

## A Constrained Eigenvalue Problem\*

Walter Gander  
*Institut für Informatik,  
ETH Zentrum  
CH-8092 Zürich  
Switzerland*

Gene H. Golub  
*Department of Computer Science  
Stanford University  
Stanford, California 94305*  
and

Urs von Matt  
*Institut für Informatik,  
ETH Zentrum  
CH-8092 Zürich  
Switzerland*

Dedicated to Alan J. Hoffman on the occasion of his 65th birthday.

Submitted by Richard A. Brualdi

---

### ABSTRACT

Let  $A$  be a real symmetric matrix. In order to find the eigenvector corresponding to the smallest eigenvalue of  $A$ , we find the minimizer of the expression  $x^T A x$  subject to  $x^T x = 1$ . In many applications, however, it is necessary to introduce linear constraints:  $N^T x = t$ . In this paper we first show how to eliminate these linear constraints. Then the minimization is tackled by employing Lagrange equations. An analysis of the solvability of the problem and the sensitivity of the solution  $x$  is given. We show how the problem can be reduced to a so-called secular equation that we solve by a conventional zero-finding process. Alternatively, we present a second method which transforms the Lagrange equations into a quadratic eigenvalue problem. The two approaches are compared to each other.

---

\*This work was in part supported by the National Science Foundation under Grant NSF CCR-8412314 and by the US Army (DAAL03-87-K-0095).

## 1. INTRODUCTION

In this paper we consider the following mathematical and computational problem. Given the quantities

- $A$ :  $(n + m)$ -by- $(n + m)$  matrix, symmetric,  $n > 0$ ,  
 $N$ :  $(n + m)$ -by- $m$  matrix with full rank,  
 $\mathbf{t}$ : vector of dimension  $m$  with  $\|(N^T)^+ \mathbf{t}\| < 1$ .

Determine an  $\mathbf{x}$  such that

$$\mathbf{x}^T A \mathbf{x} = \min \quad (1)$$

subject to the constraints

$$N^T \mathbf{x} = \mathbf{t}, \quad (2)$$

$$\mathbf{x}^T \mathbf{x} = 1. \quad (3)$$

Variants of this problem occur in many applications [1, 5, 7, 8, 11]. The problem has been studied previously when  $\mathbf{t} = \mathbf{0}$ , the null vector (cf. [4, 6]).

When  $\mathbf{t} \neq \mathbf{0}$ , then the problem becomes more complicated. We now motivate our assumptions. Suppose  $N$  does not have full rank. If  $\mathbf{t}$  is not in the range of  $N^T$ , the problem has no solution. If, however, the linear constraints are consistent, we can deflate the system until we get a submatrix of full rank. In the extreme case, where  $N = \mathbf{0}$  and  $\mathbf{t} = \mathbf{0}$ , the problem reduces to the ordinary eigenvalue problem with an eigenvector corresponding to the smallest eigenvalue of  $A$  as the solution.

Now consider the quantity  $\|(N^T)^+ \mathbf{t}\|$ . As  $(N^T)^+ \mathbf{t}$  denotes the unique solution of  $N^T \mathbf{x} = \mathbf{t}$  of minimal norm, we can make the following distinctions: when  $\|(N^T)^+ \mathbf{t}\| > 1$  there is no solution. In the case of  $\|(N^T)^+ \mathbf{t}\| = 1$  the unique solution is given by  $\mathbf{x} = (N^T)^+ \mathbf{t}$ . Therefore, the condition  $\|(N^T)^+ \mathbf{t}\| < 1$  is the only interesting case.

Thus when  $\|(N^T)^+ \mathbf{t}\| < 1$ , we can always find an  $\mathbf{x}$  that satisfies both constraints, and hence there always exists at least one solution to the problem.

## 2. PROBLEM SIMPLIFICATION

In this section we will normalize the problem in order to study its solvability. These considerations are not only of theoretical interest, but they also serve as a basis for the numerical calculations.

### 2.1. Elimination of the Linear Constraint

In the domain of the transformation  $N^T$  we can distinguish two fundamental subspaces:

$$\begin{aligned}\mathcal{N}(N^T): & \quad \text{nullity of } N^T, \\ \mathcal{N}(N^T)^\perp: & \quad \text{orthogonal complement.}\end{aligned}$$

If the transformation  $N^T$  is restricted on  $\mathcal{N}(N^T)^\perp$ , it acts as a bijection between  $\mathcal{N}(N^T)^\perp$  and  $\mathcal{R}(N^T)$  with the inverse  $(N^T)^+$ . Thus the general solution of the first constraint is given by

$$\mathbf{x} = (N^T)^+ \mathbf{t} + \boldsymbol{\xi}$$

with an arbitrary  $\boldsymbol{\xi} \in \mathcal{N}(N^T)$ . Now using the singular value decomposition of  $N^T$ ,

$$N^T = U \Sigma V^T$$

with

$$\begin{aligned}U &= \begin{bmatrix} & m \\ & U \end{bmatrix}_m \\ \Sigma &= \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_m & & \\ & & & & 0 \end{bmatrix}_{m \times n} \\ V &= \begin{bmatrix} & m & & n \\ & V_1 & & V_2 \end{bmatrix}_{m+n},\end{aligned}$$

we can write  $\mathbf{x}$  as

$$\mathbf{x} = V \Sigma^+ U^T \mathbf{t} + V_2 \mathbf{z}$$

with an arbitrary  $\mathbf{z}$ . For the subsequent minimization, we only have to consider the second constraint (3).

