

Chemical Library Subset Selection Algorithms: A Unified Derivation Using Spatial Statistics

Fred A. Hamprecht,^{*,†} Walter Thiel,[‡] and Wilfred F. van Gunsteren[§]

Interdisciplinary Center for Scientific Computing (IWR), University of Heidelberg,
Im Neuenheimer Feld 368, 69120 Heidelberg, Germany, Max-Planck-Institut für Kohlenforschung,
45470 Mülheim an der Ruhr, Germany, and Laboratory of Physical Chemistry,
Swiss Federal Institute of Technology, 8092 Zürich, Switzerland

Received July 12, 2001

If similar compounds have similar activity, rational subset selection becomes superior to random selection in screening for pharmacological lead discovery programs. Traditional approaches to this experimental design problem fall into two classes: (i) a linear or quadratic response function is assumed (ii) some space filling criterion is optimized. The assumptions underlying the first approach are clear but not always defensible; the second approach yields more intuitive designs but lacks a clear theoretical foundation. We model activity in a bioassay as realization of a stochastic process and use the best linear unbiased estimator to construct spatial sampling designs that optimize the integrated mean square prediction error, the maximum mean square prediction error, or the entropy. We argue that our approach constitutes a unifying framework encompassing most proposed techniques as limiting cases and sheds light on their underlying assumptions. In particular, vector quantization is obtained, in dimensions up to eight, in the limiting case of very smooth response surfaces for the integrated mean square error criterion. Closest packing is obtained for very rough surfaces under the integrated mean square error and entropy criteria. We suggest to use either the integrated mean square prediction error or the entropy as optimization criteria rather than approximations thereof and propose a scheme for direct iterative minimization of the integrated mean square prediction error. Finally, we discuss how the quality of chemical descriptors manifests itself and clarify the assumptions underlying the selection of diverse or representative subsets.

1. INTRODUCTION

You should call it “entropy” and for two reasons: first, the function is already in use in thermodynamics under that name; second, and more importantly, most people don’t know what entropy really is, and if you use the word “entropy” in an argument you will win every time! *John von Neumann to Claude Shannon when prompted for a suitable term; quoted after ref 1.*

This paper is concerned with optimal spatial sampling from a finite set under the assumption of a smooth response function that cannot be well modeled by low polynomial terms.

To put things in a chemical context, we are concerned with the primary steps of lead discovery (as opposed to lead optimization) for pharmaceutical or agrochemical research in a collection of N distinct chemical compounds. Once a suitable bioassay is available, exhaustive testing of all N compounds is usually impossible due to restrictions on time and availability. Consequently, a subset of $n \ll N$ must be selected that should maximize the chance of locating an active. Ignoring all chemical and biochemical knowledge, simple random selection without replacement would be the best thing one can do.

1.1. Molecular Recognition. In reality, though, the biological activity of a drug is usually the consequence of its selective binding to a specific protein or receptor. The central ingredient of this simple causal model is the concept of molecular recognition. It predicts affinity if the total strength of interactions between drug and receptor exceeds those with competing substances, notably the solvent.

Now, if two molecules “look” almost the same to a receptor as well as to the solvent, their interaction strengths should be almost the same, and the elicited activity can be expected to be almost the same. When selecting a subset, it then makes sense to avoid choosing molecules that can be expected to look very similar to the target under study—the evoked response is probably similar and the experiment would be better invested in a dissimilar molecule.

1.2. And Its Descriptors. A first problem, then, lies in determining how similar two molecules look to each other, as judged by the receptor and solvent. The difficulty is increased by the fact that two molecules may be judged similar by one receptor but dissimilar by another;² for instance, some receptors cannot discriminate between two isomers while others have a clear preference. Once molecular recognition is assumed as causal model for activity, it is near at hand to characterize the similarity of molecules to a specific receptor by a set of descriptors that can explain the affinity or absence thereof. To be useful, these descriptors should be easy to come by, that is, without performing the actual bioassay or something closely related to it. Popular choices are physicochemical characteristics such as the

* Corresponding author fax: ++49-6221-548850; e-mail: f.hamprecht@alumni.ethz.ch.

[†] University of Heidelberg.

[‡] Max-Planck-Institut für Kohlenforschung.

[§] Swiss Federal Institute of Technology.

molecular weight or the octanol/water partition coefficient, shape descriptors such as steric or electrostatic fields, geometrical descriptors of pharmacophores, or dichotomous parameters indicating the presence or absence of, say, a functional group (e.g., references in refs 3–5). Deriving suitable descriptors is a science and art requiring much chemical and problem-specific expertise, and we will not discuss it in any detail but give some general remarks instead:

1.2.1. What Are Good Descriptors? In our mind, the value of a set of descriptors and associated similarity measure can be established only after the activities of all compounds have been assayed. The combination of good descriptors with a suitable similarity measure then leads to similar compounds (as measured by their distance in descriptor space) showing similar activity. The converse need not be true: Similar activity need not imply vicinity in descriptor space, in particular the response function may feature multiple maxima. This “similar property principle”⁶ is old and has been described using various names with slightly different emphasis, such as “neighborhood behavior”,⁷ “similarity radius”,⁸ or “sampling radius”.⁹ The principle is illustrated beautifully with real data in Figure 19.1 of Martin et al.¹⁰

A compound characterized by descriptors on the interval scale (e.g., most physicochemical descriptors) can be considered as a single point in descriptor space. An embedding can also be obtained for compounds characterized on the nominal scale (by bit-strings or fingerprints) as follows: the similarity between any two compounds can be measured by one of a plethora of association coefficients such as the Tanimoto, Dice, Jaccard, etc.^{11–13} These pairwise similarities can be transformed to dissimilarities and the compounds can be embedded in a metric space using multidimensional scaling (e.g., ref 12). In the following, we will use the Euclidean distance $d(\cdot, \cdot)$ in a metric descriptor space to describe compound similarity although other choices are possible.

1.3. Beyond Random Subset Selection: Previous Work. We now turn to the second problem: once it can be assumed that the compounds have been embedded in a space such that their distances reflect their similarity, as judged by the system under study, how should one go about selecting a subset to avoid redundant experiments and gain the maximum amount of information possible? There are two distinct approaches in the literature:

1.3.1. Optimal Design or Regression with Uncorrelated Errors. The response in descriptor space is assumed to obey the form

$$Z(x) = \beta F(x) + \epsilon(x) \quad (1)$$

where β are the coefficients of first- or second-order polynomials $F(x)$ and $\epsilon(x)$ is an uncorrelated zero-mean random error. Under this assumption, the theory of optimum design of experiments¹⁴ then yields optimal sampling designs that minimize the variance of the parameter estimates for the model (D-optimality) or a related quantity (e.g., ref 15).

Such optimal designs for polynomials up to second order have been used in chemical library design.¹⁶ In lead optimization and quantitative structure–activity relationships (QSAR), the assumption that the response function can be locally approximated by low polynomial terms often works. However, we follow^{17,18} and others in arguing that polynomials, especially of low orders one and two, are not necessarily

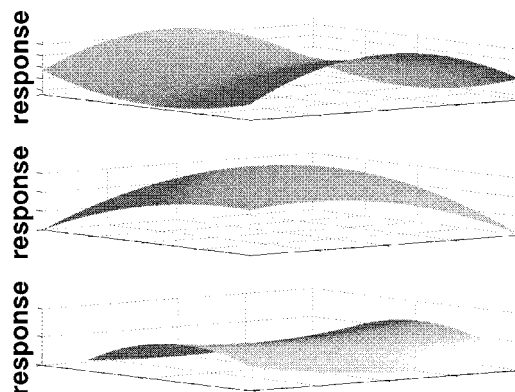


Figure 1. Polynomial models of response surfaces. Shown are some second-order polynomials with coefficients sampled uniformly from the interval $[-1, 1]$.

good basis functions for or approximations to the response functions encountered in lead discovery. Optimal designs of low orders are extreme in that the hull of the cloud of points is well sampled, while the inside is almost entirely neglected; but they are, indeed, optimal when the underlying assumption (eq 1) holds, i.e., if you can assume that the unknown “true” response surface (with the uncorrelated noise ϵ averaged out) is of the kind shown in Figure 1.

1.3.2. Space-Filling Designs. If, on the other hand, the response surface is likely to be irregular, for instance featuring multiple maxima as in Figure 2, it seems natural to sample space more evenly. This intuition has led to a quest for criteria and algorithms with more space-filling properties: the subset should sample the entire cloud of points (including its interior) evenly, and compounds should be selected in a nonredundant manner. The result has been a great number of papers on the subject, some of which will be discussed in section 6.

Summing up, D-optimal design is a strategy that rests on clear assumptions. Low-order polynomials cannot provide a good model for a response surface with a complex topography. The latter case is addressed by the intuitive concept of “space-fillingness” which, however, has remained without theoretical justification or a quantitative statement of the assumptions involved.

2. KRIGING OR BEST LINEAR UNBIASED ESTIMATION

Reiterating our assumptions, we wish to choose a subset from a number of points in a metric space so as to maximize our chance of locating an active. Nothing is known about the response function except that it is “reasonably” smooth, in the sense of proximate points featuring similar response. Experiments will be performed only after the entire subset has been selected. Intuitively, we understand that if the response surface has some smoothness, i.e., if it shows spatial correlation, performing an experiment on a given compound will also teach us something about the response in its immediate environment “for free”. How far this “immediate” environment extends depends on the degree of spatial correlation. We need a mathematical framework quantifying the total amount of “free” information we gain depending on which points we select for the subset.

