# **Computer Experiments**

J. R. Koehler and A. B. Owen

#### 1. Introduction

Deterministic computer simulations of physical phenomena are becoming widely used in science and engineering. Computers are used to describe the flow of air over an airplane wing, combustion of gases in a flame, behavior of a metal structure under stress, safety of a nuclear reactor, and so on.

9

Some of the most widely used computer models, and the ones that lead us to work in this area, arise in the design of the semiconductors used in the computers themselves. A process simulator starts with a data structure representing an unprocessed piece of silicon and simulates the steps such as oxidation, etching and ion injection that produce a semiconductor device such as a transistor. A device simulator takes a description of such a device and simulates the flow of current through it under varying conditions to determine properties of the device such as its switching speed and the critical voltage at which it switches. A circuit simulator takes a list of devices and the way that they are arranged together and computes properties of the circuit as a whole.

In each of these computer simulations, the user must specify the values of some governing variables. For example in process simulation the user might have to specify the duration and temperature of the oxidation step, and doses and energies for each of the ion implantation steps. These are continuously valued variables. There may also be discrete variables, such as whether to use wet or dry oxidation. Most of this chapter treats the case of continuous variables, but some of it is easily adaptable to discrete variables, especially those taking only two values.

Let  $X \in \mathbb{R}^p$  denote the vector of input values chosen for the computer program. We will write X as the row vector  $(X^1, \ldots, X^p)$  using superscripts to denote components of X. We assume that each component  $X^j$  is continuously adjustable between a lower and an upper limit, which after a linear transformation can be taken to be 0 and 1 respectively. (For some results where every input is dichotomous see Mitchel et al. (1990).) The computer program is denoted by f and it computes q output quantities, denoted by  $Y \in \mathbb{R}^q$ .

$$Y = f(X), \quad X \in [0,1]^p.$$
 (1)

Some important quantities describing a computer model are the number of inputs p, the number of outputs q and the speed with which f can be computed. These vary

enormously in applications. In the semiconductor problems we have considered p is usually between 4 and 10. Other computer experiments use scores or even hundreds of input variables. In our motivating applications q is usually larger than 1. For example interest might center on the switching speed of a device and also on its stability as measured by a breakdown voltage. For some problems f takes hours to evaluate on a supercomputer and for others f runs in milliseconds on a personal computer.

Equation (1) differs from the usual X - Y relationship studied by statisticians in that there is no random error term. If the program is run twice with the same X, the same Y is obtained both times. Therefore it is worth discussing why a statistical approach is called for.

These computer programs are written to calculate Y from a known value of X. The way they are often used however, is to search for good values of X according to some goals for Y. Suppose that  $X_1 = (X_1^1, \ldots, X_1^p)$  is the initial choice for X. Often  $X_1$  does not give a desirable  $Y_1 = f(X_1)$ . The engineer or scientist can often deduce why this is, from the program output, and select a new value,  $X_2$  for which  $Y_2 = f(X_2)$  is likely to be an improvement. This improvement process can be repeated until a satisfactory design is found. The disadvantage of this procedure is that it may easily miss some good designs X, because it does not fully explore the design space. It can also be slow, especially when p is large, or when improvements of  $Y^1$  say, tend to appear with worsenings of  $Y^2$  and vice versa.

A commonly used way of exploring the design space around  $X_1$  is to vary each of the  $X_1^j$  one at a time. As is well known to statisticians, this approach can be misleading if there are strong interactions among the components of X. Increasing  $X^1$  may be an improvement and increasing  $X^2$  may be an improvement, but increasing them both together might make things worse. This would usually be determined from a confirmation run in which both  $X^1$  and  $X^2$  have been increased. The greater difficulty with interactions stems from missed opportunities: the best combination might be to increase  $X^1$  while decreasing  $X^2$ , but one at a time experimentation might never lead the user to try this. Thus techniques from experimental design may be expected to help in exploring the input space.

This chapter presents and compares two statistical approaches to computer experiments. Randomness is required in order to generate probability or confidence intervals. The first approach introduces randomness by modeling the function f as a realization of a Gaussian process. The second approach does so by taking random input points (with some balance properties).

# 2. Goals in computer experiments

There are many different but related goals that arise in computer experiments. The problem described in the previous section is that of finding a good value for X according to some criterion on Y. Here are some other goals in computer experimentation: finding a simple approximation  $\hat{f}$  that is accurate enough over a region A of X values, estimating the size of the error  $\hat{f}(X_0) - f(X_0)$  for some  $X_0 \in A$ , estimating  $\int_A f \, dX$ , sensitivity analysis of Y with respect to changes in X, finding which  $X^j$  are most important for each response  $Y^k$ , finding which competing goals for Y conflict the most, visualizing the function f and uncovering bugs in the implementation of f.

# 2.1. Optimization

Many engineering design problems take the form of optimizing  $Y^1$  over allowable values of X. The problem may be to find the fastest chip, or the least expensive soda can. There is often, perhaps usually, some additional constraint on another response  $Y^2$ . The chip should be stable enough, and the can should be able to withstand a specified internal pressure.

Standard optimization methods, such as quasi-Newton or conjugate gradients (see for example Gill et al., 1981) can be unsatisfactory for computer experiments. These methods usually require first and possibly second derivatives of f, and these may be difficult to obtain or expensive to run. The standard methods also depend strongly on having good starting values. Computer experimentation as described below is useful in the early stages of optimization where one is searching for a suitable starting value. It is also useful when searching for several widely separated regions of the predictor space that might all have good Y values. Given a good starting value, the standard methods will be superior if one needs to locate the optimum precisely.

### 2.2. Visualization

As Diaconis (1988) points out, being able to compute a function f at any given value X does not necessarily imply that one "understands" the function. One might not know whether the function is continuous or bounded or unimodal, where its optimum is or whether it has asymptotes.

Computer experimentation can serve as a primitive way to visualize functions. One evaluates f at a well chosen set of points  $x_1, \ldots, x_n$  obtaining responses  $y_1, \ldots, y_n$ . Then data visualization methods may be applied to the p + q dimensional points  $(x_i, y_i)$ ,  $i = 1, \ldots, n$ . Plotting the responses versus the input variables (there are pq such plots) identifies strong dependencies, and plotting residuals from a fit can show weaker dependencies. Selecting the points with desirable values of Y and then producing histograms and plots of the corresponding X values can be used to identify the most promising subregion of X values. Sharifzadeh et al. (1989) took this approach to find that increasing a certain implant dose helped to make two different threshold voltages near their common targets and nearly equal (as they should have been). Similar exploration can identify which input combinations are likely to crash the simulator.

Roosen (1995) has used computer experiment designs for the purpose of visualizing functions fit to data.

# 2.3. Approximation

The original program f may be exceedingly expensive to evaluate. It may however be possible to approximate f by some very simple function  $\hat{f}$ , the approximation holding adequately in a region of interest, though not necessarily over the whole domain of f. If the function  $\hat{f}$  is fast to evaluate, as for instance a polynomial, neural network or a MARS model (see Friedman, 1991), then it may be feasible to make millions of  $\hat{f}$ 

evaluations. This makes possible brute force approximations for the other problems. For example, optimization could be approached by finding the best value of  $\hat{f}(x)$  over a million random runs x.

Approximation by computer experiments involves choosing where to gather  $(x_i, f(x_i))$  pairs, how to construct an approximation based on them and how to assess the accuracy of this approximation.

# 2.4. Integration

Suppose that  $X^*$  is the target value of the input vector, but in the system being modeled the actual value of X will be random with a distribution dF that hopefully is concentrated near  $X^*$ . Then one is naturally interested in  $\int f(X) dF$ , the average value of Y over this distribution. Similarly the variance of Y and the probability that Y exceeds some threshold can be expressed in terms of integrals. This sort of calculation is of interest to researchers studying nuclear safety. McKay (1995) surveys this literature.

Integration and optimization goals can appear together in the same problem. In robust design problems (Phadke, 1988), one might seek the value  $X_0$  that minimizes the variance of Y as X varies randomly in a neighborhood of  $X_0$ .

# 3. Approaches to computer experiments

There are two main statistical approaches to computer experiments, one based on Bayesian statistics and a frequentist one based on sampling techniques. It seems to be essential to introduce randomness in one or other of these ways, especially for the problem of gauging how much an estimate  $\hat{f}(X_0)$  might differ from the true value  $f(X_0)$ .

In the Bayesian framework, surveyed below in Sections 4 and 5, f is a realization of a random process. One sets a prior distribution on the space of all functions from  $[0,1]^p$  to  $\mathbb{R}^q$ . Given the values  $y_i = f(x_i)$ ,  $i = 1, \ldots, n$ , one forms a posterior distribution on f or at least on certain aspects of it such as  $f(x_0)$ . This approach is extremely elegant. The prior distribution is usually taken to be Gaussian so that any finite list of function values has a multivariate normal distribution. Then the posterior distribution, given observed function values is also multivariate normal. The posterior mean interpolates the observed values and the posterior variance may be used to give 95% posterior probability intervals. The method extends naturally to incorporate measurement and prediction of derivatives, partial derivatives and definite integrals of f.

The Bayesian framework is well developed as evidenced by all the work cited below in Sections 4 and 5. But, as is common with Bayesian methods there may be difficulty in finding an appropriate prior distribution. The simulator output might not have as many derivatives as the underlying physical reality, and assuming too much smoothness for the function can lead to Gibbs-effect overshoots. A numerical difficulty also arises: the Bayesian approach requires solving n linear equations in

*n* unknowns when there are *n* data points. The effort involved grows as  $n^3$  while the effort in computing  $f(X_1), \ldots, f(X_n)$  grows proportionally to *n*. Inevitably this limits the size of problems that can be addressed. For example, suppose that one spends an hour computing  $f(x_1), \ldots, f(x_n)$  and then one minute solving the linear equations. If one then finds it necessary to run 24 times as many function evaluations, the time to compute the  $f(x_i)$  grows from an hour to a day, while the time to solve the linear equations grows from one minute to over nine and a half days.

These difficulties with the Bayesian approach motivate a search for an alternative. The frequentist approach, surveyed in Sections 6 and 7, introduces randomness by taking function values  $x_1, \ldots, x_n$  that are partially determined by pseudo-random number generators. Then this randomness in the  $x_i$  is propagated through to randomness in  $\hat{f}(x_0)$ . This approach allows one to consider f to be deterministic, and in particular to avoid having to specify a distribution for f. The material given there expands on a proposal of Owen (1992a). There is still much more to be done.

#### 4. Bayesian prediction and inference

A Bayesian approach to modeling simulator output (Sacks et al., 1989a, b; Welch et al., 1990) can be based on a spatial model adapted from the geo-statistical Kriging model (Matheron, 1963; Journel and Huibregts, 1978; Cressie, 1986, 1993; Ripley, 1981). This approach treats the bias, or systematic departure of the response surface from a linear model, as the realization of a stationary random function. This model has exact predictions at the observed responses and predicts with increasing error variance as the prediction point moves away from all the design points.

This section introduces the Kriging (or Bayesian) approach to modeling the response surfaces of computer experiments. Several correlation families are discussed as well as their effect on prediction and error analysis. Additionally, extensions to this model are presented that allow the use and the modeling of gradient information.

# 4.1. The Kriging model

The Kriging approach uses a two component model. The first component consists of a general linear model while the second (or lack of fit) component is treated as the realization of a stationary Gaussian random function. Define  $S = [0, 1]^p$  to be the design space and let  $x \in S$  be a scaled *p*-dimensional vector of input values. The Kriging approach models the associated response as

$$Y(x) = \sum_{j=1}^{k} \beta_j h_j(x) + Z(x)$$
(2)

where the  $h_j$ 's are known fixed functions, the  $\beta_j$ 's are unknown coefficients to be estimated and Z(x) is a stationary Gaussian random function with E[Z(x)] = 0 and covariance

$$\operatorname{Cov}[Z(x_i), Z(x_j)] = \sigma^2 R(x_j - x_i).$$
(3)

For any point  $x \in S$ , the simulator output Y(x) at that point has the Gaussian distribution with mean  $\sum \beta_j h_j(x)$  and variance  $\sigma^2$ . The linear component models the drift in the response, while the systematic lack-of-fit (or bias) is modeled by the second component. The smoothness and other properties of  $Z(\cdot)$  are controlled by  $R(\cdot)$ .

