta}_k)$ is an estimate for the gradient of the penalty function. When the penalty function is analytic, this term may be replaced by $\nabla P(\hat{\theta}_k)$.

In our formulation we used a loss function based on the weighted delay throughout the system, with separate weights for delay at the gate, on the ground, and in the air (see Kleinman, Hill, and Ilenda 1998) and with the constraints as described.

## 4 MAIN RESULT

We obtained a solution that satisfied the constraints and resulted in a reduction of the air delay for the described problem. A table of critical values is given below. In each case we were able to satisfy the penalty constraint while meeting the hard constraints on gate hold times. In the revised problem the algorithm found a solution that resulted in greater gate hold times, and the bounding constraints on some of the gate holds were tight.

The decrease in the loss function was slight for all tested penalty functions and not statistically significant. We discuss the reasons for this observation. Unlike other constrained optimizations, penalty optimizations are exterior point methods (Fiacco and McCormick 1990). Consequently, the initial point, $\theta_0$, lies outside the feasible region. The penalty term causes the loss function to increase sharply until sufficient reductions in air delay occur to offset (or eliminate) this term. At that point $\theta_k$ is inside or very near the feasible region, and we have the problem in (1) with the penalty constraint satisfied, usually tightly. This occurred at about iteration 75 for our problem. The remaining iterations attempt to minimize $L(\theta)$. Our interpretation is that 300 iterations were not sufficient to gain

convergence. It is an interesting research problem to determine whether the algorithm would profit from redefining the step size $a_k$ when this switch occurs, and what this new sequence should be.

The results reported in Table 1, Figure 1, and Figure 2 are averaged over 30 Monte Carlo trials. The results of Figure 3 are from a more modest run of 15 Monte Carlo trials.

Table 1: Results from Original Problem (300 Iterations, 30 Monte Carlo Trials)

| Case | Loss Function | Total Air Delay | Total Gate Delay | Maximum Gate Hold |
|---|---|---|---|---|
| Initial | 50.81 | 4520 | 0 | 0 |
| Final ($\beta$=1) | 49.93 | 3880 | 1718 | 16.3 |
| Final ($\beta$=2) | 49.88 | 3851 | 1862 | 17.2 |

It is apparent from Figure 1 that even with 30 Monte Carlo trials, the averaged loss function is still quite noisy. Empirical testing showed that for this simulation model, Monte Carlo trial counts greater than 100 are required before sufficient smoothing is obtained. This may be caused by the constraints and the compact nature of the problem. In many cases it will be prohibitive to obtain such a high number of trials.
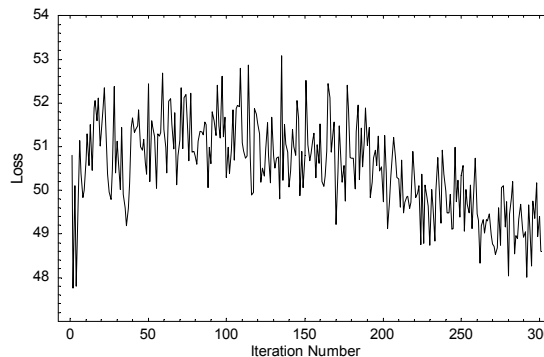


Figure 1: Trace of the Averaged Loss Function by Iteration for $\beta = 1$.

The results of a quadratic penalty function are shown in Figure 2. The trace of the loss function is similar to that in Figure 1. Statistically, the final result for $\beta = 1$ is not significantly different from the result for $\beta = 2$, meaning that no advantages were observed for one penalty function over another, though we expect that increasing the number of Monte Carlo trials and the number of iterations will allow stronger conclusions to be drawn.