Luckily, this problem has already surfaced in the exploration of ore deposits and oil fields, where each sample

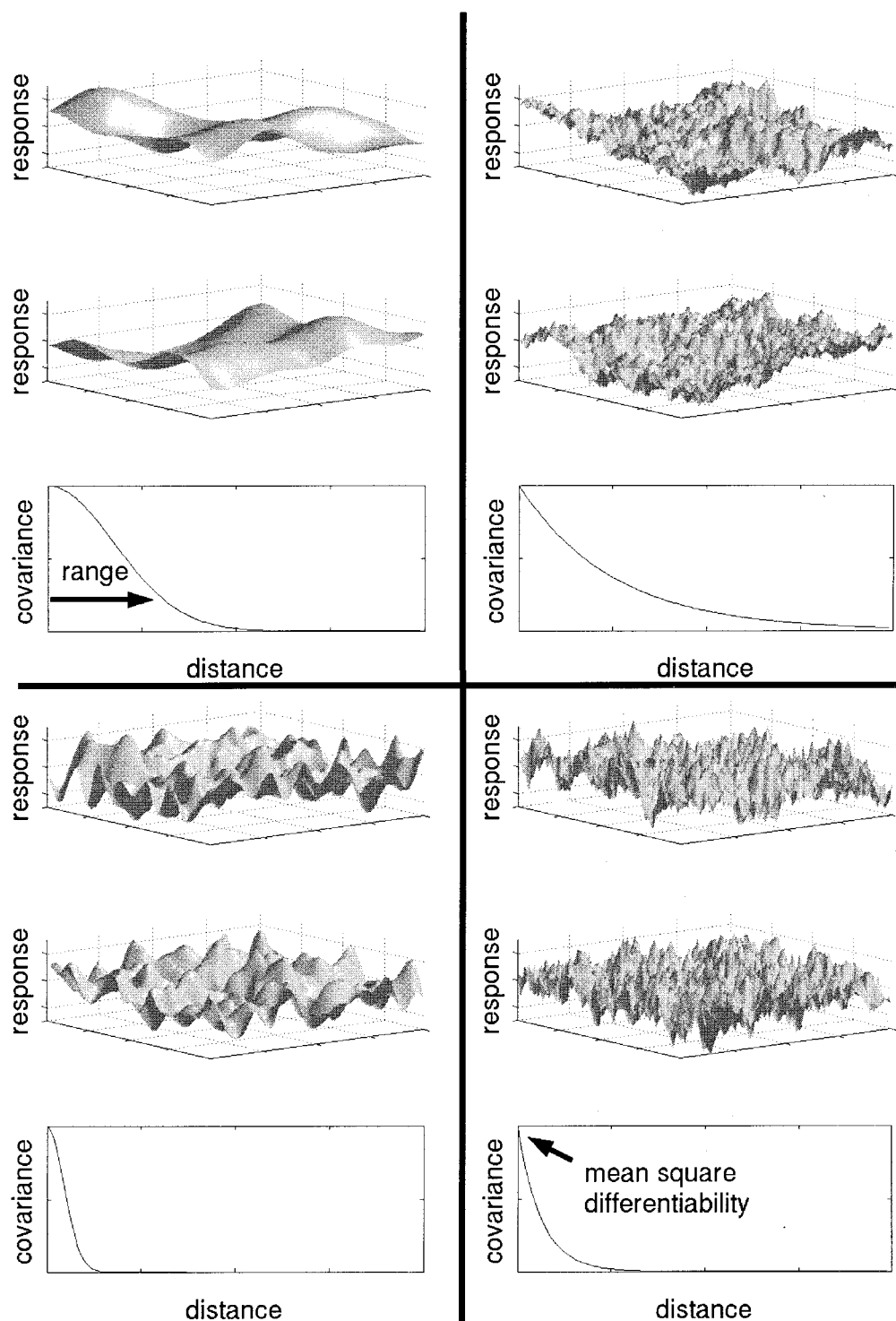


Figure 2. Stochastic models of response surfaces. Shown are two realizations each of four stochastic processes with their corresponding isotropic (circularly symmetric) covariance functions. Left: processes with a Gaussian covariance function (which have an infinite degree of mean square differentiability at the origin). Right: processes with an exponential covariance function (with zero mean square differentiability at the origin). Top: long range. Bottom: short range. All simulations performed with gstat.⁸⁶

(borehole) causes considerable costs such that a maximum of information needs to be extracted from a given number of samples. The interpolator of choice is known as the “best linear unbiased estimator” in spatial statistics and as “kriging” [named after the mining engineer D. G. Krige] in geostatistics.¹⁹ The salient feature of the kriging estimator is that it gives not only an estimate of the response between measurements (i.e., the interpolation) but also the uncertainty of that estimate. Intuitively, the uncertainty can be expected

to be small close to the experiments and to grow with distance from these. Minimizing the uncertainty maximizes the information available and thus the kriging estimator can be used to construct optimal designs. We now introduce the mathematics that allow a quantification of these ideas.

First, the application of kriging requires modeling of the activity or response surface as a particular realization of a stochastic spatial process, also known as random field. Such processes arise in the description of systems that evolve in

space under probabilistic laws. Just like a random number generator, a stochastic process can produce an arbitrary number of realizations. For a random number generator, each realization is a single number; whereas for a stochastic process, each realization is an entire surface. Four different Gaussian [A stochastic process is characterized by its finite-dimensional distributions.²⁰ For Gaussian processes, these are multivariate normal.] stochastic processes are introduced in Figure 2: each is characterized by two out of its infinite number of possible realizations, as well as by its covariance function. What distinguishes these stochastic processes is their smoothness, and these differences are reflected in their covariance functions (see below for definition).

Making no notational distinction between a random process and a realization thereof, let $Z(x)$ be a random field modeling our response surface. Assume that $Z(x)$ is second-order stationary, i.e.,

$$E[Z(x)] = m = \text{const.}$$

$$E[(Z(x_1) - m)(Z(x_2) - m)] = \text{Cov}(x_1, x_2) = \text{Cov}(x_2 - x_1)$$

E denotes expectation, i.e., an average over all possible realizations of the random process. In words, the average of all realizations at a specific coordinate x should be equal to the constant m , independent of x . [This assumption of zero drift can be relaxed, see section 4.2.] The covariance $\text{Cov}(x_1, x_2)$ at points two points x_1, x_2 should depend only on their relative distance and orientation. The covariance function then describes how similar the response is, on average, as a function of distance and orientation. To simplify the following calculations, shift the entire random process by m such that $E[Z(x)] = 0$ and $\text{Cov}(x_2 - x_1) = E[Z(x_1)Z(x_2)]$.

Consider the measured responses as samples $Z(y)$ from the realization of a random field at a set of points $y = \{y_i; i = 1, \dots, n\}$. The "simple kriging" or best linear unbiased estimator $\hat{Z}(x)$ is given²⁰ by

$$\hat{Z}(x) = \sum_{i=1}^n \lambda_i(x) Z(y_i)$$

i.e., the interpolator is a linear combination of the observations, with the weights $\lambda_i(x)$ varying with position.

The weights λ_i are selected so as to minimize the expected mean square error of the prediction, that is, the squared deviation of the true from the estimated surface, averaged over all realizations:

$$\begin{aligned} E[(\hat{Z}(x) - Z(x))^2] &= E\left[\left(\sum_{i=1}^n \lambda_i(x) Z(y_i) - Z(x)\right)^2\right] \\ &= E\left[\sum_{i=1}^n \lambda_i(x) Z(y_i) \sum_{j=1}^n \lambda_j(x) Z(y_j) - \right. \\ &\quad \left. 2 \sum_{i=1}^n \lambda_i(x) Z(y_i) Z(x) + Z(x)^2\right] \\ &= \sum_{i=1}^n \sum_{j=1}^n \lambda_i(x) \lambda_j(x) \text{Cov}(y_i, y_j) \\ &\quad - 2 \sum_{i=1}^n \lambda_i(x) \text{Cov}(y_i, x) + \text{Cov}(x, x) \quad (2) \end{aligned}$$

with derivatives

$$\frac{\partial}{\partial \lambda_i(x)} E[(\hat{Z}(x) - Z(x))^2] = 0 =$$

$$2 \sum_{j=1}^n \lambda_j(x) \text{Cov}(y_i, y_j) - 2 \text{Cov}(y_i, x)$$

That is, the weights which minimize the mean square error of the simple kriging estimator are the solution of the system of n linear equations

$$\sum_{i=1}^n \lambda_i(x) \text{Cov}(y_i, y_j) = \text{Cov}(x, y_j) \quad j = 1, \dots, n \quad (3)$$

A more compact notation which we shall use later is

$$\Sigma \lambda(x) = \sigma(x)$$

with $[\Sigma]_{ij} \equiv \text{Cov}(y_i, y_j)$ the covariance between design points y_i and y_j , $[\sigma(x)]_j \equiv \text{Cov}(x, y_j)$ the covariance between point x and design point y_j and $[\lambda(x)]_i \equiv \lambda_i(x)$ the contribution of the response at design point y_i to the response at x .

To obtain the desired estimate of the uncertainty afflicting the kriging interpolation (or prediction) of the response, eq 3 can be transformed to

$$\sum_{j=1}^n \lambda_j(x) \sum_{i=1}^n \lambda_i(x) \text{Cov}(y_i, y_j) = \sum_{j=1}^n \lambda_j(x) \text{Cov}(x, y_j) \quad (4)$$

which in turn can be plugged into eq 2 so that the simple kriging variance σ_{SK}^2 becomes

$$\begin{aligned} \sigma_{\text{SK}}^2(x) &= E[(\hat{Z}(x) - Z(x))^2] = \text{Cov}(x, x) - \sum_{j=1}^n \lambda_j(x) \text{Cov}(y_j, x) \\ &= \text{Cov}(x, x) - \sigma^T(x) \Sigma^{-1} \sigma(x) \quad (5) \end{aligned}$$

The simple kriging variance expresses the uncertainty of the prediction as a function of coordinate x and depends only on the covariance function $\text{Cov}(\cdot, \cdot)$ and, through Σ , on the coordinates of the design points, y , but not on the actual response $Z(y)$ at the design points! In other words, we now have an expression for the uncertainty which does not depend on the measurements but only on their positions. The important consequence is that once we assume a covariance structure and an initial design, we can compute the uncertainty and then optimize the design to minimize uncertainty and thus maximize information without performing any measurements.

As a numerical example, consider the one-dimensional case with a design consisting of two points, $y_1 = -0.5$ and $y_2 = 0.5$. The covariance used is plotted in Figure 3 and the simple kriging variance is shown in Figure 4; the uncertainty is zero at the locations where experiments are performed and grows with the distance from these points.

The dashed line corresponds to the uncertainty or prediction variance that results when only the information gained from a single design point is taken into account (this corresponds to the envelope used in ref 9). The solid line gives the simple kriging variance that results when all available information is taken into account.

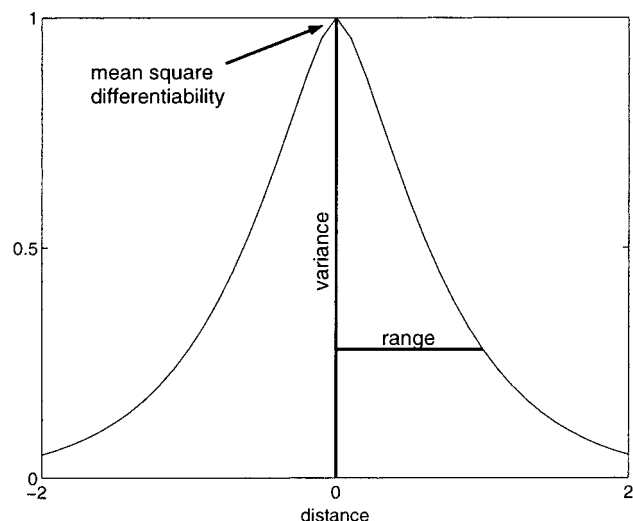


Figure 3. The Matérn⁴¹ covariance function with variance 1, range 1, and mean square differentiability 1 (cf. section 4.1).