Let design  $D = \{x_i, i = 1, ..., n\} \subset S$  yield responses  $y'_D = \{y(x_1), ..., y(x_n)\}$ and consider a linear predictor

$$\widehat{y}(x_0) = \lambda'(x_0) y_D$$

of an unobserved point  $x_0$ . The Kriging approach of Matheron (1963) treats  $\hat{y}(x_0)$  as a random variable by substituting  $Y_D$  for  $y_D$  where

$$Y'_D = (Y(x_1), \ldots, Y(x_n)).$$

The best linear unbiased predictor (BLUP) finds the  $\lambda(x_0)$  that minimizes

$$MSE[\widehat{Y}(x_0)] = E[\lambda' Y_D - Y(x_0)]^2$$

subject to the unbiasedness condition

 $\mathbf{E}[\lambda' Y_D] = \mathbf{E}[Y(x_0)].$ 

The BLUP of  $Y(x_0)$  is given by

$$\widehat{Y}(x_0) = h'(x_0)\widehat{\beta} + v'_{x_0}V_D^{-1}(Y_D - H_D\widehat{\beta})$$
(4)

where

$$\begin{aligned} h'(x_0) &= (h_1(x_0), \dots, h_k(x_0)), \\ (H_D)_{ij} &= h_j(x_i), \\ (V_D)_{ij} &= \operatorname{Cov}[Z(x_i), Z(x_j)], \\ v'_{x_0} &= (\operatorname{Cov}[Z(x_0), Z(x_1)], \dots, \operatorname{Cov}[Z(x_0), Z(x_n)]) \end{aligned}$$

and

$$\widehat{\beta} = [H'V^{-1}H]^{-1}H'V^{-1}Y_D$$

is the generalized least squares estimate of  $\beta$ . The mean square error of  $\widehat{Y}(x_0)$  is

$$MSE[\widehat{Y}(x_0)] = \sigma^2 - (h'(x_0), v'_{x_0}) \begin{pmatrix} 0 & H'_D \\ H_D & V_D \end{pmatrix}^{-1} \begin{pmatrix} h(x_0) \\ v_{x_0} \end{pmatrix}$$

The first component of equation (4) is the generalized least squares prediction at point  $x_0$  given the design covariance matrix  $V_D$ , while the second component

266



Fig. 1. A prediction example with n = 3.

"pulls" the generalized least squares response surface through the observed data points. The elasticity of the response surface "pull" is solely determined by the correlation function  $R(\cdot)$ . The predictions at the design points are exactly the corresponding observations, and the mean square error equals zero. As a prediction point  $x_0$  moves away from all of the design points, the second component of equation (4) goes to zero, yielding the generalized least squares prediction, while the mean square error at that point goes to  $\sigma^2 + h'(x_0) \left[H'V_D^{-1}H\right]^{-1} h(x_0)$ . In fact, these results are true in the wide sense if the Gaussian assumption is removed.

As an example, consider an experiment where n = 3, p = 1,  $\sigma^2 = .05$ , R(d) = $\exp(-20d^2)$  and  $D = \{.3, .5, .8\}$ . The response of the unknown function at the design is  $y'_D = (.7, .3, .5)$ . The dashed line of Figure 1 is the generalized least squares prediction surface for  $h(\cdot) \equiv 1$  where  $\hat{\beta} = .524$ . The effect of the second component of equation (4) is to pull the dashed line through the observed design points as shown by the solid line. The shape of the surface or the amount of elasticity of the "pull" is determined by the vector  $v'_x V_D^{-1}$  as a function of x and therefore is completely determined by  $R(\cdot)$ . The dotted lines are  $\pm 2\sqrt{MSE[\widehat{Y}(x)]}$  or 95% pointwise confidence envelopes around the prediction surface. The interpretation of these pointwise confidence envelopes is that for any point  $x_0$ , if the unknown function is truly generated by a random function with constant mean and correlation function  $R(d) = \exp(-20d^2)$ , then approximately 95% of the sample paths that go through the observed design points would be between these dotted lines at  $x_0$ . The predictions and confidence intervals can be very different for different  $\sigma^2$  and  $R(\cdot)$ . The effect of different correlation functions is discussed in Section 4.3. Clearly, the true function is not "generated" stochastically. The above model is used for prediction and to quantify the uncertainty of the prediction. This naturally leads to a Bayesian interpretation of this methodology.

#### 4.2. A fully Bayesian interpretation

An alternative to the above interpretation of equation (2) is the fully Bayesian interpretation which uses the model as a way of quantifying the uncertainty of the unknown function. The Bayesian approach (Currin et al., 1991; O'Hagan, 1989) uses the same model but has a different interpretation of the  $\beta_j$ 's. Here the  $\beta_j$ 's are random variables with prior distribution  $\pi_j$ . The effect of these prior distributions is to quantify the prior belief of the unknown function or to put a prior distribution on a large class of functions  $\mathcal{G}$ . Hence hopefully the true function  $y(\cdot) \in \mathcal{G}$ . The mixed convolution of the  $\pi_i$ 's and  $\pi(Z)$  yield the prior distribution  $\Pi(G)$  for subsets of functions  $G \subset \mathcal{G}$ .

Once the data  $Y_D = y_D$  has been observed, the posterior distribution  $\Pi(G \mid Y_D)$  is calculated. The mean

$$\widehat{Y}(x_0) = \int g(x_0) arPi(g \mid Y_D = y_D) \, \mathrm{d}g$$

and variance

$$\operatorname{Var}(\widehat{Y}(x_0) \mid Y_D = y_D) = \int (g(x_0) - \widehat{Y}(x_0))^2 \Pi(g \mid Y_D = y_D) \, \mathrm{d}g$$

of the posterior distribution at each input point are then used as the predictor and measure of error, respectively, at that point. In general, the Kriging and Bayesian approaches will lead to different estimators. However, if the prior distribution of  $Z(\cdot)$  is Gaussian and if the prior distribution of the  $\beta_j$ 's is diffuse, then the two approaches yield identical estimators.

As an example, consider the case where the prior distribution of the vector of  $\beta$ 's is

$$\beta \sim N_k(b, \tau^2 \Sigma)$$

and the prior distribution of  $Z(\cdot)$  is a stationary Gaussian distribution with expected value zero and covariance function given by equation (3). After the simulator function has been evaluated at the experimental design, the posterior distribution of  $\beta$  is

$$\beta \mid Y_D \sim N_k(\widetilde{\beta}, \widetilde{\Sigma})$$

where

$$\widetilde{\boldsymbol{\beta}} = \widetilde{\boldsymbol{\Sigma}} \left[ \boldsymbol{H}' \boldsymbol{V}_D^{-1} \boldsymbol{Y}_D + \boldsymbol{\tau}^{-2} \boldsymbol{\boldsymbol{\Sigma}}^{-1} \boldsymbol{b} \right]$$

and

$$\widetilde{\Sigma} = \left[ H' V_D^{-1} H + \tau^{-2} \Sigma^{-1} \right]^{-1}$$

and the posterior distribution of  $Y(x_0)$  is

$$Y(x_0) \mid Y_D \sim N \left( v'_{x_0} V_D^{-1} Y_D + c'_{x_0} \widetilde{eta}, \sigma^2 - v'_{x_0} V_D^{-1} v_{x_0} + c'_{x_0} \widetilde{\Sigma} c_{x_0} 
ight)$$

where

$$c_{x_0}' = h' - v_{x_0}' V_D^{-1} H.$$

Hence the posterior distribution is still Gaussian but it is no longer stationary. Now if  $\tau^2 \rightarrow \infty$  then

$$\begin{split} & \widetilde{\beta} \to \widehat{\beta}, \\ & \widetilde{\Sigma} \to \left[ H' V_D^{-1} H \right]^{-1} \end{split}$$

and hence the posterior variance of  $Y(x_0)$  is

$$\begin{aligned} \operatorname{Var}(Y(x_0) \mid Y_D) &= \sigma^2 - v_{x_0}' V_D^{-1} v_{x_0} + c_{x_0}' \left[ H' V_D^{-1} H \right]^{-1} c_{x_0} \\ &= \sigma^2 - v_{x_0}' V_D^{-1} v_{x_0} + h' \left[ H' V_D^{-1} H \right]^{-1} h \\ &- 2h' \left[ H' V_D^{-1} H \right]^{-1} H' V_D^{-1} v_{x_0} \\ &+ v_{x_0}' V_D^{-1} H \left[ H' V_D^{-1} H \right]^{-1} H' V_D^{-1} v_{x_0} \\ &= \sigma^2 - \left[ -h' \left[ H' V_D^{-1} H \right]^{-1} h \\ &+ 2h' \left[ H' V_D^{-1} H \right]^{-1} H' V_D^{-1} v_{x_0} \right] \\ &- \left[ v_{x_0}' \left( V_D^{-1} - V_D^{-1} H \left[ H' V_D^{-1} H \right]^{-1} H' V_D^{-1} \right) v_{x_0} \\ &= \sigma^2 - \left( h'(x_0), v_{x_0}' \right) \begin{pmatrix} 0 & H'_D \\ H_D & V_D \end{pmatrix}^{-1} \begin{pmatrix} h(x_0) \\ v_{x_0} \end{pmatrix} \end{aligned}$$

which is the same variance as the BLUP in the Kriging approach. Therefore, if  $Z(\cdot)$  has a Gaussian prior distribution and if the  $\beta$ 's have a diffuse prior, the Bayesian and the Kriging approaches yield identical estimators.

Currin et al. (1991) provide a more in depth discussion of the Bayesian approach for the model with a fixed mean  $(h \equiv 1)$ . O'Hagan (1989) discusses Bayes Linear Estimators (BLE) and their connection to equations (2) and (4). The Bayesian approach, which uses random functions as a method of quantifying the uncertainty of the unknown simulator function  $Y(\cdot)$ , is more subjective than the Kriging or frequentist approach. While both approaches require prior knowledge or an objective method of estimating the covariance function, the Bayesian approach additionally requires knowledge of parameters of the prior distribution of  $\beta$  (b and  $\Sigma$ ). For this reason, the Kriging results and Bayesian approach with diffuse prior distributions and the Gaussian assumption are widely used in computer experiments.



Fig. 2. The effects of  $\theta$  on prediction.

#### 4.3. Correlation functions

As discussed above, the selection of  $R(\cdot)$  plays a crucial role in constructing designs and in the predictive process. Consider the example of Section 4.1 where n = 3, p = 1,  $D = \{.3, .5, .8\}$ ,  $y'_d = \{.7, .3, .5\}$ ,  $R(d) = \exp\{-\theta d^2\}$  and  $\theta = 20$ . Figure 2(a) shows the effect on prediction for  $\theta = 2$ . Now  $\hat{\beta} = 1.3$  and the surface elasticity is very low. The predictions outside of the design are actually higher than the observed surface since the convex nature of the observed response indicate that the design range contains a local minimum for the total process. Eventually, the extrapolations would return to the value of  $\hat{\beta}$ . Additionally, the 95% pointwise confidence intervals are much narrower within the range of the design than in Figure 1. Figure 2(b) displays the prediction line is typically .5 with smooth curves pulling the surface through the design points. The 95% pointwise confidence intervals are wider than before.

This section presents some simplifying restrictions on  $R(\cdot)$  and four families of univariate correlation functions used in generating the simplified correlation functions. Examples of realization of these families will be shown to explain the effect on prediction by varying the parameter of these families. Furthermore, the maximum likelihood method for estimating the parameters of a correlation family along with a technique for implementation will be discussed in Section 4.4.

#### 4.3.1. Restrictions on $R(\cdot)$

Any positive definite function R with R(x, x) = 1 could be used as a correlation function, but for simplicity, it is common to restrict  $R(\cdot)$  such that for any  $x_1, x_2 \in S$ 

$$R(x_1, x_2) = R(x_1 - x_2)$$

so that the process  $Z(\cdot)$  is stationary. Some types of nonstationary behavior in the mean function of  $Y(\cdot)$  can be modeled by the linear term in equation (2). A further restriction makes the correlation function depend only on the magnitude of the distance.

$$R(x_1, x_2) = R(|x_1 - x_2|).$$

In higher dimensions  $(p \ge 2)$  a product correlation function,

$$R(x_1, x_2) = \prod_{j=1}^p R_j(|x_{1j} - x_{2j}|)$$

is often used for mathematical convenience. That is,  $R(\cdot)$  is a product of univariate correlation functions and, hence, only univariate correlation functions are of interest. The product correlation function has been used for prediction in spatial settings (Ylvisaher, 1975; Curin et al., 1991; Sacks et al., 1989a, b; Welch et al., 1990, 1992).

