

Optimal Algorithm for Linear
Problems with Gaussian Measures

CUCS-176-85

G. W. Wasilkowski

G.W. Wasilkowski

Department of Computer Science
Columbia University

February, 1984

This research was supported in part by the National Science
Foundation under Grant MCS-7823676.

•

Abstract

We study optimal algorithms for linear problems in two settings: the average case and the probabilistic case settings. We assume that the probability measure is Gaussian. This assumption enables us to consider a general class of error criteria. We prove that in both settings adaption does not help and a translated spline algorithm is optimal. We also devise optimal information under some additional assumptions concerning the error criterion.

1. Introduction

In this paper we study the optimal reduction of uncertainty for linear problems in two settings: the average case setting and the probabilistic case setting.

By a linear problem we mean the problem of approximating Sf , where S is a linear operator defined on a separable Hilbert space F_1 , when only partial information Nf on f is available. This partial information causes uncertainty. In the average case setting the intrinsic uncertainty is measured by the average size of the error of the best possible algorithm that uses N . In the probabilistic case setting it is measured by the probability that the error of the best possible algorithm is small. In this paper we assume that the probability measure on the space F_1 is Gaussian and the difference between Sf and x , the value given by an algorithm, is measured by $E(Sf-x)$, where E is an arbitrary error functional.

The average case setting has been studied in [5,7,8] for rather general class of probability measures assuming however that the error functional is of a special case. Typically it is assumed that $E(Sf-x) = \|Sf-x\|^2$ and $S(F_1)$ is a separable Hilbert space. Here restricting the class of probability measures to Gaussian measures we relax the

assumption concerning the problem and the form of the error functional E . We are able also to study the probabilistic case setting.

The following results are obtained for both the average case and the probabilistic case settings:

- 1^o For every error functional and for every adaptive information N^a there exists nonadaptive information on the same structure as N^a with uncertainty not greater than the uncertainty caused by N^a . Thus adaption does not help.
- 2^o For every error functional and for every nonadaptive information N a translated spline algorithm is optimal. A sufficient condition for the spline algorithm to be optimal is given.
- 3^o Optimal information N^* is exhibited under some additional assumptions concerning the error functional E .

We now comment on the results mentioned above. The result 1^o states that adaptive information is not more powerful than nonadaptive information in either setting. A similar result for the average case setting has been established in [5,8]. This is not merely of theoretical interest since adaptive information has several undesirable properties like eg.:

--It has more complicated structure than nonadaptive information

--It is ill-suited for parallel computation, whereas nonadaptive information can be computed very efficiently in parallel.

Since adaptive information does not decrease the uncertainty, it may be replaced in practice by nonadaptive information. We want to stress that many commonly used algorithms use adaptive information.

We comment on the result 2^o which states that in both settings a translated spline algorithm φ^* is optimal. (For a similar result for the average case setting see [5,7,9].) Since the spline algorithm is linear, the optimal algorithm φ^* is affine. Hence the cost of evaluating φ^* for given $y = Nf$ is proportional to the cost of evaluating $y = Nf$. This is a desirable property from the complexity point of view.

The result 3^o gives us the best information to be used, i.e., information which minimizes the uncertainty in two settings.

We now summarize the contents of the paper. In Section 2 we formulate the problem. In Section 3 we derive some properties of Gaussian measures. These properties will play a key role in the rest of this paper.

In Section 4 we study the average case setting, and we prove that 1° , 2° and 3° hold for that setting. In Section 5 we study the probabilistic case setting, and we prove that 1° , 2° and 3° hold for that setting. In Section 6 we prove that the spline algorithm enjoys one more optimality property. Namely, assuming that the error functional $E(Sf-x) = \|Sf-x\|^2$, the spline algorithm minimizes the variance.

2. Basic Concepts

Our aim is to approximate the solution operator S ,

$$S: F_1 \rightarrow F_2.$$

We assume that S is linear, F_1 is a separable Hilbert space and F_2 is a linear space, both F_1 and F_2 over the real field. Hence we want to construct an element $x = x(f) \in F_2$ which approximates Sf , $\forall f \in F_1$, with a small error. The error between Sf and x is measured by $E(Sf - x)$, where

$$E: F_2 \rightarrow \mathbf{R}_+,$$

is called an error functional. For example, E might be of the form $E(g) = \|g\|^P$ if F_2 is normed. Here we consider a general class of error functionals. The only assumption concerning E is that for every $g \in F_2$, $H(\cdot) \stackrel{\text{df}}{=} E(S(\cdot) - g)$ is measurable, i.e., $H^{-1}(B) \in \mathbf{B}(F_1)$ whenever $B \in \mathbf{B}(F_1)$, where $\mathbf{B}(F_1)$ stands for the σ -field of Borel sets from F_1 .

To construct $x = x(f)$ we need to know something about f . We assume that our knowledge of f is given by $N^a(f)$. Here N^a is a linear adaptive information operator (for brevity adaptive information), i.e.,

$$(2.1) \quad N^a(f) = [(f, g_1), (f, g_2(y_1)), \dots, (f, g_n(y_1, \dots, y_{n-1}))],$$

where (\cdot, \cdot) is the innerproduct in F_1 ,

$$(2.2) \quad y_1 = y_1(f) = (f, g_1), \quad y_i = y_i(f) \\ = (f, g_i(y_1, \dots, y_{i-1})), \quad i = 2, 3, \dots, n.$$

For brevity we shall write $g_i(y) = g_i(y_1, \dots, y_{i-1}) \in F_1$ for every $y = [y_1, \dots, y_n] \in \mathbb{R}^n$. We assume that $g_i(\cdot)$, as functions of y , are measurable. Without loss of generality we assume that $g_1(y), \dots, g_n(y)$ are linearly independent for every $y \in \mathbb{R}^n$. By

$$(2.3) \quad \text{card}(N^a) = n,$$

we mean the cardinality of N^a . Note that in general the i th evaluation $(f, g_i(y_1, \dots, y_{i-1}))$ depends on the previously computed information $y_1(f), \dots, y_{i-1}(f)$. That's why N^a is called adaptive. If g_i do not depend on y , $g_i(y) = g_i, \forall i, \forall y \in \mathbb{R}^n$, then N^a is called nonadaptive. To stress the nonadaptive character of N^a we often write N^{non} instead of N^a . For every adaptive information N^a , by fixing $y \in \mathbb{R}^n$ a priori and letting $g_i \stackrel{\text{df}}{=} g_i(y)$, we obtain a nonadaptive information

$$(2.4) \quad N_y^{\text{non}}(\cdot) = [(\cdot, g_1), \dots, (\cdot, g_n)]$$

which uses the same evaluations as N^a .

Knowing $N^a(f)$ we construct an approximation $x = x(f)$ by an algorithm φ ,

$$x = \varphi(N^a(f)),$$

where by an algorithm φ that uses N^a we mean any mapping

$$(2.5) \quad \varphi: N^a(F_1) = \mathbb{R}^n \rightarrow F_2.$$

We are interested in optimal algorithms, that is algorithms with minimal errors. What we mean by the error of an algorithm depends on the setting we are dealing with. In this paper we study two different settings: the average case setting and the probabilistic case setting. We begin with the average case setting.

In the average case setting the error of φ is determined by the average behavior of the error $E(Sf - \varphi(N^a f))$. More precisely, let μ be a Gaussian measure defined on $B(F_1)$. Then the average error of φ is defined by

$$(2.6) \quad e^{\text{avg}}(\varphi, N^a) = \int_{F_1} E(Sf - \varphi(N^a f)) \mu(df)$$

and an optimal algorithm φ^* that uses N^a is defined by

$$(2.7) \quad e^{\text{avg}}(\varphi^*, N^a) = r^{\text{avg}}(N^a) \stackrel{\text{df}}{=} \inf_{\varphi} e^{\text{avg}}(\varphi, N^a).$$

This means that in the average case setting we are interested in algorithms φ^* , if they exist, whose average error are minimal. In Section 4 we find φ^* for every nonadaptive information N^a and for adaptive N^a we prove that $r^{\text{avg}}(N^a) \leq r^{\text{avg}}(N_{y^*}^{\text{non}})$ for some y^* .