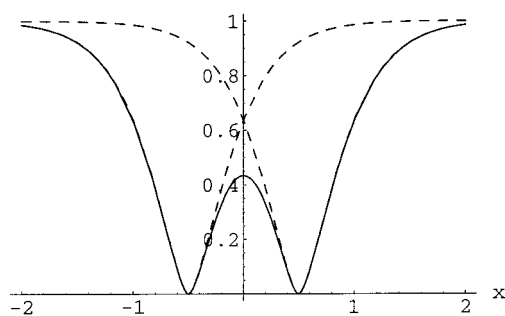


Figure 4. The expected uncertainty of a one-dimensional interpolation of two measurements located at ± 0.5 , using the covariance function from Figure 3. The solid line shows the simple kriging variance, i.e., the expected uncertainty in the interpolation. The dashed line shows the uncertainty if information from the second design point is neglected.

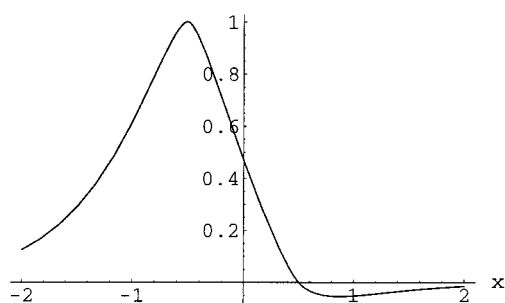


Figure 5. Weight $\lambda_1(x)$ of the left design point in Figure 4.

For further illustration, Figure 5 shows the weight $\lambda_1(x)$ of the first design point as a function of coordinate x in the simple kriging estimate $\hat{Z}(x)$. Note that it becomes identically 1 at the location of its own design point and zero at the location of the other. This feature is general and not an artifact of the symmetry present here. [This can be verified by taking $\lambda_j(y_j) = 1$, $\lambda_{i \neq j}(y_j) = 0$ to obtain $\hat{Z}(y_j)$ and noting that this is a solution of eq 3. The proof then follows from the uniqueness of the solutions which is ensured by the regularity of Σ which in turn is guaranteed by a strictly positive definite covariance function.]

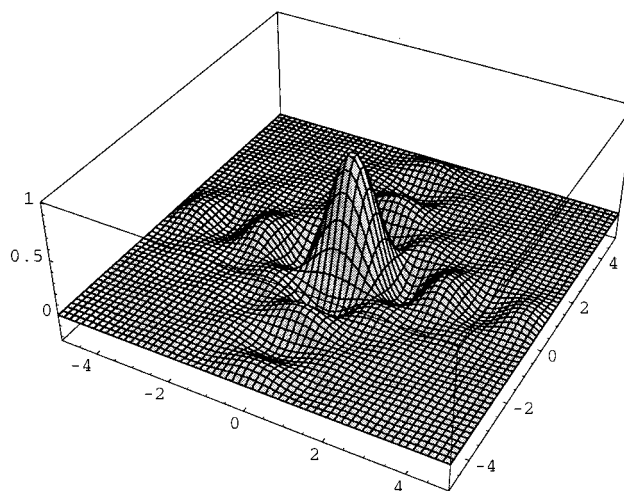


Figure 6. Simple kriging weight function for the two-dimensional square lattice, for an infinitely smooth response surface. The weight is 1 at $(0, 0)$ and zero at the coordinates of all other design points.

To provide a two-dimensional example, Figure 6 shows a simple kriging weight function for one point out of the square design $y = \{(i, j)\}$ with integer i, j .

On the technical side, the relaxation of the assumptions of perfect descriptors (eliminating the uncorrelated random error component and giving rise to a continuous response surface) and absence of a systematic drift, including a non-zero mean, will be discussed in sections 4.3 and 4.2, respectively.

3. OPTIMAL DESIGNS FOR KRIGING

We now turn to the question of the optimal design; we assume that a given number of design points have to be selected and assayed simultaneously (as opposed to sequentially). We further assume that the design region \mathcal{R} is given and for the time being also assume that the covariance function governing the smoothness of the realizations is known, e.g., from related experiments. Since the kriging variance given in eq 5 depends only on the covariance function and on the location of the design points y_i , $i = 1, \dots, n$, but not on the responses $Z(y_i)$, the question of optimal design can be discussed without knowledge of the actual response.

The question remains which design criterion to choose. Natural choices are as follows:²¹ minimization of the integrated mean square error $\int_{\mathcal{R}} \sigma_{\text{SK}}^2(x) dx$, minimization of the maximum mean square error $\max_{x \in \mathcal{R}} \sigma_{\text{SK}}^2(x)$, and maximization of the prior entropy $\int_{\Omega} g(\omega) \log g(\omega) d\omega$ where $g(\omega)$ is the density of $Z(y)$ in a probability space Ω .

Before turning to a detailed discussion of these criteria and a number of approximations to them in the ensuing sections, we wish to emphasize the fundamental distinctness between a criterion or objective function or cost function and the method or algorithm used to optimize a design according to it. Any method of local optimization (such as steepest descent) or of global optimization (such as simulated annealing, evolutionary computation, etc.) can be combined with any criterion.

3.1. Integrated Mean Square Error. 3.1.1 Direct Optimization. Allowing for a weight function $\omega(x)$ —which, for

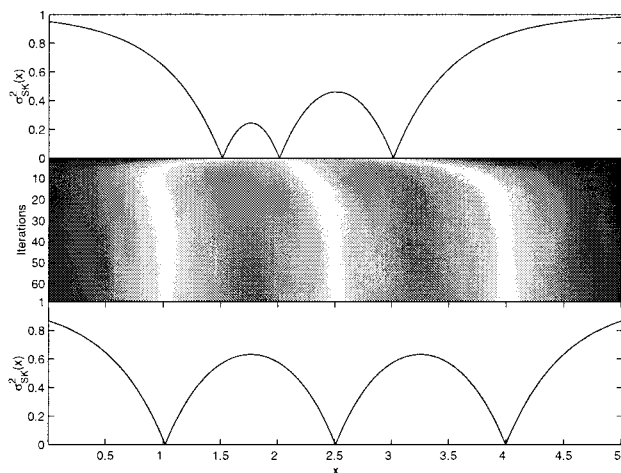


Figure 7. Finding a design with minimum integrated mean square error. Top: simple kriging variance for an arbitrary three point design using an exponential covariance function; middle: intensity plot of the simple kriging variance over 65 iterations of eq 7 with experimental region $0 \leq x \leq 5$. The first row of the intensity plot corresponds to the top of the figure, the last row corresponds to the bottom. During the first few iterations, the left two design points seek to evade each other; in the subsequent iterations, the entire design moves to the right until symmetry with respect to the experimental region is attained. Bottom: simple kriging variance of the final design that minimizes the integrated mean square prediction error.

the time being, may be set to 1—minimization of the integrated mean square error is equivalent to maximization of

$$\int_{\mathcal{R}} \omega(x) \sigma^T(x) \Sigma^{-1} \sigma(x) dx = \int_{\mathcal{R}} \omega(x) \sum_i \lambda_i(x) \text{Cov}(x, y_i) dx \quad (6)$$

which is reminiscent of an unnormalized distortion in vector quantization. In ref 21, the maximization was effected with a quasi-Newton optimizer.

In analogy to the Linde-Buzo-Gray²² and the Rose-Gurewitz-Fox²³ algorithm, we propose to minimize the contribution to the prediction error from each design point iteratively by finding new locations for the design points,

$$y_i^{(k+1)} = \frac{\int_{\mathcal{R}} x \omega(x) \lambda_i^{(k)}(x) \text{Cov}(x, y_i^{(k)}) dx}{\int_{\mathcal{R}} \omega(x) \lambda_i^{(k)}(x) \text{Cov}(x, y_i^{(k)}) dx}$$

then solving for the new $\lambda^{(k+1)}$, etc. For the case depicted in Figure 4, the first iteration leads from $y_1^{(1)} = -0.50$ to $y_1^{(2)} = -0.60$ when the experimental region is defined as $[-5, 5]$, which means the design points move apart to minimize the overall kriging variance. For illustration, another example is given in Figure 7. Qualitatively, the algorithm converges when all $\lambda_i^{(k)}(x) \text{Cov}(x, y_i^{(k)})$ become symmetric about their y_i .

Direct Optimization for Discrete Data. Integration over the entire experimental region as in the equations given above is possible but costly and unnecessary: in library screening, only a finite number of discrete points are at our disposition for testing. It is then sufficient to minimize the integrated mean square error at those points in space which are actually occupied. A replacement of the integral $\int_{\mathcal{R}} dx$ by a summation \sum_x over all compounds, in the following called data, suggests itself. However, this replacement would violate the

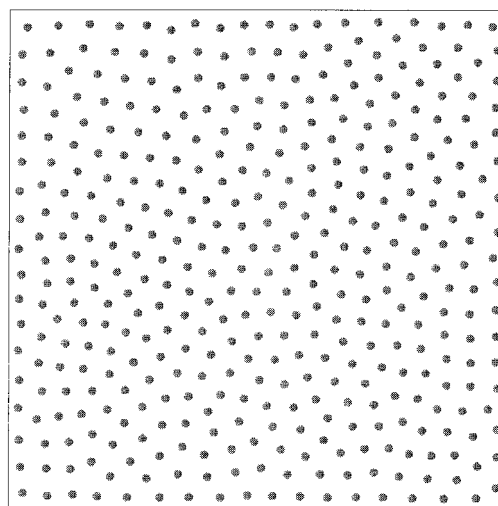


Figure 8. Four hundred point design on a square experimental region with uniform prior probability resulting from iteration of eq 7 after one single random initialization. The design is a hexagonal lattice that is perturbed by the simple cubic boundary conditions.

spirit of eq 6: the integration attempts to minimize the kriging variance in the *entire* design region, whereas a summation minimizes that variance only where data lie; in other words, the resultant design would model the density of the library (which is the very purpose of vector quantization!). This is what methods known as “representative subset selection” try to achieve, see section 5.1. We need to eliminate the effect of the clustering or density inhomogeneity as in

$$y_i^{(k+1)} = \frac{\sum_x \frac{\omega(x)}{\tilde{f}(x)} x \lambda_i^{(k)}(x) \text{Cov}(x, y_i^{(k)})}{\sum_x \frac{\omega(x)}{\tilde{f}(x)} \lambda_i^{(k)}(x) \text{Cov}(x, y_i^{(k)})} \quad (7)$$

where $\tilde{f}(x)$ is the density estimate at point x . If two points are infinitely close to each other but far from all others, their local density is twice that of a single isolated point and their contribution to the above equation is that of a single point at the same position, which is what was desired. In practice, $\tilde{f}(x)$ can be obtained from a kernel density estimate (e.g., refs 24–26), that is, a sum over kernels centered at each point. Likely candidates for these kernels are the covariance function or its square, though we currently ignore which is more appropriate. The criterion corresponding to the above iterative algorithm is

$$\sum_x \frac{\omega(x)}{\tilde{f}(x)} \sigma^T(x) \Sigma^{-1} \sigma(x) \quad (8)$$

This algorithm is a local optimizer and liable to get caught in local maxima. Multiple restarts or combination with a stochastic optimizer may be required.

