

# Estimating derivatives and integrals with Kriging

Emmanuel Vazquez

Department of Signal and Electronic Systems  
Supélec

91192 Gif-sur-Yvette, France

emmanuel.vazquez@supelec.fr

Eric Walter

Laboratoire des Signaux et Systèmes  
CNRS, Supélec, Univ. Paris-Sud

91192 Gif-sur-Yvette, France

eric.walter@lss.supelec.fr

**Abstract**—This paper formalizes a methodology based on Kriging, a technique developed by geostatisticians, for estimating derivatives and integrals of signals that are only known via possibly irregularly spaced and noisy observations. This finds direct applications, e.g., in system identification when differential algebra is used to express parameters as nonlinear functions of the inputs and outputs and their derivatives. The procedure is quite simple to implement, and allows confidence intervals on the predicted values to be derived.

## I. INTRODUCTION

Let  $f^*(\mathbf{x})$  be a real function defined on some compact set  $\mathbb{X} \subset \mathbb{R}^d$ . The problem to be considered here is the estimation of the derivative (or the integral) of  $f^*$  at any given  $\mathbf{x} \in \mathbb{X}$  from a finite set of observations  $S = \{(\mathbf{x}_i, f_{\mathbf{x}_i}^{\text{obs}}), i = 1, \dots, n\}$ . These observations may be corrupted by (not necessarily white) noise, so  $f_{\mathbf{x}_i}^{\text{obs}}$  is not equal to  $f(\mathbf{x}_i)$  in general. Moreover, the observations need not be regularly sampled. Such a problem is frequently encountered in system identification and control, for instance when algebraic differential methods are used to express parameters as nonlinear functions of the input and outputs and their derivatives, see, e.g., [1], [4].

The methodology to be presented is based upon techniques developed by Geostatisticians and known as *Kriging* and *Intrinsic Kriging*. The possibility of estimating derivatives and integrals via Kriging has already been suggested in the context of Geostatistics [2] and we would like to call the attention of the control community on its simplicity, pertinence and performances. This paper will consist of two parts. The first one will briefly recall the theory of Kriging and Intrinsic Kriging, before considering the estimation of derivatives and integrals with these methods. To the best of our knowledge, the mathematical formalization of the prediction of derivatives and integrals with intrinsic Kriging had never been published. The second part of the paper provides illustrative examples.

## II. FROM LINEAR PREDICTION TO THE ESTIMATION OF DERIVATIVES AND INTEGRALS

### A. Random Processes and Kriging

Kriging can be used to approximate or interpolate data, just as splines do. See, for instance, [2], [3], [6], [9]. As splines, it is a kernel regression method [10], [13]. Its specificity is that the kernel is chosen after a statistical analysis of the data. Let  $(\Omega, \mathcal{A}, \mathbb{P})$  be a probability space and

$L^2(\Omega, \mathcal{A}, \mathbb{P})$  be the set of second-order real-valued random variables defined on  $(\Omega, \mathcal{A}, \mathbb{P})$ .

*Definition 1:* The set of all random variables  $F(\omega, \mathbf{x}) \in L^2(\Omega, \mathcal{A}, \mathbb{P})$  obtained when  $\mathbf{x}$  runs through  $\mathbb{X}$  is called a *second-order random process* with parameter  $\mathbf{x} \in \mathbb{X}$ .

Let  $m(\mathbf{x}) = \mathbb{E}[F(\mathbf{x})]$  be the *mean* of  $F(\mathbf{x})$ , and  $k(\mathbf{x}, \mathbf{y}) = \text{cov}(F(\mathbf{x}), F(\mathbf{y})) = \mathbb{E}[(F(\mathbf{x}) - m(\mathbf{x}))(F(\mathbf{y}) - m(\mathbf{y}))]$  be its *covariance function*. The covariance function is *positive* since

$$\text{var} \left[ \sum_{i=1}^n \lambda_i F(\mathbf{x}_i) \right] = \sum_{i,j=1}^n \lambda_i \lambda_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0,$$

for all  $\lambda_i \in \mathbb{R}$ ,  $\mathbf{x}_i \in \mathbb{X}$  and  $n > 0$ .

Let  $F_1(\mathbf{x}), \dots, F_q(\mathbf{x})$  be random processes defined on the same probability space and parameter space  $\mathbb{X}$ . Let  $m_\alpha(\mathbf{x})$  be the mean of  $F_\alpha(\mathbf{x})$  and  $k_{\alpha,\beta}(\mathbf{x}, \mathbf{y})$  be the covariance  $\mathbb{E}[(F_\alpha(\mathbf{x}) - m_\alpha(\mathbf{x}))(F_\beta(\mathbf{y}) - m_\beta(\mathbf{y}))]$ ,  $\alpha, \beta \in \{1, \dots, q\}$ . Note that  $k_{\alpha,\beta}(\mathbf{x}, \mathbf{y}) = k_{\beta,\alpha}(\mathbf{y}, \mathbf{x})$  but in general  $k_{\alpha,\beta}(\mathbf{x}, \mathbf{y}) \neq k_{\alpha,\beta}(\mathbf{y}, \mathbf{x})$ .

Assume, for the time being, that  $f^*(\mathbf{x})$  is a *trajectory* of  $F(\mathbf{x})$  (i.e., there exists  $\omega \in \Omega$  such that  $F(\omega, \mathbf{x}) = f^*(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{X}$ ) and that the observations are noise-free. Each  $f_{\mathbf{x}_i}^{\text{obs}} = f^*(\mathbf{x}_i)$  is then a *realization* of the random variable  $F(\mathbf{x}_i)$ . The data  $S$  can then be interpolated by predicting  $F(\mathbf{x})$  given the random variables  $F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)$ , which amounts to finding a function  $\hat{F}(\mathbf{x})$  of  $F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)$  that minimizes  $F(\mathbf{x}) - \hat{F}(\mathbf{x})$  in some sense. A predicted value is obtained by replacing the random variables  $F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)$  by their realizations  $f_1^{\text{obs}}, \dots, f_n^{\text{obs}}$  in the expression of  $\hat{F}(\mathbf{x})$  ( $\omega$  remaining uncertain). Consider the class of *linear predictors*, which can be written as

$$\hat{F}(\mathbf{x}) = \sum_{i=1}^n \hat{\lambda}_{i,\mathbf{x}} F(\mathbf{x}_i).$$

