

# Reducing Parameter Uncertainty for Stochastic Systems\*

Szu Hui NG<sup>†</sup>

Stephen E. CHICK<sup>‡</sup>

November 21, 2005

## Abstract

The design of many production and service systems is informed by stochastic model analysis. But the parameters of statistical distributions of stochastic models are rarely known with certainty, and are often estimated from field data. Even if the mean system performance is a known function of the model's parameters, there may still be uncertainty about the mean performance because the parameters are not known precisely. Several methods have been proposed to quantify this uncertainty, but data sampling plans have not yet been provided to reduce parameter uncertainty in a way that effectively reduces uncertainty about mean performance. The optimal solution is challenging, so we use asymptotic approximations to obtain closed-form results for sampling plans. The results apply to a wide class of stochastic models, including situations where the mean performance is unknown but estimated with simulation. Analytical and empirical results for the  $M/M/1$  queue, a quadratic response-surface model, and a simulated critical care facility illustrate the ideas.

**Category:** G.3: Probability and Statistics, Probabilistic algorithms (including Monte Carlo) experimental design

**Category:** I.6: Simulation and Modeling, Simulation output analysis

**Terms:** Experimentation, Performance

**Keywords:** Stochastic simulation, uncertainty analysis, parameter estimation, Bayesian statistics

## 1 Introduction

Stochastic modeling and simulation analysis are useful approaches to evaluate the performance of production and service systems as a function of design and statistical parameters [e.g., Buzacott and Shanthikumar 1993; Law and Kelton 2000]. Models of existing systems, or variations of existing systems, can help inform system design and improvement decisions. Data may be available to help estimate statistical parameters, but estimators are subject to random variation because they are functions of random phenomena. Parameter uncertainty means that there is a risk that a planned system will not perform as expected [Cheng and Holland 1997; Chick 1997, 2001; Barton and Schruben 2001], even if the relationship between the parameters and system performance is

---

\*©ACM, (2006). This is the author's version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The definitive version was published in ACM TOMACS {Vol 16, Iss 1, Jan. 2006}. <http://www.linklings.net/tomacs/index.html>

<sup>†</sup>S.-H. Ng, Dept. Industrial and Systems Engineering, National University of Singapore, 10 Kent Ridge Crescent, Singapore 119260, [isensh@nus.edu.sg](mailto:isensh@nus.edu.sg).

<sup>‡</sup>S. E. Chick, Technology and Operations Management Area, INSEAD, Boulevard de Constance, 77305 Fontainebleau CEDEX, France, [stephen.chick@insead.edu](mailto:stephen.chick@insead.edu).

completely known. For example, the stationary mean occupancy of a stable  $M/M/1$  queue is a known function of the arrival and service rates, but the mean occupancy is unknown if those rates are not precisely known. When the system performance is an unknown function of the parameters, and is estimated by simulation, the problem is aggravated. The common practice of inputting point estimates of parameters into simulations can dramatically reduce the coverage of a supposedly  $100(1 - \alpha)\%$  confidence interval (CI) for the mean performance (Barton and Schruben 2001).

Bootstrap methods (Cheng and Holland 1997; Barton and Schruben 2001), asymptotic normality approximations (Cheng and Holland 1997; Ng and Chick 2001), and Bayesian model averaging [Draper 1995; Chick 2001; Zouaoui and Wilson 2003, 2004] can quantify the effect of input parameter uncertainty on output uncertainty.

This paper goes beyond quantifying uncertainty by suggesting how to reduce parameter uncertainty in a way that effectively reduces uncertainty about the mean performance of a system. We do so for a broad class of stochastic systems by using asymptotic variance approximations for the unknown mean system performance. Section 2 formalizes the stochastic model assumptions along with examples. Section 3 introduces the asymptotic variance approximations for input parameter uncertainty. Section 4 describes how input parameter uncertainty affects output uncertainty and formulates an optimization problem to reduce output uncertainty. When the system response is a known function of the inputs, the optimal sampling plan turns out to be very similar to known results for stratified sampling. When the system response must be estimated with simulation, a different sampling allocation is obtained. The allocation shows how to balance the cost of running additional simulations (to reduce uncertainty about the mean response and its gradient) and data collection (to reduce uncertainty about the inputs to the system). Under some general conditions, the optimal sampling allocations are shown to have the attractive property of being invariant to the coordinate systems used to represent the parameter of each input distribution. Section 5 applies the results to several examples. This paper extends preliminary work (Ng and Chick 2001) by generalizing results to a broader class of distributions, and by providing results for unknown response functions. An appendix also describes how to reduce the mean-squared error (MSE) rather than the variance. The approach taken is Bayesian, so parameter uncertainty is quantified with random variables.

## 2 Model Formulation and Examples

Consider a system with multiple statistical parameters  $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$  that influence the system performance  $g(\Theta)$ . Simulated performance realizations  $Y_l$  (“output”) with input parameter  $\Theta_l$  for replication  $l = 1, 2, \dots$ ,

$$Y_l = g(\Theta_l) + \sigma(\Theta_l)Z_l,$$

may be observed if  $g(\Theta)$  is unknown, where  $\sigma^2$  is the output variance, and the  $Z_l$  are independent zero-mean random variables with unit variance.

The system is driven by  $k$  scalar (real-valued) input processes that are mutually independent sources of randomness. Given the parameter vector  $\theta_i$  for the  $i$ th input process, the corresponding observations  $\{x_{i\ell} : \ell = 1, 2, \dots\}$  are independent and identically distributed with conditional probability density function  $p_i(x_{i\ell} | \theta_i)$ . Uncertainty about the parameter  $\theta_i$  is initially described with a probability model  $\pi_i(\theta_i)$ . In Bayesian statistics,  $\pi_i(\theta_i)$  is called a prior probability distribution. Here we presume joint independence across sources of randomness,  $\pi(\Theta) = \prod_{i=1}^k \pi_i(\theta_i)$ .

Each  $\theta_i$  is assumed to be inferrable from data. The data  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in_i})$  that has been observed so far for the  $i$ th input process can be used to get a more precise idea about the value of

the unknown parameter  $\theta_i$ , whose uncertainty is characterized by the posterior distribution

$$f_{in_i}(\theta_i | \mathbf{x}_i) \propto \pi_i(\theta_i) \prod_{\ell=1}^{n_i} p_i(x_{i\ell} | \theta_i), \quad (1)$$

for  $i = 1, 2, \dots, k$ . We will use those distributions to figure out how to collect additional data, if needed, to reduce uncertainty about input parameters. Let  $\mathbf{n} = (n_1, \dots, n_k)$  specify the number of data points observed so far for each source of randomness. Probability statements below are all conditional on all data collected so far,  $\mathcal{E}_{\mathbf{n}} = (\mathbf{x}_1^{\top}, \dots, \mathbf{x}_k^{\top})$ , unless otherwise specified.

A tilde denotes the maximum *a posteriori* (MAP) estimate of any parameter, e.g.  $\tilde{\theta}_{in_i}$  maximizes  $f_{in_i}$ . Maximum likelihood estimates (MLEs) are denoted with a hat,  $\hat{\theta}_{in_i}$ . Given the data  $\mathbf{x}_i$ , the MLE maximizes  $\prod_{\ell=1}^{n_i} p_i(x_{i\ell} | \theta_i)$ .

Parameters may be multivariate, such as  $\theta_1 = (\vartheta_1, \vartheta_2)$ . It will be convenient at times to refer to the individual components of the parameter vector  $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$  encompassing all  $k$  sources of randomness; and for this purpose we write  $\Theta = (\vartheta_1, \vartheta_2, \dots, \vartheta_K)$ , where  $K > k$  when at least one of the  $\theta_i$  are multidimensional. Decision variables and system control parameters are considered to be fixed and implicit in the definition of  $g$ , so that we focus upon the effects of input parameter uncertainty. Table 1 summarizes this notation, as well as much of the notation used in the rest of the paper.

The  $M/M/1$  queue can be described with this notation as a system with  $k = 2$  sources of randomness (arrival and service times), with arrival rate  $\theta_1 = \lambda$  and service rate  $\theta_2 = \mu$ . There are a total of  $K = 2$  components of the overall parameter vector  $\Theta$ , so  $\Theta = (\lambda, \mu)$ . The stationary mean occupancy is known in closed form if the queue is stable,  $g(\Theta) = \lambda/(\mu - \lambda)$ .

A less trivial example for which the system response is unknown is the critical care facility depicted in Figure 1. Schruben and Margolin (1978) studied that system to determine the expected number of patients per month that are denied a bed, as a function of the number of beds in the intensive care unit (ICU), coronary care unit (CCU), and intermediate care units. Patients arrive according to a Poisson process and are routed through the system depending upon their specific health condition. Their analysis presumed fixed point estimates (MLE) for the parameters of the  $k = 6$  sources of randomness: the patient arrival process (Poisson arrivals, mean  $\hat{\theta}_{1n} = 3.3/\text{day}$ ), ICU stay duration (lognormal  $\theta_2 = (\vartheta_2, \vartheta_3)$ , with mean  $\hat{\vartheta}_{2n} = 3.4$  and standard deviation  $\hat{\vartheta}_{3n} = 3.5$  days), three more bivariate input parameters  $\theta_3, \theta_4, \theta_5$  for the lognormal service times at the intermediate ICU, intermediate CCU, and CCU processes, and a sixth parameter for patient routing (multinomial,  $\hat{\theta}_{6n} = (\hat{p}_1, \hat{p}_3, \hat{p}_4) = (.2, .2, .05)$ , note that  $p_2 = 1 - p_1 - p_3 - p_4$ ). There are therefore a total of  $K = 12$  components of the overall parameter vector  $\Theta$ . The function  $g(\cdot)$  is not known, and is to be estimated from simulation near the point  $\hat{\Theta}_n$ .

This paper will use approximations that assume a constant variance,  $\sigma(\Theta) = \sigma$  (homoscedastic). If the variance actually depends upon  $\Theta$  (heteroscedastic), then we will estimate  $\sigma$  by the MAP estimator  $\tilde{\sigma}(\hat{\Theta}_n)$ , a Bayesian analog of the usual MLE for the standard deviation,  $\hat{\sigma}(\hat{\Theta}_n)$ . This approximation is a standard assumption in the design of experiments literature (Box, Hunter, and Hunter 1978; Box and Draper 1987; Myers, Khuri, and Carter 1989), and is asymptotically valid under certain conditions (Mendoza 1994, p. 176) if  $\sigma(\Theta)$  is a “nice” function of  $\Theta$  near  $\hat{\Theta}_n$ .

The approximations made below in a Bayesian context are analogous to frequentist approximations by Cheng and Holland (1997). The modification of the analysis is with a straightforward application of results of Bernardo and Smith (1994), and result in a decoupling of stochastic uncertainty from parameter uncertainty when  $g(\cdot)$  is a known function. Zouaoui and Wilson (2001) used a similar decoupling to obtain a variance reduction for estimates of  $E[g(\Theta) | \mathcal{E}_{\mathbf{n}}]$  with the Bayesian Model Average (BMA) when multiple candidate distributions are proposed for a single

Table 1: Summary Of Notation<sup>1</sup>

$k$	number of independent sources of randomness, each consisting of i.i.d. scalar (real-valued) random variables
$\boldsymbol{\theta}_i$	statistical input parameter for $i$ th source of randomness (can be multivariate with dimension $d_i \geq 1$ ), $i = 1, 2, \dots, k$
$K$	total number of components of all $k$ parameter vectors, $K = \sum_{i=1}^k d_i$
$\boldsymbol{\Theta}$	vector of all statistical parameters, $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_k)$
$\vartheta_j$	the $j$ th individual component of $\boldsymbol{\Theta}$ for $j = 1, \dots, K$ , so that we have $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_k) = (\vartheta_1, \vartheta_2, \dots, \vartheta_K)$
$g(\boldsymbol{\Theta})$	the mean response of the model, as a function of the statistical parameters
$r$	number of replications a simulated model
$\sigma$	standard deviation of output of a simulated model
$n_i$	number of observations taken from the $i$ th source of randomness and used to estimate the associated parameter vector $\boldsymbol{\theta}_i$ , for $i = 1, \dots, k$
$\mathbf{n}$	vector of number of observations, $\mathbf{n} = (n_1, n_2, \dots, n_k)$
$x_{i\ell}$	$\ell$ th observation for $i$ th statistical parameter, $i = 1, \dots, k$ ; $\ell = 1, 2, \dots, n_i$
$\mathbf{x}_i$	vector of observations $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in_i})$ for $i$ th statistical parameter
$\mathcal{E}_{\mathbf{n}}$	all data observed so far, $\mathcal{E}_{\mathbf{n}} = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_k^T)$
$\pi_i(\boldsymbol{\theta}_i)$	prior probability density of parameter $\boldsymbol{\theta}_i$ , $i = 1, \dots, k$
$p_i(x   \boldsymbol{\theta}_i)$	density function for data $x$ given parameter $\boldsymbol{\theta}_i$ , $i = 1, \dots, k$
$f_{in_i}(\boldsymbol{\theta}_i)$	posterior probability density of parameter $\boldsymbol{\theta}_i$ given data $\mathbf{x}_i$ ( $n_i$ observations)
$\tilde{\boldsymbol{\theta}}_i$	MAP estimate for $i$ th statistical parameter, $i = 1, \dots, k$
$\hat{\boldsymbol{\theta}}_i$	MLE for $i$ th statistical parameter, $i = 1, \dots, k$
$\mathbf{G}_i = \boldsymbol{\Sigma}_i^{-1}$	observed information matrix associated with the parameter vector $\boldsymbol{\theta}_i$ , given $\mathbf{x}_i$ for $i = 1, \dots, k$ . Also sometimes written $\mathbf{G}_{in_i} = \boldsymbol{\Sigma}_{in_i}^{-1}$ to emphasize the dependence on the sample size $n_i$
$\boldsymbol{\Sigma}$	Block diagonal matrix $\boldsymbol{\Sigma} = \text{diag}[\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k] = \text{diag}[\mathbf{G}_1^{-1}, \dots, \mathbf{G}_k^{-1}]$
$\mathbf{H}_i$	Expected information matrix of one observation for parameter $i$ , $i = 1, \dots, k$
$\beta_j$	Derivative of $g(\boldsymbol{\Theta})$ with respect to $j$ th component of $\boldsymbol{\Theta}$ , evaluated at $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ , for $j = 1, \dots, K$
$\boldsymbol{\beta}_i$	$\nabla_{\boldsymbol{\theta}_i} g(\boldsymbol{\Theta})  _{\boldsymbol{\theta}=\tilde{\boldsymbol{\theta}}}$ , the gradient vector of the function $g(\boldsymbol{\Theta})$ with respect to the parameter vector $\boldsymbol{\theta}_i$ for the $i$ th source of randomness, evaluated at $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$
$\boldsymbol{\beta}$	$\nabla_{\boldsymbol{\Theta}} g(\boldsymbol{\Theta})  _{\boldsymbol{\theta}=\tilde{\boldsymbol{\theta}}}$ , the gradient vector of the function $g(\boldsymbol{\Theta})$ with respect to the parameter vector $\boldsymbol{\Theta}$ , evaluated at $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$
$m_i$	number of additional data points to be collected for $i$ th statistical parameter
$\mathbf{m}$	vector of number of additional points, $\mathbf{m} = (m_1, m_2, \dots, m_k)$
$m_0$	number of additional simulation replications to be run
$\mathcal{D}$	all additional data to be collected ( $m_i$ points for $i$ th source, $i = 0, 1, \dots, k$ )