Several choices for the factors in the product correlation function are outlined below.



Fig. 3. Realizations for the cubic correlation function  $(\rho, \gamma) = (a) (.15, .03)$ , (b) (.45, .20), (c) (.70, .50), and (d) (.95, .90).

# 4.3.2. Cubic

The (univariate) cubic correlation family is parameterized by  $\rho \in [0, 1]$  and  $\gamma \in [0, 1]$ and is given for  $d \in [0, 1]$  by

$$R(d) = 1 - \frac{3(1-\rho)}{2+\gamma} d^2 + \frac{(1-\rho)(1-\gamma)}{2+\gamma} |d|^3$$

where  $\rho$  and  $\gamma$  are restricted by

$$\rho \geqslant \frac{5\gamma^2 + 8\gamma - 1}{\gamma^2 + 4\gamma + 7}$$

to ensure that the function is positive definite (see Mitchell et al., 1990). Here  $\rho = \operatorname{corr}(Y(0), Y(1))$  is the correlation between endpoint observations and  $\gamma = \operatorname{corr}(Y'(0), Y'(1))$  is the correlation between endpoints of the derivative process. The cubic correlation function implies that the derivative process has a linear correlation process with parameter  $\gamma$ .

A prediction model in one dimension for this family is a cubic spline interpolator. In two dimensions, when the correlation is a product of univariate cubic correlation functions the predictions are piece-wise cubic in each variable.

Processes generated with the cubic correlation function are once mean square differentiable. Figure 3 shows several realizations of processes with the cubic correlation function and parameter pairs (.15, .03), (.45, .20), (.70, .50), (.95, .9). Notice that the realizations are quite smooth and almost linear for parameter pair (.95, .90).

### 4.3.3. Exponential

The (univariate) exponential correlation family is parameterized by  $\theta \in (0, \infty)$  and is given by

$$R(d) = \exp(-\theta |d|)$$

for  $d \in [0, 1]$ . Processes with the exponential correlation function are Ornstein–Uhlenbeck processes (Parzen, 1962). The exponential correlation function is not mean square differentiable.

Figure 4 presents several realizations of one dimensional processes with the exponential correlation function and  $\theta = 0.5$ , 2.0, 5.0, 20. Figure 4(a) is for  $\theta = 0.5$  and these realizations have very small global trends but much local variation. Figure 4(d) is for  $\theta = 20$ , and is very jumpy. Mitchell et al. (1990) also found necessary and sufficient conditions on the correlation function so that the derivative process has an exponential correlation function. These are called smoothed exponential correlation functions.

4.3.4. Gaussian Sacks et al. (1989b) generalized the exponential correlation function by using

$$R(d) = \exp(-\theta |d|^q)$$



Fig. 4. Realizations for the exponential correlation function with  $\theta = (a) 0.5$ , (b) 2.0, (c) 5.0, and (d) 20.0.

where  $0 < q \leq 2$  and  $\theta \in (0, \infty)$ . Taking q = 1 recovers the exponential correlation function. As q increases, this correlation function produces smoother realizations. However, as long as q < 2, these processes are not mean square differentiable.

The Gaussian correlation function is the case q = 2 and the associated processes are infinitely mean square differentiable. In the Bayesian interpretation, this correlation function puts all of the prior mass on analytic functions (Currin et al., 1991). This correlation function is appropriate when the simulator output is known to be analytic. Figure 5 displays several realizations for various  $\theta$  for the Gaussian correlation function. These realizations are very smooth, even when  $\theta = 50$ .

# 4.3.5. Matérn

All of the univariate correlation functions described above are either zero, once or infinitely many times mean square differentiable. Stein (1989) recommends a more flexible family of correlation function (Matérn, 1947; Yaglom, 1987). The Matérn correlation function is parameterized by  $\theta \in (0, \infty)$  and  $\nu \in (-1, \infty)$  and is given by

$$R(d) = \frac{(\theta|d|)^{\nu}}{\Gamma(\nu)2^{\nu-1}} K_{\nu}(\theta|d|)$$



Fig. 5. Realizations for the Gaussian correlation function with  $\theta =$  (a) 0.5, (b) 2.0, (c) 10.0, and (d) 50.0.

where  $K_{\nu}(\cdot)$  is a modified Bessel function of order  $\nu$ . The associated process will be m times differentiable if and only if  $\nu > m$ . Hence, the amount of differentiability can be controlled by  $\nu$  while  $\theta$  controls the range of the correlations. This correlation family is more flexible than the other correlation families described above due to the control of the differentiability of the predictive surface.

Figure 6 displays several realizations of processes with the Matérn correlation function with  $\nu = 2.5$  and various values of  $\theta$ . For small values of  $\theta$ , the realizations are very smooth and flat while the realizations are erratic for large values of  $\theta$ .

#### 4.3.6. Summary

The correlation functions described above have been applied in computer experiments. Software for predicting with them is described in Koehler (1990). The cubic correlation function yields predictions that are cubic splines. The exponential predictions are non-differentiable while the Gaussian predictions are infinitely differentiable. The Matérn correlation function is the most flexible since the degree of differentiability and the smoothness of the predictions can be controlled. In general, enough prior information to fix the parameters of a particular correlation family and  $\sigma^2$  will not be available. A pure Bayesian approach would place a prior distribution on the parameters of a family and use the posterior-distribution of the parameter in the estimation



Fig. 6. Realizations for the Matérn correlation function with  $\nu = 2.5$  and  $\theta =$  (a) 2.0, (b) 4.0, (c) 10.0, and (d) 25.0.

process. Alternatively, an empirical Bayes approach which uses the data to estimate the parameters of a correlation family and  $\sigma^2$  is often used. The maximum likelihood estimation procedure will be presented and discussed in the next section.

### 4.4. Correlation function estimation – maximum likelihood

The previous subsections of this section presented the Kriging model, and families of correlation functions. The families of correlations are all parameterized by one or two parameters which control the range of correlation and the smoothness of the corresponding processes. This model assumes that  $\sigma^2$ , the family and parameters of  $R(\cdot)$  are known. In general, these values are not completely known a priori. The appropriate correlation family might be known from the simulator's designers experience regarding the smoothness of the function. Also, ranges for  $\sigma^2$  and the parameters of  $R(\cdot)$  might be known if a similar computer experiment has been performed. A pure Bayesian approach is to quantify this knowledge into a prior distribution on  $\sigma^2$  and  $R(\cdot)$ . How to distribute a non-informative prior across the different correlation

families and within each family is unclear. Furthermore, the calculation of the posterior distribution is generally intractable.

An alternative and more objective method of estimating these parameters is an empirical Bayes approach which finds the parameters which are most consistent with the observed data. This section presents the maximum likelihood method for estimating  $\beta$ ,  $\sigma^2$  and the parameters of a fixed correlation family when the underlying distribution of  $Z(\cdot)$  is Gaussian. The best parameter set from each correlation family can be evaluated to find the overall "best"  $\sigma^2$  and  $R(\cdot)$ .

Consider the case where the distribution of  $Z(\cdot)$  is Gaussian. Then the distribution for the response at the *n* design points  $Y_D$  is multinormal and the likelihood is given by

lik
$$(\beta, \sigma^2, R \mid Y_D) = (2\pi)^{-n/2} \sigma^{-n} |R_D|^{-1/2}$$
  
  $\times \exp\left\{-\frac{1}{2\sigma^2}(Y_D - H\beta)' R_D^{-1}(Y_D - H\beta)\right\}$ 

where  $R_D$  is the design correlation matrix. The log likelihood is

$$l_{ml}(\beta, \sigma^{2}, R_{D} \mid Y_{D}) = -\frac{n}{2} \ln (2\pi) - \frac{n}{2} \ln (\sigma^{2}) - \frac{1}{2} \ln (|R_{D}|) -\frac{1}{2\sigma^{2}} (Y_{D} - H\beta)' R_{D}^{-1} (Y_{D} - H\beta).$$
(5)

Hence

$$\frac{\partial l_{ml}(\beta,\sigma^2,R\mid Y_D)}{\partial\beta} = -\frac{1}{\sigma^2} \left( H' R_D^{-1} Y_D - H' R_D^{-1} H \beta \right)$$

which when set to zero yields the maximum likelihood estimate of  $\beta$  that is the same as the generalized least squares estimate,

$$\widehat{\beta}_{ml} = \left[ H' R_D^{-1} H \right]^{-1} H' R_D^{-1} Y_D.$$
(6)

Similarly,

$$\frac{\partial l_{ml}(\beta,\sigma^2,R_D \mid Y_D)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (Y_D - H\beta)' R_D^{-1} (Y_D - H\beta)$$

which when set to zero yields the maximum likelihood estimate of  $\sigma^2$ 

$$\widehat{\sigma}_{ml}^2 = \frac{1}{n} \left( Y_D - H\widehat{\beta} \right)' R_D^{-1} \left( Y_D - H\widehat{\beta} \right). \tag{7}$$

Therefore, if  $R_D$  is known, the maximum likelihood estimates of  $\beta$  and  $\sigma^2$  are easily calculated. However, if  $R(\cdot)$  is parameterized by  $\theta = (\theta_1, \ldots, \theta_s)$ ,

$$\frac{\partial l_{ml}(\beta, \sigma^2, R_D \mid Y_D)}{\partial \theta_i} = -\frac{1}{2} \frac{\partial |R_D|}{\partial \theta_i} - \frac{1}{2\sigma^2} (Y_D - H\beta)' \frac{\partial R_D^{-1}}{\partial \theta_i} (Y_D - H\beta)$$
$$= -\frac{1}{2} \operatorname{tr} \left\{ R_D^{-1} \frac{\partial R_D}{\partial \theta_i} \right\}$$
$$+ \frac{1}{2\sigma^2} (Y_D - H\beta)' R_D^{-1} \frac{\partial R_D}{\partial \theta_i} R_D^{-1} (Y_D - H\beta) \quad (8)$$

does not generally yield an analytic solution for  $\theta$  when set to zero for i = 1, ..., s. (Commonly s = p or 2p, but this need not be assumed.)

An alternative method to estimate  $\theta$  is to use a nonlinear optimization routine using equation (5) as the function to be optimized. For a given value of  $\theta$ , estimates of  $\beta$  and  $\sigma^2$  are calculated using equations (6) and (7), respectively. Next, equation (8) is used in calculating the partial derivatives of the objective function. See Mardia and Marshall (1984) for an overview of the maximum likelihood procedure.

### 4.5. Estimating and using derivatives

In the manufacturing sciences, deterministic simulators help describe the relationships between product design, and the manufacturing process to the product's final characteristics. This allows the product to be designed and manufactured efficiently. Equally important are the effects of uncontrollable variation in the manufacturing parameters to the end product. If the product's characteristics are sensitive to slight variations in the manufacturing process, the yield, or percentage of marketable units produced, may decrease. Furthermore, understanding the sensitivities of the product's characteristics can help design more reliable products and increase the overall quality of the product.

Many simulators need to solve differential equations and can provide the gradient of the response at a design point with little or no additional computational cost. However, some simulators require that the gradient be approximated by a difference equation. Then the cost of finding a directional derivative at a point is equal to evaluating an additional point while approximating the total gradient requires p additional runs.

Consider Figure 7 for an example in p = 1 showing the effects of including gradient information on prediction. The solid lines, Y in Figure 7(a) and Y' in Figure 7(b), are the true function and it's derivative, respectively, while the long dashed lines are Kriging predictors  $\hat{Y}_3$  and  $\hat{Y}'_3$  based on n = 3 observations. As expected  $\hat{Y}_3$  goes through the design points,  $D = \{.2, .5, .8\}$ , but  $\hat{Y}'_3$  is a poor predictor of Y'. The short dashed lines are the n = 3 predictors with derivative information  $\hat{Y}_{3'}$  and  $\hat{Y}'_{3'}$ . Notice that this predictor now matches Y' and Y at D and the interpolations are over all much better. The addition of gradient information substantially improves the fits of both Y and Y'. The dotted lines are the n = 6 predictors  $\hat{Y}_6$  and  $\hat{Y}'_6$  and is a fairer comparison if the derivative costs are equal to the response cost. The predictor  $\hat{Y}_6$  is a little better on the interior of S but  $\hat{Y}'_6$  is worse at x = 0 than  $\hat{Y}'_3$ .