We now turn to the probabilistic case setting. In this setting the goodness of φ is measured by the probability of success, i.e., by the probability that the error $E(Sf - \varphi(N^a f))$ of φ is small. More precisely, given $\epsilon \geq 0$, let

$$(2.8) \quad \text{Prob}(\varphi, N^a, \epsilon) = \mu(\{f \in F_1 : E(Sf - \varphi(N^a f)) \leq \epsilon\}),$$

where μ is a Gaussian measure defined on $B(F_1)$. Then by an optimal algorithm that uses N^a we mean an algorithm φ^* so that

$$(2.9) \quad \text{Prob}(\varphi^*, N^a, \epsilon) = \text{Prob}(N^a, \epsilon) \stackrel{\text{df}}{=} \sup_{\varphi} \text{Prob}(\varphi, N^a, \epsilon).$$

In Section 5 we find φ^* for every nonadaptive N^{non} . For adaptive N^a we prove that $\text{Prob}(N^a, \epsilon) \leq \text{Prob}(N_{y^*}^{\text{non}}, \epsilon)$ for some y^* .

We comment on the definitions (2.6) and (2.8). In order for (2.6) and/or (2.8) to be well defined, $E(S(\cdot) - \varphi(N^a(\cdot)))$, as a function of f , should be measurable.

It is shown in [6] that this assumption is not restrictive since it is possible to extend the definitions (2.6) and (2.8) for every algorithm and prove that for optimal algorithms φ^* , $E(S(\cdot) - \varphi^*(N^a(\cdot)))$ is measurable.

We now recall some basic properties of Gaussian measures. By a Gaussian measure on $B(F_1)$ we mean a measure μ such that

$$(2.10) \quad \int_{F_1} \exp\{i(f, x)\} \mu(df) = \exp\{i(a, x) - \frac{1}{2}(Ax, x)\},$$

$$\forall x \in F_1, \quad (i = \sqrt{-1}),$$

where $A: F_1 \rightarrow F_1$ is a self-adjoint nonnegative definite operator with finite trace and a is an element of F_1 .

(The left hand side of (2.10) is called the characteristic functional of μ and is denoted by $\psi_\mu(x)$.) Then the mean element m_μ of μ is given by

$$(2.11) \quad m_\mu = a$$

and the correlation operator S_μ of μ , by

$$(2.12) \quad S_\mu = A$$

(see [2,3,4]). Recall that for an arbitrary measure μ , its mean element m_μ is defined by

$$(2.13) \quad (m_{\mu}, x) = \int_{F_1} (f, x)_{\mu} (df), \quad \forall x \in F_1,$$

and its correlation operator S_{μ} by

$$(2.14) \quad (S_{\mu} g, h) = \int_{F_1} (f - m_{\mu}, g) (f - m_{\mu}, h)_{\mu} (df),$$

$$\forall g, h \in F_1.$$

Throughout the rest of this paper we shall assume without loss of generality that the mean element m_{μ} of μ is zero. $m_{\mu} = 0$, and that the correlation operator S_{μ} is positive definite.

3. Conditional measure.

In this section we exhibit an important property of the conditional measure for adaptive information. This property will be extensively used in the next sections. We begin with the definition of conditional measure (see [3]).

For an adaptive information operator N^a , let $\mu_1(\cdot, N^a)$ be the probability measure on $B(\mathbb{R}^n)$ induced by N^a , i.e.,

$$(3.1) \quad \mu_1(\cdot, A) = \mu((N^a)^{-1}(A)) = \mu(\{f \in F_1 : N^a(f) \in A\}),$$

$$\forall A \in B(\mathbb{R}^n).$$

Let $\mu_2(\cdot | Y, N^a), Y \in \mathbb{R}^n$, be a family of probability measures on $B(F_1)$ such that

$$(3.2) \quad \mu_2(F_1 | Y, N^a) = \mu_2((N^a)^{-1}(\{Y\}) | Y, N^a) = 1,$$

for almost every $y \in \mathbb{R}^n$,

$$(3.3) \quad \mu_2(B | \cdot, N^a), \text{ as a function of } y, \text{ is } \mu_1(\cdot | N^a)\text{-}$$

measurable, $\forall B \in B(F_1)$,

and

$$(3.4) \quad \mu(B) = \int_{\mathbb{R}^n} \mu_2(B | Y, N^a) \mu_1(dy, N^a), \quad \forall B \in B(F_1).$$

The family $\mu_2(\cdot | y, N^a)$ is called the conditional measure with respect to N^a and y . The existence and uniqueness of μ_2 follows from [3, Th. 8.1].

Let now G be a measurable function, $G: F_1 \rightarrow \mathbb{R}_+$.

Then

$$(3.5) \quad \int_{F_1} G(f) \mu(df) \\ = \int_{\mathbb{R}^n} \left[\int_{V(N^a, y)} G(f) \mu_2(df | y, N^a) \right] \mu_1(dy, N^a),$$

where $V(N^a, y) = (N^a)^{-1}(\{y\}) = \{f \in F_1 : N^a(f) = y\}$ is the set of elements f from F_1 which share the same information, $Nf = y$. The essence of (3.5) is that we first integrate G over all f with fixed information value y , and then over all values y from \mathbb{R}^n .

Recall that

$$(3.6) \quad N^a(f) = [(f, g_1), (f, g_2(y_1)), \dots, (f, g_n(y_1, \dots, y_{n-1}))],$$

$$y_i = y_i(f) = (f, g_i(y_1, \dots, y_{i-1})).$$

For brevity we write $g_i(y) = g_i(y_1, \dots, y_{i-1})$. Without loss of generality we assume that

$$(3.7) \quad (S_{\mu} g_i(y), g_j(y)) = \delta_{ij}, \quad \forall y \in \mathbb{R}^n.$$

Let for a fixed $y = [y_1, \dots, y_n] \in \mathbb{R}^n$,

$$(3.8) \quad m(N^a, y) = \sum_{j=1}^n y_j S_{\mu} g_j(y)$$

and

$$(3.9) \quad \sigma_{N^a, y}(\cdot) = \sum_{j=1}^n (\cdot, g_j(y)) S_{\mu} g_j(y).$$

Then $\sigma_{N^a, y} : F_1 \rightarrow F_1$ is linear and $m(N^a, y) = \sigma_{N^a, y}(g)$, for every $g \in V(N^a, y)$, and for every fixed $y \in \mathbb{R}^n$. Of course, $m(N^a, y(f))$ and $\sigma_{N^a, y(f)}$, $y(f) = N^a(f)$, need not be linear in f .

Theorem 3.1: Let N^a be an arbitrary information operator of the form (3.7).

(i) Then the induced probability measure

$$(3.10) \quad \mu_1(\cdot, N^a) = \mu_1(\cdot),$$

where μ_1 is the Gaussian measure on $B(\mathbb{R}^n)$ with mean element zero and correlation operator identity, i.e.,

$$(3.11) \quad \mu_1(A) = \frac{1}{\sqrt{(2\pi)^n}} \int_A \exp\{-\frac{1}{2}(x, x)\} d_n x.$$

(ii) The conditional measure $\mu_2(\cdot | y, N^a)$ is the Gaussian measure on $B(F_1)$ with mean element $m(N^a, y)$ given by (3.6) and correlation operator

$$(3.12) \quad S_{N^a, y} = (I - \sigma_{N^a, y}) S_{\mu} (I - \sigma_{N^a, y}^*).$$

Proof: It is shown in [8, Th. 4.2] that there exists a probability measure μ_1 on $\mathbf{B}(\mathbb{R}^n)$ such that $\mu_1(\cdot, N^a) = \mu_1(\cdot)$ for every N^a of the form (3.7). It was shown in [6, Th. 4.2 (i)] that for every nonadaptive N^{non} of the form (3.7), $\mu_1(\cdot, N^{\text{non}})$ is the Gaussian measure on $\mathbf{B}(\mathbb{R}^n)$ with mean element zero and correlation operator identity. Since $\mu_1(\cdot, N^a) = \mu_1(\cdot, N^{\text{non}}) = \mu_1$, the proof of Theorem 3.1 (i) is completed.