A design for 400 points on a square experimental region with uniform prior probability is shown in Figure 8.

To come to a slightly more realistic example, Figure 9 shows a probability density from which $N = 200\,000$ points were sampled. A number of designs with $n = 30$ (using an exponential covariance function with a range measuring 1/7th of a side of the square experimental region) are shown in Figure 10 where the generating density was used as density

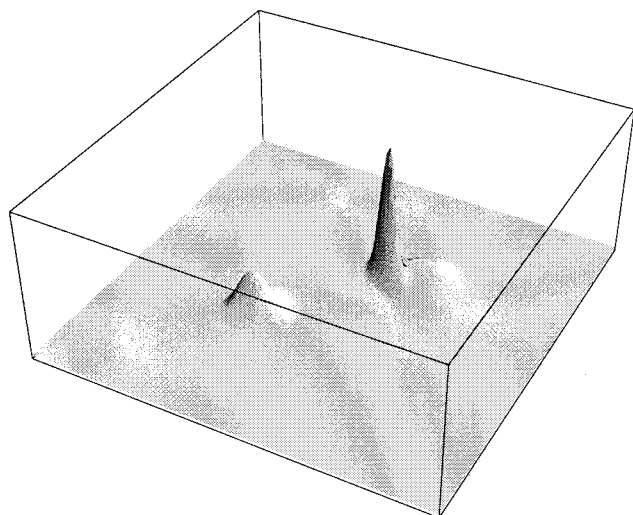
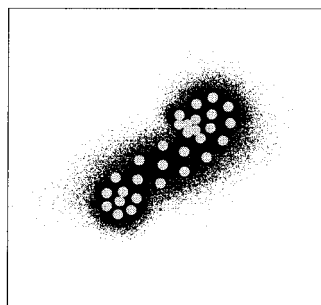
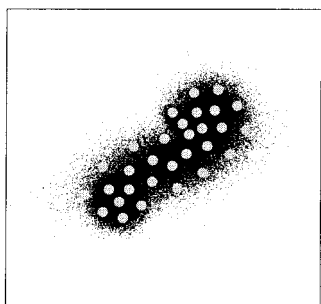


Figure 9. Probability density function (a sum of four normal distributions) from which the 200 000 points in Figure 10 have been sampled.

$\kappa=0$: subset represents density of data



$\kappa=0.5$: intermediate representation



$\kappa=1$: subset represents space covered by data

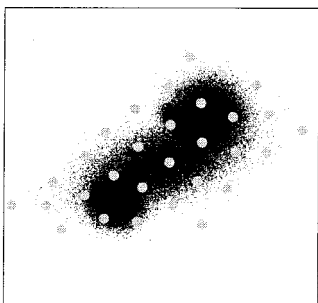


Figure 10. Designs minimizing the integrated mean square prediction error obtained from 150 iterations of eq 7 using different weight functions, see section 5.1. Top: representative subset. Bottom: diverse subset.

estimate. The design in the lower part again resembles a hexagonal lattice, as in Figure 8.

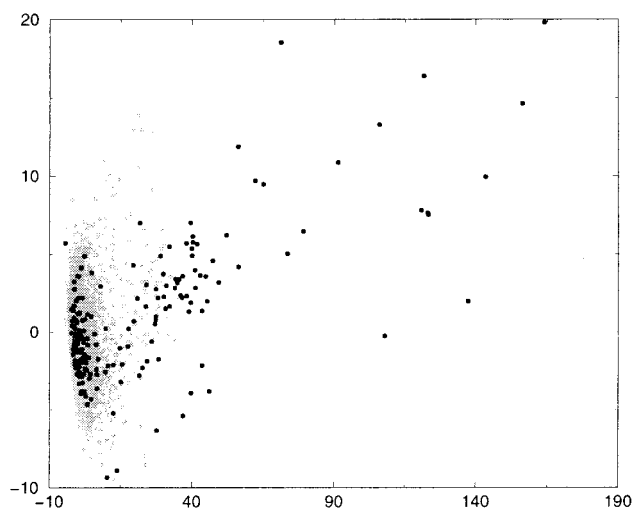


Figure 11. Subset of size 200 out of 6122, obtained with the integrated mean square error criterion. The picture shows only a two-dimensional projection of the full 10-dimensional space.⁵

A last example is shown in Figure 11: the kernel in the density estimate was an exponential of range two, which was also used as a covariance function. The weight function was set to one.

Finally, we note that the above algorithm affords easy incorporation of a preselected subset as advocated in ref 27, that is, fixed points which are to be part of the design. Their coordinates will be used in the calculation of Σ but are not updated.

3.1.2. Optimal Designs for Very Rough Stochastic Processes. When the range of the covariance function is much smaller than the separation between adjacent sampling points, the integrated mean square error criterion reduces²⁸ to the asymptotic “channel coding problem”²⁹ from communication theory. This problem cannot be solved exactly but approximately by means of the “union bound” that yields an equation closely related to sphere packing (see section 3.4.1), albeit with an additional factor taking into account the average number of design points at the minimal distance (see ref 29, chapter 3, 1.3; this is what Johnson et al.³⁰ refer to as a maximin design of lowest index). The sphere packing and two other problems are so intimately linked that section 3.4 has been devoted to them and we defer a discussion until then.

3.1.3. Optimal Designs for Very Smooth Stochastic Processes. If, on the other hand, the range and mean square differentiability at the origin go to infinity, the realizations of the stochastic process become approximately band-limited, i.e., their Fourier transform is effectively zero above a given frequency. In the limit of an infinite experimental region, the optimal sampling design is then given by the dual, or Fourier transform, of the best sphere packing lattice.^{28,31} In dimensions up to eight, the dual of the best sphere packer is the best vector quantizer, see section 3.4. Even in dimensions higher than eight, the duals of the best vector quantizers are still very good sphere packers,²⁹ such that vector quantization can be recommended in all dimensions as an approximation to the integrated mean square error criterion in the limit of very smooth response surfaces.

3.2. Maximum Mean Square Error. If the stochastic process is Gaussian and the correlation range becomes

infinitely short, a minimax (see 3.4.2) design of highest index becomes the optimal design.³⁰

3.3. Entropy. For a compound collection, the number of potential observation sites is finite, and minimization of the expected posterior entropy of the predictor at unobserved sites becomes equivalent to maximization of the prior entropy at the design points.³² If the stochastic process is Gaussian, this prior entropy is maximized by maximizing the determinant of the covariance matrix, $\det(\Sigma)$. An efficient algorithm to do so, based on the classic DETMAX,³³ is described in ref 34.

If, in addition, the correlation range becomes infinitely short, a maximin (see 3.4.1) design of lowest index becomes the optimal design.³⁰

3.4. Packing, Covering, and Quantizing. The criteria for packing, covering, and vector quantization have intimate links to number theory and geometry and are thus of fundamental mathematical interest besides their eminent practical importance in digital communication and engineering sciences. [Note that covering is used here as a technical term from the mathematical literature and not in the sense of ref 17 and others; see section 3.4.2.]

How different are these criteria really? We will first define them²⁹ and then study their optimal designs for infinite homogeneous design regions in the framework of lattices (section 3.5) before commenting on them.

3.4.1. Packing. **Packing density** is the proportion of space covered by a set of nonoverlapping spheres of equal size. This criterion tends to minimize the redundancy of representation and is related to the maximin criterion. A maximin design is one in which the minimum distance between any two design points is maximized.

Rather than integrating over space, a summation over the data analogous to eq 8 can be performed to estimate the packing density, as in

$$\frac{\sum_x \frac{\omega(x)}{\tilde{f}(x)} \sum_i M_r(x, y_i)}{\sum_x \frac{\omega(x)}{\tilde{f}(x)}} \quad (9)$$

where $M_r(x, y_i)$ is a membership function that is 1 if x lies within radius r of design point y_i and 0 otherwise. r is the packing radius, i.e., half the minimal distance between any two design points, $1/2 \min_{i,j} d(y_i, y_j)$. Note that maximization of this formulation of the packing problem avoids some of the difficulties encountered with the maximin criterion for finite design regions: the maximin criterion pushes the design points into the corners of the design region; if the number of design points is small relative to the dimensionality of the space, all design points will lie in corners of the convex hull. This is not the case for the above objective function. Again, incorporation of the density estimate $\tilde{f}(x)$ removes bias from clustering of the data.

Unfortunately, this objective function does not directly lend itself to an optimization as in the case of the integrated mean square error and intuitive heuristics have to be introduced (see section 6.1).

3.4.2. Covering. **Covering thickness** is the average number of spheres that contain a point of the space provided that all of space is covered by spheres of equal size; it is

related to the minimax criterion which minimizes the maximum distance between any datum and its closest design point.

An objective function to minimize the covering thickness can be written as

$$\frac{\sum_x \frac{\omega(x)}{\tilde{f}(x)} \sum_i M_R(x, y_i)}{\sum_x \frac{\omega(x)}{\tilde{f}(x)}} \quad (10)$$

where $M_R(x, y_i)$ is a membership function that is 1 if x lies within radius R of design point y_i and 0 otherwise. R is now the covering radius, i.e., the maximal distance between any datum and its closest design point, $\max_x \min_i d(x, y_i)$.

Again, this objective function does not directly suggest an optimization algorithm though it is clear that the most remote datum should attract its closest design point.

3.4.3. Quantizing. **Quantization error** is the average distortion, i.e., the mean of the squared distances between any datum and its closest design point; it minimizes the second moment of the Voronoi cells pertaining to the design points. [The Voronoi cell of a design point is that part of space which is closer to that design point than to any other; alternative names include area of influence, Brillouin zone, Dirichlet region, nearest neighbor region, Thiessen polygon, Wigner-Seitz cell, etc. See ref 35 or any textbook on "Computational Geometry".]

An objective function eliminating bias from the density is

$$\sum_x \frac{\omega(x)}{\tilde{f}(x)} \sum_i M_V(x, y_i) d(x, y_i)^2$$

where $M_V(x, y_i)$ is a membership function that becomes 1 if x lies within the Voronoi region of y_i , i.e., if it is closer to y_i than to any other $y_{j \neq i}$, and 0 otherwise. This criterion can be optimized using a modification of the traditional Linde-Buzo-Gray algorithm²² that takes into account the inhomogeneous density

$$y_i = \frac{\sum_x \frac{\omega(x)}{\tilde{f}(x)} x M_V(x, y_i)}{\sum_x \frac{\omega(x)}{\tilde{f}(x)} M_V(x, y_i)}$$

compare to eq 7.