<sup>1</sup> A subscript  $i$  may be dropped to denote a generic source of randomness to simplify notation. The addition of (a) a  $\hat{\cdot}$  denotes maximum likelihood estimate (MLE), (b) a  $\tilde{\cdot}$  denotes a MAP estimator, (c) an extra subscript  $n$ ,  $n_i$  or  $\mathbf{n}$  emphasizes the number of data points that determine the estimate, (d) an extra subscript  $\boldsymbol{\Theta}$ ,  $\boldsymbol{\psi}$  may be added to emphasize the coordinate system for the parameters, as in  $\boldsymbol{\Sigma}_{\boldsymbol{\Theta}}$ ,  $\boldsymbol{\Sigma}_{\boldsymbol{\psi}}$ .

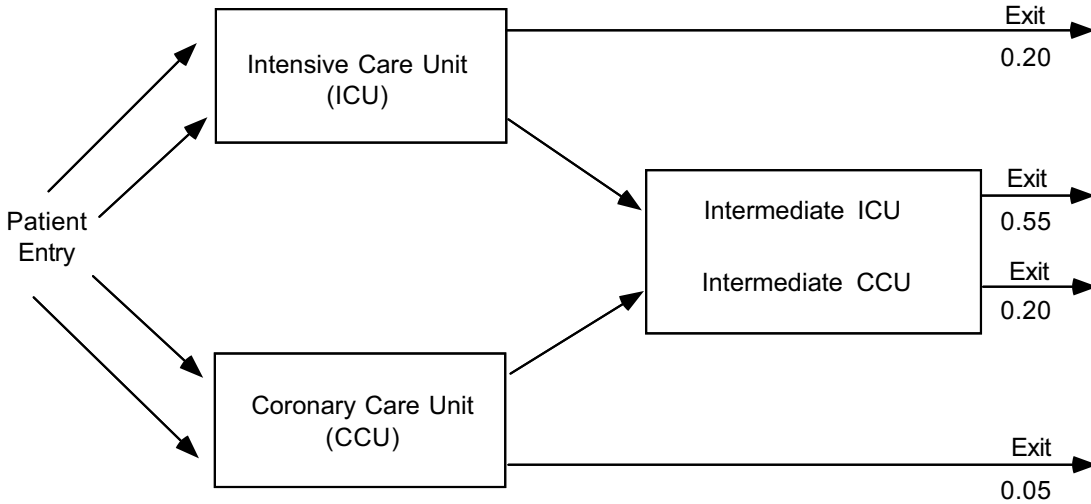


Figure 1: Fraction of patients routed through different units of a critical care facility.

source of randomness, and the response is homoscedastic. Their result appears to be extendable to the heteroscedastic case ( $n \rightarrow \infty$ ) by approximating  $\sigma$  with  $\tilde{\sigma}(\tilde{\Theta}_n)$ . The idea of decoupling stochastic and parameter uncertainty is applied to the case of an unknown  $g(\Theta)$  with an unknown gradient in Section 4.2.

### 3 Approximations for Parameter Uncertainty

Uncertainty about parameters, and the effect of data collection variability, is approximated by the variance. The mean squared error (MSE) of estimates is treated in Appendix A. The analysis is motivated by analogy with the following well-known result for confidence intervals based on asymptotic normality. Suppose that the simulation outputs  $y_1, \dots, y_r$  are observed, with sample mean  $\bar{y}_r$ , sample variance  $S_r^2$ , and  $100(1 - \alpha)\%$  CI for the mean  $\bar{y}_r \pm t_{r-1, 1-\alpha/2}(S_r^2/r)^{1/2}$ . An *approximate* expression for the number of additional observations  $m^*$  required to obtain an absolute error of  $\epsilon$  is  $m^* = \min \{m \geq 0 : t_{r+m-1, 1-\alpha/2} [S_r^2/(r+m)]^{1/2} \leq \epsilon\}$ . For large  $r$ ,  $t_{r+m-1, 1-\alpha/2}$  does not change much as a function of  $m$ . The CI width is essentially scaled by shrinking the variance for the unknown mean from  $S_r^2/r$  to

$$S_r^2/(r+m). \tag{2}$$

Analogous results exist for asymptotic normal approximations to the posterior distributions of input parameters.

For the rest of this section, we focus on *one* parameter  $\theta = (\vartheta_1, \dots, \vartheta_d)$  for a single source of randomness (dropping the subscript  $i$  for notational simplicity).

**Proposition 3.1.** *For each  $n$ , let  $f_n(\cdot | \mathbf{x}_n)$  denote the posterior p.d.f. of the  $d$ -dimensional parameter vector  $\theta_n$ ; let  $\tilde{\theta}_n$  be the associated MAP; and define the Bayesian observed information  $\mathbf{G}_n = \Sigma_n^{-1}$  by*

$$[\mathbf{G}_n]_{jl} = [\Sigma_n^{-1}]_{jl} = - \left. \frac{\partial^2 \log f_n(\theta | \mathbf{x}_n)}{\partial \vartheta_j \partial \vartheta_l} \right|_{\theta = \tilde{\theta}_n} \text{ for } j, l = 1, \dots, d. \tag{3}$$

*Then the random variable  $\Sigma_n^{-1/2}(\theta_n - \tilde{\theta}_n)$  converges in distribution to a standard (multivariate) normal random variable as  $n \rightarrow \infty$  if certain regularity conditions hold.*

*Proof.* See Bernardo and Smith (1994, Prop 5.14). The regularity conditions presume a nonzero prior probability density near the “true”  $\boldsymbol{\theta}$ , and the “steepness”, “smoothness”, and “concentration” regularity conditions imply that the data is informative for all components of  $\boldsymbol{\theta}_n$ , so the posterior becomes peaked near  $\tilde{\boldsymbol{\theta}}_n$ .  $\square$

Colloquially,  $\boldsymbol{\Sigma}_n$  is proportional to  $1/n$ , so the posterior variance shrinks as in Equation (2) as additional samples are observed. Gelman, Carlin, Stern, and Rubin (1995, Sec. 4.3) describe when the regularity conditions may not hold, including underspecified models, nonidentified parameters, cases in which the number of parameters grows with the sample size, unbounded likelihood functions, improper posterior distributions, or convergence to a boundary of the parameter space.

The expected information matrix in *one* observation is

$$[\mathbf{H}(\tilde{\boldsymbol{\theta}}_n)]_{jl} = \mathbb{E}_{X_\ell} \left[ -\frac{\partial^2 \log p(X_\ell | \boldsymbol{\theta})}{\partial \vartheta_j \partial \vartheta_l} \right] \Big|_{\boldsymbol{\theta}=\tilde{\boldsymbol{\theta}}_n} \quad \text{for } j, l = 1, \dots, d.$$

so the analog to the asymptotic approximation in Equation (2) for the variance of  $\boldsymbol{\theta}$ , given that  $m$  additional data points are to be collected, is

$$(\boldsymbol{\Sigma}_n^{-1} + m\mathbf{H}(\tilde{\boldsymbol{\theta}}_n))^{-1}. \quad (4)$$

Equation (4) simplifies under some special conditions. A sufficient condition is that the distribution for  $\boldsymbol{\theta}$  be in the regular exponential family (which includes the exponential, gamma, normal, binomial, multinomial, lognormal, Poisson, and many other distributions) and that a conjugate prior distribution be used to describe initial parameter uncertainty (see Appendix B). Use of a conjugate prior distribution implies that the posterior distribution will have the same functional form as the prior distribution.

Appendix B shows that if the prior distribution before observing any data is a conjugate prior, then  $\boldsymbol{\Sigma}_n^{-1} = (n_0 + n)\mathbf{H}(\tilde{\boldsymbol{\theta}}_n)$  for some  $n_0$  and Equation (4) simplifies:

$$(\boldsymbol{\Sigma}_n^{-1} + m\mathbf{H}(\tilde{\boldsymbol{\theta}}_n))^{-1} = \frac{\mathbf{H}^{-1}(\tilde{\boldsymbol{\theta}}_n)}{n_0 + n + m}. \quad (5)$$

Equation (4) also simplifies if the MLE  $\hat{\boldsymbol{\theta}}_n$  replaces  $\tilde{\boldsymbol{\theta}}_n$  and the expected information  $\hat{\boldsymbol{\Sigma}}_n^{-1} = n\mathbf{H}(\hat{\boldsymbol{\theta}}_n)$  replaces  $\boldsymbol{\Sigma}_n^{-1}$ .

All of the distributions for the  $M/M/1$  queue and the critical care facility are in the regular exponential family and satisfy the regularity conditions of Proposition 3.1. The specific numbers of service time and routing decision observations for the critical care facility were not published. For the analysis here, we presumed standard noninformative prior distributions for each unknown parameter (Bernardo and Smith 1994). We further assumed that the published point estimates were MLEs based on 100 patients, so that 20 patients were routed to the ICU alone, 5 were routed to the CCU alone, and so forth. This assumption uniquely determines the sufficient statistics of the data and the proper posterior distributions for the parameters. It also uniquely determines the observed information matrices  $\boldsymbol{\Sigma}_n^{-1}$  for each parameter. Ng and Chick (2001) provided details for the information matrix and input uncertainty of the Poisson and lognormal distributions. For the multinomial routing probabilities, the Dirichlet distribution is the conjugate prior,  $\pi(p_1, p_3, p_4) \propto p_1^{\alpha_{01}-1} p_3^{\alpha_{03}-1} p_4^{\alpha_{04}-1} (1 - p_1 - p_3 - p_4)^{\alpha_{02}-1}$ , with noninformative distribution  $\alpha_{0j} = 1/2$  for  $j = 1, \dots, 4$  (Bernardo and Smith 1994). Reparametrizing this into the canonical conjugate form as in Appendix B, we have  $\boldsymbol{\psi} = (\log p_1, \log p_3, \log p_4)$ . Since this mapping is bijective, the results from Appendix B also hold in the  $(p_1, p_3, p_4)$  coordinates. If  $s_j$  patients are routed to location  $j$ , then the

posterior distribution is Dirichlet with parameters  $\alpha_j = \alpha_{0j} + s_j$ . If  $\alpha_j > 1$  for all  $j$ , then the MAP is  $\tilde{p}_j = (\alpha_j - 1) / [\sum_{j=1}^4 (\alpha_j - 1)]$ , and  $\Sigma_n^{-1} = [\sum_{j=1}^4 (\alpha_j - 1)] \mathbf{H}(\tilde{\boldsymbol{\theta}}_n) = [\frac{4}{2} + n - 4] \mathbf{H}(\tilde{\boldsymbol{\theta}}_n) = (n_0 + n) \mathbf{H}(\tilde{\boldsymbol{\theta}}_n)$ , where  $n_0 = -2$ .

## 4 Parameter and Performance Uncertainty

Section 4.1 derives approximations for performance uncertainty, measured by variance, as a function of parameter uncertainty when  $g(\cdot)$  is a known function. It formulates and solves an optimization problem that allocates resources for further data collection to minimize an asymptotic approximation to performance uncertainty. Section 4.2 identifies sampling allocations when  $g(\cdot)$  is unknown. It shows whether it is more important to run more simulation replications (to improve the gradient estimate) or to collect more field data (to reduce parameter uncertainty).