To prove (ii), let $\lambda_2(\cdot | y, N^a)$ be the Gaussian measure on $\mathbf{B}(F_1)$ with mean element $m(N^a, y)$ given by (3.8) and correlation operator $S_{N^a, y}$ given by (3.12). Then, due to (2.10)

$$\begin{aligned} \int_{F_1} e^{i(f, x)} \lambda_2(df | y, N^a) \\ = \exp\{i(m(N^a, y), x) - \frac{1}{2}(S_{N^a, y} x, x)\}. \end{aligned}$$

Since $\sigma_{N^a, y}^*(x) = \sum_{j=1}^n (x, S_{\mu} g_j(y)) g_j(y)$, we have

$$\begin{aligned} (S_{N^a, y} x, x) &= (S_{\mu} (x - \sigma_{N^a, y}^*(x)), x - \sigma_{N^a, y}^*(x)) \\ &= (S_{\mu} x, x) - 2(S_{\mu} x, \sigma_{N^a, y}^*(x)) + (S_{\mu} \sigma_{N^a, y}^*(x), \sigma_{N^a, y}^*(x)) \\ &= (S_{\mu} x, x) - 2 \sum_{j=1}^n (S_{\mu} x, g_j(y))^2 + \sum_{j=1}^n (S_{\mu} x, g_j(y))^2 \\ &= (S_{\mu} x, x) - \sum_{j=1}^n (S_{\mu} x, g_j(y))^2. \end{aligned}$$

Hence

$$(3.13) \quad \int_{F_1} e^{i(f,x)} \lambda_2(df|y, N^a) = \exp\{-\frac{1}{2}(S_{\mu} x, x)\} H(x, y),$$

where

$$H(x, y) = \exp\{i(m(N^a, y), x) + \frac{1}{2} \sum_{j=1}^n (S_{\mu} x, g_j(y))^2\}.$$

Due to (3.8)

$$\begin{aligned} H(x, y) &= \exp\{\sum_{j=1}^n [iy_j (S_{\mu} x, g_j(y)) + \frac{1}{2} (S_{\mu} x, g_j(y))^2]\} \\ &= \prod_{j=1}^n \exp\{iy_j (S_{\mu} x, g_j(y)) + \frac{1}{2} (S_{\mu} x, g_j(y))^2\}. \end{aligned}$$

Recall that $g_j(y) = g_j(y_1, \dots, y_{j-1})$, and that μ_1 is the Gaussian measure. Hence

$$\begin{aligned} a \stackrel{df}{=} \int_{\mathbf{R}^n} H(x, y) \mu_1(dy) &= \frac{1}{\sqrt{(2\pi)^n}} \int_{\mathbf{R}} \cdots \int_{\mathbf{R}} \prod_{j=1}^n \exp\{iy_j (S_{\mu} x, g_j(y_1, \dots, y_{j-1})) \\ &\quad + \frac{1}{2} (S_{\mu} x, g_j(y_1, \dots, y_{j-1}))\} \exp\{-\frac{1}{2} \sum_{i=1}^n y_i^2\} d[y_1, \dots, y_n]. \end{aligned}$$

Observe that

$$\begin{aligned} &\frac{1}{\sqrt{2\pi}} \int_{\mathbf{R}} \exp\{iy_j (S_{\mu} x, g_j(y_1, \dots, y_{j-1})) \\ &\quad + \frac{1}{2} (S_{\mu} x, g_j(y_1, \dots, y_{j-1}))\} \exp\{-\frac{1}{2} y_j^2\} dy_j \\ &= \exp\{\frac{1}{2} (S_{\mu} x, g_j(y_1, \dots, y_{j-1}))\} \frac{1}{\sqrt{2\pi}} \int_{\mathbf{R}} \\ &\quad \exp\{iy_j (S_{\mu} x, g_j(y_1, \dots, y_{j-1}))\} \exp\{-\frac{1}{2} y_j^2\} dy_j \end{aligned}$$

$$\begin{aligned}
&= \exp\left\{\frac{1}{2}(S_{\mu} x, g_j(y_1, \dots, y_{j-n}))\right\} \exp\left\{-\frac{1}{2}(S_{\mu} x, g_j(y_1, \dots, y_{j-n}))\right\} \\
&= 1.
\end{aligned}$$

This yields that $a = 1$ and

$$\int_{\mathbb{R}^n} \int_{F_1} e^{i(f, x)} \lambda_2(df | y, N^a) = \exp\left\{-\frac{1}{2}(S_{\mu} x, x)\right\} = \psi_{\mu}(x),$$

where ψ_{μ} is the characteristic functional of μ . Since characteristic functional defines measure uniquely and since conditional measure is determined uniquely (up to a set of μ_1 -measure zero), $\mu_2(\cdot | y, N^a) = \lambda_2(\cdot | y, N^a)$, $\forall y \in \mathbb{R}^a$, a.e. This proves the theorem. ■

Theorem 3.1 states that the induced measure $\mu_1(\cdot, N^a)$ does not depend on N^a , it only depends on $n = \text{card}(N^a)$. From (ii) we can easily conclude that for $y \in \mathbb{R}^n$, the conditional measure $\mu_2(\cdot | y, N^a)$ is the same as the conditional measure for the nonadaptive information operator N_y^{non} ,

$$(3.14) \quad \mu_2(\cdot | y, N^a) = \mu_2(\cdot | y, N_y^{\text{non}}).$$

Furthermore, $\mu_2(\cdot | y, N^a)$ is a translated measure $\mu_2(\cdot | 0, N_y^{\text{non}})$, i.e.,

$$(3.15) \quad \mu_2(B | y, N^a) = \mu_2(B - m(N^a, y) | 0, N_y^{\text{non}}), \quad \forall B \in \mathcal{B}(F_1).$$

In particular, if N^{non} is nonadaptive then

$$(3.16) \quad \mu_2(B|y, N^{\text{non}}) = \mu_2(B - m(N^{\text{non}}, y) | 0, N^{\text{non}}), \quad \forall B \in \mathcal{B}(F_1).$$

We end this section by two lemmas whose proofs, because of their length, are presented in the Appendix.

Lemma 3.1: For every Gaussian measure λ with mean element zero and for every balanced and convex set B ,

$$(3.17) \quad \lambda(B) \geq \lambda(B+h), \quad \forall h \in F_1. \quad \blacksquare$$

Lemma 3.2: Let λ_1, λ_2 be two Gaussian measures on a separable Hilbert space with mean elements zero and correlation operators S_{λ_1} and S_{λ_2} respectively. Let $\alpha_{1,i}, \alpha_{2,i}, \dots, (\alpha_{j,i} \geq \alpha_{j+1,i})$ be the eigenvalues of operators S_{λ_i} , $i = 1, 2$. If

$$\alpha_{j,1} \leq \alpha_{j,2}, \quad \forall j = 1, 2, \dots$$

then

$$\lambda_1(J(0, \epsilon)) \geq \lambda_2(J(0, \epsilon)), \quad \forall \epsilon \geq 0,$$

where $J(0, \epsilon)$ stands for the ball with center zero and radius ϵ . \blacksquare

4. Spline algorithm and adaptive information on the average.

In this section we prove that for every error functional E and for every nonadaptive information, a translated spline algorithm is optimal. We also prove that for every adaptive information N^a there exists nonadaptive information of the same cardinality and whose radius is not greater than the radius of N^a .

Let N^a and φ be given. Recall that the (global) average error of φ is defined by

$$(4.1) \quad e^{\text{avg}}(\varphi, N^a) = \int_{F_1} E(Sf - \varphi(N^a f))_{\perp}(df)$$

and the (global) average radius of N^a , by

$$(4.2) \quad r^{\text{avg}}(N^a) = \inf_{\varphi} e^{\text{avg}}(\varphi, N^a).$$

Hence the global average radius of N^a is the minimal global average error made by any algorithm φ that uses N^a , and the optimal algorithm φ^* that uses N^a is defined so that its error is minimal, i.e.,

$$(4.3) \quad e^{\text{avg}}(\varphi^*, N^a) = r^{\text{avg}}(N^a).$$

We now define the concept of the local average error as studied in [6]. Due to (3.5) and Theorem 3.1(i)

$$(4.4) \quad e_{avg}(\omega, N_a) = \int_{\mathbb{R}^n} e_{avg}(\omega, N_a, Y) \mu_1(dy),$$

where the local average error $e_{avg}(\omega, N_a, Y)$ is given by

$$(4.5) \quad e_{avg}(\omega, N_a, Y) = \int E(Sf - \omega)(Y) \mu_2(dY | Y, N_a).$$

Theorem 4.1: For every nonadaptive information N_{non} of

the form (3.7) the average radius is given by

$$(4.6) \quad r_{avg}(N_{non}) = \inf_{g \in F_2} \int_{F_1} E(Sf - g) \mu_2(dY | 0, N_{non}).$$

Let

$$(4.7) \quad P = \{g^* \in F_2 : \int_{F_1} E(Sf - g^*) \mu_2(dY | 0, N_{non})\}.$$

$$= r_{avg}(N_{non}).$$

An algorithm ω^* that uses N_{non} is optimal iff

$$(4.8) \quad g(Y) \stackrel{df}{=} \omega^*(Y) - Sm(N_{non}, Y) \in P, \text{ for almost}$$

every $Y \in \mathbb{R}^n$.