Fig. 7. (a) An example of a response (Y) and three predictors  $(\widehat{Y}_3, \widehat{Y}_{3'}, \widehat{Y}_6)$ . (b) An example of a derivative (Y') and three predictors  $(\widehat{Y}'_3, \widehat{Y}'_3, \widehat{Y}'_6)$ .

The Kriging methodology easily extends to model gradients. To see this for p = 1, let  $E[Y(\cdot)] = \mu$  and  $d = t_2 - t_1$ , then

Cov 
$$[Y(t_1), Y'(t_2)] = E[Y(t_1)Y'(t_2)] - E[Y(t_1)]E[Y'(t_2)].$$

Now due to the stationarity of  $Y(\cdot)$ ,  $\mathrm{E}[Y'(\cdot)]=0$  and

$$Cov [Y(t_1), Y'(t_2)] = E [Y(t_1)Y'(t_2)]$$
$$= E \left[Y(t_1) \lim_{\delta \to 0} \frac{Y(t_2 + \delta) - Y(t_2)}{\delta}\right]$$
$$= E \left[\lim_{\delta \to 0} \frac{Y(t_1)Y(t_2 + \delta) - Y(t_1)Y(t_2)}{\delta}\right]$$
$$= \sigma^2 \lim_{\delta \to 0} \frac{R(d + \delta) - R(d)}{\delta}$$
$$= \sigma^2 R'(d)$$

for differentiable  $R(\cdot)$ . Similarly,

$$\operatorname{Cov}\left[Y'(t_1), Y(t_2)\right] = -\sigma^2 R'(d)$$

and

$$\operatorname{Cov}\left[Y'(t_1), Y'(t_2)\right] = -\sigma^2 R''(d)$$

For more general p and for higher derivatives, following Morris et al. (1993) let

$$Y^{(a_1,\ldots,a_p)}(t) = \frac{\partial^a}{\partial t_1^{(a_1)}\cdots \partial t_p^{(a_p)}}Y(t)$$

where  $a = \sum_{j=1}^{p} a_j$  and  $t_j$  is the *j*th component of *t*. Then  $E[Y^{(a_1,...,a_p)}] = 0$  and

$$\operatorname{Cov}\left[Y^{(a_1,\ldots,a_p)}(t_1),Y^{(b_1,\ldots,b_p)}(t_2)\right] = (-1)^a \sigma^2 \prod_{j=1}^p R_j^{(a_j+b_j)}(t_{2j}-t_{1j})$$

for  $R(d) = \prod_{j=1}^{p} R_j(d_j)$ . Furthermore, for directional derivatives, let  $Y'_{\nu}(t)$  be the directional derivative of Y(t) in the direction  $\nu = (\nu_1, \dots, \nu_p)'$ ,  $\sum_{j=1}^{p} \nu_j^2 = 1$ ,

$$Y'_{
u}(t) = \sum_{j=1}^{p} \frac{\partial Y(t)}{\partial t_j} \nu_j = \langle \bigtriangledown Y(t), \nu \rangle \,.$$

Then  $E[Y'_{\nu}(t)] = 0$  and for d = t - s,

$$\operatorname{Cov}\left[Y(s), Y_{\nu}'(t)\right] = \operatorname{E}\left[Y(s)Y_{\nu}'(t)\right]$$
$$= \sum_{j=1}^{p} \operatorname{E}\left[Y(s)\frac{\partial Y(t)}{\partial t_{j}}\nu_{j}\right]$$
$$= \sum_{j=1}^{p} \operatorname{Cov}\left[Y(s), \frac{\partial Y(t)}{\partial t_{j}}\right]\nu_{j}$$
$$= \sigma^{2} \sum_{j=1}^{p} \frac{\partial \dot{R}(d)}{\partial d_{j}}\nu_{j}$$
$$= \sigma^{2} \langle \dot{R}(d), \nu \rangle$$
(9)

where  $\dot{R}(d) = [\partial R(d)/\partial d_1, \dots, \partial R(d)/\partial d_p]'$ . Similarly,

$$\operatorname{Cov}\left[Y_{\nu}'(s), Y(t)\right] = -\sigma^{2} \left\langle \dot{R}(d), \nu \right\rangle \tag{10}$$

J. R. Koehler and A. B. Owen

and

$$\operatorname{Cov}\left[Y_{\nu_{s}}'(s), Y_{\nu_{t}}'(t)\right] = -\sigma^{2}\nu_{s}'\ddot{R}(d)\nu_{t}$$
(11)

where

$$\left(\ddot{R}(d)\right)_{kl} = \frac{\partial^2 R(d)}{\partial d_k \partial d_l}$$

is the matrix of 2nd partial derivatives evaluated at d.

The Kriging methodology is modified to model gradient information by letting

$$y_D^* = [y(x_1), \dots, y(x_n), y'_{\nu_{11}}(x_1), y'_{\nu_{12}}(x_1), \dots, y'_{\nu_{nm}}(x_n)]'$$

where  $\nu_{il}$  is the direction of the *l*th directional derivative at  $x_i$ . Also let

$$\mu^{*}=(\mu,\mu,\ldots,\mu,0,0,\ldots,0)^{\prime}$$

with  $n \mu s$  and mn 0s and let  $V^*$  be the combined covariance matrix for the design responses and derivatives with the entries as prescribed above (equations (9), (10), and (11)). Then

$$\widehat{Y}(x_0) = \mu + {v_{x_0}'}^* {V^*}^{-1} \left( y_D^* - \mu^* 
ight)$$

and

$$\widehat{Y}'_{\nu}(x_0) = v'^*_{x_0,\nu} V^{*^{-1}} \left( y^*_D - \mu^* \right)$$

where  $v_{x_0}^{\prime *} = \operatorname{Cov}[Y(x_0), Y_D^*]$ , and  $v_{x_0,\nu}^{\prime *} = \operatorname{Cov}[Y_{\nu}^{\prime}(x_0), Y_D^*]$ .

Notice that once differentiable random functions need twice differentiable correlation functions. One problem with using the total gradient information is the rapid increase in the covariance matrix. For each additional design point,  $V^*$  increases by p + 1 rows and columns. Fortunately, these new rows and columns generally have lower correlations than the corresponding rows and columns for an equal number of response. The inversion of  $V^*$  is more computationally stable than for an equally sized  $V_D$ . More research is needed to provide general guidelines for using gradient information efficiently.

# 4.6. Complexity of computer experiments

Recent progress in complexity theory, a branch of theoretical computer science, has shed some light on computer experiments. The dissertation of Ritter (1995) contains an excellent summary of this area. Consider the case where Y(x) = Z(x), that is where there is no regression function. If for  $r \ge 1$  all of the r'th order partial derivatives of Z(x) exist in the mean square sense and obey a Holder condition of order  $\beta$ , then it

280

is possible (see Ritter et al., 1993) to approximate Z(x) with an  $L^2$  error that decays as  $O(n^{-(r+\beta)/p})$ . This error is a root mean square average over randomly generated functions Z.

When the covariance has a tensor product form, like those considered here, one can do even better. Ritter et al. (1995) show that the error rate for approximation in this case is  $n^{-r-1/2}(\log n)^{(p-1)(r+1)}$  for products of covariances satisfying Sacks-Ylvisaker conditions of order  $r \ge 0$ . When Z is a p dimensional Wiener sheet process, for which r = 0, the result is  $n^{-1/2}(\log n)^{(p-1)}$  which was first established by Wozniakowski (1991).

In the general case, the rate for integration is  $n^{-1/2}$  times the rate for approximation. A theorem of Wasilkowski (1994) shows that a rate  $n^{-d}$  for approximation can usually be turned into a rate  $n^{-d-1/2}$  for integration by the simple device of fitting an approximation with n/2 function evaluations, integrating the approximation, and then adjusting the result by the average approximation error on n/2 more Monte Carlo function evaluations. For tensor product kernels the rate for integration is  $n^{-r-1}(\log n)^{(p-1)/2}$  (see Paskov, 1993), which has a more favorable power of  $\log n$  than would arise via Wasilkowski's theorem.

The fact that much better rates are possible under tensor product models than for general covariances suggests that the tensor product assumption may be a very strong one. The tensor product assumption is at least strong enough that under it, there is no average case curse of dimensionality for approximation.

# 5. Bayesian designs

Selecting an experimental design, D, is a key issue in building an efficient and informative Kriging model. Since there is no random error in this model, we wish to find designs that minimize squared-bias. While some experimental design theories (Box and Draper, 1959; Steinberg, 1985) do investigate the case where bias rather than solely variance plays a crucial role in the error of the fitted model, how good these designs are for the pure bias problem of computer experiments is unclear. Box and Draper (1959) studied the effect of scaling factorial designs by using a first order polynomial model when the true function is a quadratic polynomial. Box and Draper (1983) extended the results to using a quadratic polynomial model when the true response surface is a cubic polynomial. They found that mean squared-error optimal designs are close to bias optimal designs. Steinberg (1985) extended these ideas further by using a prior model proposed by Young (1977) that puts prior distributions on the coefficients of a sufficiently large polynomial. However, model (2) is more flexible than high ordered polynomials and therefore better designs are needed.

This section introduces four design optimality criteria for use with computer experiments: entropy, mean squared-error, maximin and minimax designs. Entropy designs maximize the amount of information expected for the design while mean squared-error designs minimize the expected mean squared-error. Both these designs require a priori knowledge of the correlation function  $R(\cdot)$ . The design criteria described below are for the case of fixed design size n. Simple sequential designs, where the location of



Fig. 8(a). Maximum entropy designs for p = 2, n = 1-16, and the Gaussian correlation function with  $\theta = (0.5, 0.5)$ .

the *n*th design point is determined after the first n-1 points have been evaluated, will not be presented due to their tendencies to replicate (Sacks et al., 1989b). However, sequential block strategies could be used where the above designs could be used as starting blocks. Depending upon the ultimate goal of the computer experiment, the first design block might be utilized to refine the design and reduce the design space.

#### 5.1. Entropy designs

Lindley (1956) introduced a measure, based upon Shannon's entropy (Shannon, 1948), of the amount of information provided by an experiment. This Bayesian measure uses the expected reduction in entropy as a design criterion. This criterion has been used in Box and Hill (1967) and Borth (1975) for model discrimination. Shewry and Wynn (1987) showed that, if the design space is discrete (i.e., a lattice in  $[0, 1]^p$ ), then minimizing the expected posterior entropy is equivalent to maximizing the prior entropy.



Fig. 8(b). Maximum entropy designs for p = 2, n = 1-16, and the Gaussian correlation function with  $\theta = (2, 2)$ .

DEFINITION 1. A design  $D_E$  is a Maximum Entropy Design if

$$E_{Y}\left[-\ln P(Y_{D_{E}})\right] = \min_{D} E_{Y}\left[-\ln P(Y_{D})\right]$$

where  $P(Y_D)$  is the density of  $Y_D$ .

In the Gaussian case, this is equivalent to finding a design that maximizes the determinant of the variance of  $Y_D$ . In the Gaussian prior case, where  $\beta \sim N_k(b, \tau^2 \Sigma)$ , the determinant of the unconditioned covariance matrix is

$$\begin{aligned} \left| V_D + \tau^2 H \Sigma H' \right| &= \begin{vmatrix} V_D + \tau^2 H \Sigma H' & H \\ 0 & I \end{vmatrix} \\ &= \begin{vmatrix} \begin{pmatrix} V_D & H \\ -\tau^2 \Sigma H' & I \end{pmatrix} \begin{pmatrix} I & 0 \\ \tau^2 \Sigma H' & I \end{pmatrix} \end{aligned}$$



Fig. 8(c). Maximum entropy designs for p = 2, n = 1-16, and the Gaussian correlation function with  $\theta = (10, 10)$ .

$$= \begin{vmatrix} V_D & H \\ -\tau^2 \Sigma H' & I \end{vmatrix}$$

$$= \begin{vmatrix} \begin{pmatrix} I & 0 \\ \tau^2 \Sigma H' V_D^{-1} & I \end{pmatrix} \begin{pmatrix} V_D & H \\ -\tau^2 \Sigma H' & I \end{pmatrix} \end{vmatrix}$$

$$= \begin{vmatrix} V_D & H \\ 0 & \tau^2 \Sigma H' V_D^{-1} H + I \end{vmatrix}$$

$$= |V_D| \left| \tau^2 \Sigma H' V_D^{-1} H + \tau^{-2} \Sigma^{-1} \right| \left| \tau^2 \Sigma \right|.$$

Since  $\tau^2 \Sigma$  is fixed, the maximum entropy criterion is equivalent to finding the design  $D_E$  that maximizes

$$|V_D| |H'V_D^{-1}H + \tau^{-2}\Sigma^{-1}|.$$

If the prior distribution is diffuse,  $\tau^2 \rightarrow \infty$ , the maximum entropy criterion is equivalent to

$$|V_D| \left| H' V_D^{-1} H \right|$$

and if  $\beta$  is treated as fixed, then the maximum entropy criterion is equivalent to  $|V_D|$ .