3.5. Lattices. Consider an infinite homogeneous isotropic design region with uniform prior probability of finding an active that is to be represented by an infinite number of design points according to one of the above criteria. It is intuitive (though notoriously difficult to prove) that the optimum designs for packing, covering, and quantizing will then be very regular themselves. [A book-sized proof for the Kepler conjecture, that there is no packing denser than the face-centered cubic in three dimensions, was accomplished only in 1998 and relied heavily on computational resources.^{36,37}] In fact, the best arrangements known in all dimensions of practical interest are tessellations (very regular structures indeed) defined in the following. Consider an

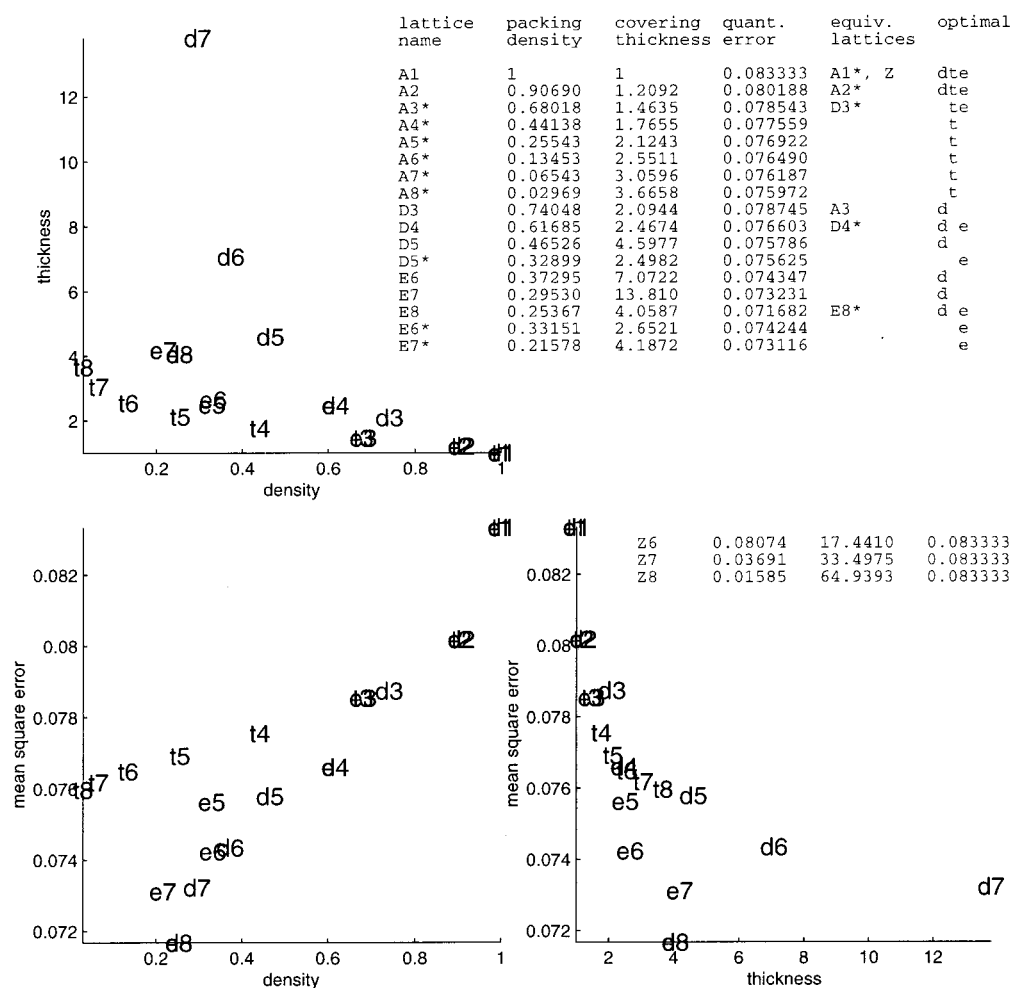


Figure 12. Packing density, covering thickness, and quantization error for lattices that are optimal for either of those in up to eight dimensions. In the three plots, dn , tn , and en indicate a lattice that is optimal for packing, covering, and quantizing, respectively, in dimension n . The table gives the same information in numbers. In the first column, the letters A, D, E, and Z indicate the “family” that a lattice belongs to, the number gives their dimension, and a ★ indicates duality; dual lattices are each other’s Fourier transform. Columns two through four characterize these lattices. The fifth column shows some equivalences between lattices. The last column indicates by which criterion a lattice is optimal. For more details on lattices and notation, refer to ref 29. All numbers shown result from calculus, not from numerical optimization.

infinite set of points and the associated Voronoi regions. If any two regions can be superimposed by translation, rotation, and reflection only, the set of points may be called a tessellation.³⁸

A special case are lattices²⁹ whose Voronoi regions can be superimposed by translation only. These have been studied more extensively and the following discussion is restricted to members from that family. Since all Voronoi regions are congruent, it is sufficient to consider a single lattice point and its Voronoi region only. It then turns out²⁹ that packing attempts to maximize the inradius of the Voronoi region, while covering attempts to minimize its circumradius and quantizing seeks to minimize its second moment, all under the constraint of keeping the volume of the Voronoi cell constant and the requirement of retaining its tessellating property. In other words, all three criteria (plus the channel coding problem) are biased toward sphericity of the lattice’s Voronoi cell. Does that make them equivalent? Indeed, in two dimensions, the optimal design for all three criteria is the hexagonal lattice which may have led people to believe they are essentially the same. However, they are not. In dimensions greater than two, the optimal lattices start to diverge strongly, as illustrated in Figure 12.

Note that the optimal lattices for packing and quantizing are similar with respect to these measures. We hypothesize that this may be a consequence of Conway and Sloane’s conjecture that optimal packing and quantizing lattices are always each other’s duals. However, a counter example has been found³⁸ for the best tessellations in the ninth and tenth dimension and the numerical values supplied therein indicate that this similarity with respect to each other’s measure may be lost in dimensions greater than eight. Also, lattices which pack optimally feature a bad covering thickness.

For reference, the characteristics of the simple cubic lattice in dimensions six through eight (Z6, Z7, Z8) have been supplied. They are bad according to all criteria!

3.6. Discussion of the Primary Criteria. In the absence of further information, we currently know of no compelling objective reason to prefer one criterion over the other. However, we would like to give our subjective view on the matter.

Often six and more descriptors are used to span the chemical space, meaning that the “curse of dimensionality” (e.g., refs 24 and 25) makes itself felt. One of its symptoms is that the Voronoi cells corresponding to the design points become extremely “spiky”. That is, the vertices of the

Voronoi polyhedra (called the “deep holes” in lattice theory) become very remote from the design points, while the volume or probability mass of the tips of the Voronoi cells becomes small.³⁹ An example is provided by the hypercubic Voronoi cells corresponding to a simple cubic lattice of design points (see section 3.5): in dimensions greater than nine, the corner of a cube is more remote from its own design point than some points in the second nearest neighbor cell⁴⁰!

In other words, the maximum mean square error criterion leads to a design that caters optimally for a negligible volume of the entire experimental region, at the cost of the rest of the volume. This problem is alleviated when the integration is replaced by a summation over the data as in eq 10: for the very reason that the volume of these regions is low, a datum is not likely to lie in it. Moreover, only optimal placements lead to a satisfaction of requirement 2 of ref 9, namely that each additional design point should lead to a reduction of the uncertainty as long as it does not coincide with a previous design point.

On these grounds, we tend to dissuade from using the maximum mean square error in higher dimensions.

The integrated mean square error criterion is the only one directly yielding an iterative algorithm for its optimization.

Besides its information theoretic justification, a prime advantage of the entropy criterion has already been quoted in the motto. A drawback is its requirement for the assumption of a specific distribution for the stochastic process before an explicitly calculable objective function is obtained.

4. KRIGING TECHNICALITIES

This section is not central to our exposition and can be omitted at first reading.

4.1. The Covariance Function. We refer to ref 41 for an eloquent account of the prime importance of the degree of differentiability around zero of the covariance function for the behavior of a stochastic process. The impact of the degree of the differentiability is also visible in Figure 2: the markedly different appearances of surfaces on the left and right-hand side are mainly due to their varying differentiabilitys. A consequence is that covariance functions with adjustable mean square differentiability are to be preferred and Stein⁴¹ and others recommend usage of the Matérn covariance.⁴² The parametrization employed here largely decouples the range ρ and the degree of mean square differentiability, given by $[\nu]$ ⁴¹ [p 50]. In one dimension

$$\text{Cov}(d) = \frac{\varsigma}{2^{\nu-1}\Gamma(\nu)} \left(\frac{2\nu^{1/2}d}{\rho} \right)^{\nu} \mathcal{K}_{\nu} \left(\frac{2\nu^{1/2}d}{\rho} \right)$$

The variance ς is the value of the covariance at zero distance d ; Γ is the gamma function; and \mathcal{K}_{ν} is the modified Bessel function of the second kind of order ν . This formulation encompasses several well-known covariance functions: $\nu = 1/2$ yields⁴³ an exponential covariance function, $\nu = 1$ gives the Whittle covariance function (see ref 44, section 9 for a motivation), and $\nu \rightarrow \infty$ affords a Gaussian covariance function (which often leads to oversmoothing and underestimation of the prediction variance, especially when setting the nugget effect to zero, see section 4.3).

4.1.1. Isotropy and Estimation. The bad news is that the covariance function is usually unknown before any assays

have been performed. When available knowledge has already been used to rescale the axes spanning the space, isotropy has to be assumed (that is, the covariance between two points is a function of their distance only, not of the relative orientation). Once responses have been measured, a parametric nonisotropic covariance function from an admissible class, guaranteed to be positive definite, can be fitted to the data manually, as geostatisticians usually propose, or the parameters can be estimated using maximum likelihood⁴⁵ (which requires the assumption of normality) or cross-validation,⁴⁶ an approach more favored by statisticians; for Bayesian formulations, see refs 43, 47, 48. The parameters describing the anisotropy can be considered as describing the optimal scaling of the axes for the response under study.

If the region of high activity is well localized in chemical space, the assumption of stationarity breaks down: there will be a large number of inactives at short and long distances and these data bias the range to large and the variance ς to small values. These data need to be weighed appropriately when estimating a covariance function.

Moreover, the kriging variance is always underestimated because the derivation of the kriging equations was based on the assumption of a known exact covariance function. Progress is ongoing⁴⁹ toward a quantification of a more realistic prediction variance, mostly in Bayesian settings (e.g., ref 43 estimating the posterior probability density of the parameters of Matérn covariances that have been fit to conditional simulations given the actual observations).

The good news is that the actual designs obtained are robust against badly estimated covariance functions. Especially in dimensions 1, 2, 4, and 8 in which the best packing lattice is also the best vector quantizer, the optimal designs are the same for a wide range of covariance functions, at least for the integrated mean square error criterion.