Proof: Let ω be an algorithm that uses N_{non} . Consider

the local error $e_{avg}(\omega, N_{non}, Y)$. Due to (3.16) and

linearity of S ,

$$e_{avg}(\omega, N_{non}, Y) = \int_{F_1} E(S(f+m)(N_{non}, Y))$$

$$- \omega(Y) \mu_2(dY | 0, N_{non})$$

$$= \int_{F_1} E(Sf - \omega)(N_{non}, Y) \mu_2(dY | 0, N_{non})$$

$$\geq \inf_{g \in F_2} \int_{F_1} E(Sf-g) \mu_2(df|0, N^{\text{non}}) \stackrel{\text{df}}{=} H.$$

This proves that

$$r^{\text{avg}}(N^{\text{non}}) \geq H.$$

To prove that $r^{\text{avg}}(N^{\text{non}}) = H$ we can assume that H is finite. Then for every $\delta > 0$, there exists $g_\sigma \in F_2$ such that $\int_{F_1} E(Sf-g_\sigma) \mu_2(df|0, N^{\text{non}}) \leq H + \delta$. Define $\varphi_\delta(y) = Sm(N^{\text{non}}, y) + g_\sigma$. Then

$$e^{\text{avg}}(\varphi_\delta, N^{\text{non}}) = \int_{\mathbb{R}^n} e^{\text{avg}}(\varphi_\delta, N^{\text{non}}, y) \mu_1(dy) \leq H + \delta.$$

Since δ is arbitrary, $r^{\text{avg}}(N^{\text{non}}) \leq H$ and consequently $r^{\text{avg}}(N^{\text{non}}) = H$. This proves (4.6). To complete the proof observe that if $H = +\infty$ then every algorithm is optimal and $P = F_2$. Therefore we can assume that $H < +\infty$. If $\varphi^*(y) = Sm(N^{\text{non}}, y) + g^*(y)$ with $g^*(y) \in P$ for almost every y , then, obviously, $\varphi^*(y) = Sm(N^{\text{non}}, y) + g^*(y)$ with $g^*(y) \in P$ for almost every y , then, obviously, φ^* is optimal. On the other hand, take an arbitrary algorithm φ . Define

$$Y = \{y \in \mathbb{R}^n : g(y) = \varphi(y) - Sm(N^{\text{non}}, y) \notin P\}.$$

If Y has a positive μ_1 measure, then

$$\begin{aligned}
e^{\text{avg}}_{(\varphi, N^{\text{non}})} &= \int_Y \int_{F_1} E(Sf-g(y))_{u_2} (df|0, N^{\text{non}})_{u_1} (dy) \\
&\quad + \int_{\mathbb{R}^n \setminus Y} \int_{F_1} E(Sf-g(y))_{u_2} (df|0, N^{\text{non}})_{u_1} (dy) \\
&> u_1(Y) r^{\text{avg}}_{(N^{\text{non}})} + u_1(\mathbb{R}^n \setminus Y) r^{\text{avg}}_{(N^{\text{non}})} \\
&= r^{\text{avg}}_{(N^{\text{non}})}.
\end{aligned}$$

Hence φ is not optimal. This completes the proof of Theorem 4.1. ■

Theorem 4.1 states that there exists an optimal algorithm iff the infimum in (4.8) is attained by some element g^* . Of course g^* need not be unique, but taking any g^* satisfying (4.8), the algorithm

$$\varphi^*(\cdot) = \varphi^S(\cdot) + g^*$$

is optimal where

$$\varphi^S(N^{\text{non}}f) = \sum_{i=1}^n (f, g_i) S_{\mu_i} g_i = S m(N^{\text{non}}, N^{\text{non}}f).$$

The algorithm φ^S , called the spline algorithm, is linear.

Hence φ^* is an affine mapping, which is a desirable property from the complexity point of view. On the other hand if the infimum in (4.8) is not attained, i.e., $P = \emptyset$, then there is no optimal algorithm. In this case taking g^* so that $\int_{F_1} E(Sf-g^*)_{u_2} (df|0, N^{\text{non}})$ is sufficiently

close to $r^{\text{avg}}(N^{\text{non}})$, say is not greater than $r^{\text{avg}}(N^{\text{non}}) + \delta$,
the following affine algorithm

$$\varpi^*(\cdot) = \varpi^s(\cdot) + g^*$$

is almost optimal, since

$$e^{\text{avg}}(\varpi^*, N^{\text{non}}) \leq r^{\text{avg}}(N^{\text{non}}) + \delta.$$

We now prove that adaption does not help on the average.
Let N^a be adaptive information of the form (3.7), and let

$$H(y) = \inf_{g \in F_2} \int_{F_2} E(Sf - g) \mu_2(df | 0, N_y^{\text{non}}).$$

Then, due to (3.15) and (4.6),

$$(4.9) \quad r^{\text{avg}}(N^a) = \int_{\mathbb{R}^n} H(y) \mu_1(dy) = \int_{\mathbb{R}^n} r^{\text{avg}}(N_y^{\text{non}})^2 \mu_1(dy).$$

Let $y^*, y^* \in \mathbb{R}^n$, be such that

$$(4.10) \quad r^{\text{avg}}(N_{y^*}^{\text{non}}) \leq r^{\text{avg}}(N^a).$$

Observe that such y^* exists. Indeed, $r^{\text{avg}}(N_y^{\text{non}}) > r^{\text{avg}}(N^a)$
for every y would contradict (4.9). Hence we have proven

Theorem 4.2: For every adaptive information N^a there
exists $y^* \in \mathbb{R}^n$ such that

$$r^{\text{avg}}(N_{y^*}^{\text{non}}) \leq r^{\text{avg}}(N^a).$$

We now give a sufficient condition on the error

functional E so that the spline algorithm φ^S is optimal. Technically, this means that $0 \in P$.

Theorem 4.3: If E is convex and symmetric (with respect to zero) then for every nonadaptive information N^{non} the spline algorithm φ^S is optimal. ■

Proof: Although Theorem 4.3 follows immediately from [6], we present its proof for the completeness. Take $g \in F_2$. Then, due to the symmetricity of $\mu_2(\cdot) = \mu_2(\cdot | 0, N^{\text{non}})$ (i.e., $\mu_2(B) = \mu_2(-B)$, $\forall B \in \mathcal{B}(F_1)$),

$$\int_{F_1} E(Sf-g)\mu_2(df) = \frac{1}{2} \int_{F_1} [E(Sf-g) + E(-Sf-g)]\mu_2(df).$$

Since E is symmetric and convex,

$$\frac{1}{2}[E(Sf-g)+E(Sf+g)] = \frac{1}{2}[E(Sf-g)+E(Sf+g)] \geq E(Sf).$$

Hence

$$\int_{F_1} E(Sf-g)\mu_2(df) \geq \int_{F_1} E(Sf)\mu_2(df), \quad \forall g \in F_2.$$

This proves that $g^* = 0 \in P$ and completes the proof of Theorem 4.2. ■

Remark 4.1: Optimality of the spline algorithm on the average has been established in [7,8] for orthogonally invariant measures assuming that F_2 is a separable Hilbert space and $E(g) = \|g\|^2$. The same result was obtained in

[5] assuming that F_1 is a finite dimensional space and $E(Sf - \varrho(N^a f)) = \int \|Sf - \varrho(N^a f)\|_D^2 ((S f, f))$ for some function ϱ .

In this paper, restricting the class of probability measures to Gaussian measures we relax the assumptions concerning E and the spaces F_1 and F_2 .