### 4.1 Known Response Function

Denote gradient information of the known response surface,  $g(\boldsymbol{\Theta})$ , evaluated at the MAP estimator of  $\boldsymbol{\Theta}$ , by

$$\begin{aligned} \beta_j &= \left. \frac{\partial g(\boldsymbol{\Theta})}{\partial \vartheta_j} \right|_{\boldsymbol{\Theta}=\tilde{\boldsymbol{\theta}}_n} \quad \text{for } j = 1, \dots, K \\ \boldsymbol{\beta}_i &= \left. \nabla_{\boldsymbol{\theta}_i} g(\boldsymbol{\Theta}) \right|_{\boldsymbol{\Theta}=\tilde{\boldsymbol{\theta}}_n} \quad \text{for } i = 1, \dots, k \\ \boldsymbol{\beta} &= (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_k). \end{aligned} \tag{6}$$

Suppose that the parameters for each source of randomness are estimated with  $n$  observations each, and that each satisfies the asymptotic normality properties of Proposition 3.1. Proposition 4.1 states that  $g(\boldsymbol{\Theta})$  is also asymptotically normal as more observations are collected from each of the sources of randomness.

**Proposition 4.1.** *Use the assumptions, notation and technical conditions of Proposition 3.1, with the modification that  $\boldsymbol{\Theta}_n = (\boldsymbol{\theta}_{1n}, \dots, \boldsymbol{\theta}_{kn})$  is the  $n$ th vector of parameters, so that  $f_{in}(\cdot)$  is the posterior density function of the random variable  $\boldsymbol{\theta}_{in}$ , the MAP estimator of  $\boldsymbol{\Theta}$  is  $\tilde{\boldsymbol{\Theta}}_n = (\tilde{\boldsymbol{\theta}}_{1n}, \dots, \tilde{\boldsymbol{\theta}}_{kn})$ , and  $\Sigma_n = \text{diag}[\Sigma_{1n} \dots \Sigma_{kn}]$  is the  $n$ th block diagonal matrix with appropriate submatrices for parameters  $i = 1, \dots, k$ , assuming  $n$  observations are available from each source of randomness. Suppose that  $\boldsymbol{\Theta}$  has an asymptotic Normal  $(\tilde{\boldsymbol{\Theta}}_n, \Sigma_n)$  distribution, as in the conclusion of Proposition 3.1. Suppose that  $g$  is continuously differentiable near  $\boldsymbol{\Theta}_0$ , that  $\bar{\sigma}_n^2 \rightarrow 0$  and  $\tilde{\boldsymbol{\Theta}}_n \rightarrow \boldsymbol{\Theta}_0$  in probability, and  $\bar{\sigma}_n^2 = O(\underline{\sigma}_n^2)$ , where  $\bar{\sigma}_n^2$  and  $\underline{\sigma}_n^2$  respectively denote the largest and smallest eigenvalues of  $\Sigma_n$ .*

*The random variable  $(\boldsymbol{\beta}_n \Sigma_n \boldsymbol{\beta}_n^T)^{-1/2} [g(\boldsymbol{\Theta}) - g(\tilde{\boldsymbol{\Theta}}_n)]$  converges in distribution to the standard normal distribution. Colloquially,  $g(\boldsymbol{\Theta})$  is asymptotically distributed*

$$\text{Normal} \left( g(\tilde{\boldsymbol{\Theta}}_n), \boldsymbol{\beta}_n \Sigma_n \boldsymbol{\beta}_n^T \right).$$

*Proof.* Mendoza (1994) and Bernardo and Smith (1994, Prop. 5.17) provide an analogous result for bijective functions  $g(\boldsymbol{\Theta})$  when  $k = 1$ . The result holds for univariate  $g(\boldsymbol{\Theta})$  because the marginal distribution of a single dimension of a multivariate joint Gaussian is also Gaussian. That result generalizes directly to  $k \geq 1$ , as required here, with straightforward algebra. This is a Bayesian analog for classical results from Serfling (1980, Sec. 3.3).  $\square$

Block diagonality reflects an assumption that the parameters for different sources of randomness are *a priori* independent, and that samples are conditionally independent, given the value of  $\Theta$ .

The plug-in estimator  $g(\hat{\Theta}_n)$  converges to  $g(\Theta_0)$  but is biased for nonlinear functions. Appendix A discusses bias and the mean squared error of estimates, using a first-order bias approximation. For now we focus on the variance of the estimators.

The number of observations for each of the  $k$  sources of randomness may differ, as may the number of additional data points to collect. Let  $n_i$  be the number of data points for source  $i$ , suppose that  $m_i$  additional points are to be collected for each source of randomness  $i$ , and set  $\mathbf{m} = (m_1, \dots, m_k)$ . Using the approximation of Equation (4), we see that the variance in the unknown performance after collecting additional information is approximately

$$V_{\text{par}}(\mathbf{m}) = \sum_{i=1}^k \beta_{i\mathbf{n}_i} (\Sigma_{i\mathbf{n}_i}^{-1} + m_i \mathbf{H}_i(\tilde{\theta}_{i\mathbf{n}_i}))^{-1} \beta_{i\mathbf{n}_i}^T \quad (7)$$

If each observed information matrix is proportional to the corresponding expected information matrix, as when conjugate priors are used, then Equation (7) simplifies:

$$V_{\text{par}}(\mathbf{m}) = \sum_{i=1}^k \frac{\xi_i}{n_{0i} + n_i + m_i}, \quad \text{where } \xi_i = \beta_{i\mathbf{n}_i} \mathbf{H}_i^{-1}(\tilde{\theta}_{i\mathbf{n}_i}) \beta_{i\mathbf{n}_i}^T. \quad (8)$$

The optimal sampling plan for reducing this asymptotic approximation to the variance in system performance, assuming that sampling costs for source  $i$  is linear in the number of samples, is therefore the solution to Problem (9).

$$\begin{aligned} \min \quad & V_{\text{par}}(\mathbf{m}) \\ \text{s.t.} \quad & m_i \geq 0 \quad \text{for } i = 1, \dots, k \\ & \sum c_i m_i = B \end{aligned} \quad (9)$$

Cost differences may arise if data collection can be automated for some but not all sources of randomness, if suppliers differ in willingness to share data, and so forth.

**Proposition 4.2.** *If  $V_{\text{par}}(\mathbf{m})$  simplifies to Equation (8), the integer restriction is relaxed (continuous  $m_i$ ), and  $B$  is large, then the solution to Problem (9) is:*

$$m_i^* = \frac{B + \sum_{\ell=1}^k (n_{0\ell} + n_\ell) c_\ell}{\sum_{j=1}^k \left( \frac{\xi_j c_j}{\xi_i} \right)^{1/2}} - (n_{0i} + n_i) \quad \text{for } i = 1, \dots, k \quad (10)$$

*Proof.* See Appendix C. □

For small  $B$ , the nonnegativity constraints for the  $m_i$  need consideration. Special features of the data collection process can be handled by adding constraints, e.g. set a subset of  $m_i$ 's to be equal for simultaneously reported data.

**Example:** The mean occupancy of an  $M/M/1$  queue is  $g(\lambda, \mu) = \lambda/(\mu - \lambda)$  when  $\mu > \lambda$ . The gamma distribution is the conjugate prior for an unknown rate of samples with an exponential distribution. Suppose that a gamma prior distribution  $\text{Gamma}(\alpha, \nu)$  models uncertainty for the unknown arrival rate, with mean  $\alpha/\nu$  and density  $f_\Lambda(w) = \nu^\alpha w^{\alpha-1} \exp(-w\nu)/\Gamma(\alpha)$  for all  $w \geq 0$ . If  $n_a$  arrivals are observed, the posterior distribution is also Gamma with parameters  $\alpha_\lambda = \alpha + n_a$ ,  $\nu_\lambda = \nu + \sum_{i=1}^{n_a} x_{il}$ . The MAP estimate is  $\tilde{\lambda} = (\alpha_\lambda - 1)/\nu_\lambda$ , and  $\mathbf{H}_1(\tilde{\lambda}) = 1/\tilde{\lambda}^2 = \nu_\lambda^2/(\alpha_\lambda - 1)^2$ .



With calculus,  $\Sigma_{1n_a}^{-1} = (\alpha_\lambda - 1)/\tilde{\lambda}^2 = \nu_\lambda^2/(\alpha_\lambda - 1)$ . Similarly assume the service rate has a conjugate gamma prior, that  $n_s$  service times are observed, resulting in a Gamma  $(\alpha_\mu, \nu_\mu)$  posterior distribution. If  $\tilde{\mu} > \tilde{\lambda}$  (otherwise, we would consider a multiserver system), then Equation (7) implies

$$V_{\text{par}}(\mathbf{m}) = \frac{\tilde{\lambda}^2 \tilde{\mu}^2}{(\tilde{\mu} - \tilde{\lambda})^4} \left( \frac{1}{\alpha_\lambda - 1 + m_a} + \frac{1}{\alpha_\mu - 1 + m_s} \right).$$

As conjugate priors are used,  $\Sigma_{1n_a}^{-1} = (\alpha_\lambda - 1)\mathbf{H}_1(\tilde{\lambda}_{n_a})$  and  $\Sigma_{2n_s}^{-1} = (\alpha_\mu - 1)\mathbf{H}_2(\tilde{\mu}_{n_s})$ . This reduces to Equation (8) where  $n_{0a} + n_a = \alpha_\lambda - 1$  and  $n_{0s} + n_s = \alpha_\mu - 1$ .

If  $m_a$  interarrival times can be collected at cost  $c_a$  each, and  $m_s$  service times can be collected at cost  $c_s$  each, then by Proposition 4.2, it is optimal to sample with  $m_s^* = (B - c_a m_a^*)/c_s$  and

$$m_a^* = \frac{B + (\alpha_\lambda - 1)c_a + (\alpha_\mu - 1)c_s}{c_a + (c_a c_s)^{1/2}} - (\alpha_\lambda - 1)$$

for large  $B$ . If  $c_a = c_s$ , this choice minimizes the difference between  $\alpha_\lambda + m_a$  and  $\alpha_\mu + m_s$  (the same effective number of data points is desired for both parameters).

If the arrival and service rates are believed to have a gamma distribution, then the queue is unstable with positive probability. In fact, the expectation of the average queue length does not exist even after conditioning on stability,  $E_{\lambda, \mu}[\lambda/(\mu - \lambda) \mid \lambda < \mu] = \infty$  (Chick 2001). This means that using Bayesian models for unknown statistical input parameters may lead to simulation output processes whose expectation is not defined. At least two alternative approaches can be taken. First, a similar analysis can be done for a model that not only has finite moments but is also more realistic (say, by considering a finite waiting capacity or by modeling a finite time horizon rather than an infinite horizon time average). Second, one can use the asymptotic normality approximation in spite of the lack of moments as a starting point for reducing input uncertainty. Section 5.2 illustrates by example when this second alternative might be expected to work well or not.

The optimal allocation for Problem (9) does not depend on the parametrization  $\Theta$  under certain conditions. Let  $\mathbf{T} : \mathbb{R}^K \mapsto \mathbb{R}^K$  be a one-to-one continuously differentiable transformation  $\boldsymbol{\psi} = (\psi_1, \dots, \psi_K) = \mathbf{T}(\Theta) = [T_1(\Theta), \dots, T_K(\Theta)]$  with differentiable inverse  $\Theta = \mathbf{T}^{-1}(\boldsymbol{\psi})$ , and nonsingular matrix of derivatives  $\mathbf{J}_{\mathbf{T}}(\Theta)$ , where  $[\mathbf{J}_{\mathbf{T}}(\Theta)]_{j\ell} = \partial T_j / \partial \vartheta_\ell$ . We write  $\mathbf{J}_{\mathbf{T}^{-1}}(\boldsymbol{\psi})$ , where  $[\mathbf{J}_{\mathbf{T}^{-1}}(\boldsymbol{\psi})]_{j\ell} = \partial T_j^{-1} / \partial \psi_\ell$ . Let  $\tilde{\boldsymbol{\psi}}_{\mathbf{n}}$  be the MAP estimator for  $\boldsymbol{\psi}$  given  $\mathcal{E}_{\mathbf{n}}$ . Denote the response gradient with respect to  $\boldsymbol{\psi}$  by  $\boldsymbol{\beta}_{\boldsymbol{\psi}}$  and the gradient in Equation (6) with respect to  $\Theta$  by  $\boldsymbol{\beta}_{\Theta}$ . Similarly denote the observed information matrix and expected information in one observation (measured in  $\boldsymbol{\psi}$  coordinates) by  $\Sigma_{\boldsymbol{\psi}, \mathbf{n}}$  and  $\mathbf{H}_{\boldsymbol{\psi}}$  to emphasize the coordinate system. The expression  $V_{\text{par}}(\mathbf{m})$  in Equation (8) explicitly depends upon  $\Theta$ . We write  $V_{\text{par}}(\mathbf{m}; \Theta)$  to denote that dependence explicitly. Let  $V_{\text{par}}(\mathbf{m}; \boldsymbol{\psi})$  denote the analogous approximation approximation when computed in  $\boldsymbol{\psi}$  coordinates.