Assume, moreover, that  $F(\mathbf{x})$  is a zero-mean process ( $m(\mathbf{x}) = 0, \forall \mathbf{x} \in \mathbb{X}$ ), an hypothesis that will be relaxed in Section II-B. The *best predictor* or *Kriging predictor*  $\hat{F}(\mathbf{x})$  of  $F(\mathbf{x})$  is the orthogonal projection of  $F(\mathbf{x})$  onto the subspace

$$\mathcal{H}_S = \text{span}\{F(\mathbf{x}_i), i = 1, \dots, n\}.$$

Since

$$(F(\mathbf{x}) - \hat{F}(\mathbf{x}), F(\mathbf{x}_i))_{L^2(\Omega, \mathcal{A}, \mathbb{P})} = 0 \Rightarrow$$

$$k(\mathbf{x}, \mathbf{x}_i) - \sum_{j=1}^n \hat{\lambda}_{j, \mathbf{x}} k(\mathbf{x}_j, \mathbf{x}_i) = 0, \quad (1)$$

for  $i = 1, \dots, n$ , the scalars  $\hat{\lambda}_{i, \mathbf{x}}$ ,  $i = 1, \dots, n$ , are obtained by solving a linear system of equations.

*Proposition 1 (Kriging):* Let  $F_1(\mathbf{x}), \dots, F_q(\mathbf{x})$  be second-order random processes, with zero mean, and covariance functions  $k_{\alpha, \beta}(\mathbf{x}, \mathbf{y})$ . The best predictor of  $F_\alpha(\mathbf{x})$  from  $F_{\alpha_i}(\mathbf{x}_i)$ ,  $i = 1, \dots, n$ , is the orthogonal projection of  $F_\alpha(\mathbf{x})$  onto  $\mathcal{H}_S = \text{span}\{F_{\alpha_i}(\mathbf{x}_i), i = 1, \dots, n\}$ , written as

$$\hat{F}_\alpha(\mathbf{x}) = \sum_{i=1}^n \hat{\lambda}_{i, \mathbf{x}} F_{\alpha_i}(\mathbf{x}_i), \quad (2)$$

where the  $\hat{\lambda}_{i, \mathbf{x}}$ s are the solution of the linear system<sup>1</sup> of equations

$$\mathbf{K} \hat{\boldsymbol{\lambda}}_{\mathbf{x}} = \mathbf{k}_{\mathbf{x}}. \quad (3)$$

In (3),  $\mathbf{K}$  is the  $n \times n$  covariance matrix of the random vector  $\mathbf{F}_S = (F_{\alpha_1}(\mathbf{x}_1), \dots, F_{\alpha_n}(\mathbf{x}_n))^T$ ,  $\hat{\boldsymbol{\lambda}}_{\mathbf{x}} = (\hat{\lambda}_{1, \mathbf{x}}, \dots, \hat{\lambda}_{n, \mathbf{x}})^T$  is the vector of the Kriging coefficients, and  $\mathbf{k}_{\mathbf{x}} = (k_{\alpha, \alpha_1}(\mathbf{x}, \mathbf{x}_1), \dots, k_{\alpha, \alpha_n}(\mathbf{x}, \mathbf{x}_n))^T$  is the covariance vector of  $\mathbf{F}_S$  and  $F_\alpha(\mathbf{x})$ . Confidence intervals are obtained by evaluating the variance of  $F_\alpha(\mathbf{x}) - \hat{F}_\alpha(\mathbf{x})$ .

### B. Intrinsic Kriging and Intrinsic Random Functions

*Intrinsic Kriging (IK)* [8] extends linear prediction to the case where the mean of  $F(\mathbf{x})$  is unknown. In this framework, the function  $f^*$  generating the data is assumed to fluctuate around  $m(\mathbf{x})$ , which can be written as a linear parametric function  $\mathbf{b}^T \mathbf{r}(\mathbf{x})$ , where  $\mathbf{b}$  and  $\mathbf{r}(\mathbf{x})$  are  $l$ -dimensional vectors. Let  $\mathcal{N}$  be the vector space  $\{\mathbf{b}^T \mathbf{r}(\mathbf{x}), \mathbf{b} \in \mathbb{R}^l\}$  and  $F(\mathbf{x})$  be a random process with mean  $m(\mathbf{x}) \in \mathcal{N}$ . The main idea of IK is to find some linear transformations of  $F(\mathbf{x})$  *filtering out* the mean so as to consider a zero-mean process again.

We first recall the notion of generalized random processes. Let  $\tilde{\Lambda}$  be the vector space of *finite-support measures*, i.e. the space of linear combinations  $\sum_{i=1}^n \lambda_i \delta_{\mathbf{x}_i}$ , where  $\delta_{\mathbf{x}}$  stands for the Dirac measure, such that for any  $B \subset \mathbb{X}$ ,  $\delta_{\mathbf{x}}(B)$  equals one if  $\mathbf{x} \in B$  and zero otherwise. Let  $\tilde{\Lambda}_{\mathcal{N}^\perp}$  be the subset of the elements of  $\tilde{\Lambda}$  that vanish on  $\mathcal{N}$ . Thus,  $\lambda \in \tilde{\Lambda}$  implies

$$\langle \lambda, f \rangle := \sum_{i=1}^n \lambda_i f(\mathbf{x}_i) = 0, \quad \forall f \in \mathcal{N}.$$

*Definition 2:* A symmetric function  $k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$  is *conditionally positive* with respect to  $\mathcal{N}$  if, for all  $\lambda \in \tilde{\Lambda}_{\mathcal{N}^\perp}$ ,  $k(\lambda, \lambda) \geq 0$ , where  $k(\lambda, \mu)$ ,  $\lambda, \mu \in \tilde{\Lambda}_{\mathcal{N}^\perp}$ , is defined by

$$k(\lambda, \mu) := \sum_{i, j=1}^{n, m} \lambda_i \mu_j k(\mathbf{x}_i, \mathbf{y}_j).$$

<sup>1</sup> $\mathbf{K}$  in (3) is generally a full rank matrix since covariances are most often *positive definite* functions. However, adapted solving techniques such as rank reduction must be used when the condition number of  $\mathbf{K}$  is large, which may happen for instance when two observations are close in the space of factors.