We now exhibit an n -th optimal information operator N^* of $\text{card}(N^*) = n$, i.e., N^* satisfying

$$r^{\text{avg}}(N^*) \leq r^{\text{avg}}(N^a), \quad \forall N^a, \quad \text{card}(N^a) = n.$$

We find N^* under some additional assumptions on F_2, S and E . Namely, we assume that F_2 is a separable Hilbert space, S is continuous and

$$(4.11) \quad E(g) = H(\|g\|)$$

for some function $H: \mathbb{R}_+ \rightarrow \mathbb{R}_+$ which is convex and nondecreasing. Observe that then E is convex and symmetric and therefore the spline algorithm is optimal for every nonadaptive information.

To find N^* we proceed as follows. Let

$$(4.12) \quad R = \int S S^*: F_2 \rightarrow F_2.$$

Since S is continuous, R is a nonnegative definite

operator with finite trace. Let $\zeta_1^*, \zeta_2^*, \dots$ be eigenelements of R corresponding to the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$, i.e.,

$$R\zeta_i^* = \lambda_i \zeta_i^*, \quad (\zeta_i^*, \zeta_j^*) = \delta_{i,j}.$$

Take

$$(4.13) \quad g_i^* = \frac{1}{\sqrt{\lambda_i}} S^* \zeta_i^*, \quad i = 1, 2, \dots$$

Remark 4.2: The optimal information for the average case setting studied in [7] is derived from the operator K defined by

$$K = S_{\mu}^{1/2} S^* S S_{\mu}^{1/2} : F_1 \rightarrow F_1.$$

Observe that if η is an eigenvector of K corresponding to an eigenvalue β , $K\eta = \beta\eta$, then letting

$$\zeta = S S_{\mu}^{1/2} \eta$$

we get

$$R\zeta = S S_{\mu} S^* S S_{\mu}^{1/2} \eta = S S_{\mu}^{1/2} K\eta = \beta S S_{\mu}^{1/2} \eta = \beta\zeta.$$

Hence the operators K and R have the same eigenvalues. Furthermore η is an eigenvector of K iff $S S_{\mu}^{1/2}$ is an eigenvector of R . ■

Define the nonadaptive information operator

$$(4.14) \quad N^*(f) = [(f, g_1^*), \dots, (f, g_n^*)].$$

Note that N^* satisfies (3.7).

Theorem 4.4: The information operator N^* defined by

(4.14) is n th optimal.

Proof: Due to Theorem 4.2, we need only to prove that

$$r^{\text{avg}}(N^*) \leq r^{\text{avg}}(N^{\text{non}}),$$

for every N^{non} of the form (3.7). Due to Theorems 4.1 and 4.3,

$$r^{\text{avg}}(N^{\text{non}}) = \int_{F_1} H(\|Sf\|) \mu_2(df | 0, N^{\text{non}}).$$

If H is constant, $H(x) \equiv c$, then $r^{\text{avg}}(N^{\text{non}}) = r^{\text{avg}}(N^*) = c$ for every N^{non} . Hence without loss of generality we can assume that H is not constant. Then $H(\mathbb{R}_+) = [H(0), +\infty)$.

Indeed, convexity of H yields $2H(x) \leq H(0) + H(x)$,

$\forall x \in \mathbb{R}_+$. Since H is nondecreasing, $\sup\{H(x) : x \in \mathbb{R}_+\}$

$= \lim_{x \rightarrow \infty} H(x) = c$. Note that $H(0) < c$. If $c < +\infty$ then

$2c \leq H(0) + c < 2c$ which is a contradiction. Hence

$H(\mathbb{R}_+) = [H(0), +\infty)$ as claimed. Define

$$(4.15) \quad \gamma(B, N^{\text{non}}) = \mu_2(\{f \in F_1 : H(\|Sf\|) \in B\} | 0, N^{\text{non}}),$$

$$\forall B \in \mathcal{B}(H(\mathbb{R}_+)).$$

Then $\gamma(\cdot, N^{\text{non}})$ is a probability measure on $\mathbb{B}(H(\mathbb{R}_+))$ and

$$(4.16) \quad r^{\text{avg}}(N^{\text{non}}) = \int_{H(0)}^{+\infty} t \gamma(dt, N^{\text{non}}).$$

Let $D(\cdot, N^{\text{non}})$ be the distribution function for $\gamma(\cdot, N^{\text{non}})$, i.e.,

$$(4.17) \quad D(x, N^{\text{non}}) = \int_{H(0)}^x \gamma(dt, N^{\text{non}}), \quad \forall x \in H(\mathbb{R}_+).$$

We shall prove that

$$(4.18) \quad D(x, N^{\text{non}}) \leq D(x, N^*), \quad \forall x \in H(\mathbb{R}_+), \quad \forall N^{\text{non}}.$$

Before presenting the proof of (4.18), we show that (4.18) will complete the proof of Theorem 4.4. For this end, observe that

$$t = \lim_k \sum_{i=1}^k t_{i,k} \chi_{(a_{i,k}, a_{i+1,k}]}(t), \quad \forall t > H(0),$$

for some numbers $H(0) = a_{1,k} < a_{2,k} < \dots < a_{k,k} < a_{k+1,k} = +\infty$ and $t_i \in (a_{i,k}, a_{i+1,k}]$. Hence for every N^{non} ,

$$\begin{aligned} r^{\text{avg}}(N^{\text{non}}) &= \lim_k \sum_{i=1}^k t_{i,k} \gamma((a_{i,k}, a_{i+1,k}], N^{\text{non}}) \\ &= \lim_k \sum_{i=1}^k t_{i,k} [D(a_{i+1,k}, N^{\text{non}}) - D(a_{i,k}, N^{\text{non}})] \\ &= \lim_k [\sum_{i=1}^{k-1} (t_{i,k} - t_{i+1,k}) D(a_{i+1,k}, N^{\text{non}}) + t_{k,k}] \end{aligned}$$

since $D(a_{1,k}, N^{\text{non}}) = 0$ and $D(a_{k+1,k}, N^{\text{non}}) = D(+\infty, N^{\text{non}}) = 1$.

Hence,

$$\begin{aligned}
& r^{\text{avg}}(N^{\text{non}}) - r^{\text{avg}}(N^*) \\
&= \lim_k \sum_{i=1}^k (t_{i,k} - t_{i+1,k}) [D(a_{i+1,k}, N^{\text{non}}) - D(a_{i+1,k}, N^*)]
\end{aligned}$$

and $t_i - t_{i+1} < 0$ imply that

$$r^{\text{avg}}(N^{\text{non}}) - r^{\text{avg}}(N^*) \geq 0.$$

Hence to complete the proof of Theorem 4.4 it is enough to show that (4.18) holds. Observe that

$$\begin{aligned}
(4.19) \quad D(x, N^{\text{non}}) &= \mu_2(\{f: H(\|Sf\|) \leq x\} | 0, N^{\text{non}}) \\
&= \mu_2(\{f: \|Sf\| \leq H^{-1}(x)\} | 0, N^{\text{non}}).
\end{aligned}$$

Define

$$\lambda(B, N^{\text{non}}) = \mu_2(\{f \in F_1: Sf \in B\} | 0, N^{\text{non}}),$$

$$\forall B \in \mathcal{B}(F_2).$$

Then $\lambda(\cdot, N^{\text{non}})$ is a probability measure on $\mathcal{B}(F_2)$ and, due to (4.19),

$$(4.20) \quad D(x, N^{\text{non}}) = \lambda(J(0, z), N^{\text{non}}),$$

where now $z = H^{-1}(x)$ and $J(0, z)$ is the ball in F_2 with center zero and radius z . We need the following two lemmas.

Lemma 4.1: For every N^{non} , $\lambda(\cdot, N^{\text{non}})$ is the Gaussian measure with mean element zero and correlation operator

$$R_{N^{\text{non}}} = S(I - \sigma_{N^{\text{non}}}) S_{\mu} (I - \sigma_{N^{\text{non}}}^*) S^* \dots \quad \blacksquare$$

Proof: Observe that for the characteristic functional

$\psi_{N^{\text{non}}}$ of $\lambda(\cdot, N^{\text{non}})$ we have

$$\psi_{N^{\text{non}}}(h) = \int_{F_2} \exp\{i(g, h)\} \lambda(dg, N^{\text{non}}).$$

Change variables by setting $f = Sg$. Then

$$\begin{aligned} \psi_{N^{\text{non}}}(y) &= \int_{F_1} \exp\{i(f, S^*h)\} \mu_2(df | 0, N^{\text{non}}) \\ &= \exp\left\{-\frac{1}{2} \left((I - \sigma_{N^{\text{non}}}) S_{\mu} (I - \sigma_{N^{\text{non}}}^*) S^*h, S^*h \right)\right\} \\ &= \exp\left\{-\frac{1}{2} (R_{N^{\text{non}}} h, h)\right\}, \quad \forall h \in F_2. \end{aligned}$$

This completes the proof of Lemma 4.1. \blacksquare

Let $\gamma_1, \gamma_2, \dots$ ($\gamma_i \geq \gamma_{i+1} \geq 0$) be the eigenvalues of $R_{N^{\text{non}}}$. It is easy to check that for N^* , $\lambda_{n+1}, \lambda_{n+2}, \dots$ are the dominating eigenvalues of R_{N^*} .