Figure 3 shows the relative performance of the various penalty functions on air delay. No attempt has been made to distinguish these on the graph as there are no significant differences between them. The heavy bar represents the penalty constraint (20% reduction in total system air delay time). It is apparent that this goal is satisfied in about 75
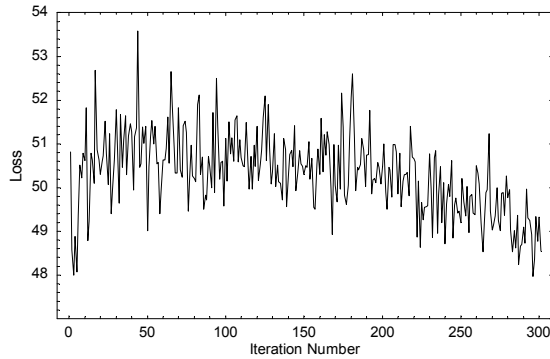
Figure 2: Trace of the Averaged Loss Function by Iteration for β = 2
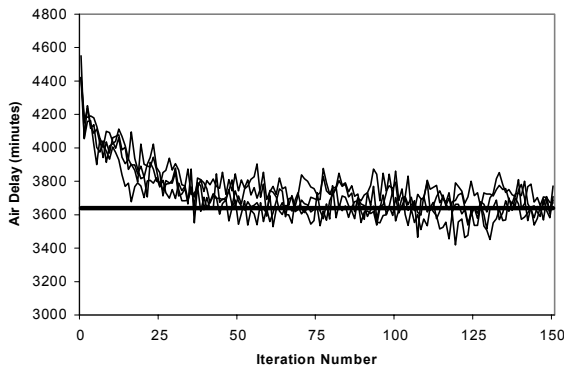


Figure 3: Effect of Penalty Function on Air Delay

iterations of the SPSA algorithm for all penalty functions studied. As mentioned, the results are averaged over 15 Monte Carlo trials, and even when this smoothing is applied, there is sufficient remnant noise that we see no discernable or statistical differences between the penalty functions over the set of betas chosen.

There are several reasons why this could be true. First, the problem may be too noisy to allow such distinctions. Penalty terms are notoriously sensitive to noise (Fiacco and McCormick 1990) and in this problem the noise was fairly high with loss function values of about 50 and a standard deviation of 2.78 for the error. If differences do, indeed, exist, they may be discernable with a higher number of Monte Carlo trials. This suggests an area for future work.

## 5 CONCLUSIONS

Our study suggests the viability of SPSA for solving high dimensional constrained optimization problems efficiently. For the problem at hand, we were able to reduce air delays to a stated goal, while minimizing costs or delays elsewhere in the system subject to simple constraints. It is perhaps significant that we obtained a solution at all. It is not possible in the general case to guarantee *a priori* that the solution space is not empty.

We looked at minimizing the total cost of delay using a weighted loss function comprising gate delays, ground delays, and air delays. We were able to demonstrate that penalty functions are a useful way to perform simulation optimization in the presence of constraints that are nonlinear, implicit, and noisy. Unfortunately we were not able to discern differences in the relative efficiency of the various penalty functions studied. The lack of a discernable difference may be attributed to the extremely noisy nature of the simulation.

These results suggest several interesting avenues for further research. A more systematic examination of these and other classes of penalty functions would be useful to identify conditions that argue for one penalty function over another.

## REFERENCES

Andradottir, S. 1998. Simulation optimization, in *Handbook of Simulation*, 307-333. J. Banks, ed. Wiley.

ATAC Corporation. 1995. *SIMMOD Reference Manual*, ATAC Corporation, Sunnyvale, CA.

Bazaraa, M S., H. D. Sherali, and C. M. Shetty. 1993. *Nonlinear Programming*, John Wiley and Sons, Inc.

Fiacco, A. V. and G. P. McCormick. 1990. *Nonlinear Programming*, Society for Industrial and Applied Mathematics, Philadelphia.

Fu, M. C. 1994. A tutorial review of techniques for simulation optimization, *Proceedings of the 1994 Winter Simulation Conference*, 149-156.

Fu, M. C. and S. D. Hill. 1997. Optimization of discrete event systems via simultaneous perturbation stochastic approximation, *IIE Transactions*, 29:233-243.

Gilbo, E. P. 1993. Airport capacity: representation, estimation, optimization, *IEEE Transactions on Control Systems Technology*, 1:144-154.