If, moreover,  $k(\lambda, \lambda) = 0$  implies  $\lambda = 0$ , for all  $\lambda \in \tilde{\Lambda}_{\mathcal{N}^\perp}$ , then  $k(\mathbf{x}, \mathbf{y})$  is *conditionally positive definite*.

Let  $F_G(\lambda)$  be a linear application defined on  $\tilde{\Lambda}_{\mathcal{N}^\perp}$ , with values in  $L^2(\Omega, \mathcal{A}, \mathbb{P})$ . Assume that  $F_G(\lambda)$  is zero-mean for all  $\lambda$  and that  $\text{cov}[F_G(\lambda), F_G(\mu)] = k(\lambda, \mu)$ , where  $k(\mathbf{x}, \mathbf{y})$  is a conditionally positive definite function. Then,  $F_G(\lambda)$  is called a *generalized random process*. Such a random process is no longer defined on  $\mathbb{X}$  but on a space of measures, and  $k(\mathbf{x}, \mathbf{y})$  is called a *generalized covariance* (see Section II-E). Let  $\tilde{\mathcal{H}}_{\mathcal{N}^\perp}$  be the subspace of  $L^2(\Omega, \mathcal{A}, \mathbb{P})$  generated by  $F_G(\lambda)$ ,  $\lambda \in \tilde{\Lambda}_{\mathcal{N}^\perp}$ . Since random variables in  $\tilde{\mathcal{H}}_{\mathcal{N}^\perp}$  are zero-mean, the inner product of  $L^2(\Omega, \mathcal{A}, \mathbb{P})$  can be expressed in  $\tilde{\mathcal{H}}_{\mathcal{N}^\perp}$  as

$$(F_G(\lambda), F_G(\mu))_{L^2(\Omega, \mathcal{A}, \mathbb{P})} = k(\lambda, \mu) = \sum_{i, j} \lambda_i \mu_j k(\mathbf{x}_i, \mathbf{y}_j),$$

where  $\lambda = \sum_i \lambda_i \delta_{\mathbf{x}_i}$  and  $\mu = \sum_j \mu_j \delta_{\mathbf{y}_j}$  are in  $\tilde{\Lambda}_{\mathcal{N}^\perp}$ . Thus, the bilinear form  $k(\lambda, \mu)$  endows  $\tilde{\Lambda}_{\mathcal{N}^\perp}$  and  $\tilde{\mathcal{H}}_{\mathcal{N}^\perp}$  with a structure of pre-Hilbert space. The completions  $\mathcal{H}_{\mathcal{N}^\perp}$  and  $\Lambda_{\mathcal{N}^\perp}$  of  $\tilde{\mathcal{H}}_{\mathcal{N}^\perp}$  and  $\tilde{\Lambda}_{\mathcal{N}^\perp}$  under this inner product define isomorphic Hilbert spaces.  $F_G(\lambda)$  can be extended on  $\Lambda_{\mathcal{N}^\perp}$  by continuity. The generalized random process  $F_G(\lambda)$  is used as a random model. Simplifying hypotheses are introduced in the next paragraph.

*Intrinsic Random Functions (IRF)* are obtained when generalized random processes are endowed with a stationarity property. IRF are flexible models to use since unknown means can be conveniently dealt with and stationarity makes the inference of the parameters of their (generalized) covariance function feasible. Let  $\tau_{\mathbf{h}} : \tilde{\Lambda}_{\mathcal{N}^\perp} \rightarrow \tilde{\Lambda}$  be the translation operator such that for  $\lambda = \sum_i \lambda_i \delta_{\mathbf{x}_i} \in \tilde{\Lambda}_{\mathcal{N}^\perp}$ ,  $\tau_{\mathbf{h}} \lambda = \sum_i \lambda_i \delta_{\mathbf{x}_i + \mathbf{h}}$ . Assume that  $\tilde{\Lambda}_{\mathcal{N}^\perp}$  is stable under translation.  $\mathcal{N}$  must therefore be itself a stable space of functions under  $\tau_{\mathbf{h}}$ . Assume further that the generalized covariance  $k(\mathbf{x}, \mathbf{y})$  is invariant by translation. In the following, we shall write  $k(\mathbf{h})$  with  $\mathbf{h} = \mathbf{x} - \mathbf{y}$  instead of  $k(\mathbf{x}, \mathbf{y})$ , when the covariance is assumed to be stationary. Then  $\tau_{\mathbf{h}}$  is continuous and can be uniquely extended on  $\Lambda_{\mathcal{N}^\perp}$ .

*Definition 3:* Let  $F_G(\lambda)$  be a zero-mean generalized random process defined on  $\Lambda_{\mathcal{N}^\perp}$ , with stationary generalized covariance  $k(\mathbf{h})$ . The random process  $\mathbf{h} \mapsto F(\tau_{\mathbf{h}} \lambda)$ ,  $\lambda \in \Lambda_{\mathcal{N}^\perp}$ , is therefore weakly stationary.  $F_G(\lambda)$ ,  $\lambda \in \Lambda_{\mathcal{N}^\perp}$ , is then an *Intrinsic Random Function*.

The stability of  $\mathcal{N}$  under the group of translations implies that  $\mathcal{N}$  is necessarily a vector space of *exponential-polynomial* functions [7]. Such a space is generated by functions that can be written as  $\mathbf{x}^l e^{\mathbf{a}^T \mathbf{x}}$ , where  $\mathbf{a}$  is a real or complex vector,  $l$  is the vector-valued index  $(l_1, \dots, l_d)$  and  $\mathbf{x}^l = x_{[1]}^{l_1} \dots x_{[d]}^{l_d}$ . (For a vector-valued index  $l$ , we shall write  $|l| = l_1 + \dots + l_d$ .) For  $\mathcal{N}$  to be stable by linear combinations and translations, the monomials  $\mathbf{x}^l$  must form a complete basis. We restrict ourselves to the case where  $\mathcal{N}$  is a vector space of polynomials of degree at most equal to  $l$ . Let  $\mathcal{N}_l = \text{span}\{\mathbf{x}^l, \forall l \text{ such that } |l| \leq l\}$  and  $\tilde{\Lambda}_l = \tilde{\Lambda}_{\mathcal{N}_l^\perp}$ . Let  $\Lambda_l$  be a completion of  $\tilde{\Lambda}_l$  under the inner product  $k(\lambda, \mu)$ . If the IRF  $F_G(\lambda)$  is defined on  $\Lambda_l$ ,  $F_G(\lambda)$  is an IRF of order  $l$ ,

or IRF( $l$ ) in short.