Lemma 4.2:

$$(4.21) \quad \lambda_{n+k} \leq \gamma_k, \quad \forall k = 1, 2, \dots \quad \blacksquare$$

Proof (induction on k): For $k = 1$, (4.21) holds trivially.

Suppose therefore that (4.21) holds for every $k \leq k_0$. We prove that (4.21) also holds for $k = k_0 + 1$.

For this end, let $\eta_1, \eta_2, \dots, \eta_k$ be eigenelements of

$K_{N^{\text{non}}}$ corresponding to $\gamma_1, \gamma_2, \dots, \gamma_k$. Take

$$g = \sum_{i=1}^{n+k} x_i \zeta_i^* \in F_2$$

Such that

$$(4.22) \quad \|g\|^2 = \sum_{i=1}^{n+k} x_i^2 = 1,$$

$$(4.23) \quad \sigma_{N^{\text{non}}}^*(S^*g) = 0,$$

and

$$(4.24) \quad (g, \eta_i) = 0, \quad i = 1, 2, \dots, k_0 = k-1.$$

Since (4.23) and (4.24) are equivalent to a homogeneous system of $n+k-1$ linear equations with $n+k$ unknowns, such g exists. Furthermore, (4.22) and (4.24) yield that $\gamma_k \geq (R_{N^{\text{non}}} g, g)$. Hence, due to (4.23), we get

$$\begin{aligned} \gamma_k \geq (R_{N^{\text{non}}} g, g) &= (Rg, g) = \sum_{i=1}^{n+k} \lambda_i x_i^2 \geq \lambda_{n+k} \sum_{i=1}^{n+k} x_i^2 \\ &= \lambda_{n+k}, \end{aligned}$$

which completes the proof of Lemma 4.2. ■

We are ready to complete the proof of Theorem 4.4. Due to Lemmas 4.1, 4.2 and 3.2,

$$\lambda(J(0, z), N^*) \geq \lambda(J(0, z), N^{\text{non}}), \quad \forall N^{\text{non}}, \quad \forall z \in \mathbb{R}_+.$$

Hence (4.20) yields that

$$D(x, N^*) \geq D(x, N^{\text{non}}), \quad \forall N^{\text{non}}, \forall x \in \mathbb{R}_+.$$

This completes the proof of (4.18) as well as the proof of Theorem 4.4. ■

5. Spline algorithm and adaptive information in the probabilistic setting.

In this section we prove that for every error functional E and for every nonadaptive information N , the probability of the fact that the error does not exceed ϵ , is maximized by a translated spline algorithm. We also prove that adaption does not help in this setting.

Recall that for given $\epsilon \geq 0$, N^a and φ ,

$$(5.1) \quad \text{Prob}(\varphi, N^a, \epsilon) = \mu(\{f \in F_1 : E(Sf - \varphi(N^a f)) \leq \epsilon\})$$

is the probability of the fact that the error $E(Sf - \varphi(N^a f))$ made by φ is not greater than ϵ , and

$$(5.2) \quad \text{Prob}(N^a, \epsilon) = \inf_{\varphi} \text{Prob}(\varphi, N^a, \epsilon).$$

Then $\text{Prob}(N^a, \epsilon)$ is the maximal probability among all algorithms that use N^a , and the optimal algorithm φ^* that uses N^a is defined so that

$$(5.3) \quad \text{Prob}(\varphi^*, N^a, \epsilon) = \text{Prob}(N^a, \epsilon).$$

Theorem 5.1: For every nonadaptive information N^{non} of the form (3.7)

$$(5.4) \quad \text{Prob}(N^{\text{non}}, \epsilon) = \sup_{g \in F_2} \mu_2(\{f \in F_1 : E(Sf - g) \leq \epsilon\} | 0, N^{\text{non}}).$$

Let

$$(5.5) \quad P = \{g^* \in F_2: \mu_2(\{f \in F_1: E(Sf-g^*) \leq \epsilon\} | 0, N^{\text{non}})\} \\ = \text{Prob}(N^{\text{non}}, \epsilon).$$

An algorithm φ^* that uses N^{non} is optimal iff

$$(5.6) \quad g(y) \stackrel{\text{df}}{=} \varphi^*(y) - S m(N^{\text{non}}, y) \in P, \text{ for almost} \\ \text{every } y \in \mathbb{R}^n. \quad \blacksquare$$

Proof: The proof of this theorem differs from the proof of Theorem 4.1 only at the beginning. Observe that for every algorithm φ that uses N^{non} we have, due to (3.16) and linearity of S ,

$$\text{Prob}(\varphi, N^{\text{non}}, \epsilon) = \int_{\mathbb{R}^n} \mu_2(\{f \in F_1: E(Sf - \varphi(y)) \leq \epsilon\} | y, N^{\text{non}}) \mu_1(dy) \\ = \int_{\mathbb{R}^n} \mu_2(\{f \in F_1: E(Sf - (\varphi(y) - S m(N^{\text{non}}, y))) \\ \leq \epsilon\} | 0, N^{\text{non}}) \mu_1(dy) \\ \leq \sup_{g \in F_2} \mu_2(\{f \in F_1: E(Sf - g) \leq \epsilon\} | 0, N^{\text{non}}).$$

Hence using the same reasoning as in the proof of Theorem 4.1 one can easily complete the proof of Theorem 5.1.

Therefore we skip this part. \blacksquare

Let N^a be adaptive. Similar as in (4.10), let $y^*, y^* \in \mathbb{R}^n$, be such that

$$(5.7) \quad \text{Prob}(N^a, \epsilon) = \int_{\mathbb{R}^n} \text{Prob}(N_y^{\text{non}}, \epsilon) \mu_1(dy) \leq \text{Prob}(N_{y^*}^{\text{non}}, \epsilon).$$

Of course, such y^* exists.

Theorem 5.2: For every adaptive information N^a there exists $y^* \in \mathbb{R}^n$ such that

$$\text{Prob}(N_{y^*}^{\text{non}}, \epsilon) \leq \text{Prob}(N^a, \epsilon). \quad \blacksquare$$

As in Section 4 we give a sufficient condition on E for the spline algorithm to be optimal, i.e., $g^* = 0 \in P$.

Theorem 5.3: If E is convex and symmetric (with respect to zero) and if $F_2 = S(F_1)$ then for every nonadaptive information N^{non} the spline algorithm φ^S is optimal. \blacksquare

Proof: To prove this theorem it is enough to show that

$$(5.8) \quad \mu_2(\{f \in F_1: E(Sf) \leq \epsilon\}) \geq \mu_2(\{f \in F_1: E(Sf-g) \leq \epsilon\}),$$

$$\forall g \in F_2,$$

where $\mu_2(\cdot) = \mu_2(\cdot | 0, N^{\text{non}})$. Let

$$B(g) = \{f \in F_1: E(Sf-g) \leq \epsilon\}$$

and

$$B = B(0) = \{f \in F_1: E(Sf) \leq \epsilon\}.$$

Since $F_2 = S(F_1)$, there exists an element $h \in F_1$ such that $Sh = g$. Observe that

$$B(g) \subseteq B + h.$$

Indeed, for $f \in B(g)$ let $\tilde{f} = f - h$. Since $E(S\tilde{f}) = E(S(f-h)) = E(Sf-g) \leq \epsilon$. Thus $\tilde{f} \in B$ and $f = \tilde{f} + h \in B + h$ as claimed. This means that

$$\mu_2(\{f \in F_1 : E(Sf-g) \leq \epsilon\}) = \mu_2(B(g)) \leq \mu_2(B + h).$$