**Proposition 4.3.** *Suppose the setup of Proposition 4.1, that the distribution of  $\Theta$  is in the regular exponential family, and that a conjugate prior is adopted. Let  $\mathbf{T} : \mathbb{R}^K \mapsto \mathbb{R}^K$  be a one-to-one continuously differentiable transformation  $\boldsymbol{\psi} = \mathbf{T}(\Theta) = [T_1(\Theta), \dots, T_K(\Theta)]$  with differentiable inverse  $\Theta = \mathbf{T}^{-1}(\boldsymbol{\psi})$ , and nonsingular Jacobian  $\mathbf{J}_{\mathbf{T}}(\Theta)$ . If conditions (a) for each  $i$ ,  $\psi_i = T_i(\Theta)$  depends upon only one  $\theta_j$ , and (b)  $\tilde{\boldsymbol{\psi}}_{\mathbf{n}} = \mathbf{T}(\tilde{\Theta}_{\mathbf{n}})$  are true, then  $V_{\text{par}}(\mathbf{m}; \Theta) = V_{\text{par}}(\mathbf{m}; \boldsymbol{\psi})$ .*

*Proof.* Presume the setup of Proposition 4.1, define  $\mathbf{T}$  and the other notation as in the problem statement, and suppose conditions (a) and (b) hold. Condition (a) means that  $\mathbf{J}_{\mathbf{T}}(\Theta)$  has the same block diagonal structure as  $\Sigma$ , and that no  $\psi_i$  is a function of parameter vectors for multiple sources of randomness.

Suppose first that there is only  $k = 1$  source of randomness with parameter  $\theta$  and drop the subscript  $i$  for notational simplicity. In this notation,  $\xi_i$  can be rewritten  $\xi = \beta_{\Theta} \mathbf{H}_{\Theta}^{-1}(\tilde{\Theta}_n) \beta_{\Theta}^{\text{T}}$  in  $\Theta$  coordinates. By the chain rule for partial differentiation, and the assumption  $\tilde{\psi}_n = \mathbf{T}(\tilde{\Theta}_n)$ ,

$$\begin{aligned} \beta_{\psi} &= \nabla_{\psi} g(\psi) \Big|_{\psi=\tilde{\psi}_n} \\ &= \left( \nabla_{\Theta} g(\Theta) \Big|_{\Theta=\tilde{\Theta}_n} \right) \left( \mathbf{J}_{\mathbf{T}^{-1}}(\tilde{\psi}_n) \right). \end{aligned}$$

By Schervish (1995, p. 115),  $\mathbf{H}_{\psi}(\psi) = \mathbf{J}_{\mathbf{T}^{-1}}^{\text{T}}(\psi) \mathbf{H}_{\Theta}(\Theta) \mathbf{J}_{\mathbf{T}^{-1}}(\psi)$ , so that with additional calculus  $\mathbf{H}_{\psi}^{-1}(\psi) = (\mathbf{J}_{\mathbf{T}^{-1}}^{\text{T}}(\psi) \mathbf{H}_{\Theta}(\Theta) \mathbf{J}_{\mathbf{T}^{-1}}(\psi))^{-1} = \mathbf{J}_{\mathbf{T}}(\Theta) \mathbf{H}_{\Theta}^{-1}(\Theta) \mathbf{J}_{\mathbf{T}}^{\text{T}}(\Theta)$ . Therefore, in  $\psi$  coordinates,

$$\begin{aligned} \xi &= \beta_{\psi} \mathbf{H}_{\psi}^{-1}(\tilde{\psi}_n) \beta_{\psi}^{\text{T}} \\ &= [\nabla_{\Theta} g(\tilde{\Theta}_n) \mathbf{J}_{\mathbf{T}^{-1}}(\tilde{\psi}_n)] [\mathbf{J}_{\mathbf{T}}(\tilde{\Theta}_n) \mathbf{H}_{\Theta}^{-1}(\tilde{\Theta}_n) \mathbf{J}_{\mathbf{T}}^{\text{T}}(\tilde{\Theta}_n)] [\mathbf{J}_{\mathbf{T}^{-1}}^{\text{T}}(\tilde{\psi}_n) \nabla_{\Theta}^{\text{T}} g(\tilde{\Theta}_n)] \\ &= \nabla_{\Theta} g(\tilde{\Theta}_n) \mathbf{H}_{\Theta}^{-1}(\tilde{\Theta}_n) \nabla_{\Theta}^{\text{T}} g(\tilde{\Theta}_n) \\ &= \beta_{\Theta} \mathbf{H}_{\Theta}^{-1}(\tilde{\Theta}_n) \beta_{\Theta}^{\text{T}}. \end{aligned}$$

The third equality follows because  $\mathbf{J}_{\mathbf{T}}(\tilde{\Theta}_n) \mathbf{J}_{\mathbf{T}^{-1}}(\tilde{\psi}_n) = \mathbf{J}_{\mathbf{T}^{-1}}(\tilde{\psi}_n) \mathbf{J}_{\mathbf{T}}(\tilde{\Theta}_n)$  is the  $K \times K$  identity matrix, when  $\tilde{\psi}_n = \mathbf{T}(\tilde{\Theta}_n)$  holds.

For  $k > 1$ , the block diagonality assumption of Condition (a) assures that the  $\xi_i$  decouple, one for each of the  $k$  independent sources of randomness.  $\square$

The MAP estimator is not invariant to coordinate changes. That is, if  $\tilde{\Theta}_n$  is the MAP in  $\Theta$  coordinates, and  $\tilde{\psi}_n$  is the MAP in  $\psi$  coordinates, then it may not be true that  $\tilde{\psi}_n = \mathbf{T}(\tilde{\Theta}_n)$ . For example, if the unknown rate  $\psi$  of an exponential distribution has a Gamma  $(\alpha_n, \nu_n)$  distribution, then  $\tilde{\psi}_n = (\alpha_n - 1)/\nu_n$  and the mean  $\theta = 1/\psi$  has an inverted gamma distribution, with MAP  $\tilde{\theta}_n = \nu_n/(\alpha_n + 1) \neq 1/\tilde{\psi}_n$ . The multivariate generalization of Proposition 4.1 implies that  $\tilde{\psi}_n$  and  $\mathbf{T}(\tilde{\Theta}_n)$  are asymptotically equal. For a Bayesian, then, the allocations are asymptotically invariant to parametrization, in this sense.

For a frequentist, the allocations are invariant for finite  $n$ , as the MLE is invariant, given mild regularity conditions (Edwards 1984), and expected information matrices behave as required for a development parallel to Proposition 4.3.

## 4.2 Unknown Response Function

If the response function  $g(\Theta)$  is unknown, an approximation is needed to estimate the mean of some future output  $Y_{r+1}$ . The asymptotic approximations above “hit their limit” in some sense with linear approximations, so we approximate  $g(\Theta)$  near  $\tilde{\Theta}_n$  with a first-order regression model that may depend on past simulation output.

$$\mathbb{E}[Y_{r+1} \mid \Theta] \approx h(\Theta, \beta) = \beta_0 + \sum_{j=1}^K \beta_j (\vartheta_j - \tilde{\vartheta}_j) \quad (11)$$

The notation  $h(\Theta, \beta)$  replaces the notation  $g(\Theta)$  to make it explicit that the estimated response depends upon both  $\Theta$  and  $\beta$ . Note that  $\beta_0$  takes the role of  $g(\tilde{\Theta})$  and  $\beta_j$  takes the role of  $\partial g(\Theta)/\partial \vartheta_j$ . If the derivatives  $\beta_j = \partial h/\partial \vartheta_j$  are estimated from simulations, for  $j = 1, 2, \dots, K$ , then uncertainty about both  $\Theta$  and  $\beta$  must be considered. We take a Bayesian approach and presume  $\beta$  to be a random vector, whose posterior distribution is updated when simulation replications are run.

Suppose that  $r$  replications have been run, with  $r \times (K+1)$  design matrix  $\mathbf{M} = [(\Theta_1 - \tilde{\Theta}_n), (\Theta_2 - \tilde{\Theta}_n), \dots, (\Theta_r - \tilde{\Theta}_n)]^T$ , resulting in outputs  $\mathbf{y} = [y_1, \dots, y_r]^T$  with sample mean  $\bar{y} = r^{-1} \sum_{i=1}^r y_i$ . The  $r$  runs may have come from  $q$  replications at each of  $t$  design points selected with classical design of experiments (DOE) techniques, with  $r = qt$ . Or the  $\Theta_i$  may have come from using a BMA to sample  $t$  values of  $\Theta_i$  independently from  $f(\Theta | \mathbf{x})$ , with  $q$  replications for each  $\Theta_i$  (Chick (2001) used  $q = 1$ , Zouaoui and Wilson (2003) extended to  $q > 1$  to assess stochastic uncertainty, also see Zouaoui and Wilson (2004)). In either case, the posterior distribution of  $\beta$ , assuming the standard linear model with a noninformative prior distribution  $p(\beta, \sigma^2) \propto \sigma^{-2}$ , and that  $r > (K+1)$  (Gelman, Carlin, Stern, and Rubin 1995, Sec. 8.3), is

$$\begin{aligned} p(\beta | \mathbf{y}, \mathbf{M}, \sigma^2) &\sim \text{Normal} \left( \hat{\beta}_r, (\mathbf{M}^T \mathbf{M})^{-1} \sigma^2 \right), \\ p(\sigma^2 | \mathbf{y}, \mathbf{M}) &\sim \text{Inverse-Chi}^2 (r - (K + 1), s^2), \\ \hat{\beta}_r &= (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k) = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T \mathbf{y} \\ s^2 &= \frac{1}{r - (K + 1)} (\mathbf{y} - \mathbf{M} \hat{\beta}_r)^T (\mathbf{y} - \mathbf{M} \hat{\beta}_r). \end{aligned} \quad (12)$$

The standard classical estimate for  $\beta$  is  $\hat{\beta}_r$ , and for  $\sigma^2$  is  $s^2$ .

The posterior marginal distribution of  $\beta$  is actually multivariate  $t$ , not normal. We approximate this distribution by estimating  $\sigma^2$  by  $s^2$  and presuming the variance is known, as we did in Equation (2), to arrive at a normal distribution approximation for the unknown response vector that has mean  $\hat{\beta}_r$  and variance of the form  $\Sigma_\beta / (r_0 + r)$ . Similarly, we presume that the posterior variance of  $\beta$  is of the form  $\Sigma_\beta / (r_0 + r + m_0)$  when  $m_0$  additional replications are run (cf. Equation (8)). Amongst other scenarios, this approximation makes sense when  $m_0 = q't$ , with  $q'$  additional replications are run at each of the  $t$  originally selected  $\Theta_i$  as above.

Then the approximation  $V_{\text{tot}} = \text{Var}(h(\Theta, \beta) | \mathcal{E}_n, \mathbf{y}, \mathbf{M}, \sigma^2 = s^2)$  for overall performance uncertainty  $\text{Var}(E[Y_{r+1}] | \mathcal{E}_n, \mathbf{y}, \mathbf{M})$  asymptotically consists of components that reflect parameter uncertainty ( $V_{\text{par}}$ ), response surface/gradient uncertainty ( $V_{\text{resp}}$ ), and stochastic noise ( $\sigma^2 = s^2$ ). The following expectations, through Equation (13), implicitly condition on all observations so far,  $\mathcal{E}_n, \mathbf{y}, \mathbf{M}$ .

$$\begin{aligned} V_{\text{tot}} &\approx \frac{s^2}{r_0 + r} + \sum_{i=1}^K \sum_{j=1}^K \text{E} \left\{ \left[ \beta_i (\vartheta_i - \tilde{\vartheta}_i) (\vartheta_j - \tilde{\vartheta}_j)^T \beta_j \right] \right. \\ &\quad \left. - \text{E} \left[ \beta_i (\vartheta_i - \tilde{\vartheta}_i) \right] \text{E} \left[ \beta_j (\vartheta_j - \tilde{\vartheta}_j) \right] \right\} \\ &= \frac{s^2}{r_0 + r} + \sum_{i=1}^K \sum_{j=1}^K \text{E} [\beta_i \beta_j] \text{Cov}(\vartheta_i, \vartheta_j) \\ &= \frac{s^2}{r_0 + r} + \sum_{i=1}^K \sum_{j=1}^K \text{Cov}(\vartheta_i, \vartheta_j) [\text{E}(\beta_i) \text{E}(\beta_j) + \text{Cov}(\beta_i, \beta_j)] \\ &\approx \frac{s^2}{r} + \sum_{i=1}^k \nabla_{\theta_i} h(\tilde{\Theta}_n, \hat{\beta}_r) \Sigma_{in_i} \nabla_{\theta_i} h(\tilde{\Theta}_n, \hat{\beta}_r)^T \\ &\quad + \sum_{i=1}^K \sum_{j=1}^K \frac{[\Sigma_\beta]_{ij}}{r} [\Sigma_n]_{ij} \\ &\stackrel{\text{def}}{=} V_{\text{stoch}} + V_{\text{par}} + V_{\text{resp}} \end{aligned} \quad (13)$$

The last equality defines  $V_{\text{resp}}$ . The preceding equalities depend upon independence and asymptotic normality assumptions. The term  $V_{\text{stoch}} = s^2/r$  in the second approximation (with  $r_0 = 0$ ) follows directly from plugging in  $s^2$  for  $\sigma^2$  and examining the posterior conditional variance of  $\beta_0$  in Equation (12). Cheng and Holland (2004) also found an estimator of  $V_{\text{stoch}}$  that is proportional to  $1/r$ .