*Proposition 2:* Any IRF( $l$ ) is an IRF( $l + 1$ ).

*Proof:* This follows from the fact that the spaces  $\Lambda_l$  are nested:

$$\Lambda_{l+1} \subset \Lambda_l,$$

and that any IRF  $F_G(\lambda)$  defined on  $\Lambda_l$  will be stationary on  $\Lambda_{l+1}$ . ■

Remark:  $\Lambda_l$  can be viewed as a set of finite-difference (increment) type operators. The condition for  $\lambda = \sum_{i=1}^n \lambda_i \delta_{\mathbf{x}_i}$  to be in  $\Lambda_0$  can be expressed as  $\sum_{i=1}^n \lambda_i = 0$ . Thus,  $\lambda = \sum_{i=1}^n \lambda_i (\delta_{\mathbf{x}_i} - \delta_{\mathbf{x}_1})$ , so  $\lambda$  is a linear combination of increment measures  $\delta_{\mathbf{x}_i} - \delta_{\mathbf{x}_1}$ . For  $l > 0$ , generalized increments are obtained. Note that if  $\mathbb{X}$  is a space of dimension  $d$ , a minimum of  $\binom{d+l}{l}$  points have to be taken to fully specify an element of  $\Lambda_l$ .

A generalized random process can be viewed as a class of equivalence of random processes defined on  $\mathbb{X}$  with mean in  $\mathcal{N}$ . If  $F(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{X}$ , is a second-order random process, with mean in  $\mathcal{N}$  and covariance  $k(\mathbf{x}, \mathbf{y})$ , the linear application

$$\begin{aligned} F : \tilde{\Lambda}_{\mathcal{N}^\perp} &\rightarrow \mathcal{H} \\ \lambda = \sum_{i=1}^n \lambda_i \delta_{\mathbf{x}_i} &\mapsto F(\lambda) = \sum_{i=1}^n \lambda_i F(\mathbf{x}_i), \end{aligned}$$

extends  $F(\mathbf{x})$  on  $\tilde{\Lambda}_{\mathcal{N}^\perp}$ , where  $\mathcal{H}$  stands for the Hilbert space generated by  $F(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{X}$ . Since the mean of  $F(\mathbf{x})$  is in  $\mathcal{N}$ ,  $F(\lambda)$ ,  $\lambda \in \Lambda_{\mathcal{N}^\perp}$ , is a zero-mean random variable, as  $\lambda$  filters out any function of  $\mathcal{N}$ . Assume that  $k(\mathbf{x}, \mathbf{y})$  is positive definite. Then  $(\lambda, \mu)_{\tilde{\Lambda}_{\mathcal{N}^\perp}} := (F(\lambda), F(\mu))_{\mathcal{H}}$  defines an inner product on  $\tilde{\Lambda}_{\mathcal{N}^\perp}$ . Let  $\Lambda_{\mathcal{N}^\perp}$  be the completion of  $\tilde{\Lambda}_{\mathcal{N}^\perp}$  under this inner product and extend  $F(\lambda)$  on  $\Lambda_{\mathcal{N}^\perp}$  by continuity. A generalized random process is thus obtained. The next paragraph indicates how the procedure may be reversed.

*Definition 4:* Let  $F_G(\lambda)$  be a generalized random process defined on  $\Lambda_{\mathcal{N}^\perp}$ . A second-order random process  $F(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{X}$ , is a *representation* of  $F_G(\lambda)$  if

$$F_G(\lambda) = F(\lambda), \quad \forall \lambda \in \Lambda_{\mathcal{N}^\perp}.$$

If  $F_G(\lambda)$  is an IRF( $l$ ), its representations can be written explicitly by taking appropriate measures  $\lambda \in \Lambda_l$  [8]. If  $F_0(\mathbf{x})$  is any representation of  $F_G(\lambda)$ , other representations of  $F_G(\lambda)$  can be written as

$$F(\mathbf{x}) = F_0(\mathbf{x}) + \sum_{i=1}^q B_i p_i(\mathbf{x}), \quad (4)$$

where the  $p_i$ s form a basis of  $\mathcal{N}_l$  and the  $B_i$ s are any second-order random variables [8]. Thus, the representations of an IRF( $l$ ) constitutes a class of random processes with mean in  $\mathcal{N}_l$ .

### C. Derivation

In this section the notion of derivative of an IRF is developed. The aim is to estimate the derivative of a function  $f^*$  modeled by an IRF( $l$ ), which means that  $f^*$  comprises an unknown polynomial drift of degree at most equal to  $l$ . To simplify presentation, assume that  $x \in \mathbb{R}$ . Extension to the multi-dimensional case is straightforward. ■

Recall that a zero-mean stationary second-order random process  $F(x)$  with covariance function  $k(h)$  is mean-square differentiable at  $x$  if

$$F_h(x) = \frac{1}{h}(F(x+h) - F(x)) \quad (5)$$

converges in mean square when  $h \rightarrow 0$ . The limit exists if and only if  $k^{(2)}(0)$  exists, and  $F(x)$  is then mean-square differentiable at all  $x$ . The limit process is called the *derivative process* and denoted by  $F^{(1)}(x)$ . Higher-order derivatives are obtained by iteration, and it is straightforward to check that

$$\text{cov}[F^{(q)}(x), F^{(r)}(y)] = (-1)^{(r)} k^{(q+r)}(x-y). \quad (6)$$

Let  $F_G(\lambda)$  be an IRF( $l$ ), with generalized covariance  $k(h)$ . The difficulty for defining the derivative of  $F_G(\lambda)$  lies in the fact that neither  $F_G(\lambda)$  nor its derivatives can be defined point-wise. Thus, the notion of differentiability cannot be defined using the variance of an expression such as (5).

To define a derivative, we must use elements of  $\Lambda_l$ . Since for  $\lambda = \sum_i \lambda_i \delta_{\mathbf{x}_i} \in \Lambda_l$ ,  $\tau_h \lambda \in \Lambda_l$ ,  $\forall h \in \mathbb{R}$ , define

$$\lambda_h = \frac{1}{h}(\tau_h \lambda - \lambda) \in \Lambda_l.$$

*Definition 5:* An IRF( $l$ )  $F_G(\lambda)$  is *mean-square differentiable* at  $\lambda \in \Lambda_l$  if  $F_G(\lambda_h)$  converges in mean square as  $h \rightarrow 0$ . When the limit exists, it is denoted by  $F_G^{(1)}(\lambda)$ .