Shewry and Wynn (1987, 1988) applied this measure in designs for spatial models. Currin et al. (1991) and Mitchell and Scott (1987) have applied the entropy measure to finding designs for computer experiments. By this measure, the amount of information in experimental design is dependent on the prior knowledge of  $Z(\cdot)$  through  $R(\cdot)$ . In general,  $R(\cdot)$  will not be known a priori. Additionally, these optimal designs are difficult to construct due to the required  $n \times p$  dimensional optimization of the *n* design point locations. Currin et al. (1991) describe an algorithm adopted from DETMAX (Mitchell, 1974) which successively removes and adds points to improve the design.

Figure 8(a) shows the optimal entropy designs for p = 2, n = 1, ..., 16,  $R(d) = \exp\{-\theta \sum d_j^2\}$  where  $\theta = 0.5, 2, 10$ . The entropy designs tend to spread the points out in the plane and favor the edge of the design space over the interior. For example, the n = 16 designs displayed in Figure 8(a) have 12 points on the edge and only 4 points in the interior. Furthermore, most of the designs are similar across the different correlation functions although there are some differences. Generally, the ratio of the edge to interior points are constant. The entropy criterion appears to be insensitive to changes in the location of the interior points. Johnson et al. (1990) indicate that entropy designs (see Section 5.3).

# 5.2. Mean squared-error designs

Box and Draper (1959) proposed minimizing the normalized integrated mean squarederror (IMSE) of  $\hat{Y}(x)$  over  $[0,1]^p$ . Welch (1983) extended this measure to the case when the bias is more complicated. Sacks and Schiller (1988) and Sacks et al. (1989a) discuss in more detail IMSE designs for computer experiments.

DEFINITION 2. A design  $D_I$  is an Integrated Mean Squared-Error (IMSE) design if

$$J(D_I) = \min_D J(D)$$

where

$$J(D) = \frac{1}{\sigma^2} \int_{[0,1]^p} \mathbb{E} \big[ Y(x) - \widehat{Y}(x) \big]^2 \, \mathrm{d}x.$$

J(D) is dependent on  $R(\cdot)$  through Y(x). For any design, J(D) can be expressed as

$$J(D) = \sigma^2 - \operatorname{trace}\left\{ \begin{bmatrix} 0 & H' \\ H & V_D \end{bmatrix}^{-1} \int \begin{bmatrix} h(x)h'(x) & h(x)v'_x \\ v_xh'(x) & v_xv'_x \end{bmatrix} dx \right\}$$

and, as pointed out by Sacks et al. (1989a), if the elements of h(x) and  $V_x$  are products of functions of a single input variable, the multidimensional integral simplifies to products of one-dimensional integrals. As in the entropy design criterion, the minimization of J(D) is a optimization in  $n \times p$  dimensions and is also dependent on  $R(\cdot)$ .

Sacks and Schiller (1988) describe the use of a simulated annealing method for constructing IMSE designs for bounded and discrete design spaces. Sacks et al. (1989b) use a quasi-Newton optimizer on a Cray X-MP48. They found that optimizing a n = 16, p = 6 design with  $\theta_1 = \cdots = \theta_6 = 2$  took 11 minutes. The PACE program (Koehler, 1990) uses the optimization program NPSOL (Gill et al., 1986) to solve the IMSE optimization for a continuous design space. For n = 16, p = 6, this optimization requires 13 minutes on a DEC3100, a much less powerful machine than the Cray. Generally, these algorithms can find only local minima and therefore many random



Fig. 9(a). Minimum integrated mean square error designs for p = 2, n = 1-9, and the Gaussian correlation function with  $\theta = (.5, .5)$ .

starts are required.

Since J(D) is dependent on  $R(\cdot)$ , robust designs need to be found for general  $R(\cdot)$ . Sacks et al. (1989a) found that for n = 9, p = 2 and  $R(d) = \exp\{-\theta \sum_{j=1}^{2} d_{j}^{2}\}$  (see Section 4.3.4 for details on the Gaussian correlation function) the IMSE design for  $\theta = 1$  is robust in terms of relative efficiency. However, this analysis used a quadratic polynomial model and the results may not extend to higher dimensions nor different linear model components. Sacks et al. (1989b) used the optimal design for the Gaussian correlation function with  $\theta = 2$  for design efficiency-robustness.

Figure 9(a) displays IMSE designs for p = 2 and n = 1, ..., 9 for  $\theta = .5, 2, 10$ . The designs, in general lie in the interior of S. For fixed design size n, the designs usually are similar geometrically for different  $\theta$  values with the scale decreasing as  $\theta$  increases. They have much symmetry for some values of n, particularly n = 12. Notice that for the case when n = 5 that the design only takes on three unique values for each of the input variables. These designs tend to have clumped projections onto



Fig. 9(b). Minimum integrated mean square error designs for p = 2, n = 1-16, and the Gaussian correlation function with  $\theta = (2, 2)$ .



Fig. 9(c). Minimum integrated mean square error designs for p = 2, n = 1-16, and the Gaussian correlation function with  $\theta = (10, 10)$ .

lower dimension marginals of the input space. Better projection properties are needed when the true function is only dependent on a subset of the input variables.

# 5.3. Maximin and minimax designs

Johnson et al. (1990) developed the idea of minimax and maximin designs. These designs are dependent on a distance measure or metric. Let  $d(\cdot, \cdot)$  be a metric on  $[0, 1]^p$ . Hence  $\forall x_1, x_2, x_3 \in [0, 1]^p$ ,

$$\begin{aligned} &d(x_1, x_2) = d(x_2, x_1), \\ &d(x_1, x_2) \geqslant 0, \\ &d(x_1, x_2) = 0 \Leftrightarrow x_1 = x_2, \\ &d(x_1, x_2) \leqslant d(x_1, x_3) + d(x_3, x_2). \end{aligned}$$



Fig. 10. (a) Minimax and (b) Maximin designs for n = 6 and p = 2 with Euclidean distance.

DEFINITION 3. Design  $D_{MI}$  is a Minimax Distance Design if

$$\min_{D} \max_{x} d(x, D) = \max_{x} d(x, D_{MI})$$

where

$$d(x,D) = \min_{x_0 \in D} d(x,x_0).$$

Minimax distance designs ensure that all points in  $[0, 1]^p$  are not too far from a design point. Let  $d(\cdot, \cdot)$  be Euclidean distance and consider placing a *p*-dimensional sphere with radius r around each design point. The idea of a minimax design is to place the n points so that the design space is covered by the spheres with minimal r. As an illustration, consider the owner of a petroleum corporation who wants to open some franchise gas stations. The gas company would like to locate the stations in the most convenient sites for the customers. A minimax strategy of placing gas stations. would ensure that no customer is too far from one of the company's stations.

Figure 10(a) shows a minimax design for p = 2 and n = 6 with  $d(\cdot, \cdot)$  being Euclidean distance. The maximum distance to a design point is .318. For small n, minimax designs will generally lie in the interior of the design space.

DEFINITION 4. A design  $D_{MA}$  is a Maximin Distance Design if

$$\max_{D} \min_{x_1, x_2 \in D} d(x_1, x_2) = \min_{x_1, x_2 \in D_{MA}} d(x_1, x_2)$$

Again, let  $d(\cdot, \cdot)$  be Euclidean distance. Maximin designs pack the *n* design points, with their associated spheres, into the design space, *S*, with maximum radius. Parts of the sphere may be out of *S* but the design points must be in *S*. Analogous to the minimax illustration above is the position of the owners the gas station franchises. They wish to minimize the competition from each other by locating the stations as far apart as possible. A maximin strategy for placing the franchises would ensure that no two stations are too close to each other.

Figure 10(b) shows a maximin design for p = 2, n = 6 and  $d(\cdot, \cdot)$  Euclidean distance. For small n, maximin designs will generally lie on the exterior of S and fill in the interior as n becomes large.

#### 5.4. Hyperbolic cross points

Under the tensor product covariance models, it is possible to approximate and integrate functions with greater accuracy than in the general case. One gets the same rates of convergence as in univariate problems, apart from a multiplicative penalty that is some power of log n. Hyperbolic cross point designs, also known as sparse grids have been shown to achieve optimal rates in these cases. See Ritter (1995). These point sets were first developed by Smolyak (1963). They were used in interpolation by Wahba (1978) and Gordon (1971) and by Paskov (1993) for integration. Chapter 4 of Ritter (1995) gives a good description of the construction of these points and lists other references.

### 6. Frequentist prediction and inference

The frequentist approach to prediction and inference in computer experiments is based on numerical integration. For a scalar function Y = f(X), consider a regression model of the form

$$Y = f(X) \doteq Z(X)\beta \tag{12}$$

where Z(X) is a row vector of predictor functions and  $\beta$  is a vector of parameters. Suitable functions Z might include low order polynomials, trigonometric polynomials wavelets, or some functions specifically geared to the application. Ordinarily Z(X)includes a component that is always equal to 1 in order to introduce an intercept term into equation (12).

It is unrealistic to expect that the function f will be exactly representable as the finite linear combination given by (12), and it is also unrealistic to expect that the residual will be a random variable with mean zero at every fixed  $X_0$ . This is why we only write  $f \doteq Z\beta$ . There are many ways to define the best value of  $\beta$ , but an especially natural approach is to choose  $\beta$  to minimize the mean squared error of the approximation, with respect to some distribution F on  $[0, 1]^p$ . Then the optimal value for  $\beta$  is

$$\beta_{LS} = \left(\int Z(X)'Z(X) \,\mathrm{d}F\right)^{-1} \int Z(X)'f(X) \,\mathrm{d}F.$$

So if one can integrate over the domain of X then one can fit regression approximations there.

The quality of the approximation may be assessed globally by the integrated mean squared error

$$\int (Y-Z(X)\beta)^2 \,\mathrm{d} F.$$

For simplicity we take the distribution F to be uniform on  $[0, 1]^p$ . Also for simplicity the integration schemes to be considered usually estimate  $\int g(X) dF$  by

$$\frac{1}{n}\sum_{i=1}^n g(x_i)$$

for well chosen points  $x_1, \ldots, x_n$ . Then  $\beta_{LS}$  may be estimated by linear regression

$$\widehat{\beta} = \left(\frac{1}{n} \sum_{i=1}^{n} Z(x_i)' Z(x_i)\right)^{-1} \frac{1}{n} \sum_{i=1}^{n} Z(x_i)' f(x_i),$$

or when the integrals of squares and cross products of Z's are known by

$$\widetilde{\beta} = \left(\int Z(X)'Z(X)\,\mathrm{d}F\right)^{-1}\frac{1}{n}\sum_{i=1}^{n}Z(x_i)'f(x_i).$$
(13)

Choosing the components of Z to be an orthogonal basis, such as tensor products of orthogonal polynomials, multivariate Fourier series or wavelets, equation (13) simplifies to

$$\widetilde{\beta} = \frac{1}{n} \sum_{i=1}^{n} Z(x_i)' f(x_i)$$
(14)

and one can avoid the cost of matrix inversion. The computation required by equation (14) grows proportionally to nr not  $n^3$ , where r = r(n) is the number of regression variables in Z. If r = O(n) then the computations grow as  $n^2$ . Then, in the example from Section 3, an hour of function evaluation followed by a minute of algebra would scale into a day of function evaluation followed by 9.6 hours of algebra, instead of the 9.6 days that an  $n^3$  algorithm would require. If the  $Z(x_i)$  exhibit some sparsity then it may be possible to reduce the algebra to order n or order  $n \log n$ .