Cheng and Holland (2004) describe several methods to estimate the gradients of a stochastic simulation model with an unknown response model, and use those results to approximate  $V_{\text{stoch}}$  and  $V_{\text{par}}$ . A full study of the merits of the many ways to estimate the gradient  $\beta$ , and their influence upon  $V_{\text{tot}}$ , is an area for further study.

Since  $\text{Cov}(\beta_i, \beta_j)$  is inversely proportional to the number of runs, the asymptotic approximation for all uncertainty if  $m_0$  more replications and  $m_i$  data samples for each input are collected can be approximated (cf. Equation (8)) by

$$V_{\text{tot}}(\mathbf{m}, m_0) = \frac{s^2}{r + m_0} + \sum_{i=1}^k \frac{\xi_i}{n_{0i} + n_i + m_i} + \frac{a_i}{(r + m_0)(n_{0i} + n_i + m_i)}, \quad (14)$$

where  $a_i$  represents the sum of the subset of the terms in  $\sum_{l=1}^K \sum_{j=1}^K [\Sigma_{\beta}]_{lj} [\mathbf{H}_{\mathbf{n}}^{-1}]_{lj}$  that have both  $l$  and  $j$  indexing components of the overall parameter vector that deal with  $\theta_i$ , for  $i = 1, \dots, k$ . The double sum decouples into one term per input parameter vector due to the block diagonality of  $\Sigma_{\mathbf{n}}$ . Derivatives confirm that  $V_{\text{tot}}(\mathbf{m}, m_0)$  is convex and decreasing in  $m_0$  and each  $m_i$ .

Suppose there is a sampling constraint for more field data and simulation runs,

$$B = c_0 m_0 + \sum_{i=1}^k c_i m_i.$$

By Proposition 4.2, the  $\mathbf{m}$  that minimizes  $V_{\text{tot}}(\mathbf{m}, m_0)$  for a fixed  $m_0 < B/c_0$ , subject to the sampling constraint, can be found by substituting  $B$  with  $B - c_0 m_0$ , and replacing  $\xi_i$  with  $\xi_i + a_i/(r + m_0)$  for  $i = 1, 2, \dots, k$  in Equation (10). The optimal  $(\mathbf{m}, m_0)$  can then be found by univariate search for the optimal  $m_0$ . In practice, one might run more replications if  $V_{\text{stoch}} + V_{\text{resp}}$  is large relative to  $V_{\text{par}}(\mathbf{m})$ , and do more data collection if  $V_{\text{par}}(\mathbf{m})$  is unacceptably large.

## 5 Experimental Analysis

This section presents numerical studies to assess how close the asymptotic approximations are to the actual variance. It also assesses the variance, bias, and root mean squared error (RMSE) of estimators of the unknown mean performance as a function of both the number of initially collected samples of data and the total number of additional samples to be collected. We evaluate the influence of response surface uncertainty for a quadratic response in Section 5.1, consider a known nonlinear function with unknown parameters in Section 5.2, and apply the ideas to the ICU simulation in Section 5.3. In all three examples, we assume noninformative prior distributions for the unknown parameters, and initial data collection is complemented with a follow-up data collection plan based upon Equation (7) or Equation (14).

### 5.1 Unknown Nonlinear Response

The first numerical experiment assesses the reduction in uncertainty about the mean performance assuming that the response surface is unknown. Assume that the unknown parameters are the mean  $\mu_i$  and precision  $\lambda_i$  of  $k = 3$  sources of randomness for which data  $x_{i\ell}$  has been observed

( $i = 1, 2, 3; \ell = 1, 2, \dots, n_i$ ). The response function is quadratic in  $\Theta = (\mu_1, \lambda_1, \mu_2, \lambda_2, \mu_3, \lambda_3)$ , with known functional form but unknown response parameters.

Figure 2 summarizes the experiment. Based on initial data collection, the input parameters of the  $k$  normal distributions are estimated. The “true” response is

$$Y = g(\Theta) + \sigma Z = \beta_0 + \left( \sum_{i=1}^k \beta_{2i-1} \mu_i + \beta_{2i} \lambda_i \right) + \beta_7 \mu_1 \mu_2 + \beta_8 \mu_1^2 + \sigma Z$$

where  $\beta = (0, -2, 1, -2, 1, 1, 1, 1, 1)$ ,  $\mu = (2, 10, 2)$ ,  $\lambda = (1, 2, 1)$ ,  $\sigma^2 = 1/2$ . The vector  $(\beta, \sigma^2)$  is estimated from independent model output tested with  $q = 1$  replication at  $t = 29$  design points from a central composite design (CCD, Box and Draper 1987, with a  $2^{6-2}$  design for the factorial portion,  $2 \times 6$  star points with  $\alpha = 1$ , and one center point, for  $16+12+1=29$  points total). Design points were taken to be the MAP and endpoints of a 95% credible set for the six parameters (using a normal distribution approximation). Follow-up data collection and simulation run allocations from Equation (14) are used to reduce input and response parameter uncertainty.

For the inference process, we used a noninformative prior distribution  $\pi_i(\mu_i, \lambda_i) \propto \lambda_i^{-1}$  for  $i = 1, 2, 3$ . After  $n$  observations for a given source of randomness, and dropping the  $i$  momentarily for notational clarity, this results in normal gamma posterior pdf  $f_n(\mu, \lambda \mid \mathcal{E}_n) = \frac{(n\lambda)^{\frac{1}{2}}}{\sqrt{2\pi}} e^{-\frac{n\lambda}{2}(\mu - \tilde{\mu}_n)^2} \frac{\nu_n^{\alpha_n}}{\Gamma(\alpha_n)} \lambda^{\alpha_n-1} e^{-\nu_n \lambda}$ , where  $\tilde{\mu}_n = \sum_{\ell=1}^n x_{\ell}/n$ ,  $\alpha_n = (n-1)/2$ , and  $\nu_n = \sum_{\ell=1}^n (x_{\ell} - \tilde{\mu}_n)^2/2$  (Bernardo and Smith 1994). This implies that the posterior means of an unknown mean  $\mu$  is the sample average  $\tilde{\mu}_n$ , and the posterior marginal distribution of  $\lambda$  is a Gamma  $(\alpha_n, \nu_n)$  distribution, with MAP  $\hat{\lambda} = (\alpha_n - 1)/\nu_n = (n-3)/(2\nu_n)$  and posterior variance  $\text{Var}[\lambda \mid \mathcal{E}_n] = \alpha_n/\nu_n^2$ .

We estimated the mean of the variance approximation,  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$ , the expectation taken over the sampling distribution of the initial data samples, and with  $\mathbf{m}^*$  computed for  $B = 0, 10, \dots, 100$ ; and the variance  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[h(\tilde{\Theta}_{\mathbf{n}+\mathbf{m}^*}, \tilde{\beta}_{\mathbf{n}+\mathbf{m}^*})]$  and bias  $\kappa = E_{\mathcal{E}_n, \mathcal{D}}[h(\tilde{\Theta}_{\mathbf{n}+\mathbf{m}^*}, \tilde{\beta}_{\mathbf{n}+\mathbf{m}^*})] - g(\Theta)$  of the overall estimator of the mean, determined with respect to the sampling distributions of both rounds of sampling. Those samples determine the estimates  $\tilde{\Theta}_{\mathbf{n}+\mathbf{m}^*}$ ,  $\tilde{\beta}_{\mathbf{n}+\mathbf{m}^*}$  and best local linear approximation  $h(\Theta, \tilde{\beta}_{\mathbf{n}+\mathbf{m}^*})$  to  $g(\Theta)$  near  $\tilde{\Theta}_{\mathbf{n}+\mathbf{m}^*}$ . Estimates are based on 1000 macroreplications of data collection experiments for each combination of the number of initial samples ( $n_i = n = 10, 20, \dots, 100$  for  $i = 1, \dots, k$ ).

Figure 3 shows that  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$  and  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[h(\tilde{\Theta}_{\mathbf{n}+\mathbf{m}^*}, \tilde{\beta}_{\mathbf{n}+\mathbf{m}^*})]$  track each other reasonably well, which is expected since the response is a low order polynomial. The allocations tended to allocate more samples from inputs with a larger variance and/or larger influence on the output (intuitive, so data not shown). In spite of the relatively close match of variance, Figure 4 indicates the fraction  $m_1^*/B$  of samples allocated to learning about the major quadratic term  $(\mu_1, \lambda_1)$  varies as a function of  $B$  and  $n$  more significantly when there is less information about the response parameters (small  $r = qt$ ). This indicates that variability in the regression estimates may affect the optimal data collection allocations.

Figure 5 illustrates how significantly deleterious ignoring parameter uncertainty can be for the coverage probability and the mean half-width size of confidence intervals. For this example, constructing 95% CI with stochastic uncertainty alone (using  $z_{97.5}\sqrt{V_{\text{stoch}}}$ ) results in an extremely poor coverage probability (about 0 to 2% rather than 95%) because input parameters are poorly identified. The wider CI half-widths based on  $z_{97.5}\sqrt{V_{\text{tot}}}$  result in much better coverage (empirical coverage of 93-95%, very close to the nominal 95% level, when  $n_i = n = 20, 40, 60, 80$  data points were initially collected and when  $B = 0$ ). For a fixed  $B$ , the coverage is worst when the least amount of initial data is available (smaller  $n$ ). For a fixed  $n$ , the coverage declined somewhat with increasing  $B$ , dropping to 87-89% when  $B = 1000$  and down further to 77-82% when  $B = 4000$ . For

1. For  $i = 1, 2, \dots, 1000$  macroreplications
  - (a) Sample  $n$  conditionally independent, normally distributed samples ( $k = 3$  sources of randomness),  $x_{j\ell} \sim \text{Normal}(\mu_j, \lambda_j^{-1})$ , for  $j = 1, \dots, k; \ell = 1, 2, \dots, n$ .
  - (b) Determine MAP estimators  $\tilde{\mu}_i = \sum_{\ell=1}^{n_i} x_{i\ell}/n_i$ ,  $\tilde{\lambda}_i$  given the samples from Step (a).
  - (c) Generate  $q$  independent samples of  $y_{(l-1)q+j} = g(\Theta_{(l-1)q+j})$  for  $j = 1, \dots, q$ , at each design point  $\Theta_l$ , for  $l = 1, 2, \dots, t$ . The design points were on the lattice generated by the MAP and  $\text{MAP} \pm z_{0.95} \text{SE}$  for each of the  $2k$  components of the parameter vector (the unknown  $\mu_i, \lambda_i$ ).
  - (d) Estimate the response surface parameters  $\beta, \sigma^2$ , using Equation (12).
  - (e) Compute the optimal sampling allocations  $\mathbf{m}_{j_i}^*, m_{0_i}^*$  for the  $i$ th macroreplication to minimize the approximation  $V_{\text{tot}}(\mathbf{m}, m_0)$  in Equation (14) (given sampling budget  $B$  with costs  $c_0 = c_j = 1$ , for  $j = 1, \dots, k$ ).
  - (f) Generate  $m_{0_i}^*$  more additional runs at the simulation design points and reestimate the response surface (update  $\tilde{\beta}, \tilde{\sigma}$ ).
  - (g) Observe  $\mathbf{m}_i^*$  additional data points then update the estimates  $\tilde{\mu}_j, \tilde{\lambda}_j$  for  $j = 1, 2, 3$ , and  $\tilde{y}_i = h(\tilde{\Theta}, \tilde{\beta})$  (i.e. the estimate of the response function  $g(\Theta)$ ).
2. Compute the sample mean  $\bar{y}_n = \sum_{i=1}^{1000} \tilde{y}_i/1000$ , sample bias  $\bar{\kappa}_n$ , sample variance  $S_n^2$  of the  $\tilde{y}_i$ , and average number of additional samples  $\bar{m}$ .

Figure 2: Algorithm to assess variance, bias, and RMSE for the response surface example in Section 5.1, assuming  $\Theta_r = (\mu_1, \lambda_1, \dots, \mu_k, \lambda_k)$  and  $Y_r = g(\Theta_r) + \sigma Z_r$ .

fixed  $n$ , we hypothesize that the coverage degrades somewhat as  $B$  increases because the allocations and response surface estimates may be influenced by errors in the gradient estimators. In spite of the degradation in coverage for large  $B$ , the coverage is still orders of magnitude better than the coverage that considers stochastic uncertainty alone. Note that  $z_{97.5} \sqrt{V_{\text{tot}}}$  does not depend strongly upon  $n$  for large  $B$  in this example because  $V_{\text{resp}}$  remains a factor in the uncertainty. That term could be reduced by running additional replications at each design point.

Similar results were obtained for the first-order bias corrected estimators, for a different values of the response gradients, and for the 90% confidence level.

## 5.2 $M/M/1$ Queue

We assessed the influence of input uncertainty about the arrival rate  $\lambda$  and service rate  $\mu$  on the stationary mean  $g(\lambda, \mu) = \lambda/(\mu - \lambda)$  for an  $M/M/1$  queue. Noninformative prior distributions are used for the unknown arrival and service parameters so that  $\pi_\Lambda(\lambda) \propto \lambda^{-1}$  and  $\pi_\mu(\mu) \propto \mu^{-1}$ .