4.2. Relaxing the Assumptions of Zero Mean and Drift.

When the mean is constant but unknown, “ordinary kriging” can be used. If, in addition, a drift is present, the “universal kriging” estimator that represents the drift by a linear combination of fixed basis functions is available. The equations include simple kriging as a special case, see e.g. ref 20. In passing, note that linear and cubic splines are special cases of kriging with particular covariance functions.⁵⁰

4.3. Relaxing the Assumption of Zero Random Error.

In practice, the descriptors will not ensure a continuous response function. Uncorrelated random error can be incorporated by instilling the covariance function with a “nugget effect”. The name reveals its origin in the geosciences. In gold mining, nuggets have vanishing volume on the scale of distances used. A consequence is that the response, the gold grade, can vary in a jerky manner when going from zero to the smallest nonzero distance. Mathematically, this is modeled by a discontinuity of the covariance function at the origin.²⁰

The kriging estimator is an exact interpolator for zero nugget effect only; if there is a nugget effect, some uncertainty remains even at the locations at which measurements have been performed, and the best linear unbiased estimator does not necessarily go through the measurements any more.

4.4. On Lead Optimization. The framework developed here can not only be applied to the first step where a subset needs to be selected simultaneously with no prior knowledge but also is amenable to optimization. Kriging can suggest

which experiments offer the largest potential for high response. The basic idea⁵¹ is to assay those data for which $\hat{Z}(x) + \sigma_{SK}^2$ is maximal, that is, the compounds for which the estimated activity plus the prediction uncertainty become maximal. In the case of a discrete library, the computational burden is much lighter than in the original continuous implementation.⁵¹

4.5. Computational Cost. If a diverse subset is desired, an estimate of the local density $\tilde{f}(x)$ is required. The cost of a brute-force kernel density estimation is $\mathcal{O}(N^2)$.

The computational complexity of direct optimization of the integrated mean square error is of $\mathcal{O}((\text{no. of iterations}) \cdot (n^3 + N(n^2)))$ with n the number of compounds in the design and N the number of compounds in the library. The inversion of Σ gives rise to the n^3 while the matrix multiplications for each out of N compounds account for the second summand. The largest matrices stored permanently are $n \times n$ for Σ and $N \times \text{dim}$ for the coordinates of all compounds, with dim the dimensionality of the descriptor space.

Overall, the algorithm is linear in N and its memory and CPU requirements are modest; in the two-dimensional example with $N = 200\,000$ compounds in the library and $n = 30$ in the subset, only 5 MB of RAM were required when using double precision. These calculations take of the order of an hour on a 500 MHz Pentium III, depending on the desired convergence.

In the examples shown, design point coordinates were initialized uniformly random within a hypercube containing all data and one single optimization was performed, that is, no use was made of multiple restarts or other stochastic optimization techniques.

In our view, the importance of computational cost tends to be overemphasized in subset selection; considering duration and cost of an entire drug or agrochemical development process, CPU time seems well spent if it can raise even marginally the probability of finding a good lead.

5. REMARKS ON SUBSET SELECTION STRATEGIES AND THEIR EVALUATION

5.1. Representativeness vs Diversity and Prior Beliefs.

Traditionally, methods have been classified into those selecting a subset that is representative of the library and those selecting a subset that retains the diversity thereof. In the latter context, the subset has often been required to “span” or “fill” the space explored by the library well.

We argue that, implicitly or explicitly, the first class of methods attempts to represent the *density* of compounds in the collection (as in vector quantization), while the second class attempts to represent the *volume* occupied by the library.

Selecting a representative subset then expresses the prior belief that an active compound is most likely to be found in regions that are densely populated by existing compounds. This implies a bias by previous research. The validity of the assumption varies from case to case and is beyond the realm of a general theoretical study. We do emphasize, however, that it is only one out of an infinite number of possible priors that can be incorporated into the objective function through the weight function $\omega(x)$ used in sections 3.1 and 3.4 (as an aside, we note that Bayesian formulations^{34,52,53} of kriging allow a consistent incorporation of prior beliefs). It is through

this weight function that existing chemical knowledge can be incorporated; in subset selection, the term “filter” is often used instead (e.g., ref 54).

Specifically, many methods yielding “representative” subsets are recovered by letting $\omega(x) = \tilde{f}(x)$, that is, by admitting bias from the density! In other words, we interpret diversity as optimal representation of a uniform probability density.

Researchers have tried to construct implicitly and explicitly⁵⁵ algorithms that compromise between “representativeness” and “diversity”. A continuous interpolation between these goals becomes trivial in the framework described here: taking as a weight function $\omega(x) = \tilde{f}(x)^\kappa$ with κ varying from 0 to 1, the first factors in the numerators and denominators of the equations in sections 3.1 and 3.4 change from $1/\tilde{f}(x)$ to $\tilde{f}(x)/\tilde{f}(x) = 1$, thus giving diverse or representative subsets. Of course, the specific transformation used to convert $\tilde{f}(x)$ to 1 is arbitrary. The principle is illustrated in Figure 10 for various values of κ .

On a more technical note, we have made reference to the density of a library and the volume occupied by it. These are continuous notions, whereas the library per se is just a collection of points with zero total volume. The key parameter that determines how much volume a single point covers or how many other points it “sees” is the kernel bandwidth (also called window width or smoothing parameter) in a kernel density estimate.^{24,25} In the present application, it must be related to the range of the covariance function.

5.2. Calculating Overlap of Libraries. The question of which additional library should be acquired to supplement an existing compound collection can be addressed in the mathematical framework described here.

After all compounds have been embedded in a chemical space as discussed in section 1.2 and a covariance function has been assumed, the required additional ingredients are the specification of a continuous experimental region (for instance, a hyperbox containing all points) and a weight function on that experimental space (which may equal one). Then the integrated mean square error or the entropy on the continuous experimental region may be calculated by making each point a design point. [Integration over a high-dimensional continuous space may be replaced by summation over a large number of quadrature points. If the resulting covariance matrices Σ are too large to invert, the so-called “screening effect” [ref 20, chapter 3.6] may be exploited to obtain good local approximations to the global solution.] The resulting library giving the lower integrated mean square error or higher prior entropy should be preferred.

5.3. Rational vs Random Design. Randomly selecting compounds while ignoring their spatial arrangement in chemical space obviously leads to a subset that represents the density (since the process is stochastic, exceptions are possible but unlikely). If density representation is what is desired, techniques such as vector quantization are much safer if the sample size of selected points is small. In a single subset selection experiment, a random selection may outperform a rational design. On average, however, we do not doubt the superiority of rational design and dismiss the arguments brought forward against it.

5.4. Experimental Comparison of Subset Selection Algorithms. Numerous experimental studies have been

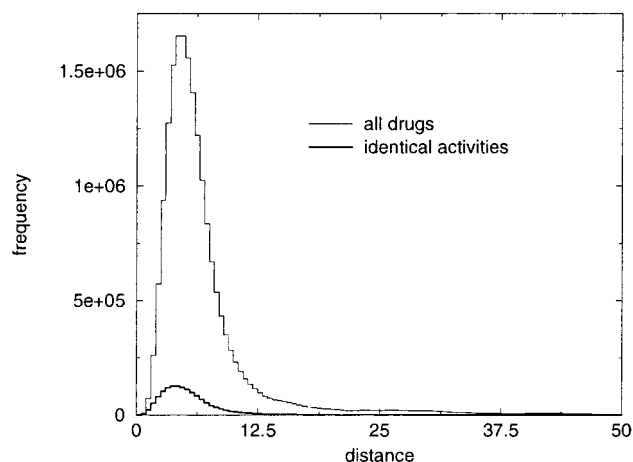


Figure 13. Lack of similar property principle; the histogram shows distances up to 50 units between all drugs in the library⁵ and between drugs sharing at least one activity. Drugs with the same activity are not sufficiently tightly clustered in descriptor space (see section 5.4).

performed, comparing different algorithms or criteria using one library and one or few different assays (e.g., refs 8 and 56–58). While we believe in the ultimate superiority of a single criterion, we consider much of the data presented insufficient for establishing a general quality ranking. We consider investigation of one compound collection with one assay as one experiment (or one realization of a stochastic process), prohibiting the determination of confidence intervals. In addition, most published implementations of different criteria are liable to get caught in local minima, and it is not sufficient to let a method choose one single subset (corresponding to one local minimum) to rank it against others; in a library composed of actives only which Bayada et al.⁵ were kind enough to share, 130 optimizations of around 200 design points led, under the integrated mean square error criterion, to subsets containing between 225 and 297 distinct activities (compare to Figure 3 in Bayada et al.). This deviation is of the order of the differences that the authors found between the methods.

In this application, the integrated mean square error criterion does not perform significantly better than random selection. This may be attributed to the failure of the “PC 10” descriptors⁵ to lead to an embedding that features similar activities for similar molecules, see Figure 13. Part of the reason may be that the activities given sometimes refer to broad classes and not to specific mechanisms involving similar or identical receptors.

We believe that in the limit of a large number of experiments which are conducted independently over the years in the entire world, it does matter which criterion is used, but in the absence of more experimental data (which is difficult to obtain due to the economic interests involved), theory should be conceded more authority than has been the case hitherto.

6. PREVIOUS WORK ON SPACE-FILLING ALGORITHMS

This is not a formal and comprehensive review of the field (for recent developments, see e.g., ref 59), and we ask forgiveness of the authors whose work is not cited or who are not given the credit of originality in case we only cite follow-up work.

6.1. Packing Type Algorithms. The following are direct implementations of or approximations to the packing problem which we recover in sections 3.1.2 and 3.3 for rough response surfaces under the integrated mean square error and the entropy criterion, respectively.

6.1.1. Maximin. The maximin criterion has been optimized using simulated annealing,^{60,61} Monte Carlo techniques,¹⁸ and exchange algorithms.^{5,62}

6.1.2. Soft Maximin. The maximin problem is “hard” because one single collision between two spheres can determine overall performance. There are several ways to make it “softer” and thus more amenable to optimization: for instance by using a soft-core repulsion (a term we borrow from molecular dynamics) as in ref 60 where simulated annealing was used as an optimizer: the objective function is a sum of repulsive potentials between the design points which is zero at large distances and rising up to a constant when decreasing the distance.

An alternative approximation which seems to perform well and is apparently unknown in the chemical community might be described as a model for a cluster of noble gas atoms with attraction, repulsion, and an exaggerated gravitational term where the repulsion potential becomes harder during optimization.⁶³

Alternative approaches to softening the maximin problem are obtained with

$$\left(\sum_i (\min_j d(y_i, y_j))\right)^{1/p} p$$

which includes the ordinary maximin problem for $p \rightarrow \infty$. It has been used with $p = 1$ in refs 60 and 61. Similar formulas corresponding to the harmonic and geometric mean of all squared distances between design points have been used (under different names) in ref 18.