Thus the idea of turning the function into data and making exploratory plots can be extended to turning the function into data and applying regression techniques. The theoretically simplest technique is to take  $X_i$  iid  $U[0, 1]^p$ . Then  $(X_i, Y_i)$  are iid pairs with the complication that Y has zero variance given X. The variance matrix of  $\tilde{\beta}$  is then

$$\frac{1}{n} \left( \int Z' Z \, \mathrm{d}F \right)^{-1} \operatorname{Var} \left( Z(X)' Y(X) \right) \left( \int Z' Z \, \mathrm{d}F \right)^{-1}$$
(15)

and for orthogonal predictors this simplifies further to

$$\frac{1}{n} \operatorname{Var} \left( Z(X)' Y(X) \right). \tag{16}$$

Thus any integration scheme that allows one to estimate variances and covariances of averages of Y times components of Z allows one to estimate the sampling variance matrix of the regression coefficients  $\tilde{\beta}$ . For iid sampling one can estimate this variance matrix by

$$\frac{1}{n-r-1}\sum_{i=1}^n \left(Z(x_i)Y(x_i)-\widetilde{\beta}\right)' \left(Z(x_i)Y(x_i)-\widetilde{\beta}\right)$$

when the row vector Z comprises an intercept and r additional regression coefficients.

This approach to computer experimentation should improve if more accurate integration techniques are substituted for the iid sampling. Owen (1992a) investigates the case of Latin hypercube sampling for which a central limit theorem also holds.

Clearly more work is needed to make this method practical. For instance a scheme for deciding how many predictors should be in Z, or otherwise for regularizing  $\tilde{\beta}$  is required.

# 7. Frequentist experimental designs

The frequentist approach proposed in the previous section requires a set of points  $x_1, \ldots, x_n$  that are good for numerical integration and also allow one to estimate the sampling variance of the corresponding integrals. These two goals are somewhat at odds. Using an iid sample makes variance estimation easier while more complicated schemes described below improve accuracy but make variance estimation harder.

The more basic goal of getting points  $x_i$  into "interesting corners" of the input space, so that important features are likely to be found is usually well served by point sets that are good for numerical integration.

We assume that the region of interest is the unit cube  $[0, 1]^p$ , and that the integrals of interest are with respect to the uniform distribution over this cube. Other regions of interest can usually be reduced to the unit cube and other distributions can be changed to the uniform by a change of variable that can be subsumed into f.

Throughout this section we consider an example with p = 5, and plot the design points  $x_i$ .



Fig. 11. 25 distinct points among 625 points in a 5<sup>5</sup> grid.

# 7.1. Grids

Since varying one coordinate at a time can cause one to miss important aspects of f, it is natural to consider instead sampling f on a regular grid. One chooses k different values for each of  $X^1$  through  $X^p$  and then runs all  $k^p$  combinations. This works well for small values of p, perhaps 2 or 3, but for larger p it becomes completely impractical because the number of runs required grows explosively.

Figure 11 shows a projection of  $5^5 = 625$  points from a uniform grid in  $[0, 1]^5$  onto two of the input variables. Notice that with 625 runs, only 25 distinct values appear in the plane, each representing 25 input settings in the other three variables. Only 5 distinct values appear for each of input variable taken singly. In situations where one of the responses  $Y^k$  depends very strongly on only one or two of the inputs  $X^j$  the grid design leads to much wasteful duplication.

The grid design does not lend itself to variance estimation since averages over the grid are not random. The accuracy of a grid based integral is typically that of a univariate integral based on  $k = n^{1/p}$  evaluations. (See Davis and Rabinowitz, 1984.) For large p this is a severe disadvantage.



Fig. 12. A 34 point Fibonacci lattice in  $[0, 1]^2$ .

#### 7.2. Good lattice points

A significant improvement on grids may be obtained in integration by the method of good lattice points. (See Sloan and Joe (1994) and Niederreiter (1992) for background and Fang and Wang (1994) for applications to statistics.)

For good lattice points

$$x_i^j = \left\{\frac{h_j(i-1) + 0.5}{n}\right\}$$

where  $\{z\}$  is z modulo 1, that is, z minus the greatest integer less than or equal to z and  $h_j$  are integers with  $h_1 = 1$ . The points  $v_i$  with  $v_i^j = ih_j/n$  for integer *i* form a lattice in  $\mathbb{R}^p$ . The points  $x_i$  are versions of these lattice points confined to the unit cube, and the term "good" refers to a careful choice of *n* and  $h_j$  usually based on number theory.

Figure 12 shows the Fibonacci lattice for p = 2 and n = 34. For more details see Sloan and Joe (1994). Here  $h_1 = 1$  and  $h_2 = 21$ . The Fibonacci lattice is only available in 2 dimensions. Appendix A of Fang and Wang (1994) lists several other choices for good lattice points, but the smallest value of n there for p = 5 is 1069. Hickernell (1996) discusses greedy algorithms for finding good lattice points with smaller n.

The recent text (Sloan and Joe, 1994) discusses lattice rules for integration, which generalize the method of good lattice points. Cranley and Patterson (1976) consider randomly perturbing the good lattice points by adding, modulo 1, a random vector uniform over  $[0, 1]^p$  to all the  $x_i$ . Taking r such random offsets for each of the n data points gives nr observations with r - 1 degrees of freedom for estimating variance.

Lattice integration rules can be extraordinarily accurate on smooth periodic integrands and thus an approach to computer experiments based on Cranley and Patterson's method might be expected to work well when both f(x) and Z(x) are smooth and periodic. Bates et al. (1996) have explored the use of lattice rules as designs for computer experiments.

#### 7.3. Latin hypercubes

While good lattice points start by improving the low dimensional projections of grids, Latin hypercube sampling starts with iid samples. A Latin hypercube sample has

$$X_{i}^{j} = \frac{\pi^{j}(i) - U_{j}^{i}}{n}$$
(17)

where the  $\pi^{j}$  are independent uniform random permutations of the integers 1 through n, and the  $U_{i}^{j}$  are independent U[0, 1] random variables independent of the  $\pi_{j}$ .

Latin hypercube sampling was introduced by McKay et al. (1979) in what is widely considered to be the first paper on computer experiments. The sample points are stratified on each of p input axes. A common variant of Latin hypercube sampling has centered points

$$X_i^j = \frac{\pi^j(i) - 0.5}{n}.$$
 (18)

Point sets of this type were studied by Patterson (1954) who called them lattice samples.

Figure 13 shows a projection of 25 points from a (centered) Latin hypercube sample over 5 variables onto two of the coordinate axes. Each input variable gets explored in each of 25 equally spaced bins.

The stratification in Latin hypercube sampling usually reduces the variance of estimated integrals. Stein (1987) finds an expression for the variance of a sample mean under Latin hypercube sampling. Assuming that  $\int f(X)^2 dF < \infty$  write

$$f(X) = \mu + \sum_{j=1}^{p} \alpha_j(X^j) + e(X)$$
(19)



Fig. 13. 25 points of a Latin hypercube sample. The range of each input variable may be partitioned into 25 bins of equal width, drawn here with horizontal and vertical dotted lines, and each such bin contains one of the points.

where  $\mu = \int f(X) dF$  and  $\alpha_j(x) = \int_{X:X^j=x} (f(X) - \mu) dF_{-j}$  in which  $dF_{-j} = \prod_{k \neq j} dX^k$  is the uniform distribution over all input variables except the j'th. Equation (19) expresses f as the sum of a grand mean  $\mu$ , univariate main effects  $\alpha_j$  and a residual from additivity e(X).

Stein shows that under Latin hypercube sampling

$$\operatorname{Var}\left(\frac{1}{n}\sum_{i=1}^{n}f(x_{i})\right) = \frac{1}{n}\int e(X)^{2}\,\mathrm{d}F + \operatorname{o}\left(\frac{1}{n}\right) \tag{20}$$

whereas under iid sampling

$$\operatorname{Var}\left(\frac{1}{n}\sum_{i=1}^{n}f(x_{i})\right) = \frac{1}{n}\left(\int e(X)^{2}\,\mathrm{d}F + \sum_{j=1}^{p}\int\alpha_{j}(X^{j})^{2}\,\mathrm{d}F\right).$$
 (21)

By balancing the univariate margins, Latin hypercube sampling has removed the main effects of the function f from the error variance.

Owen (1992a) proves a central limit theorem for Latin hypercube sampling of bounded functions and Loh (1993) proves a central limit theorem under weaker conditions. For variance estimation in Latin hypercube sampling see (Stein, 1987; Owen, 1992a).

#### 7.4. Better Latin hypercubes

Latin hypercube samples look like random scatter in any bivariate plot, though they are quite regular in each univariate plot. Some effort has been made to find especially good Latin hypercube samples.

One approach has been to find Latin hypercube samples in which the input variables have small correlations. Iman and Conover (1982) perturbed Latin hypercube samples in a way that reduces off diagonal correlation. Owen (1994b) showed that the technique in Iman and Conover (1982) typically reduces off diagonal correlations by a factor of 3, and presented a method that empirically seemed to reduce the off diagonal correlations by a factor of order n from  $O(n^{-1/2})$  to  $O(n^{-3/2})$ . This removes certain bilinear terms from the lead term in the error. Dandekar (1993) found that iterating the method in Iman and Conover (1982) can lead to large improvements.

Small correlations are desirable but not sufficient, because one can construct centered Latin hypercube samples with zero correlation (unless n is equal to 2 modulo 4) which are nonetheless highly structured. For example the points could be arranged in a diamond shape in the plane, thus missing the center and corners of the input space.

Some researchers have looked for Latin hypercube samples having good properties when considered as designs for Bayesian prediction. Park (1994) studies the IMSE criterion and Morris and Mitchell (1995) consider entropy.

# 7.5. Randomized orthogonal arrays

An orthogonal array A is an n by p matrix of integers  $0 \leq A_i^j \leq b-1$ . The array has strength  $t \leq p$  if in every n by t submatrix of A all of the  $b^t$  possible rows appear the same number  $\lambda$  of times. Of course  $n = \lambda b^t$ .

Independently Owen (1992b, 1994a) and Tang (1992, 1993) considered using orthogonal arrays to improve upon Latin hypercube samples.

A randomized orthogonal array (Owen, 1992b) has two versions,

$$X_{i}^{j} = \frac{\pi_{j}(A_{i}^{j}) + U_{i}^{j}}{b}$$
(22)

and

$$X_i^j = \frac{\pi_j(A_i^j) + 0.5}{b}$$
(23)

just as Latin hypercube sampling has two versions. Indeed Latin hypercube sampling corresponds to strength t = 1, with  $\lambda = 1$ . Here the  $\pi_j$  are independent uniform



Fig. 14. 25 points of a randomly centered randomized orthogonal array. For whichever two (of five) variables that are plotted, there is one point in each reference square.

permutations of  $0, \ldots, b - 1$ . Patterson (1954) considered some schemes like the centered version.

If one were to plot the points of a randomized orthogonal array in t or fewer of the coordinates, the result would be a regular grid. The points of a randomized orthogonal array of strength 2 appear to be randomly scattered in 3 dimensions.

Figure 14 shows a projection of 25 points from a randomly centered randomized orthogonal array over 5 variables onto two of the coordinate points. Each pair of variables gets explored in each of 25 square bins. The plot for the centered version of a randomized orthogonal array is identical to that for a grid as shown in Figure 11.

The analysis of variance decomposition used above for Latin hypercube sampling can be extended to include interactions among 2 or more factors. See Efron and Stein (1981), Owen (1992b) and Wahba (1990) for details. Gu and Wahba (1993) describe how to estimate and form confidence intervals for these main effects in noisy data.

Owen (1992b) shows that main effects and interactions of t or fewer variables do not contribute to the asymptotic variance of a mean over a randomized orthogonal array, and Owen (1994a) shows that the variance is approximately  $n^{-1}$  times the sum of integrals of squares of interactions among more than t inputs.

Computer experiments



Fig. 15. 25 points of an orthogonal array based Latin hypercube sample. For whichever two (of five) variables that are plotted, there is one point in each reference square bounded by solid lines. Each variable is sampled once within each of 25 horizontal or vertical bins.

Tang (1993) introduced orthogonal array based Latin hypercube samples. The points of these designs are Latin hypercube samples  $X_i^j$ , such that  $\lfloor bX_i^j \rfloor$  is an orthogonal array. Here b is an integer and  $\lfloor z \rfloor$  is the smallest integer less than or equal to z. Tang (1993) shows that for a strength 2 array the main effects and two variable interactions do not contribute to the integration variance.

Figure 15 shows a projection of 25 points from an orthogonal array based Latin hypercube sample over 5 variables onto two of the coordinate points. Each variable individually gets explored in each of 25 equal bins and each pair of variables gets explored in each of 25 squares.