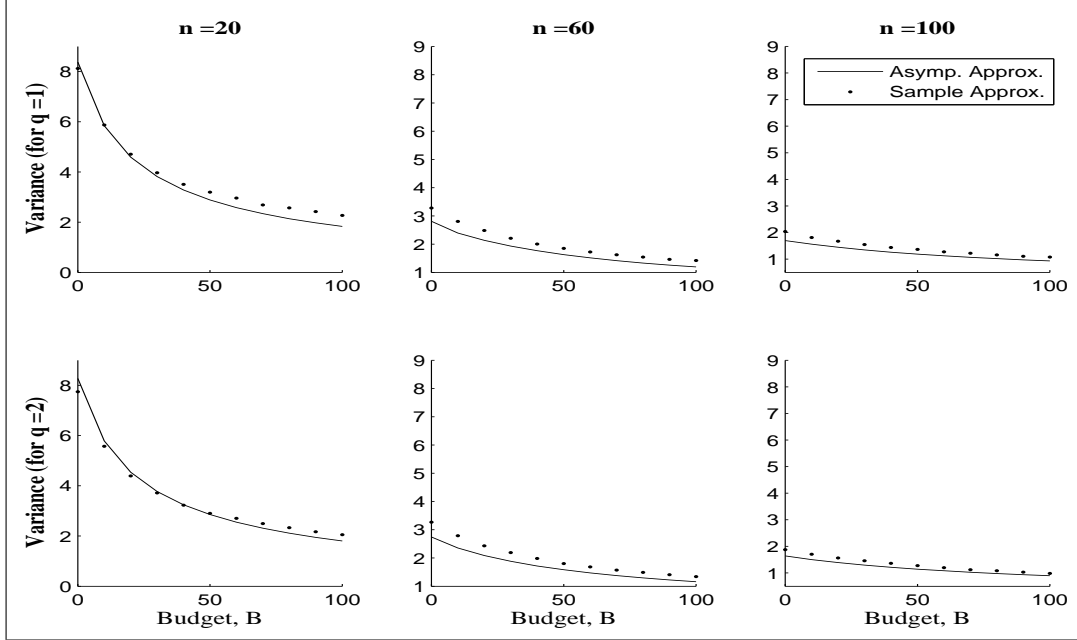


Figure 3: Comparison of mean asymptotic variance approximation  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$  (labeled “Asymp. Approx.”) with realized sample variance estimate of  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[h(\tilde{\Theta}_{n+m^*}, \tilde{\beta}_{n+m^*})]$  (labeled “Sample Approx.”), if  $q$  iterations are taken from each of  $t = 29$  design points in a CCD.

After observing  $n$  simulated arrival times and service times,  $\mathcal{E}_n = (a_1, a_2, \dots, a_n, s_1, s_2, \dots, s_n)$ , the modeler’s posterior distribution for  $\Lambda$  is Gamma( $n, \sum_{\ell=1}^n a_{\ell}$ ) (Bernardo and Smith 1994). A similar posterior distribution holds for  $\mu$ . Interarrival times and service times were simulated from the “true” exponential( $\lambda$ ) and exponential( $\mu$ ) distributions respectively.

Additional data  $\mathcal{D}$  collected from  $\mathbf{m}^* = (m_a^*, m_s^*)$  optimally-allocated follow-up samples were then simulated. We compared the mean of the asymptotic variance approximation  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$  predicted after the initial samples were collected but before the follow-up samples were collected (averaged over realizations of the initial  $n$  data observations) with the variance of the estimated mean  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[g(\tilde{\Theta}_{n+m^*})]$  (expectation with respect to both stages of sampling that estimate  $\Theta = (\lambda, \mu)$ ) and computed the overall bias  $\kappa = E_{\mathcal{E}_n, \mathcal{D}}[g(\tilde{\Theta}_{n+m^*})] - g(\Theta)$  after both rounds of sampling. Because expectations may not exist, we conditioned on the utilization being less than 0.995. Estimates are based on 1000 macroreplications of data collection experiments as in Figure 6, for each combination of the number of initial samples ( $n = 20, 40, \dots, 100$ ), follow-up samples ( $B = m_a^* + m_s^* = 0, 10, \dots, 100, 200, \dots, 1000$ ), and “true” utilizations ( $\lambda/\mu = 0.2, 0.4, 0.6, 0.8$  with  $\mu = 1$ ).

Several qualitative points that can be made. One, the mean asymptotic variance approximation  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$  better matches the sample variance estimate of  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[g(\tilde{\Theta}_{n+m^*})]$  in areas where  $g$  is more linear (low utilization) and/or the number of initial samples  $n$  is larger (Figure 7). The approximation can be considered very useful when  $g$  is relatively linear in the area where most of the posterior probability of the parameters is found. For higher utilization levels and fewer data points, the asymptotic approximation for variance and for bias was higher than the empirically observed variance and bias.

Two, the variance, bias, and MSE of the plug-in estimator  $g(\tilde{\Theta}_{n+m^*})$  was still reduced even when the approximation was poor. For example, the solid line in Figure 8 predicts how much first-order bias is expected, given that  $n$  observations of arrival and service times have been ob-

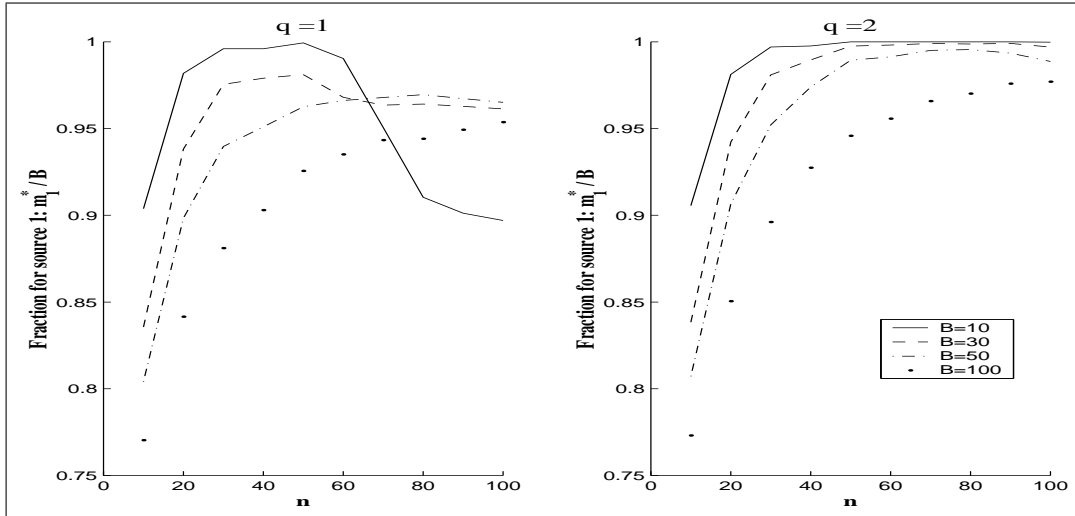


Figure 4: Fraction of budget for sampling from source of randomness 1, if  $q$  iterations are taken from each of  $t = 29$  design points in a CCD.

served, and another  $m_a^*, m_s^*$  observations are expected (based on the variance-reducing solution in Proposition 4.2). The first-order bias estimate is given in Appendix A. The dotted line represents the empirically observed bias after those additional observations were collected, and the updated plug-in estimator is used. The bias is low and the predicted bias is proportional to the estimated bias when the response surface is linear (note the scale in the row for  $\rho = 0.2$ ). The approximation is worse where the response surface is nonlinear, but an improvement in bias is still observed with more data.

Three, the optimal allocation for variance reduction samples both service times and arrival times equally, but the optimal MSE allocation tends to allocate more samples for service times, particularly with a small budget for additional samples or a high utilization (Figure 9).

Four, the MSE allocation does slightly better than the variance allocation in terms of bias and MSE, but not significantly so. Including the bias term in the optimal sampling allocation therefore added to the computing time to determine the optimal allocation, but did not significantly improve the statistics of any estimates (data not shown, graphs are so similar to the variance allocation graphs).

Five, the function  $g$  is particularly nonlinear for larger values of  $\hat{\rho}$ , so the estimator  $g(\hat{\lambda}, \hat{\mu})$  is biased. We also checked the variance, bias, and MSE of the first-order bias-corrected estimator (the plug-in estimator  $g(\hat{\lambda}, \hat{\mu})$  minus the bias correction from Equation (15) in Section A) and found that the bias correction actually hurt the variance, bias and MSE of  $g(\hat{\lambda}, \hat{\mu})$  because the first-order bias correction fared particularly poorly when  $\hat{\rho} > 0.92$ . We therefore recommend the plug-in estimator and asymptotic variance sampling allocation for the  $M/M/1$  queue, and by extension to other functions where the first-order bias correction is poor. We do not replicate the poor confidence interval coverage for the  $M/M/1$  queue already demonstrated (Barton and Schruben 2001) in related work.

### 5.3 Critical Care Facility

The critical care facility described in Section 2 can be used to illustrate the above analysis on a practical problem when the response surface is unknown. Suppose initial field data is collected,



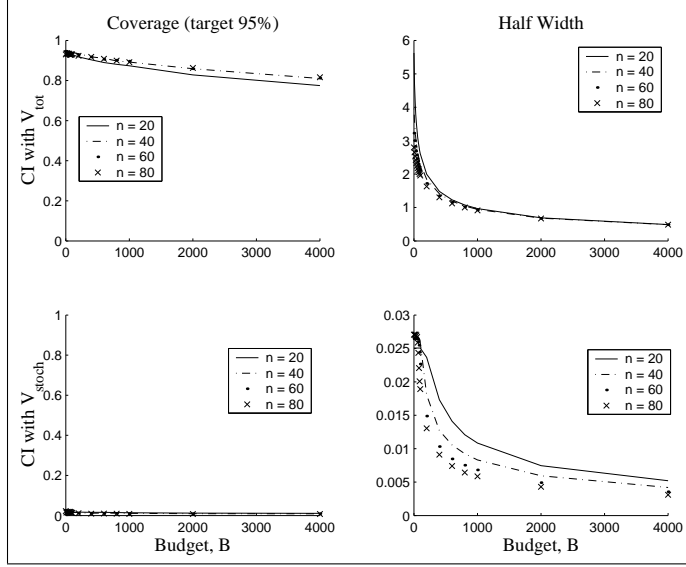


Figure 5: Empirical coverage and half-width of CIs constructed with  $V_{\text{tot}}$  (accounts for input and output uncertainty) and  $V_{\text{stoch}}$  (output uncertainty only).

resulting in asymptotically normal approximations  $\theta_i \sim \text{Normal}(\hat{\theta}_{in_i}, \Sigma_{in_i})$  for each of the  $k = 6$  input parameters. For the unknown arrival rate  $\lambda$  for a Poisson number of arrivals per day, we used a noninformative prior distribution,  $\pi(\lambda) \propto \lambda^{-1/2}$ , resulting in a Gamma( $1/2 + \sum_{\ell=1}^n x_{\ell}, n$ ) posterior distribution after observing  $n$  daily counts of arrivals (Bernardo and Smith 1994). So  $\bar{\lambda} = (-1/2 + \sum_{\ell=1}^n x_{\ell})/n$ . The prior distribution for the unknown routing probabilities is given in Section 3. The prior distributions for the lognormal service time parameters were chosen identical to those for the unknown normal distribution parameters in Section 5.1.

We used a BMA to sample 32 independent sets of input parameters for the critical care simulation. For each set of input parameters, we ran  $q = 4$  independent replications of the critical care unit, for a total of  $r = 128$  replications. Each replication simulated 50 months of operation (after a 10 month warm-up period).

Assuming the linear model in Equation (11), the simulation output implies that  $\tilde{\beta} = (28.9, 1.52, -0.29, 29.4, -0.57)$  and  $\tilde{\sigma}^2 = 1.078$ . This results in  $V_{\text{par}} \approx 98.1$  and approximations for the variance in the performance due to parameter uncertainty of  $V_{\text{resp}} \approx 0.99$ . The uncertainty due to unknown parameters greatly outweighs the response parameter uncertainty, which in turn outweighs the stochastic uncertainty from the random output  $V_{\text{stoch}} \approx \tilde{\sigma}^2/128 = 0.0084$ .

The mean allocations for further data collection suggested by minimizing Equation (13) subject to nonnegative sampling constraints for a variety of sampling budgets are shown in Table 2. The table presumes that a total of  $B$  data points can be collected, that the collection cost for each area is the same ( $c_i = 1$  for  $i = 1, \dots, k$ ), and that each “real” data sample is 4 times as expensive as running one more replication at each of the original 32 design points ( $c_0 = 1/(4 \times 32)$ ). The most important area for additional data collection, based on this analysis, is the arrival rate. Uncertainty about routing probabilities and intermediate ICU service times are of secondary importance. Simulation replications may help reduce performance uncertainty due to response surface uncertainty. Uncertainty about the other input parameters play a much less important role for output uncertainty.