6.1.3. Maximum Dissimilarity. In contrast to the previous section where the objective function is given and the optimization method subject to choice, the maximum dissimilarity method is specified by a deterministic algorithm which contains the objective function (optimal packing) only implicitly. An early reference is ref 64; it has been used in refs 57, 58, 65, 66, and 67 and constitutes one of the limiting cases in ref 68.

Maximum dissimilarity is a sequential design strategy; the location of the next design point depends on the coordinates of all previous design points, which remain fixed. A sequential strategy is optimal if assays are performed one by one because the results can then be used in deciding on the next design point.^{21,51} In a “single stage” experiment, however, simultaneous optimization of the entire design must be superior.

6.1.4. Sphere Exclusion. Sphere exclusion algorithms could, depending on the precise formulation, arguably be included with the class of clustering algorithms, because they sample high-density regions first in most formulations. We consider them as yet another approximation to the packing problem. They have been used, for instance, in refs 5, 58, and 69 and as one of the limiting cases in ref 68.

6.2. Channel Coding Type Algorithms. The primary objective function in ref 9 coincides with the channel coding problem. The figures in that article are instructive, and the requirements set forth for a sensible selection criterion are

intuitive. The authors seek to optimize the channel coding problem using a first-order approximation to the “union bound”. If only they had considered the full union bound²⁹ [chapter 3, section 1.3], they would have arrived at the packing and vector quantizing problems for infinitely short and long ranges of the covariance function, respectively.

6.3. Vector Quantization Type Algorithms. We consider the different clustering algorithms in use as approximations to vector quantization (indeed, k-means is a popular clustering technique). Applications include refs 5 and 67. Vector quantization was stated to give the optimal design in the limiting case of very smooth processes under the integrated mean square error criterion in section 3.1.3.

6.4. Partitioning of Descriptor Space. Partitioning algorithms are computationally cheap and have a number of favorable properties.⁷⁰ However, we discourage the use of the simple cubic lattice as basis for partitions on the grounds discussed in section 3.5 and recommend instead use of the more favorable lattices in Figure 12. When the lattice is not simple cubic, finding the closest lattice point for any point in space cannot be accomplished with a simple modulus operation, but efficient searches are nonetheless possible [chapter 20²⁹]⁷¹ (for a highly readable account, see the introduction of ref 35).

6.5. Unrelated Criteria. Methods which we cannot easily derive from kriging include the cosine criterion,⁷² criticized in ref 61; the weight of the minimum spanning tree,⁵⁶ which is, as noted by the authors themselves, not monotonic with respect to the addition or deletion of compounds; the Kohonen map^{5,67} which arbitrarily forces a two-dimensional topology onto the design, which is useful for visualization but not justifiable otherwise; and some of the algorithms in ref 73 that treat the different dimensions on a different footing, which seems hard to justify if space is assumed isotropic.

6.6. Previous Work on Optimal Design under Correlated Errors. Apart from the series of papers on “Designs for Computer Experiments”^{21,34,45} the majority of work in the area has concentrated on the estimation of the mean or the integral of a stochastic process (yielding designs that are useful for quadrature) and much of it was formulated in the domain of time series⁷⁴ or used coefficients that are much simpler than the optimal kriging ones^{42,75} in the sense that the weight of one observation is independent of all others, thus avoiding inversion of Σ ; for an overview, see ref 76 [chapter 5.6].

Under specific covariance functions, superiority of the hexagonal lattice in two dimensions has already been established by refs 77 and 78.

7. DISCUSSION

It is a pleasure to credit Lutz and Kenakin for the first explicit mention that we have found in a chemical screening context of space-filling designs based on the notion of a correlated stochastic process⁷⁹ [p 315]. Our connections between kriging and chemical spatial sampling design algorithms rest on the work of Johnson, Moore, and Ylvisaker.³⁰

The integrated mean square error and entropy criteria satisfy all requirements put forward in ref 9. Only the integrated mean square error criterion directly affords a local

optimization algorithm without recurring to approximations, and we recommend its use.

In the limit of very rough response surfaces, the integrated mean square error and entropy criteria are optimized by a closest sphere packing of design points (sections 3.1.2, 3.3). In the limit of very smooth surfaces, the integrated mean square error criterion is well approximated by the vector quantization problem (section 3.1.3).

In dimensions greater than one, the cubic lattice is suboptimal. This inefficiency grows with dimensionality, thus simple cubic lattices are to be avoided.

We are aware that the modeling of activity as realization of a stochastic process might raise violent objections. However, this is a battle which has already been waged, sometimes too emotionally, but often at a high intellectual level, in a setting which is entirely equivalent to the one discussed here, namely in mining. Before easily dismissing our ansatz, the reader is encouraged to read up on the debate (mainly in *Mathematical Geology*) circling around the issues of ergodicity, stationarity, the modeling of a unique natural phenomenon as realization of a stochastic process, etc. That community has, after two decades of discussion, reached something like a consensus: that kriging is permissible and does make sense. The scepticism has been addressed with both theoretical arguments (for a short summary of transitive theory which avoids the “epistemological problems associated with the uniqueness of phenomena”, see ref 20; for the original, cf. ref 80 and see also ref 81) and scores of experimental comparisons (e.g., refs 82–84) which have mostly demonstrated the superiority of kriging to alternative schemes.

In summary, we have been able to set up a unifying framework for spatial experimental design comprising most established subset selection methods as limiting cases and have spelled out the assumptions underlying these methods. We propose to find designs that are optimal under the framework’s primary criteria rather than under the restrictive approximations leading to known techniques, especially given that the primary criteria can be computationally more favorable. Some additional aspects of high-dimensional experimental spaces and lattices are discussed in ref 39.

ACKNOWLEDGMENT

It is a pleasure to thank M. Mächler (Seminar for Statistics, ETH Zürich), E. Agrell (Chalmers University, Göteborg), and R. Tobias (SAS, Cary) for substantial advice. T. Hansson and J. Pitera from our laboratory have contributed through helpful discussions. D. M. Bayada and V. J. van Geerestein (Organon, Oss) have generously shared the descriptor values and activities for a library of drugs investigated in ref 5. M. Waldman (Molecular Simulation, San Diego) has kindly provided articles prior to publication and J. Pilz (University of Klagenfurt) and J. Holland (Oxford Molecular, Oxford) have shared their slides from conference presentations. We thank the developers of the GNU scientific library for providing this valuable tool.⁸⁵

REFERENCES AND NOTES

- (1) *The Maximum entropy formalism: a conference*; Levine, R. D., Tribus, M., Eds.; MIT Press: Cambridge, MA, 1978.