### 7.6. Scrambled nets

Orthogonal arrays were developed to balance discrete experimental factors. As seen above they can be embedded into the unit cube and randomized with the result that sampling variance is reduced. But numerical analysts and algebraists have developed some integration techniques directly adapted to balancing in a continuous space. Here we describe (t, m, s)-nets and their randomizations. A full account of (t, m, s)-nets



Fig. 16. 25 points of a scrambled (0, 2, 5)-net in base 5. For whichever two (of five) variables that are plotted, there is one point in each reference square. Each variable is sampled once within each of 25 equal bins.

is given by Niederreiter (1992). Their randomization is described by Owen (1995, 1996a).

Let  $p = s \ge 1$  and  $b \ge 2$  be integers. An elementary subcube in base b is of the form

$$E = \prod_{j=1}^{s} \left[ \frac{c_j}{b^{k_j}}, \frac{c_j+1}{b^{k_j}} \right)$$

for integers  $k_j$ ,  $c_j$  with  $k_j \ge 0$  and  $0 \le c_j < b^{k_j}$ .

Let  $m \ge 0$  be an integer. A set of points  $X_i$ ,  $i = 1, ..., b^m$ , of from  $[0, 1)^s$  is a (0, m, s)-net in base b if every elementary subcube E in base b of volume  $b^{-m}$ has exactly 1 of the points. That is, every cell that "should" have one point of the sequence does have one point of the sequence.

This is a very strong form of equidistribution and by weakening it somewhat, constructions for more values of s and b become available. Let  $t \leq m$  be a nonnegative integer. A finite set of  $b^m$  points from  $[0,1)^s$  is a (t,m,s)-net in base b if every elementary subcube in base b of volume  $b^{t-m}$  contains exactly  $b^t$  points of the sequence.



Fig. 17. The 125 points of a scrambled (0, 3, 5)-net in base 5. For whichever two (of five) variables that are plotted, the result is a 5 by 5 grid of 5 point Latin hypercube samples. Each variable is sampled once within each of 125 equal bins. Each triple of variables can be partitioned into 125 congruent cubes, each of which has one point.

Cells that "should" have  $b^t$  points do have  $b^t$  points, though cells that "should" have 1 point might not.

By common usage the name (t, m, s)-net assumes that the letter s is used to denote the dimension of the input space, though one could speak of (t, m, p)-nets. Another convention to note is that the subcubes are half-open. This makes it convenient to partition the input space into congruent subcubes.

The balance properties of a (t, m, s)-net are greater than those of an orthogonal array. If  $X_i^j$  is a (t, m, s)-net in base b then  $\lfloor bX_i^j \rfloor$  is an orthogonal array of strength  $\min\{s, m-t\}$ . But the net also has balance properties when rounded to different powers of b on all axes, so long as the powers sum to no more than m-t. Thus the net combines aspects of orthogonal arrays and multi-level orthogonal arrays all in one point set.

In the case of a (0, 4, 5)-net in base 5, one has 625 points in  $[0, 1)^5$  and one can count that there are 43750 elementary subcubes of volume 1/625 of varying aspect ratios each of which has one of the 625 points.



Fig. 18. The 625 points of a scrambled (0, 4, 5)-net in base 5. For whichever two (of five) variables that are plotted, the square can be divided into 625 squares of side 1/25 or into 625 rectangles of side 1/5 by 1/125 or into 625 rectangles of side 1/125 by 1/5 and each such rectangle has one of the points. Each variable is sampled once within each of 625 equal bins. Each triple of variables can be partitioned into 625 hyperrectangles in three different ways and each such hyperrectangle has one of the points. Each quadruple of variables can be partitioned into 625 congruent hypercubes of side 1/5, each of which has one point.

For  $t \ge 0$ , an infinite sequence  $(X_i)_{i\ge 1}$  of points from  $[0,1)^s$  is a (t,s)-sequence in base b if for all  $k \ge 0$  and  $m \ge t$  the finite sequence  $(X_i)_{i=kb^m+1}^{(k+1)b^m}$  is a (t,m,s)-net in base b.

The advantage of a (t, s)-sequence is that if one finds that the first  $b^m$  points are not sufficient for an integration problem, one can find another  $b^m$  points that also form a (t, m, s)-net and tend to fill in places not occupied by the first set. If one continues to the point of having b such (t, m, s)-nets, then the complete set of points comprises a (t, m + 1, s)-net.

The theory of (t, m, s)-nets and (t, s)-sequences is given in Niederreiter (1992). A famous result of the theory is that integration over a (t, m, s)-net can attain an accuracy of order  $O(\log(n)^{s-1}/n)$  while restricting to (t, s)-sequences raises this slightly to  $O(\log(n)^s/n)$ . These results require that the integrand be of bounded variation in the sense of Hardy and Krause. For large s, it takes unrealistically large n for these rates

to be clearly better than  $n^{-1/2}$  but in examples they seem to outperform simple Monte Carlo.

The construction of (t, m, s)-nets and (t, s)-sequences is also described in Niederreiter (1992). Here we remark that for prime numbers s a construction by Faure (1982) gives (0, s)-nets in base s and Niederreiter extended the method to prime powers s. (See Niederreiter, 1992.) Thus one can choose b to be the smallest prime power greater than or equal to s and use the first s variables of the corresponding (0, b)-sequence in base b.

Owen (1995) describes a scheme to randomize (t, m, s)-nets and (t, s)-sequences. The points are written in a base b expansion and certain random permutations are applied to the coefficients in the expansion. The result is to make each permuted  $X_i$ uniformly distributed over  $[0, 1)^s$  while preserving the (t, m, s)-net or (t, s)-sequence structure of the ensemble of  $X_i$ . Thus the sample estimate  $n^{-1} \sum_{i=1}^n f(X_i)$  is unbiased for  $\int f(X) dF$  and the variance of it may be estimated by replication. On some test integrands in (Owen, 1995) the randomized nets outperformed their unrandomized counterparts. It appears that the unscrambled nets have considerable structure, stemming from the algebra underlying them, and that this structure is a liability in integration.

Figure 16 shows the 25 points of a scrambled (0, 2, 5)-net in base 5 projected onto two of the five input coordinates. These points are the initial 25 points of a (0, 5)sequence in base 5. This design has the equidistribution properties of an orthogonal array based Latin hypercube sample. Moreover every consecutive 25 points in the sequence  $X_{25a+1}, X_{25a+2}, \ldots, X_{25(a+1)}$  has these equidistribution properties. The first 125 points, shown in Figure 17 have still more equidistribution properties: any triple of the input variables can be split into 125 subcubes each with one of the  $X_i$ , in any pair of variables the points appear as a 5 by 5 grid of 5 point Latin hypercube samples and each individual input variable can be split into 125 cells each having one point. The first 625 points, are shown in Figure 18.

Owen (1996a) finds a variance formula for means over randomized (t, m, s)-nets and (t, s)-sequences. The formula involves a wavelet-like anova combining nested terms on each coordinate, all crossed against each other. It turns out that for any square integrable integrand, the resulting variance is  $o(n^{-1})$  and it therefore beats any of the usual variance reduction techniques, which typically only reduce the asymptotic coefficient of  $n^{-1}$ .

For smooth integrands with s = 1, the variance is in fact  $O(n^{-3})$  and in the general case Owen (1996b) shows that the variance is  $O(n^{-3}(\log n)^{s-1})$ .

### 8. Selected applications

One of the largest fields using and developing deterministic simulators is in the designing and manufacturing of VLSI circuits. Alvarez et al. (1988) describe the use of SUPREM-III (Ho et al., 1984) and SEDAN-II (Yu et al., 1982) in designing BIMOS devices for manufacturability. Aoki et al. (1987), use CADDETH a two dimensional device simulator, for optimizing devices and for accurate prediction of device sensitivities. Sharifzadeh et al. (1989) use SUPREME-III and PISCES-II (Pinto et al., 1984) to compute CMOS device characteristics as a function of the designable technology parameters. Nasif et al. (1984) describe the use of FABRICS-II to estimate circuit delay times in integrated circuits.

The input variables for the above work are generally device sizes, metal concentrations, implant doses and gate oxide temperatures. The multiple responses are threshold voltages, subthreshold slopes, saturation currents and linear transconductance although the output variables of concern depend on the technology under investigation. The engineers use the physical/numerical simulators to assist them in optimizing process, device, and circuit design before the costly step of building prototype devices. They are also concerned with minimizing transmitted variability as this can significantly reduce the performance of the devices and hence reduce yield. For example, Welch et al. (1990), Currin et al. (1991) and Sacks et al. (1989b) discuss the use of simulators to investigate the effect of transistor dimensions on the asynchronization of two clocks. They want to find the combination of transistor widths that produce zero clock skews with very small transmitted variability due to uncontrollable manufacturing variability in the transistors.

TIMS, a simulator developed by T. Osswald and C. L. Tucker III, helps in optimizing a compression mold filling process for manufacturing automobiles (Church et al., 1988). In this process a sheet of molding compound is cut and placed in a heated mold. The mold is slowly closed and a constant force is applied during the curing reaction. The controlling variables of the process are the geometry and thickness of the part, the compound viscosity, shape and location within the charge, and the mold closing speed. The simulator then predicts the position of the flow front as a function of time.

Miller and Frenklach (1983) discuss the use of computers to solve systems of differential equations describing chemical kinetic models. In their work, the inputs to the simulator are vectors of possibly unknown combustion rate constants and the outputs are induction-delay times and concentrations of chemical species at specified reaction times. The objectives of their investigations are to find values of the rate constants that agree with experimental data and to find the most important rate constant to the process. Sacks et al. (1989a) explore some of the design issues and applications to this field.

TWOLAYER, a thermal energy storage model developed by Alan Solomon and his colleagues at the Oak Ridge National Laboratory, simulates heat transfer through a wall containing two layers of different phase change material. Currin et al. (1991) utilize TWOLAYER in a computer experiment. The inputs into TWOLAYER are the layers dimensions, the thermal properties of the materials and the characteristics of the heat source. The object of interest was finding the configuration of the input variables that produce the highest value of a heat storage utility index.

FOAM (Bartell et al., 1981) models the transport of polycyclic aromatic hydrocarbon spills in streams using structure activity relationships. Bartell et al. (1983) modified this model to predict the fate of anthracene when introduced into ponds. This model tracks the "evaporation and dissolution of anthracene from a surface slick of synthetic oil, volatilization and photolytic degradation of dissolved anthracene, sorption to suspended particulate matter and sediments and accumulation by pond biota" (Bartell, 1983). They used Monte Carlo error analyses to assess the effect of the uncertainty in model parameters on their results.