*Proposition 3:* Let  $F_G(\lambda)$  be an IRF( $l$ ), with generalized covariance  $k(h)$ . If  $k^{(2)}(h)$  exists for all  $h$ , then  $F_G(\lambda)$  is mean-square differentiable for all  $\lambda$  in  $\Lambda_l$ , in which case  $F_G^{(1)}(\lambda)$  is an IRF( $l$ ) with generalized covariance  $-k^{(2)}(h)$ .

*Proof:* Let  $\lambda = \sum_{i=1}^n \lambda_i \delta_{\mathbf{x}_i}$  be in  $\Lambda_l$ . Then

$$\begin{aligned} \|F_G(\lambda_h)\|^2 &= \frac{1}{h^2} \left\| F_G \left( \sum_{i=1}^n \lambda_i (\delta_{\mathbf{x}_i+h} - \delta_{\mathbf{x}_i}) \right) \right\|^2 \\ &= \frac{1}{h^2} \sum_{i,j=1}^n \lambda_i \lambda_j (2k(x_i - x_j) \\ &\quad - k(x_i - x_j + h) - k(x_i - x_j - h)). \end{aligned}$$

Assume further that  $k(h)$  is twice differentiable for all  $h \in \mathbb{R}$ . Then  $\|F_G(\lambda_h)\|$  converges when  $h \rightarrow 0$  and

$$\lim_{h \rightarrow 0} \|F_G(\lambda_h)\|^2 = - \sum_{i,j=1}^n \lambda_i \lambda_j k^{(2)}(x_i - x_j).$$

It follows that  $F_G^{(1)}(\lambda)$  is a generalized random process on  $\Lambda_l$  with zero-mean and generalized covariance  $-k^{(2)}(h)$ . ■

Remark that the convergence of  $F_G(\lambda_h)$  in  $L^2(\Omega, \mathcal{A}, \mathbb{P})$  when  $h \rightarrow 0$  is equivalent to the convergence of  $\lambda_h$  in  $\Lambda_l$ . Let  $\lambda^{(1)} = \lim_{h \rightarrow 0} \lambda_h$ , if allowable. We shall then identify  $F_G^{(1)}(\lambda)$  and  $F_G(\lambda^{(1)})$ .

*Proposition 4:* Let  $F_G(\lambda)$  be an IRF( $l$ ) and  $F(\mathbf{x})$  be a representation. Then,  $F^{(1)}(\mathbf{x})$  is a representation of  $F_G^{(1)}(\lambda)$ .

*Proof:* For all  $\lambda \in \Lambda_l$ ,

$$F_G^{(1)}(\lambda) = \lim_{h \rightarrow 0} F_G(\lambda_h) = \lim_{h \rightarrow 0} F(\lambda_h) = F(\lambda^{(1)}) = F^{(1)}(\lambda). \quad \blacksquare$$

Order  $r$  derivatives are denoted by  $F_G^{(r)}(\lambda)$  and  $\lambda^{(r)}$ . Given  $\lambda = \sum_i \lambda_i \delta_{x_i}^{(q_i)}$  and  $\mu = \sum_j \mu_j \delta_{y_j}^{(r_j)}$  in  $\Lambda_l$ , it is easy to check that

$$\text{cov}[F_G(\lambda), F_G(\mu)] = \sum_{i,j} (-1)^{r_j} \lambda_i \mu_j k^{(q_i+r_j)}(x_i - y_j).$$

It becomes now possible to predict the derivatives of a representation of an IRF( $l$ ). The case where observations are corrupted by additive noise is directly studied. Observed values then correspond to realizations of the random variables  $F^{\text{obs}}(x_i) = F(x_i) + N_i$ ,  $i = 1, \dots, n$ , where  $F(x)$  is an unknown representation of  $F_G(\lambda)$ , and the  $N_i$ s are zero-mean random variables independent of  $F(x)$ , with covariance matrix  $\mathbf{K}_N$ . When the noise is white,  $\mathbf{K}_N = \sigma_N^2 \mathbf{I}_n$ . A linear predictor  $\widehat{F}^{(r)}(x)$  of  $F^{(r)}(x)$  can be written as

$$\widehat{F}^{(r)}(x) = \sum_i \widehat{\lambda}_{i,x} F^{\text{obs}}(x_i).$$

In IK the prediction error  $F(x) - \widehat{F}(x)$  is minimized under the constraint  $\delta_x - \sum \widehat{\lambda}_{i,x} \delta_{x_i} \in \Lambda_l$ . To deal with derivatives, we similarly minimize  $\text{var}[F^{(r)}(x) - \widehat{F}^{(r)}(x)]$  under the constraint

$$\delta_x^{(r)} - \sum_i \widehat{\lambda}_{i,x} \delta_{x_i} \in \Lambda_l. \quad (7)$$

The solution can be obtained using  $\text{var}[F(\delta_x^{(r)} - \sum_i \widehat{\lambda}_{i,x} \delta_{x_i})] = \text{var}[F_G(\delta_x^{(r)} - \sum_i \widehat{\lambda}_{i,x} \delta_{x_i})]$ . One can then check that the coefficients  $\widehat{\lambda}_{i,x}$ ,  $i = 1, \dots, n$ , are solutions of a system of linear equations, which can be written in matrix form as

$$\begin{pmatrix} \mathbf{K} + \mathbf{K}_N & \mathbf{P}^\top \\ \mathbf{P} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \widehat{\boldsymbol{\lambda}}_x \\ \boldsymbol{\mu} \end{pmatrix} = \begin{pmatrix} \mathbf{k}_x^{(r)} \\ \mathbf{p}_x^{(r)} \end{pmatrix}, \quad (8)$$

where  $\mathbf{K}$  is the  $n \times n$  matrix of generalized covariances  $k(x_i - x_j)$ ,  $\mathbf{P} = (x_j^i)_{i=0, j=1}^{l, n}$  is a  $(l+1) \times n$  matrix,  $\boldsymbol{\mu}$  is a vector of Lagrange coefficients,  $\mathbf{k}_x^{(r)}$  is a vector of size  $n$  with entries  $k^{(r)}(x - x_i)$  and  $\mathbf{p}_x^{(r)}$  is a vector of size  $l+1$  with entries  $(x^i)^{(r)}$ ,  $i = 0, \dots, l$ . Note that

$$(x^i)^{(r)} = \begin{cases} 0 & \text{if } i < r \\ \frac{i!}{(i-r)!} x^{i-r} & \text{if } i \geq r \end{cases}$$