Hence to prove (5.8) we need only to show that

$$(5.9) \quad \mu_2(B) \geq \mu_2(B + h), \quad \forall h \in F_1.$$

Observe that B is convex and balanced. Indeed, if $f_1, f_2 \in B$ then $E(tf_1 + (1-t)f_2) \leq tE(f_1) + (1-t)E(f_2) \leq \epsilon$, i.e., $tf_1 + (1-t)f_2 \in B$, and if $f \in B$ then $E(-f) = E(f) \leq \epsilon$, i.e., $-f \in B$. Since μ_2 is a Gaussian measure with mean element zero, Lemma 3.1 completes the proof of Theorem 5.3. ■

The next theorem is about n-th optimal information

N^* . The information N^* of cardinality n is optimal iff

$$\text{Prob}(N^*, \epsilon) \geq \text{Prob}(N^a, \epsilon), \quad \forall N^a, \text{card}(N^a) = n.$$

Theorem 5.4: Let E be of the form (4.11) and let S be continuous. Then the information N^* defined by (4.14)

is n -th optimal for every $\epsilon \geq 0$. ■

Proof: This theorem follows immediately from (4.18).
Indeed, due to Theorem 5.2, we need only to consider
nonadaptive information N^{non} . But then for every N^{non}
and every $\epsilon \geq 0$,

$$\text{Prob}(N^{\text{non}}, \epsilon) = \mu_2(\{f \in F_1 : H(\|Sf\|) \leq \epsilon\}) = D(H^{-1}(\epsilon), N^{\text{non}}).$$

Hence (4.18) implies that

$$\text{Prob}(N^*, \epsilon) \geq \text{Prob}(N^{\text{non}}, \epsilon), \quad \forall N^{\text{non}}, \quad \forall \epsilon \geq 0,$$

which completes the proof. ■

We end this section by the following problem. For
a given set $A \in \mathcal{B}(\mathbb{R}^n)$ let

$$(5.10) \quad \text{Prob}(\varphi, N^a, \epsilon, A) = \mu(\{f \in F_1 : E(Sf - \varphi(N^a f)) \\ \leq \epsilon \wedge N^a f \in A\}).$$

We want to find φ^* such that

$$(5.11) \quad \text{Prob}(N^a, \epsilon, A) \stackrel{\text{df}}{=} \sup_{\varphi} \text{Prob}(\varphi, N^a, \epsilon, A) \\ = \text{Prob}(\varphi^*, N^a, \epsilon, A).$$

Observe that $\text{Prob}(\varphi, N^a, \epsilon, A)$ is the probability that
 $E(Sf - \varphi(N^a f)) \leq \epsilon$ under the condition that $N^a(f) \in A$. Of

course, for $A = \mathbb{R}^n$, $\text{Prob}(\varphi, N^a, \epsilon, A) = \text{Prob}(\varphi, N^a, \epsilon)$.

For every adaptive information N^a ,

$$\begin{aligned} \text{Prob}(\varphi, N^a, \epsilon, A) &= \int_A \mu_2(\{f \in F_1 : E(Sf - \varphi(Y)) \leq \epsilon\} | Y, N_Y^{\text{non}}) \mu_1(dy) \\ &\leq \int_A \sup_{g \in F_2} \mu_2(\{f \in F_1 : E(Sf - g) \leq \epsilon\} | 0, N_Y^{\text{non}}) \mu_1(dy) \\ &= \int_A \text{Prob}(N_Y^{\text{non}}, \epsilon) \mu_1(dy). \end{aligned}$$

From this we can conclude

Theorem 5.5:

(i) For every adaptive information N^a there exists $y^* \in \mathbb{R}^n$ such that

$$\text{Prob}(N^a, \epsilon, A) \leq \text{Prob}(N_{y^*}^{\text{non}}, \epsilon, A), \quad \forall \epsilon, \quad \forall \epsilon \geq 0, \quad \forall A \in \mathcal{B}(\mathbb{R}^n).$$

(ii) For every nonadaptive information N^{non}

$$\text{Prob}(N^{\text{non}}, \epsilon, A) = \text{Prob}(N^{\text{non}}, \epsilon) \mu_1(A), \quad \forall \epsilon, \quad \forall \epsilon \geq 0$$

$$\forall A \in \mathcal{B}(\mathbb{R}^n).$$

In particular, φ^* is optimal independently of A .

(iii) If $F_2 = S(F_1)$ and E is convex and symmetric (with respect to zero) then the spline algorithm φ^S is optimal for every N^{non} , every $\epsilon \geq 0$ and every $A \in \mathcal{B}(\mathbb{R}^n)$.

(iv) If F_2 is a separable Hilbert space, S is continuous and E is of the form (4.11), then N^* defined

by (4.14) is optimal for every $\epsilon \geq 0$ and every $A \in \mathcal{B}(\mathbb{R}^n)$. ■

Theorem 5.5 states that the probability of a small error does not depend on the value N^{non} of information. This result will be used in a future paper for studying optimal stopping criteria.

5. Variance of spline algorithm.

In previous sections we showed when the spline algorithm φ^s is optimal. Here we exhibit another optimality property of φ^s showing that it minimizes the variance whenever F_2 is a Hilbert space and $E(g) = \|g\|^2$.

Let N^{non} be nonadaptive information and let φ be an algorithm that uses N^{non} . By the variance of φ we mean

$$(6.1) \quad \text{var}(\varphi) = \int_{F_1} \{ \|Sf - \varphi(N^{\text{non}}f)\|^2 - e^{\text{avg}}(\varphi, N^{\text{non}}) \}^2 \mu(df),$$

where

$$e^{\text{avg}}(\varphi, N^{\text{non}}) = \int_{F_1} \|Sf - \varphi(N^{\text{non}}f)\|^2 \mu(df).$$

Theorem 6.1:

$$(6.2) \quad \text{var}(\varphi^s) = \inf_{\varphi} \text{var}(\varphi). \quad \blacksquare$$

Proof: Let φ be an algorithm. Define $h = \varphi - \varphi^s$, i.e., $\varphi(N^{\text{non}}f) = \varphi^s(N^{\text{non}}f) + h(N^{\text{non}}f)$. Then due to (3.18),

$$\begin{aligned} \text{var}(\varphi) &= \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf - \varphi^s(y) - h(y)\|^2 \\ &\quad - e^{\text{avg}}(\varphi, N^{\text{non}}) \}^2 \mu_2(df|y, N^{\text{non}}) \mu_1(dy) \\ &= \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf - h(y)\|^2 - e^{\text{avg}}(\varphi, N^{\text{non}}) \}^2 \mu_2(df) \mu_1(dy), \end{aligned}$$

where $\mu_2(\cdot) = \mu_2(\cdot | 0, N^{\text{non}})$. Observe that $\|Sf - h(y)\|^2$
 $= \|Sf\|^2 - 2(f, S^*h(y)) + \|h(y)\|^2$. Since mean element of
 μ_2 is zero, $\int_{F_1} (f, S^*h(y)) \mu_2(df) = 0$, and

$$\begin{aligned} e^{\text{avg}}(\varphi, N^{\text{non}}) &= \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf\|^2 - 2(f, S^*h(y)) \\ &\quad + \|h(y)\|^2 \} \mu_2(df) \mu_1(dy) \\ &= \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf\|^2 + \|h(y)\|^2 \} \mu_2(df) \mu_1(dy) \\ &= e^{\text{avg}}(\varphi^s, N^{\text{non}}) + \int_{\mathbf{R}^n} \|h(y)\|^2 \mu_1(dy). \end{aligned}$$

Hence

$$\begin{aligned} \text{var}(\varphi) &= \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) - 2(f, S^*h(y)) + \|h(y)\|^2 \\ &\quad - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \mu_2(df) \mu_1(dy). \end{aligned}$$