Gilbo, E. P. 1997a. Optimization of air traffic management strategies at airports with uncertainty in airport capacity, *Proceedings of the 8th IFAC/IFIP/IFORS Symposium*, Chania, Greece, 35-40.

Gilbo, E. P. 1997b. Optimizing airport capacity utilization in air traffic flow management subject to constraints at arrival and departure fixes, *IEEE Transactions on Control Systems Technology*, 490-503.

Kiefer, J. and J. Wolfowitz. 1952. Stochastic Estimation of a Regression Function. *Annals of Mathematical Statistics*, 43:462- 466.

Kleinman, N. L., S. D. Hill, and V. A. Ilenda. 1998. Simulation optimization of air traffic delay cost, *Proceed-*

*ings of the 1998 Winter Simulation Conference*, Washington, DC, 1177-1181.

Kushner, H. J. and G. G. Yin. 1997. *Stochastic Approximation Algorithms and Applications*. Springer-Verlag: New York.

L'Ecuyer, P., N. Giroux, and P. W. Glynn. 1994. Stochastic optimization by simulation: numerical experiments with the M/M/1 queue in steady-state, *Management Science*, 40:1245-1261.

Odoni, A. R. 1987. The flow management problem in air traffic control, in *Flow Control of Congested Networks*, A. R. Odoni, ed., Springer-Verlag, 270-288.

Pflug, G. C. 1981. On the convergence of a penalty-type stochastic optimization procedure, *Journal of Information and Optimization Sciences*, 2:249-258.

Robbins, H. and S. Monro. 1951. A Stochastic Approximation Method, *Annals of Mathematical Statistics*, 22:400-407.

Sadegh, P. 1997. Constrained optimization via stochastic approximation with a simultaneous perturbation gradient approximation, *Automatica*, 33:889-892.

Sadegh, P., and J. C. Spall. 1998. Optimal random perturbations for stochastic approximation using a simultaneous perturbation gradient approximation, *IEEE Transactions on Automatic Control*, 43:1480-1484.

Shapiro, A. 1996. Simulation-based optimization – convergence analysis and statistical inference, *Communications in Statistics – Stochastic Models*, 12:425-453.

Spall, J. C. 1987. A stochastic approximation technique for generating maximum likelihood parameter estimates. *Proceedings of the American Control Conference*, Minneapolis, MN, 1161-1167.

Spall, J. C. 1992. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation, *IEEE Transactions on Automatic Control*, 37:332-341.

Spall, J. C. 1998. Implementation of the simultaneous perturbation algorithm for stochastic optimization, *IEEE Transactions on Aerospace and Electronic Systems*, 34:817-823.

Spall, J. C. 2000. Adaptive stochastic approximation by the simultaneous perturbation method, *IEEE Transactions on Automatic Control*, 45:1839-1853.

Wang, I.-J., and J. C. Spall. 1999. A constrained simultaneous perturbation stochastic approximation algorithm based on penalty functions, *Proceedings of the American Control Conference*, San Diego, CA, 393-399.

## AUTHOR BIRGRAPHIES

**DAVID W. HUTCHISON,** Lieutenant Colonel (U.S. Army Retired), is a graduate of the U.S. Military Academy. He received an M.S. degree in applied mathematics from the Massachusetts Institute of Technology in 1983. He has been a consultant to the RAND Corporation and Science Applications International Corporation. He is currently a Ph.D. candidate in the Department of Mathematical Sciences, Johns Hopkins University. His research interests are in stochastic optimization and modeling. His email address is <hutch@brutus.mts.jhu.edu>.

**STACY D. HILL** received his B.S. and M.S. degrees from Howard University in 1975 and 1977, respectively, and the D.Sc. degree in control systems engineering and applied mathematics from Washington University in 1983. He is currently on the Senior Professional Staff of The Johns Hoepkins University Applied Physics Laboratory where he has been project and technical lead in developing, testing, and applying statistical techniques and software, and has led systems analysis and modeling projects. He has published papers on diverse topics in statistics and engineering, including subjects such as simulation, optimization, and parameter estimation. His email address is <stacy.hill@jhuapl.edu>.