- (2) Pearlman, R. S.; Smith, K. M. Metric validation and the receptor-relevant subspace concept. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 28–35.
- (3) Agrafiotis, D. K. Diversity of chemical libraries. In *Encyclopedia of Computational Chemistry*; Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III, Schreiner, P. R., Eds.; Wiley: Chichester, 1998; Vol. 1, pp 742–761.
- (4) Gillet, V. J. Background theory of molecular diversity. In *Molecular Diversity in Drug Design*; Dean, P. M., Lewis, R. A., Eds.; Kluwer: 1999.
- (5) Bayada, D. M.; Hamersma, H.; van Geerestein, V. J. Molecular diversity and representativity in chemical databases. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 1–10.
- (6) *Concepts and applications of molecular similarity*; Johnson, M. A., Maggiora, G. M., Eds.; Wiley: 1990.
- (7) Patterson, D. E.; Cramer, R. D.; Ferguson, A. M.; Clark, R. D.; Weinberger, L. E. Neighborhood behavior: a useful concept for validation of “molecular diversity” descriptors. *J. Med. Chem.* **1996**, 39, 3049–3059.
- (8) Matter, H. Selecting optimally diverse compounds from structure databases: a validation study of two-dimensional and three-dimensional molecular descriptors. *J. Med. Chem.* **1997**, 40, 1219–1229.
- (9) Waldman, M.; Li, H.; Hassan, M. Novel algorithms for the optimization of molecular diversity of combinatorial libraries. *J. Mol. Graph. Mod.* **2000**, 18, 412–426.
- (10) Martin, Y. C.; Brown, R. D.; Bures, M. G. Quantifying diversity. In *Combinatorial Chemistry and Molecular Diversity in Drug Discovery*; Gordon, E. M., Kerwin, J. F., Jr., Eds.; Wiley: 1998.
- (11) Gower, J. C.; Legendre, P. Metric and Euclidean properties of dissimilarity coefficients. *J. Classif.* **1986**, 3, 1–48.
- (12) Cox, T. F.; Cox, M. A. *Multidimensional Scaling, Monographs on Statistics and Applied Probability*; Chapman & Hall: 1995.
- (13) Willett, P. Chemical similarity searching. *J. Chem. Inf. Comput. Sci.* **1998**, 38, 983–996.
- (14) Fedorov, V. V. *Theory of Optimal Experiments*; Academic: London, 1972.
- (15) Heiberger, R. M.; Bhaumik, D. K.; Holland, B. Optimal data augmentation strategies for additive models. *J. Am. Stat. Assoc.* **1993**, 88, 926–938.
- (16) Martin, E. J.; Blaney, J. M.; Siani, M. A.; Spellmeyer, D. C.; Wong, A. K.; Moos, W. H. Measuring diversity – experimental-design of combinatorial libraries for drug discovery. *J. Med. Chem.* **1995**, 38, 1431–1436.
- (17) Higgs, R. E.; Bemis, K. G.; Watson, I. A.; Wikel, J. H. Experimental design for selecting molecules from large chemical databases. *J. Chem. Inf. Comput. Sci.* **1997**, 37, 861–870.
- (18) Hassan, M.; Bielawski, J. P.; C., H. J.; Waldman, M. Optimization and visualization of molecular diversity of combinatorial libraries. *Molecular Diversity* **1996**, 2, 64–74.
- (19) Cressie, N. A. C. The origins of kriging. *Math. Geol.* **1990**, 22, 239–252.
- (20) Chilès, J.-P.; Delfiner, P. *Geostatistics: Modeling Spatial Uncertainty*; Wiley series in probability and statistics; Wiley: New York, 1999.
- (21) Sacks, J.; Welch, W. J.; Mitchell, T. J.; Wynn, H. P. Design and analysis of computer experiments (with discussion). *Stat. Sci.* **1989**, 4, 409–435.
- (22) Linde, Y.; Buzo, A.; Gray, R. M. An algorithm for vector quantizer design. *IEEE Trans. Commun.* **1980**, COM-28, 84–95.
- (23) Rose, K.; Gurewitz, E.; Fox, G. C. Statistical mechanics and phase transitions in clustering. *Phys. Rev. Lett.* **1990**, 65, 945–948.
- (24) Silverman, B. W. *Density Estimation for Statistics and Data Analysis. Monographs on Statistics and Applied Probability*; Chapman & Hall: 1986.
- (25) Scott, D. W. *Multivariate Density Estimation*; Wiley: New York, 1992.
- (26) Wand, M. P.; Jones, M. C. *Kernel Smoothing. Monographs on Statistics and Applied Probability*; Chapman & Hall: London, 1995.
- (27) Martin, E.; Wong, A. Sensitivity analysis and other improvements to tailored combinatorial library design. *J. Chem. Inf. Comput. Sci.* **1999**, 40, 215–220.
- (28) Agrell, E.; Hamprecht, F. A.; Künsch, H. Optimal sampling and interpolation of non wavenumber limited signals in higher dimensions. Manuscript in preparation.
- (29) Conway, J. H.; Sloane, N. J. A. *Sphere Packings, Lattices and Groups*, 2nd ed.; Springer: New York, 1988; Vol. 290 of *Grundlehren der Mathematischen Wissenschaften*.
- (30) Johnson, M. E.; Moore, L. M.; Ylvisaker, D. Minimax and maximin distance designs. *J. Stat. Plann. Inf.* **1990**, 26, 131–148.
- (31) Petersen, D. P.; Middleton, D. Sampling and reconstruction of wavenumber-limited functions in *N*-dimensional Euclidean spaces. *Information Control* **1962**, 5, 279–323.
- (32) Shewry, M. C.; Wynn, H. P. Maximum entropy sampling. *J. Appl. Stat.* **1987**, 14, 165–170.
- (33) Mitchell, T. J. An algorithm for the construction of “D-optimal” experimental designs. *Technometrics* **2000**, 42, 48–54. Reprint from 1974.
- (34) Currin, C.; Mitchell, T.; Morris, M.; Ylvisaker, D. Bayesian prediction of deterministic functions, with applications to the design and analysis of computer experiments. *J. Am. Stat. Assoc.* **1991**, 86, 953–963.
- (35) Agrell, E. Voronoi-Based Coding, Ph.D. Thesis, Chalmers University of Technology, Göteborg, Sweden, 1997.
- (36) Hales, T. C. *The Kepler conjecture*; 1998; <http://www.math.lsa.umich.edu/~hales>.
- (37) Sloane, N. J. A. Kepler’s conjecture confirmed. *Nature* **1998**, 395, 435–436.
- (38) Agrell, E.; Eriksson, T. Optimization of lattices for quantization. *IEEE Trans. Inform. Theory* **1998**, 44, 1814–1828.
- (39) Hamprecht, F. A.; Agrell, E. Exploring a space of materials: spatial sampling design and subset selection. In *Experimental Design for Combinatorial and High Throughput Materials Development*; Cawse, J. N., Ed.; Wiley: New York, 2002.
- (40) Golbraikh, A. Molecular dataset diversity indices and their applications to comparison of chemical databases and QSAR analysis. *J. Chem. Inf. Comput. Sci.* **2000**, 40, 414–425.
- (41) Stein, M. L. *Interpolation of data: some theory for kriging*; Springer Series in Statistics; Springer: New York, 1999.
- (42) Matérn, B. *Spatial Variation*, volume 36 of *Lecture Notes in Statistics*, 2nd ed.; Springer: Berlin, 1986; 1st ed. published by Meddelanden från Statens Skogsforskningsinstitut, Band 49, No. 5, 1960.
- (43) Pilz, J. Bayesian spatial prediction using the Matérn class of covariance functions. Presented at the Int. Conference on Spatial Statistics in the Agro-, Bio- and Geosciences, Freiberg, 2000, and at International Data Analysis, Innsbruck, 2000; manuscript in preparation.
- (44) Whittle, P. On stationary processes in the plane. *Biometrika* **1954**, 41, 434–449.
- (45) Welch, W. J.; Buck, R. J.; Sacks, J.; Wynn, H. P.; Mitchell, T. J.; Morris, M. D. Screening, predicting and computer experiments. *Technometrics* **1992**, 34, 15–25.
- (46) Davis, B. M. Uses and abuses of cross-validation in Geostatistics. *Math. Geol.* **1987**, 19, 241–248.
- (47) Ecker, M. D.; Gelfand, A. E. Bayesian modeling and inference for geometrically anisotropic spatial data. *Math. Geol.* **1999**, 31, 67–83.
- (48) Pardo-Igúzquiza, E. Bayesian inference of spatial covariance parameters. *Math. Geol.* **1999**, 31, 47–65.
- (49) Abt, M. Estimating the prediction mean squared error in Gaussian stochastic processes with exponential correlation structure. *Scand. J. Stat.* **1999**, 26, 563–578.
- (50) Wahba, G. *Spline models for observational data*. SIAM, Philadelphia, **1990**.
- (51) Jones, D. R.; Schonlau, M.; Welch, W. J. Efficient global optimization of expensive black-box functions. *J. Glob. Opt.* **1998**, 13, 455–492.
- (52) Omre, H. Bayesian kriging – merging observations and qualified guesses in kriging. *Math. Geol.* **1987**, 19, 25–39.
- (53) Handcock, M. S.; Stein, M. L. A Bayesian analysis of kriging. *Technometrics* **1993**, 35, 403–410.
- (54) Stanton, R. V. Combinatorial library design: maximizing model-fitting compounds within matrix synthesis constraints. *J. Chem. Inf. Comput. Sci.* **2000**, 40, 701–705.
- (55) Clark, R. D.; Langton, W. J. Balancing representativeness against diversity using optimizable k-dissimilarity and hierarchical clustering. *J. Chem. Inf. Comput. Sci.* **1998**, 38, 1079, 1086.
- (56) Mount, J.; Ruppert, J.; Welch, W.; Jain, A. N. IcePick: a flexible surface-based system for molecular diversity. *J. Med. Chem.* **1999**, 42, 60–66.
- (57) Pötter, T.; Matter, H. Random or rational design? Evaluation of diverse compound subsets from chemical structure databases. *J. Med. Chem.* **1998**, 41, 478–488.
- (58) Snarey, M.; Terrett, N. K.; Willett, P.; Wilton, D. J. Comparison of algorithms for dissimilarity-based compound selection. *J. Mol. Graph. Mod.* **1997**, 15, 372–385.
- (59) Willett, P. Chemoinformatics – similarity and diversity in chemical libraries. *Curr. Opin. Biotech.* **2000**, 11, 85–88.
- (60) Agrafiotis, D. K. Stochastic algorithms for maximizing molecular diversity. *J. Chem. Inf. Comput. Sci.* **1997**, 37, 841–851.
- (61) Agrafiotis, D. K. An efficient implementation of distance-based diversity measures based on k-d trees. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 51–58.
- (62) Marengo, E.; Todeschini, R. A new algorithm for optimal, distance-based experimental design. *Chemometrics Intell. Lab. Syst.* **1992**, 16, 37–44.
- (63) Hardin, R. H.; Sloane, N. J. A. *Operating manual for Gosset: a general-purpose program for constructing experimental designs*, 2nd ed.; AT&T Bell Labs: Murray Hill, NJ, Dec 1994.
- (64) Kennard, R. W. Computer aided design of experiments. *Technometrics* **1969**, 11, 137–148.

- (65) Lajiness, M. S. An evaluation of the performance of dissimilarity selection. In *QSAR: Rational Approaches to the Design of Bioactive Compounds*; Silipo, C., Vittoria, A., Eds.; Elsevier: Amsterdam, 1991; pp 201–204.
- (66) Flower, D. R. DISSIM: A program for the analysis of chemical diversity. *J. Mol. Graph. Mod.* **1998**, *16*, 239–253.
- (67) Tominaga, Y. Data structure using box counting analysis. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 867–875.
- (68) Clark, R. D. Optimisim: An extended dissimilarity selection method for finding diverse representative subsets. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 1181–1188.
- (69) Hudson, B. D.; Hyde, R. M.; Wood, J. Parameter based methods for compound selection from chemical databases. *Quant. Struct.-Act. Relat.* **1996**, *15*, 285–289.
- (70) Schnur, D. Design and diversity analysis of large combinatorial libraries using cell-based methods. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 36–45.
- (71) Agrell, E.; Eriksson, T.; Vardy, A.; Zeger, K. Closest point search in lattices. submitted to *IEEE Trans. Inform. Theory* **2000**.
- (72) Turner, D. B.; Tyrrell, S. M.; Willett, P. Rapid quantification of molecular diversity for selective database acquisition. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 18–22.
- (73) Bayley, M. J.; Willett, P. Binning schemes for partition-based compound selection. *J. Mol. Graph. Mod.* **1999**, *17*, 10–18.
- (74) Cambanis, S. Sampling designs for time series. In *Handbook of Statistics*; Hannan, E. J., Krishnaiah, P. R., Rao, M. M., Eds.; Elsevier: Amsterdam, 1985; Vol. 5, pp 337–362.
- (75) Bellhouse, R. D. Some optimal designs for sampling in two dimensions. *Biometrika* **1977**, *64*, 605–611.
- (76) Cressie, N. A. C. *Statistics for spatial data*; Wiley series in probability and mathematical statistics; Wiley: New York, 1991.
- (77) Olea, R. A. Sampling design optimization for spatial functions. *Math. Geol.* **1984**, *16*, 369–392.
- (78) Yfantis, E. A.; Flatmann, G. T.; Behar, J. V. Efficiency of kriging estimation for square, triangular and hexagonal grids. *Math. Geol.* **1987**, *19*, 183–205. The hexagonal grid is denoted “triangular” in this work; the authors use “hexagonal” to denote another structure.
- (79) Lutz, M.; Kenakin, T. *Quantitative molecular pharmacology and informatics in drug discovery*; Wiley: Chichester, 1999.
- (80) Mathéron, G. *Les variables régionalisées et leur estimation. Une application de la théorie des fonctions aléatoires aux Sciences de la Nature*; Masson: Paris, 1965.
- (81) Journel, A. G. The deterministic side of kriging. *Math. Geol.* **1985**, *17*, 1–15.
- (82) Zimmermann, D.; Pavlik, C.; Ruggles, A.; Armstrong, M. P. An experimental comparison of ordinary and universal kriging and inverse distance weighting. *Math. Geol.* **1999**, *31*, 375–390.
- (83) Haaland, P.; McMillan, N.; Nychka, D.; Welch, W. Analysis of space-filling designs. In *Computationally intensive statistical methods: proceedings of the 26th Symposium on the Interface: Interface '94, volume 26 of Proceedings of Computer Science and Statistics: Annual Symposium on the Interface*; Sall, J., Lehman, A., Eds.; Interface Foundation of North America: 1994; pp 111–120.
- (84) Laslett, G. M. Kriging and splines: an empirical comparison of their predictive performance in some applications (with discussion). *J. Am. Stat. Assoc.* **1994**, *89*, 391–400.
- (85) GNU scientific library, 2001; <http://sources.redhat.com/gsl/>.
- (86) Pebesma, E. J.; Heuveling, G. B. M. Latin hypercube sampling of Gaussian random fields. *Technometrics* **1999**, *41*, 303–312. <http://www.gstat.org>.

CI010376B