The specific mechanism for estimating the gradient in this example did not significantly affect

1. For  $i = 1, 2, \dots, 1000$  macroreplications
  - (a) Sample  $n$  independent exponential( $\lambda$ ) inter-arrival times and  $n$  independent exponential( $\mu$ ) service times to obtain MAP estimates  $\tilde{\mu}_i, \tilde{\lambda}_i$  (resampling until  $\tilde{\lambda}_i/\tilde{\mu}_i < .995$ ), assuming a noninformative prior distribution  $p(\lambda, \mu) \propto \lambda^{-1}\mu^{-1}$ .
  - (b) Given a sampling budget  $B$  with costs  $c_i = 1$ , compute the optimal sampling allocations  $m_{a,v}^*, m_{s,v}^*$  to minimize the variance approximation  $V_{\text{par}}(\mathbf{m})$  and optimal allocations  $m_{a,\text{mse}}^*, m_{s,\text{mse}}^*$  to minimize the MSE approximation (see Appendix A).
  - (c) Collect additional input samples to obtain updated estimates  $\tilde{\lambda}_{i,v}, \tilde{\mu}_{i,v}$ , and  $y_{i,v} = g(\tilde{\lambda}_{i,v}, \tilde{\mu}_{i,v})$  for the variance allocation, and  $\tilde{\lambda}_{i,\text{mse}}, \tilde{\mu}_{i,\text{mse}}$ , and  $y_{i,\text{mse}} = g(\tilde{\lambda}_{i,\text{mse}}, \tilde{\mu}_{i,\text{mse}})$  for the MSE allocation.
2. Compute the sample mean  $\bar{y}_{n,v}$  to estimate  $g(\lambda, \mu)$ , the sample bias  $\bar{\kappa}_{n,v}$ , and the sample variance  $S_{n,v}^2$  of the output, as well as the average number of additional samples  $\bar{m}_{a,v}^*, \bar{m}_{s,v}^*$  for the variance-based allocation (subscript  $v$ ), and analogous statistics  $\bar{y}_{n,\text{mse}}, \bar{\kappa}_{n,\text{mse}}, S_{n,\text{mse}}^2, \bar{m}_{a,\text{mse}}^*, \bar{m}_{s,\text{mse}}^*$  for the MSE-based allocation (subscript  $\text{mse}$ ).

Figure 6: Algorithm to assess variance, bias, and RMSE for the  $M/M/1$  queue example in Section 5.2 for a given  $\lambda, \mu$  and sampling parameters  $n, B$ .

the data collection allocation that was ultimately suggested. Kleijnen (2001) notes that spurious regressors in parametric regression can deteriorate the predictor, so a parsimonious metamodel is desirable. We repeated the experiment with another technique for response surface modeling due to Raftery, Madigan, and Hoeting (1997) that identifies the most important response parameters (effectively screening out  $\beta_j$  that are relatively close to 0). That analysis determined essentially the same sampling plans in Table 2. We did not check the CI coverage for this example since the true performance is unknown.

## 6 Comments and Conclusions

Stochastic models are useful system design and analysis tools. The mean system performance depends on statistical parameters that describe the model. There are several ways to handle parameter uncertainty. One is to design a system to handle “worst case” scenarios, but this may be cost ineffective. Another way is to develop robust designs that perform well over a range of parameter values. This begs the question of what range to consider. A reduction of parameter uncertainty can reduce the range of parameters for which the design must perform well. When possible, a third way is to learn about parameters as the system operates and adapt the operation of the system to the extent possible as parameters are inferred, given the existing design of a system (e.g. as did Scarf (1959) for the news vendor problem).

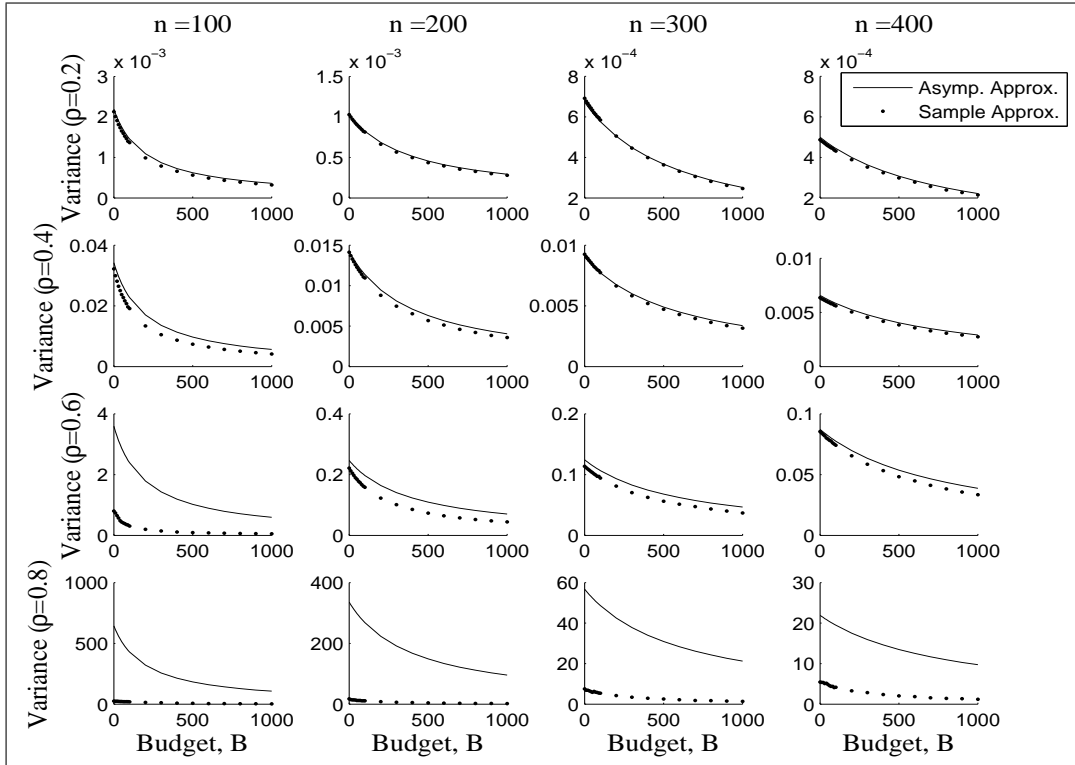


Figure 7: Comparison of mean asymptotic variance approximation  $E_{\mathcal{E}_n}[V_{\text{par}}(\mathbf{m}^*)]$  (“Asymp. Approx.”) with realized sample variance estimate of  $\text{Var}_{\mathcal{E}_n, \mathcal{D}}[g(\hat{\Theta}_{n+\mathbf{m}^*})]$  (“Sample Approx.”).

Another approach, taken by this paper, assumes that a number of parameters may need estimation, and that data sampling is possible to better infer their values before a final design is chosen. Asymptotic approximations simplify the problem, and closed form solutions can be found for a broad class of input distribution models. One advantage of this approach is its generality for input models, and the ability to use a variety of gradient estimation tools, including but not limited to regression with the BMA, a natural tool for exploring input parameter uncertainty. Our approach was Bayesian, but some frequentist results emerge as corollaries.

Several issues warrant further exploration. The asymptotic variance approximation is poor for some systems, including the  $M/M/1$  queue where even the posterior mean of the system performance may not exist. The nonexistence of moments might be avoided by making the model more realistic (assume capacitated queues, run transient rather than steady state simulations), but nonnormal approximations may be of use. Also, some Bayesian representations of parameter uncertainty may have no obvious sampling mechanism for further inference. Nonetheless, the formulas in this paper present an implementable mechanism for guiding data collection plans to reduce parameter uncertainty in a way that in some sense effectively reduces performance uncertainty.

## APPENDIX

### A Reducing MSE with Bias Estimates

The paper principally discusses sampling plans to reduce an approximation to the variance of the unknown performance, based on parameter variance and gradient information. The linear,

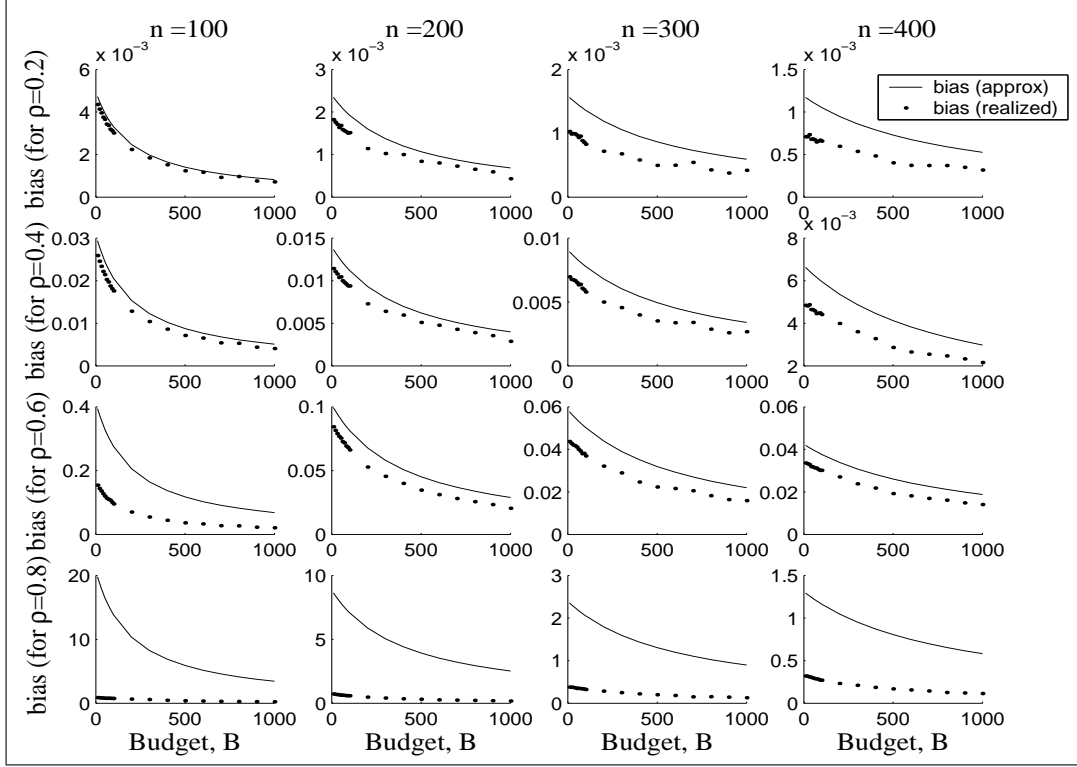


Figure 8: Bias comparison using the optimal allocation for variance, measuring the predicted bias and the empirical bias.

gradient-based approximation of the response function ignores bias introduced by naive estimates of the mean due to any nonlinearity of  $g(\cdot)$ . Bias increases the mean squared error. In particular, if  $\boldsymbol{\theta} = (\vartheta_1, \dots, \vartheta_d)$  is a single parameter for which  $n$  data points have been observed, and  $\boldsymbol{\Sigma}_n$  is as in Equation (3), then the bias  $E[g(\tilde{\boldsymbol{\theta}}_n)] - g(\boldsymbol{\theta})$  is approximately

$$\gamma_n = \frac{1}{2} \sum_{j=1}^d \sum_{l=1}^d \frac{\partial^2 g(\tilde{\boldsymbol{\theta}}_n)}{\partial \vartheta_j \partial \vartheta_l} \text{Cov}(\vartheta_j, \vartheta_l) = \frac{1}{2} \text{tr} \left( \nabla_{\boldsymbol{\theta}}^2 g(\tilde{\boldsymbol{\theta}}_n) \boldsymbol{\Sigma}_n \right), \quad (15)$$

a well-known result obtained with a Taylor series expansion with linear and quadratic terms. The MSE is (asymptotically) approximately  $V_{\text{par}} + \gamma_n^2$ . Informally, as  $n$  increases,  $V_{\text{par}}$  decreases like  $1/n$  and the bias-squared term decreases like  $1/n^2$ .

If  $V_{\text{par}}(\mathbf{m})$  simplifies to Equation (8) with  $k \geq 1$  input parameters then the asymptotic estimate for MSE after  $\mathbf{m}$  samples are collected simplifies.

$$\begin{aligned} \text{MSE}_n(\mathbf{m}) &= V_{\text{par}}(\mathbf{m}) + \gamma_n^2(\mathbf{m}) \\ \gamma_n(\mathbf{m}) &= \frac{1}{2} \left( \text{tr} \sum_{i=1}^k \nabla_{\boldsymbol{\theta}_i}^2 g(\tilde{\boldsymbol{\Theta}}_n) \left( \boldsymbol{\Sigma}_{in_i}^{-1} + (m_i \mathbf{H}_i(\tilde{\boldsymbol{\theta}}_{in_i}))^{-1} \right) \right) \\ &= \sum_{i=1}^k \frac{\zeta_i}{n_{0i} + n_i + m_i} \end{aligned} \quad (16)$$

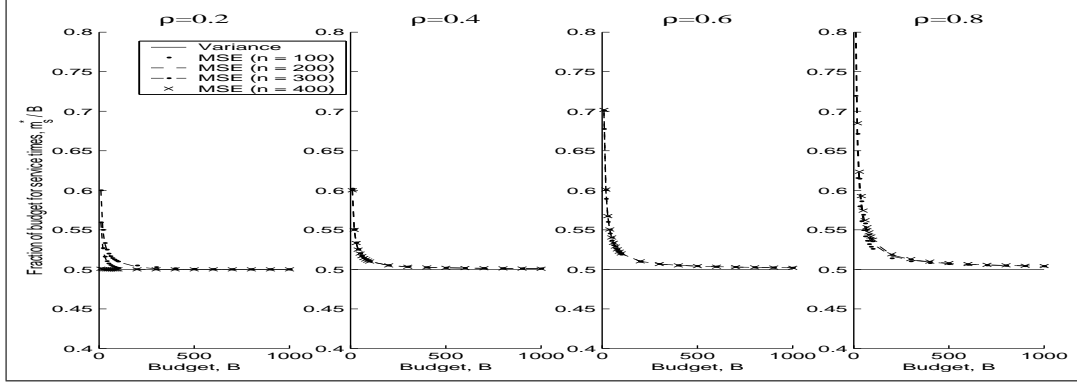


Figure 9: Average fraction of samples  $\bar{m}_s^*/B$  allocated for service times for the variance allocation (always 50%) and the MSE optimal allocation for various  $n$ .