Change the variables by letting $f = -f$. Then

$$\begin{aligned} \text{var}(\varphi) &= \frac{1}{2} \int_{\mathbf{R}^n} \int_{F_1} \{ (\|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) - 2(f, S^*h(y)) + \|h(y)\|^2 \\ &\quad - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \\ &\quad + (\|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) + 2(f, S^*h(y)) + \|h(y)\|^2 \\ &\quad - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \} \mu_2(df) \mu_1(dy) \\ &\geq \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) + \|h(y)\|^2 \\ &\quad - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \mu_2(df) \mu_1(dy), \end{aligned}$$

since $\frac{1}{2}[(a+b)^2 + (a-b)^2] \geq a^2$. Hence

$$(6.3) \quad \text{var}(\varphi) \geq \text{var}(\varphi^s) + 2H_1 + H_2,$$

where

$$H_1 = \int_{\mathbf{R}^n} \int_{F_1} \{ \|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) \} \{ \|h(y)\|^2 - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \mu_2(df) \mu_1(dy)$$

and

$$H_2 = \int_{\mathbf{R}^n} \int_{F_1} \{ \|h(y)\|^2 - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \mu_2(df) \mu_1(dy).$$

Of course, $H_2 \geq 0$ and therefore

$$(6.4) \quad \text{var}(\varphi) \geq \text{var}(\varphi^s) + 2H_1.$$

We now prove that $H_1 = 0$. Indeed,

$$H_1 = \int_{\mathbf{R}^n} \{ \{ \|h(y)\|^2 - \int_{\mathbf{R}^n} \|h(z)\|^2 \mu_1(dz) \}^2 \cdot \int_{F_1} \{ \|Sf\|^2 - e^{\text{avg}}(\varphi^s, N^{\text{non}}) \} \mu_2(df) \} \mu_1(dy)$$

and since $e^{\text{avg}}(\varphi^s, N^{\text{non}}) = \int_{F_1} \|Sf\|^2 \mu_2(df)$, see Theorem 4.3 and (4.8), $H_1 = 0$ as claimed. Hence

$$\text{var}(\varphi) \geq \text{var}(\varphi^s), \quad \forall \varphi,$$

which completes the proof. ■

We want to stress that the minimal variance of the

spline algorithm strongly depends on the form of E ,
i.e., $E(g) = \|g\|^2$. For arbitrary E (even convex
and symmetric) the spline algorithm need not minimize
the variance.

7. Appendix.

We prove Lemmas 3.1 and 3.2. Since these lemmas are well known for finite dimensional spaces, the proofs are mainly to show that the infinite dimensional case can be reduced to a finite dimensional one.

Proof of Lemma 3.1: We prove that (3.17) can be reduced to a problem with a finite dimensional Gaussian measure. Then the Anderson's inequality will complete the proof.

Let ζ_1, ζ_2, \dots be eigenvalues of the covariance operator S_λ , $S_\lambda \zeta_i = \alpha_i \zeta_i$ and $(\zeta_i, \zeta_j) = \delta_{ij}$. Let $X = \ker S_\lambda$ and let X^\perp be the orthogonal complement of X , $F_1 = X^\perp \oplus X$. Then for every $f \in F_1$, $f = f_1 + f_2$, where $f_1 \in X$ and $f_2 \in X^\perp$, and for every $C \in \mathcal{B}(F_1)$

$$(A.1) \quad \lambda(C) = \lambda^\perp(C \cap X^\perp),$$

where λ^\perp is the Gaussian measure on $\mathcal{B}(X^\perp)$ with mean element zero and covariance operator $S_{\lambda^\perp} = S_\lambda|_{X^\perp}$ (see [4]). Observe that $B \cap X^\perp$ is convex and balanced and that $(B + h) \cap X^\perp \subset (B \cap X^\perp) + h_2$ ($h = h_1 + h_2$, $h_1 \in X$ and $h_2 \in X^\perp$). Hence, due to (A.1),

$$\lambda(B) = \lambda^\perp(B \cap X^\perp) \text{ and } \lambda(B+h) = \lambda^\perp((B+h) \cap X^\perp) \leq \lambda^\perp((B \cap X^\perp) + h_2).$$

This means that to prove (3.17) we can assume without

loss of generality that $X^\perp = F_1$ and $\lambda = \lambda^\perp$, i.e., that all eigenvalues of S_λ are positive.

For $k = 1, 2, \dots$ define $P_k: F_1 \rightarrow \mathbb{R}^k$,

$$(A.2) \quad P_k(f) = \left[\left(f, \frac{\zeta_1}{\sqrt{\alpha_1}} \right), \dots, \left(f, \frac{\zeta_k}{\sqrt{\alpha_k}} \right) \right].$$

Observe that for every set $C \in \mathcal{B}(F_1)$, $P_k^{-1}(P_k(C)) \supset P_{k+1}^{-1}(P_{k+1}(C))$ and $\bigcap_{k=1}^{\infty} P_k^{-1}(P_k(C)) = C$. Hence

$$(A.3) \quad \lambda(C) = \lim_k \lambda(P_k^{-1}(P_k(C))), \quad \forall C \in \mathcal{B}(F_1).$$

Let λ_k be the probability measure on $\mathcal{B}(\mathbb{R}^k)$ induced by P_k , i.e.,

$$\lambda_k(A) = \lambda(P_k^{-1}(A)), \quad \forall A \in \mathcal{B}(\mathbb{R}^k).$$

Then (A.3) can be rewritten as

$$(A.4) \quad \lambda(C) = \lim_k \lambda_k(P_k(C)).$$

Since for every $k = 1, 2, \dots$ the operator P_k is of the form (3.7) then, due to (3.10), λ_k is the Gaussian measure on $\mathcal{B}(\mathbb{R}^k)$ with mean element zero and correlation operator identity. Observe also, that $P_k(B)$ is convex and balanced and that $P_k(B+h) = P_k(B) + P_k(h)$. Hence the Anderson's inequality [1] yields that

$$(A.5) \quad \lambda_k(P_k(B)) \geq \lambda_k(P_k(B+h)), \quad \forall k = 1, 2, \dots$$

This and (A.4) implies that

$$\lambda(B) \geq \lambda(B+h)$$

which completes the proof of Lemma 3.1. ■

Proof of Lemma 3.2: Let $\alpha_{j,i}$ be the eigenvalues of S_{λ_i} ($i = 1, 2$), and

$$(A.6) \quad \alpha_{j,1} \leq \alpha_{j,2}, \quad \forall j = 1, 2, \dots$$

Similarly as in the proof of Lemma 3.1 we can assume that $\alpha_{j,i} > 0$. Then

$$(A.7) \quad \lambda_i(J(0, \epsilon)) = \lim_k A_{i,k}, \quad i = 1, 2,$$

where

$$A_{i,k} = \frac{1}{\sqrt{(2\pi)^k}} \int_{B_{i,k}} \exp\left[-\frac{1}{2} \sum_{j=1}^k y_j^2\right] d(y_1, \dots, y_k)$$

and

$$B_{i,k} = \{y \in \mathbb{R}^k : \sum_{j=1}^k \alpha_{j,i} y_j^2 \leq \epsilon^2\}.$$

Since $\alpha_{j,1} \leq \alpha_{j,2}$, $\forall j = 1, 2, \dots$, then $B_{2,k} \subset B_{1,k}$ which implies that $A_{1,k} \geq A_{2,k}$, $k = 1, 2, \dots$. This and (A.7) complete the proof of Lemma 3.2. ■

Acknowledgments

I wish to thank J.F. Traub and H. Woźniakowski for their valuable suggestions concerning this paper.

References

- [1] Anderson, T.W., "The Integral of Symmetric Unimodal Function over a Symmetric Convex Set and some Probability Inequalities." Proc. Amer. Math. Soc., Vol. 6 (1955), pp. 170-175.
- [2] Kuo, Hui-Hsuing, Gaussian Measures in Banach Spaces, Lecture Notes in Mathematics 463, Springer-Verlag, Berlin, 1975.
- [3] Parthasarathy, K.R., Probability Measures on Metric Spaces, Academic Press, New York, 1967.
- [4] Skorohad, A.V., Integration in Hilbert Space, Springer-Verlag, New York, 1979.
- [5] Traub, J.F., Wasilkowski, G.W. and Woźniakowski, H., "Average Case Optimality for Linear Problems," Dept. of Computer Science Report, Columbia University, 1981. To appear in J.T.C.S. 28 (1984).
- [6] Wasilkowski, G.W., "Local Average Error," Dept. of Computer Science Report, Columbia University, 1983.
- [7] Wasilkowski, G.W. and Woźniakowski, H., "Average Case Optimal Algorithm in Hilbert Spaces," Dept. of Computer Science Report, Columbia University, 1982.
- [8] Wasilkowski, G.W. and Woźniakowski, H., "Can Adaption Help on the Average?" 1983. To appear in Num. Math.