The variance of the prediction error is given by

$$\text{var}[F^{(r)}(x) - \widehat{F}^{(r)}(x)] = -k^{(2r)}(0) - \begin{pmatrix} \widehat{\boldsymbol{\lambda}} \\ \boldsymbol{\mu} \end{pmatrix}^\top \begin{pmatrix} \mathbf{k}_x^{(r)} \\ \mathbf{p}_x^{(r)} \end{pmatrix}.$$

It can be used to assess confidence intervals, as illustrated in Section III.

#### D. Integration

This problem can be viewed as the prediction of a function  $f(x)$  from observations of its derivative. Formally, this is equivalent to the previous problem, with straightforward adaptation.

#### E. Choice and estimation of the covariance

Once the covariance function is chosen, the procedure of estimating a derivative is quite simple to implement, since it boils down to solving a linear system. The question of the choice of the covariance is now considered. Asymptotic results [11], [14] suggest that satisfactory performance may be obtained even if the covariance is chosen incorrectly. It could therefore be argued that covariance choice is not an important issue. However, for satisfactory performance with a finite and often relatively small number of samples, a proper choice of the covariance turns out to be very important.

Any classical parametrized covariance can be used as a generalized covariance. For instance, [11] advocates the *Matérn covariance*, which can be written as

$$k(\mathbf{h}) = \frac{\sigma^2}{2^{\nu-1} \Gamma(\nu)} \left( \frac{2\nu^{1/2} \|\mathbf{h}\|}{\rho} \right)^\nu \mathcal{K}_\nu \left( \frac{2\nu^{1/2} \|\mathbf{h}\|}{\rho} \right),$$

where  $\mathcal{K}_\nu$  is the modified Bessel function of the second kind,  $\nu > 0$  controls the regularity (the differentiability) of the covariance at the origin,  $\sigma^2 > 0$  is the variance ( $k(0) = \sigma^2$ ), and  $\rho > 0$  is a scale parameter.

*Polynomial covariances* are also a useful class of generalized covariances [8]. Given an order  $l$ , they can be written as

$$k(\mathbf{h}) = \sum_{p=0}^l (-1)^{p+1} a_p \|\mathbf{h}\|^{2p+1} \quad \text{with } a_p > 0, \forall p.$$

Note that this expression is linear in its parameters  $a_p$ . For example, intrinsic Kriging based on a covariance written as  $-a_0 \|\mathbf{h}\|$  gives a piecewise-linear interpolation.

We use *Maximum Likelihood* to estimate the vector  $\boldsymbol{\theta}$  of the parameters of a covariance  $k_\theta(\mathbf{x}, \mathbf{y})$  when the mean of the covariance is known [5]. Let  $F(\mathbf{x})$  be a zero-mean Gaussian random process with covariance  $k_\theta(\mathbf{x}, \mathbf{y})$ . Assume also that the observation noise is Gaussian. Let  $\mathbf{K}(\boldsymbol{\theta})$  be the covariance matrix of  $\mathbf{F}_S = [F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)]^\top$  and  $\mathbf{K}_N(\boldsymbol{\theta}')$  be the covariance matrix of the random vector  $\mathbf{N}$  of the measurement noise, assumed Gaussian, with zero mean and a covariance depending on some parameter vector  $\boldsymbol{\theta}'$ .

To simplify presentation, take  $\bar{\boldsymbol{\theta}} = [\boldsymbol{\theta}^\top, \boldsymbol{\theta}'^\top]^\top$  and  $\mathbf{K}(\bar{\boldsymbol{\theta}}) = \mathbf{K}(\boldsymbol{\theta}) + \mathbf{K}_N(\boldsymbol{\theta}')$ . The log-likelihood of the data can then be written as

$$\begin{aligned} L(\mathbf{f}^{\text{obs}} | \bar{\boldsymbol{\theta}}) &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log \det \mathbf{K}(\bar{\boldsymbol{\theta}}) \\ &\quad - \frac{1}{2} \mathbf{f}^{\text{obs}^\top} \mathbf{K}(\bar{\boldsymbol{\theta}})^{-1} \mathbf{f}^{\text{obs}}. \end{aligned} \quad (9)$$

In the following paragraph we recall the Restricted Maximum Likelihood (REML) approach to estimating the covariance parameters of a random process with unknown mean. Instead of the likelihood function of the data, REML maximizes that of the increments (or generalized increments) of these data [11].

Let  $F_G(\lambda)$  be a Gaussian IRF( $l$ ). Let  $\mathbf{F}^{\text{obs}}$  be the random observation vector, the sum of  $\mathbf{F}_S = [F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)]^\top$ , with  $F(\mathbf{x})$  a representation of  $F_G(\lambda)$ , and some zero-mean

measurement noise vector. Let  $\mathbf{P} = (\mathbf{x}_j^{l_i})_{i,j=1}^{q,n}$  be the  $q \times n$  matrix of basis functions of  $\mathcal{N}_l$  evaluated on  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . Since the dimension of  $\mathcal{N}_l$  is  $q$ , the dimension of the space of the measures with support  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  that cancel out the functions of  $\mathcal{N}_l$  is  $n - q$ . Assume an  $n \times (n - q)$  matrix  $\mathbf{W}$  with rank  $n - q$  has been found, such that