Table 2: Predicted Input Uncertainty  $V_{\text{tot}}(\mathbf{m}, m_0)$  as a Function of Sampling Budget  $B$  and Resulting Optimal Sampling Plan ( $m_1$ =Days of Arrival Data,  $m_2$ =# ICU Service Times,  $m_3$ =# Intermediate ICU Service Times,  $m_4$ =# Intermediate CCU Service Times,  $m_5$ =# CCU Service Times,  $m_6$ =# Routing Decisions,  $q' = m_0/32$ =# Replications per Design Point)

$B$	$m_1$	$m_2$	$m_3$	$m_4$	$m_5$	$m_6$	$q'$	$V_{\text{par}}(\mathbf{m})$	$V_{\text{tot}}(\mathbf{m}, r)$
0	0	0	0	0	0	0	0	98.1	99.1
100	99	0	0	0	0	0	1	53.8	54.3
200	198	0	6	0	0	0	8	38.8	39.2
500	440	0	25	0	0	30	20	23.2	23.3
1000	794	0	78	6	0	115	28	14.0	14.1

where  $\zeta_i = \text{tr}[\nabla_{\theta_i}^2 g(\tilde{\Theta}_n) \mathbf{H}_i^{-1}(\tilde{\theta}_{in_i})]/2$ . For the  $M/M/1$  queue in Section 4.1,

$$\gamma_n(\mathbf{m}) = \frac{\tilde{\lambda}}{(\tilde{\mu} - \tilde{\lambda})^3} \frac{\tilde{\mu}^2}{(\alpha_\mu - 1 + m_s)} + \frac{\tilde{\mu}}{(\tilde{\mu} - \tilde{\lambda})^3} \frac{\tilde{\lambda}^2}{(\alpha_\lambda - 1 + m_a)}.$$

## B Canonical conjugate prior distribution.

This section justifies the claimed simplification of Equation (4) to Equation (5) when (i) the likelihood is in the regular exponential family and has a finite dimensional sufficient statistic, (ii) the prior distribution for the unknown parameter is a canonical conjugate distribution, and (iii) some differentiability conditions hold. This section fixes typographical errors in Bernardo and Smith (1994, Prop. 5.16).

If  $X$  has a likelihood model in the regular exponential family, with  $d$ -dimensional parameter  $\theta = (\vartheta_1, \dots, \vartheta_d)$ , the likelihood function can be written

$$p(x | \theta) = a_1(x)g(\theta) \exp \left[ \sum_{j=1}^d c_j \phi_j(\theta) h_j(x) \right] \quad (17)$$

for some  $a_1(\cdot)$ ,  $g(\cdot)$ ,  $c_j$ ,  $\phi_j(\cdot)$ , and  $h_j(\cdot)$ . The conjugate prior distribution for  $\theta$  with parameter

$n_0, \mathbf{t} = (t_1, \dots, t_d)$  is

$$p(\boldsymbol{\theta} \mid n_0, \mathbf{t}) = [K(n_0, \mathbf{t})]^{-1} g(\boldsymbol{\theta})^{n_0} \exp \left[ \sum_{j=1}^d c_j \phi_j(\boldsymbol{\theta}) t_j \right], \quad (18)$$

where  $K(n_0, \mathbf{t}) = \int g(\boldsymbol{\theta})^{n_0} \exp[\sum_{j=1}^d c_j \phi_j(\boldsymbol{\theta}) t_j] d\boldsymbol{\theta} < \infty$ .

Suppose that the data  $x_1, \dots, x_n$  are observed. Set  $v_{\ell j} = h_j(x_\ell)$  and  $\mathbf{v}_\ell = (v_{\ell 1}, \dots, v_{\ell d})$  with sample average  $\bar{\mathbf{v}}_n = \sum_{\ell=1}^n \mathbf{v}_\ell / n$ . Then  $n\bar{\mathbf{v}}_n$  is a sufficient statistic for  $\boldsymbol{\theta}$  and the posterior distribution of  $\boldsymbol{\theta}$  is  $p(\boldsymbol{\theta} \mid n_0 + n, \mathbf{t} + n\bar{\mathbf{v}}_n)$ .

A final reparametrization is useful. Set  $\psi_j = c_j \phi_j(\boldsymbol{\theta})$  and  $\boldsymbol{\psi} = (\psi_1, \dots, \psi_d)$ . Equation (17) can be rewritten in what is known as canonical conjugate form (Bernardo and Smith 1994, Definition 4.12),

$$p(\mathbf{v}_\ell \mid \boldsymbol{\psi}) = a_2(\mathbf{v}) \exp [\mathbf{v}_\ell \boldsymbol{\psi}^\top - b(\boldsymbol{\psi})]$$

for some real-valued  $a_2(\mathbf{v})$ ,  $b(\boldsymbol{\psi})$ . The canonical conjugate prior distribution is

$$p(\boldsymbol{\psi} \mid n_0, \mathbf{t}) = c(n_0, \mathbf{t}) \exp [n_0 \mathbf{t} \boldsymbol{\psi}^\top - n_0 b(\boldsymbol{\psi})],$$

where  $n_0, \mathbf{t} = (t_1, \dots, t_d)$  are parameters of the prior distribution.

Bernardo and Smith (1994, Prop. 5.16) note that the posterior distribution has the same functional form as the prior distribution,

$$p(\boldsymbol{\psi} \mid n_0, \mathbf{t}, \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n) = p(\boldsymbol{\psi} \mid n_0 + n, n_0 \mathbf{t} + n\bar{\mathbf{v}}_n).$$

Further, if  $\tilde{\boldsymbol{\psi}}$  is the MAP estimator of  $\boldsymbol{\psi}$ , then  $(\boldsymbol{\Sigma}_n^{-1})_{jl} = (n_0 + n) \frac{\partial^2 b(\boldsymbol{\psi})}{\partial \psi_j \partial \psi_l} \Big|_{\boldsymbol{\psi}=\tilde{\boldsymbol{\psi}}}$ . The claim is proven by noting that the expected information in one observation is

$$(\mathbf{H}(\boldsymbol{\psi}))_{jl} = \mathbb{E} \left[ \frac{\partial^2 b(\boldsymbol{\psi})}{\partial \psi_j \partial \psi_l} \right] \Big|_{\tilde{\boldsymbol{\psi}}} = \frac{\partial^2 b(\boldsymbol{\psi})}{\partial \psi_j \partial \psi_l} \Big|_{\tilde{\boldsymbol{\psi}}}.$$

The relationship also holds in the  $\boldsymbol{\theta}$  coordinate system if the map from  $\boldsymbol{\psi}$  to  $\boldsymbol{\theta}$  is bijective in a neighborhood of  $\tilde{\boldsymbol{\psi}}$ .

## C Proof of Proposition 4.2.

Using the method of Lagrange multipliers, we obtain the Lagrangian function

$$\begin{aligned} L(\mathbf{m}, \tau) &= V_{\text{par}}(\mathbf{m}) - \sum_{i=1}^k \tau (c_i m_i - B) \\ &= \sum_{i=1}^k \frac{\xi_i}{n_{0i} + n_i + m_i} - \sum_{i=1}^k \tau (c_i m_i - B), \end{aligned}$$

where  $\tau$  is the Lagrange multiplier. Set the  $(k+1)$  partial derivatives to 0 to get

$$-\frac{\xi_i}{(n_{0i} + n_i + m_i)^2} - \tau c_i = 0, \text{ for } i = 1, \dots, k, \text{ and} \quad (19)$$

$$\sum_{i=1}^k c_i m_i - B = 0. \quad (20)$$

Equation (19) implies  $\frac{\xi_i/c_i}{(n_{0i}+n_i+m_i)^2} = \frac{\xi_j/c_j}{(n_{0j}+n_j+m_j)^2}$ . After some algebra we get

$$c_j(n_{0j} + n_j + m_j) = \left( \frac{\xi_j c_i c_j}{\xi_i} \right)^{1/2} (n_{0i} + n_i + m_i). \quad (21)$$

Equation (20) can be rewritten  $\sum_{j=1}^k c_j(n_{0j} + n_j + m_j) = B + \sum_{\ell=1}^k c_\ell(n_{0\ell} + n_\ell)$ . Substitute Equation (21) into this expression to obtain

$$(n_{0i} + n_i + m_i) \sum_{j=1}^k \left( \frac{\xi_j c_i c_j}{\xi_i} \right)^{1/2} = B + \sum_{\ell=1}^k c_\ell(n_{0\ell} + n_\ell). \quad (22)$$

Algebraic rearranging gives the claimed result for the optimal  $m_i^*$ ,

$$m_i^* = \frac{B + \sum_{\ell=1}^k c_\ell(n_{0\ell} + n_\ell)}{\sum_{j=1}^k \left( \frac{\xi_j c_i c_j}{\xi_i} \right)^{1/2}} - (n_{0i} + n_i).$$

When  $B$  is sufficiently large, the constraint  $m_i^* \geq 0$  is satisfied as desired.

**Acknowledgments:** We thank Jim Wilson for thorough and helpful feedback.

## Bibliography

- Barton, R. R. and L. W. Schruben (2001). Simulating real systems. in submission.
- Bernardo, J. M. and A. F. M. Smith (1994). *Bayesian Theory*. Chichester, UK: Wiley.
- Box, G. and N. Draper (1987). *Empirical Model-Building and Response Surfaces*. New York: Wiley.
- Box, G., W. Hunter, and J. Hunter (1978). *Statistics for Experimenters. An Introduction to Design, Data Analysis and Model Building*. New York: Wiley.
- Buzacott, J. A. and J. G. Shanthikumar (1993). *Stochastic Models of Manufacturing Systems*. Englewood Cliffs, New Jersey: Prentice-Hall, Inc.
- Cheng, R. C. H. and W. Holland (1997). Sensitivity of computer simulation experiments to errors in input data. *J. Statistical Computing and Simulation* 57, 219–241.
- Cheng, R. C. H. and W. Holland (2004). Calculation of confidence intervals for simulation output. *ACM Transactions on Modeling and Computer Simulation* 14(4), 344–362.
- Chick, S. E. (1997). Bayesian analysis for simulation input and output. In S. Andradóttir, K. Healy, D. Withers, and B. Nelson (Eds.), *Proc. 1997 Winter Simulation Conference*, Piscataway, NJ, pp. 253–260. IEEE, Inc.
- Chick, S. E. (2001). Input distribution selection for simulation experiments: Accounting for input uncertainty. *Operations Research* 49(5), 744–758.
- Draper, D. (1995). Assessment and propagation of model uncertainty (with discussion). *J. Royal Statistical Society, Series B* 57(1), 45–97.
- Edwards, A. W. F. (1984). *Likelihood*. Cambridge University Press.
- Gelman, A., J. B. Carlin, H. S. Stern, and D. R. Rubin (1995). *Bayesian Data Analysis*. London: Chapman & Hall.

- Kleijnen, J. P. C. (2001). Comments on M.C. Kennedy & A. O'Hagan's 'Bayesian calibration of computer models'. *Journal Royal Statistical Society, Series B* 63(3), in press.
- Law, A. M. and W. D. Kelton (2000). *Simulation Modeling & Analysis* (3rd ed.). New York: McGraw-Hill, Inc.
- Mendoza, M. (1994). Asymptotic normality under transformations. A result with Bayesian applications. *TEST* 3(2), 173–180.
- Myers, R., A. Khuri, and W. Carter (1989). Response surface methodology: 1966-1988. *Technometrics* 31(2), 137–157.
- Ng, S.-H. and S. E. Chick (2001). Reducing input distribution uncertainty for simulations. In B. Peters, J. Smith, M. Rohrer, and D. Madeiros (Eds.), *Proc. 2001 Winter Simulation Conference*, Piscataway, NJ, pp. 364–371. IEEE, Inc.
- Raftery, A. E., D. Madigan, and J. A. Hoeting (1997). Bayesian model averaging for linear regression models. *Journal of the American Statistical Association* 92(437), 179–191.
- Scarf, H. (1959). Bayes solutions of the statistical inventory problem. *Annals of Mathematical Statistics* 30, 490–508.
- Schervish, M. J. (1995). *Theory of Statistics*. New York: Springer.
- Schruben, L. W. and B. H. Margolin (1978). Pseudorandom number assignment in statistically designed simulation and distribution sampling experiments. *Journal of the American Statistical Association* 73(363), 504–525.
- Serfling, R. J. (1980). *Approximation Theorems of Mathematical Statistics*. New York: Wiley.
- Zouaoui, F. and J. R. Wilson (2001). Accounting for input model and parameter uncertainty in simulation input modeling. In B. Peters, J. Smith, M. Rohrer, and D. Madeiros (Eds.), *Proc. 2001 Winter Simulation Conference*, Piscataway, NJ, pp. 290–299. IEEE, Inc.
- Zouaoui, F. and J. R. Wilson (2003). Accounting for parameter uncertainty in simulation input modeling. *IIE Transactions* 35, 781–792.
- Zouaoui, F. and J. R. Wilson (2004). Accounting for input-model and input-parameter uncertainties in simulation. *IIE Transactions* 36(11), 1135–1151.