# References

- Alvarez, A. R., B. L. Abdi, D. L. Young, H. D. Weed, J. Teplik and E. Herald (1988). Application of statistical design and response surface methods to computer-aided VLSI device design. *IEEE Trans. Comput. Aided Design* 7(2), 271–288.
- Aoki, Y., H. Masuda, S. Shimada and S. Sato (1987). A new design-centering methodology for VLSI device development. *IEEE Trans. Comput. Aided Design* 6(3), 452–461.
- Bartell, S. M., R. H. Gardner, R. V. O'Neill and J. M. Giddings (1983). Error analysis of predicted fate of anthracene in a simulated pond. *Environ. Toxicol. Chem.* 2, 19–28.
- Bartell, S. M., J. P. Landrum, J. P. Giesy and G. J. Leversee (1981). Simulated transport of polycyclic aromatic hydrocarbons in artificial streams. In: W. J. Mitch, R. W. Bosserman and J. M. Klopatek, eds., *Energy and Ecological Modelling*. Elsevier, New York, 133–143.
- Bates, R. A., R. J. Buck, E. Riccomagno and H. P. Wynn (1996). Experimental design and observation for large systems (with discussion). J. Roy. Statist. Soc. Ser. B 58(1), 77–94.
- Borth, D. M. (1975). A total entropy criterion for the dual problem of model discrimination and parameter estimation. J. Roy. Statist. Soc. Ser. B 37, 77–87.
- Box, G. E. P. and N. R. Draper (1959). A basis for the selection of a response surface design. J. Amer. Statist. Assoc. 54, 622–654.
- Box, G. E. P. and N. R. Draper (1963). The choice of a second order rotatable design. *Biometrika* 50, 335-352.
- Box, G. E. P. and W. J. Hill (1967). Discrimination among mechanistic models. Technometrics 9, 57-70.
- Church, A., T. Mitchell and D. Fleming (1988). Computer experiments to optimize a compression mold filling process. Talk given at the Workshop on Design for Computer Experiments in Oak Ridge, TN, November.
- Cranley, R. and T. N. L. Patterson (1976). Randomization of number theoretic methods for multiple integration. SIAM J. Numer. Anal. 23, 904–914.
- Cressie, N. A. C. (1986). Kriging nonstationary data. J. Amer. Statist. Assoc. 81, 625-634.
- Cressie, N. A. C. (1993). Statistics for Spatial Data (Revised edition). Wiley, New York.
- Currin, C., M. Mitchell, M. Morris and D. Ylvisaker (1991). Bayesian prediction of deterministic functions, with applications to the design and analysis of computer experiments. J. Amer. Statist. Assoc. 86, 953–963.
- Dandekar, R. (1993). Performance improvement of restricted pairing algorithm for Latin hypercube sampling Draft Report, Energy Information Administration, U.S.D.O.E.
- Davis, P. J. and P. Rabinowitz (1984). *Methods of Numerical Integration*, 2nd. edn. Academic Press, San Diego.
- Diaconis, P. (1988). Bayesian numerical analysis In: S. S. Gupta and J. O. Berger, eds., Statistical Decision Theory and Related Topics IV, Vol. 1. Springer, New York, 163–176.
- Efron, B. and C. Stein (1981). The jackknife estimate of variance. Ann. Statist. 9, 586-596.
- Fang, K. T. and Y. Wang (1994). Number-theoretic Methods in Statistics. Chapman and Hall, London.
- Faure, H. (1982). Discrépances des suites associées à un système de numération (en dimension s). Acta Arithmetica 41, 337-351.
- Friedman, J. H. (1991). Multivariate adaptive regression splines (with Discussion). Ann. Statist. 19, 1-67.
- Gill, P. E., W. Murray, M. A. Saunders and M. H. Wright (1986). User's guide for npsol (version 4.0): A Fortran package for nonlinear programming. SOL 86-2, Stanford Optimization Laboratory, Dept. of Operations Research, Stanford University, California, 94305, January.
- Gill, P. E., W. Murray and M. H. Wright (1981). Practical Optimization. Academic Press, London.
- Gordon, W. J. (1971). Blending function methods of bivariate and multivariate interpolation and approximation. SIAM J. Numer. Anal. 8, 158-177.
- Gu, C. and G. Wahba (1993). Smoothing spline ANOVA with component-wise Bayesian "confidence intervals". J. Comp. Graph. Statist. 2, 97–117.
- Hickernell, F. J. (1996). Quadrature error bounds with applications to lattice rules. SIAM J. Numer. Anal. 33 (in press).
- Ho, S. P., S. E. Hansen and P. M. Fahey (1984). Suprem III a program for integrated circuit process modeling and simulation. TR-SEL84 1, Stanford Electronics Laboratories.

- Iman, R. L. and W. J. Conover (1982). A distributon-free approach to inducing rank correlation among input variables. Comm. Statist. B11(3), 311-334.
- Johnson, M. E., L. M. Moore and D. Ylvisaker (1990). Minimax and maximin distance designs. J. Statist. Plann. Inference 26, 131-148.
- Journel, A. G. and C. J. Huijbregts (1978). Mining Geostatistics. Academic Press, London.
- Koehler, J. R. (1990). Design and estimation issues in computer experiments. Dissertation, Dept. of Statistics, Stanford University.
- Lindley, D. V. (1956). On a measure of the information provided by an experiment. Ann. Math. Statist. 27, 986-1005.
- Loh, W.-L. (1993). On Latin hypercube sampling. Tech. Report No. 93-52, Dept. of Statistics, Purdue University.
- Loh, W.-L. (1994). A combinatorial central limit theorem for randomized orthogonal array sampling designs. Tech. Report No. 94-4, Dept. of Statistics, Purdue University.
- Mardia, K. V. and R. J. Marshall (1984). Maximum likelihood estimation of models for residual covariance in spatial regression. *Biometrika* 71(1), 135–146.
- Matérn, B. (1947). Method of estimating the accuracy of line and sample plot surveys. *Medd. Skogsforskn Inst.* 36(1).
- Matheron, G. (1963). Principles of geostatistics. Econom. Geol. 58, 1246-1266.
- McKay, M. (1995). Evaluating prediction uncertainty. Report NUREG/CR-6311, Los Alamos National Laboratory.
- McKay, M., R. Beckman and W. Conover (1979). A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* **21**(2), 239–245.
- Miller, D. and M. Frenklach (1983). Sensitivity analysis and parameter estimation in dynamic modeling of chemical kinetics. *Internat. J. Chem. Kinetics* 15, 677–696.
- Mitchell, T. J. (1974). An algorithm for the construction of 'D-optimal' experimental designs. *Technometrics* **16**, 203–210.
- Mitchell, T., M. Morris and D. Ylvisaker (1990). Existence of smoothed stationary processes on an interval. Stochastic Process. Appl. 35, 109–119.
- Mitchell, T., M. Morris and D. Ylvisaker (1995). Two-level fractional factorials and Bayesian prediction. Statist. Sinica 5, 559–573.
- Mitchell, T. J. and D. S. Scott (1987). A computer program for the design of group testing experiments. Comm. Statist. Theory Methods 16, 2943–2955.
- Morris, M. D. and T. J. Mitchell (1995). Exploratory designs for computational experiments. J. Statist. Plann. Inference 43, 381–402.
- Morris, M. D., T. J. Mitchell and D. Ylvisaker (1993). Bayesian design and analysis of computer experiments: Use of derivative in surface prediction. *Technometrics* 35(3), 243–255.
- Nassif, S. R., A. J. Strojwas and S. W. Director (1984). FABRICS II: A statistically based IC fabrication process simulator. *IEEE Trans. Comput. Aided Design* 3, 40–46.
- Niederreiter, H. (1992). Random Number Generation and Quasi-Monte Carlo Methods. SIAM, Philadelphia, PA.
- O'Hagan, A. (1989). Comment: Design and analysis of computer experiments. Statist. Sci. 4(4), 430-432.
- Owen, A. B. (1992a). A central limit theorem for Latin hypercube sampling. J. Roy. Statist. Soc. Ser. B 54, 541-551.
- Owen, A. B. (1992b). Orthogonal arrays for computer experiments, integration and visualization. *Statist. Sinica* 2, 439–452.
- Owen, A. B. (1994a). Lattice sampling revisited: Monte Carlo variance of means over randomized orthogonal arrays. Ann. Statist. 22, 930–945.
- Owen, A. B. (1994b). Controlling correlations in latin hypercube samples. J. Amer. Statist. Assoc. 89, 1517-1522.
- Owen, A. B. (1995). Randomly permuted (t, m, s)-nets and (t, s)-sequences. In: H. Niederreiter and P. J.-S. Shiue, eds., Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing. Springer, New York, 299–317.
- Owen, A. B. (1996a). Monte Carlo variance of scrambled net quadrature. SIAM J. Numer. Anal., to appear.

- Owen, A. B. (1996b). Scrambled net variance for integrals of smooth functions. Tech. Report Number 493, Department of Statistics, Stanford University.
- Paskov, S. H. (1993). Average case complexity of multivariate integration for smooth functions. J. Complexity 9, 291-312.
- Park, J.-S. (1994) Optimal Latin-hypercube designs for computer experiments. J. Statist. Plann. Inference **39**, 95–111.
- Parzen, A. B. (1962). Stochastic Processes. Holden-Day, San Francisco, CA.
- Patterson, H. D. (1954). The errors of lattice sampling. J. Roy. Statist. Soc. Ser. B 16, 140-149.
- Phadke, M. (1988). Quality Engineering Using Robust Design. Prentice-Hall, Englewood Cliffs, NJ.
- Pinto, M. R., C. S. Rafferty and R. W. Dutton (1984). PISCES-II-posson and continuity equation solver. DAGG-29-83-k 0125, Stanford Electron. Lab.
- Ripley, B. (1981). Spatial Statistics. Wiley, New York.
- Ritter, K. (1995). Average case analysis of numerical problems. Dissertation, University of Erlangen.
- Ritter, K., G. Wasilkowski and H. Wozniakowski (1993). On multivariate integration for stochastic processes. In: H. Brass and G. Hammerlin, eds., *Numerical Integration*, Birkhauser, Basel, 331–347.
- Ritter, K., G. Wasilkowski and H. Wozniakowski (1995). Multivariate integration and approximation for random fields satisfying Sacks-Ylvisaker conditions. Ann. Appl. Prob. 5, 518-540.
- Roosen, C. B. (1995). Visualization and exploration of high-dimensional functions using the functional ANOVA decomposition. Dissertation, Dept. of Statistics, Stanford University.
- Sacks, J. and S. Schiller (1988). Spatial designs. In: S. S. Gupta and J. O. Berger, eds., Statistical Decision Theory and Related Topics IV, Vol. 2. Springer, New York, 385–399.
- Sacks, J., S. B. Schiller and W. J. Welch (1989). Designs for computer experiments. *Technometrics* **31**(1), 41-47.
- Sacks, J., W. J. Welch, T. J. Mitchell and H. P. Wynn (1989). Design and analysis of computer experiments. Statist. Sci. 4(4), 409–423.
- Shannon, C. E. (1948). A mathematical theory of communication. Bell Syst. Tech. J. 27, 379-423, 623-656.
- Sharifzadeh, S., J. R. Koehler, A. B. Owen and J. D. Shott (1989). Using simulators to model transmitted variability in IC manufacturing. *IEEE Trans. Semicond. Manufact.* 2(3), 82–93.
- Shewry, M. C. and H. P. Wynn (1987). Maximum entropy sampling. J. Appl. Statist. 14, 165-170.
- Shewry, M. C. and H. P. Wynn (1988). Maximum entropy sampling and simulation codes. In: Proc. 12th World Congress on Scientific Computation, Vol. 2, IMAC88, 517–519.
- Sloan, I. H. and S. Joe (1994). Lattice Methods for Multiple Integration. Oxford Science Publications, Oxford.
- Smolyak, S. A. (1963). Quadrature and interpolation formulas for tensor products of certain classes of functions. Soviet Math. Dokl. 4, 240–243.
- Stein, M. L. (1987). Large sample properties of simulations using Latin hypercube sampling. *Technometrics* **29**(2), 143–151.
- Stein, M. L. (1989). Comment: Design and analysis of computer experiments. Statist. Sci. 4(4), 432-433.
- Steinberg, D. M. (1985). Model robust response surface designs: Scaling two-level factorials. *Biometrika* **72**, 513–26.
- Tang, B. (1992). Latin hypercubes and supersaturated designs. Dissertation, Dept. of Statistics and Actuarial Science, University of Waterloo.
- Tang, B. (1993). Orthogonal array-based Latin hypercubes. J. Amer. Statist. Assoc. 88, 1392-1397.
- Wahba, G. (1978). Interpolating surfaces: High order convergence rates and their associated designs, with applications to X-ray image reconstruction. Tech. report 523, Statistics Department, University of Wisconsin, Madison.
- Wahba, G. (1990). Spline Models for Observational Data. CBMS-NSF Regional Conference Series in Applied Mathematics, Vol. 59. SIAM, Philadelphia, PA.
- Wasilkowski, G. (1993). Integration and approximation of multivariate functions: Average case complexity with Wiener measure. Bull. Amer. Math. Soc. (N. S.) 28, 308–314. Full version J. Approx. Theory 77, 212–227.
- Wozniakowski H. (1991). Average case complexity of multivariate integration. Bull. Amer. Math. Soc. (N. S.) 24, 185–194.

- Welch, W. J. (1983). A mean squared error criterion for the design of experiments. *Biometrika* 70(1), 201-213.
- Welch, W. Yu, T. Kang and J. Sacks (1990). Computer experiments for quality control by parameter design. J. Quality Technol. 22, 15–22.
- Welch, W. J., J. R. Buck, J. Sacks, H. P. Wynn, T. J. Mitchell and M. D. Morris. Screening, prediction, and computer experiments. *Technometrics* 34(1), 15–25.
- Yaglom, A. M. (1987). Correlation Theory of Stationary and Related Random Functions, Vol. 1. Springer, New York.
- Ylvisaker, D. (1975). Designs on random fields. In: J. N. Srivastava, ed., A Survey of Statistical Design and Linear Models. North-Holland, Amsterdam, 593-607.
- Young, A. S. (1977). A Bayesian approach to prediction using polynomials. Biometrika 64, 309-317.
- Yu, Z., G. G. Y. Chang and R. W. Dutton (1982). Supplementary report on sedan II. TR-G201 12, Stanford Electronics Laboratories.