$$\mathbf{P}\mathbf{W} = \mathbf{0}.$$

(The columns of  $\mathbf{W}$  are in the kernel of  $\mathbf{P}$ .) The columns of  $\mathbf{W}$  are therefore the coefficients of measures with support  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ ,  $\sum_{j=1}^n w_{i,j} \delta_{\mathbf{x}_j} \in \Lambda_l$ . Then  $\mathbf{Z} = \mathbf{W}^\top \mathbf{F}^{\text{obs}}$  is a Gaussian random vector taking its values in  $\mathbb{R}^{n-q}$ , with zero mean and covariance matrix  $\mathbf{W}^\top (\mathbf{K}(\boldsymbol{\theta}) + \mathbf{K}_N(\boldsymbol{\theta}')) \mathbf{W} = \mathbf{W}^\top \mathbf{K}(\bar{\boldsymbol{\theta}}) \mathbf{W}$ , where  $\mathbf{K}(\boldsymbol{\theta})$  is the generalized covariance matrix with entries  $k_{\boldsymbol{\theta}}(\mathbf{x}_i - \mathbf{x}_j)$  and where  $\mathbf{K}_N(\boldsymbol{\theta}')$  is the covariance matrix of the observation noise. The random vector  $\mathbf{Z}$  is a *contrast vector*. The log-likelihood of the contrasts is given by

$$\begin{aligned} L(\mathbf{z} | \bar{\boldsymbol{\theta}}) &= -\frac{n-q}{2} \log 2\pi - \frac{1}{2} \log \det(\mathbf{W}^\top \mathbf{K}(\bar{\boldsymbol{\theta}}) \mathbf{W}) \\ &\quad - \frac{1}{2} \mathbf{z}^\top (\mathbf{W}^\top \mathbf{K}(\bar{\boldsymbol{\theta}}) \mathbf{W})^{-1} \mathbf{z}. \end{aligned} \quad (10)$$

Various methods may be employed to compute the matrix  $\mathbf{W}$ . We favor the QR decomposition of  $\mathbf{P}^\top$

$$\mathbf{P}^\top = (\mathbf{Q}_1 | \mathbf{Q}_2) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix},$$

where  $(\mathbf{Q}_1 | \mathbf{Q}_2)$  is an  $n \times n$  orthogonal matrix and  $\mathbf{R}$  is a  $q \times q$  upper triangular matrix. It is trivial to check that the columns of  $\mathbf{Q}_2$  form a basis of the kernel of  $\mathbf{P}$ , so we may chose  $\mathbf{W} = \mathbf{Q}_2$ . Note that  $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_{n-q}$ .

### III. EXAMPLES

#### A. Estimation of derivatives

Figure 1 represents a system output and its derivative. We see that the prediction of the derivative from a number of irregularly spaced noise-free observations of this output is satisfactory. Confidence intervals for this prediction can be provided, which is one of the advantages of the methodology advocated here. As could be expected, the uncertainty intervals are narrower close to the locations of the observations but, may be more surprisingly, when two observations are close enough, prediction is best *between* these observations.

The experiment is repeated in Figure 2, now with the addition of noise on the observations. Again the prediction is quite satisfactory, and the potential applications of such predictors are numerous.

#### B. Black-box model with prior information on derivatives

The previous examples illustrate only partially the possibilities offered by Kriging for the prediction of the derivatives of a signal. In Figure 3, observations of both the function and its derivative are used. This makes it possible to improve prediction, for instance by taking into account the knowledge that the value of the derivative at time zero is zero. This is an opportunity for introducing some prior information in a black-box model [12].

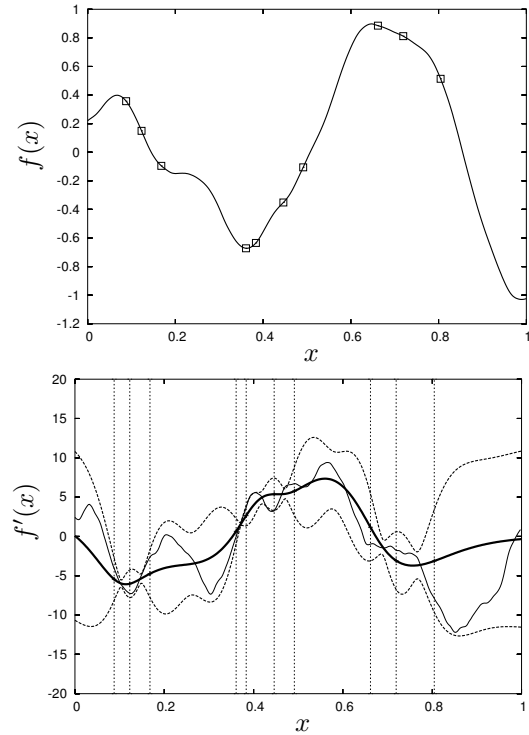


Fig. 1. Top:  $f(x)$ ,  $x \in [0, 1]$ , and 10 noise-free irregularly sampled observations (squares). Bottom: predicted value of the derivative of  $f(x)$  (bold solid line) from previous observations. True derivative (plain solid line) and 95% confidence intervals (dashed lines) are shown. Vertical bars indicate the positions of the observations.

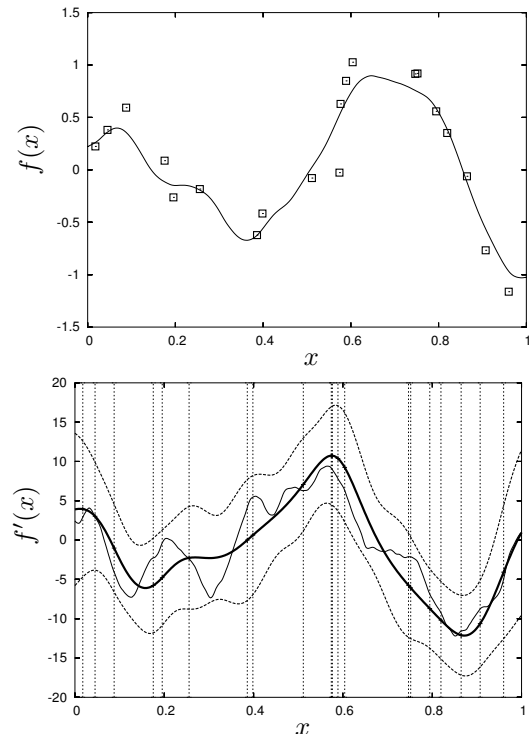


Fig. 2. Top:  $f(x)$ ,  $x \in [0, 1]$ , and 20 noisy irregularly sampled observations (squares). Standard deviation of the noise is 0.2. Bottom: predicted value of the derivative of  $f(x)$  (bold solid line) from previous observations. True derivative (plain solid line) and 95% confidence intervals (dashed lines) are shown. Vertical bars indicate the positions of the observations.

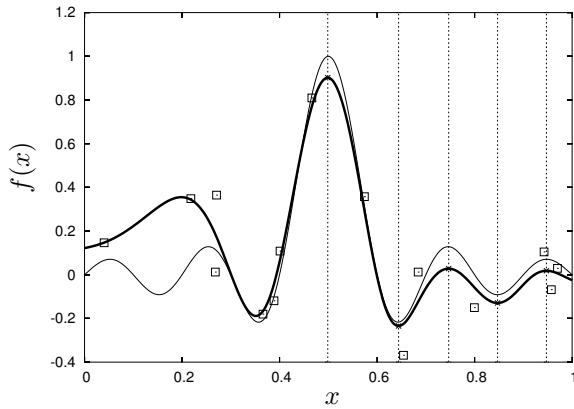


Fig. 3. Approximation (bold solid line) of the sinc function (plain solid line) from noisy observations (squares) taking into account the derivative of the sinc function at positions indicated by the vertical bars.

### C. Approximation of a function from observations of its derivative — integration.

To conclude this section, let us apply Kriging to integration (Figure 4). Note that at least one observation of  $f$  is required in order to specify an initial condition. This approach may be extended to the numerical integration of partial differential equations. It may also be used to model nonlinear dynamical systems described by the state equation  $\dot{x} = f(x, u)$ , which constitutes a promising perspective.

## IV. CONCLUSIONS AND PERSPECTIVES

The methodology presented and formalized in this paper provides tools for solving two basic problems in control and signal processing, namely the differentiation and the integration of possibly multivariable signals that are only known via possibly noisy and irregularly sampled observations. It is based on intrinsic Kriging, a general-purpose technique for black-box modeling developed by geostatisticians during more than 50 years but still relatively ignored by the control community. As demonstrated by the examples treated, the resulting methodology is quite versatile and allows prior information, e.g., on boundary values of the derivative to be taken into account. Another distinctive feature of this statistically-based approach is that it allows confidence intervals on the predicted values of the derivative or integral to be provided. If the covariance function is chosen a priori, the technique requires only the solution of linear systems of equations. Better results, however, are to be expected if this covariance function is estimated from the data, for instance by maximum likelihood or Bayesian estimation. Numerical differentiation finds direct applications, e.g., in the framework of algebraic differential methods for parameter and state estimation, while the possibilities offered by numerical integration in the context of dynamical system modeling look very promising.

## REFERENCES

[1] M. Braci and S. Diop. On numerical algorithms for nonlinear estimation. In *Proc. IEEE CDC 2003*, pages 2896–2901, Hawaii, USA, 2003.

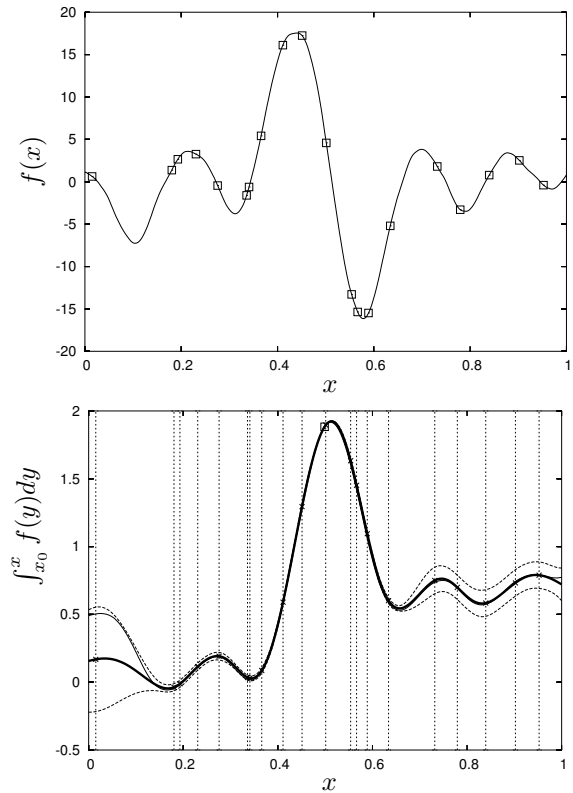


Fig. 4. Integration. Top:  $f(x)$  and 20 noisy irregularly sampled observations (squares). Bottom: approximation of the integral of  $f(x)$  (bold solid line) from previous observations knowing the value of the primitive at  $x_0 = 0.5$  (indicated by a square). The true integral (plain solid line) and confidence intervals (dashed lines) are shown. Vertical bars indicate the positions of the observations.

[2] J.-P. Chilès and P. Delfiner. *Geostatistics: Modeling Spatial Uncertainty*. Wiley, New York, 1999.

[3] N. Cressie. *Statistics for Spatial Data*. Wiley, New York, 1993.

[4] S. Diop, V. Fromion, and J. W. Grizzle. A resettable extended Kalman filter based on numerical differentiation. In *Proc. IEEE CDC 2001*, New York, 2001.

[5] K. V. Mardia and R. J. Marshall. Maximum likelihood estimation of models for residual covariance in spatial regression. *Biometrika*, 73:135–146, 1984.

[6] G. Matheron. Principles of geostatistics. *Economic Geology*, 58:1246–1266, 1963.

[7] G. Matheron. La théorie des fonctions aléatoires intrinsèques généralisées. Note Géostatistiques 117, Centre de Géostatistique de l’Ecole des Mines, Paris, 1971.

[8] G. Matheron. The intrinsic random functions, and their applications. *Adv. Appl. Prob.*, 5:439–468, 1973.

[9] J. Sacks, W. J. Welch, T. J. Mitchell, and H. P. Wynn. Design and analysis of computer experiments. *Statistical Science*, 4(4):409–435, 1989.

[10] B. Schölkopf and A. Smola. *Learning with Kernels*. MIT Press, Cambridge, 2002.

[11] M. L. Stein. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer, New York, 1999.

[12] E. Vazquez and E. Walter. Intrinsic Kriging and prior information. *Applied Stochastic Models in Business and Industry*, 2 (to appear), 2005.

[13] G. Wahba. *Spline Models for Observational Data*, volume 59 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. SIAM, Philadelphia, 1990.

[14] S. J. Yakowitz and F. Szidarovszky. A comparison of Kriging with nonparametric regression methods. *J. Multivariate Analysis*, 16:21–53, 1985.