For practical reasons, we can use the  $QR$  decomposition of  $N$  instead of the singular value decomposition. To simplify the first constraint (2), we write

$$P^T N = \begin{bmatrix} R \\ 0 \end{bmatrix}, \quad (4)$$

where  $P$  denotes an orthogonal matrix, and  $R$  is a  $m$ -by- $m$  upper triangular

matrix. Now the problem can be formulated as

$$\mathbf{x}^T A \mathbf{x} = \mathbf{x}^T P P^T A P P^T \mathbf{x} = \min,$$

$$N^T \mathbf{x} = \begin{bmatrix} R^T & 0 \end{bmatrix} P^T \mathbf{x} = \mathbf{t},$$

$$\mathbf{x}^T \mathbf{x} = \mathbf{x}^T P P^T \mathbf{x} = 1.$$

We now make the definitions

$$P^T A P =: \begin{bmatrix} B & \Gamma^T \\ \Gamma & C \end{bmatrix}_{\substack{m \\ n}}^{\substack{m \\ n}} \quad (5)$$

$$P^T \mathbf{x} =: \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix}_{\substack{m \\ n}}. \quad (6)$$

Note that  $C^T = C$ . Now

$$\begin{aligned} \mathbf{x}^T A \mathbf{x} &= \begin{bmatrix} \mathbf{y}^T & \mathbf{z}^T \end{bmatrix} \begin{bmatrix} B & \Gamma^T \\ \Gamma & C \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} \\ &= \mathbf{y}^T B \mathbf{y} + \mathbf{y}^T \Gamma^T \mathbf{z} + \mathbf{z}^T \Gamma \mathbf{y} + \mathbf{z}^T C \mathbf{z} \\ &= \mathbf{y}^T B \mathbf{y} + 2\mathbf{z}^T \Gamma \mathbf{y} + \mathbf{z}^T C \mathbf{z}, \\ N^T \mathbf{x} &= \begin{bmatrix} R^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = R^T \mathbf{y} = \mathbf{t}, \\ \mathbf{x}^T \mathbf{x} &= \begin{bmatrix} \mathbf{y}^T & \mathbf{z}^T \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = \mathbf{y}^T \mathbf{y} + \mathbf{z}^T \mathbf{z} = 1. \end{aligned}$$

So we have reduced the constraint (2) to an ordinary linear system with an upper triangular matrix. Now

$$\mathbf{y} = R^{-T} \mathbf{t} \quad (7)$$

and with the help of the definitions

$$s^2 := 1 - \mathbf{y}^T \mathbf{y} > 0, \quad (8)$$

$$\mathbf{b} := -\Gamma \mathbf{y}, \quad (9)$$

we get the simplified problem

$$\begin{aligned} \mathbf{z}^T C \mathbf{z} - 2\mathbf{b}^T \mathbf{z} &= \min, \\ \mathbf{z}^T \mathbf{z} &= s^2. \end{aligned} \quad (10)$$

This problem has been extensively studied in the literature (cf. [2, 7, 8, 10]).

## 2.2. Stationary Points

In order to calculate the stationary points, we set up the so-called *Lagrange principal function*:

$$\Phi(\mathbf{z}, \lambda) := \mathbf{z}^T C \mathbf{z} - 2\mathbf{b}^T \mathbf{z} - \lambda(\mathbf{z}^T \mathbf{z} - s^2). \quad (11)$$

Differentiating  $\Phi$  with respect to  $\mathbf{z}$  and  $\lambda$  yields the equations

$$\begin{aligned} 2C\mathbf{z} - 2\mathbf{b} - 2\lambda\mathbf{z} &= 0, \\ \mathbf{z}^T \mathbf{z} - s^2 &= 0, \end{aligned}$$

or, normalized,

$$\begin{aligned} C\mathbf{z} &= \lambda\mathbf{z} + \mathbf{b}, \\ \mathbf{z}^T \mathbf{z} &= s^2. \end{aligned} \quad (12)$$

Now let us compare the values  $\mathbf{z}^T C \mathbf{z} - 2\mathbf{b}^T \mathbf{z}$  of different tuples  $(\lambda, \mathbf{z})$ . Following the proof given in [3, 12], a short calculation shows that the smallest  $\lambda$  is needed in order to minimize the value  $\mathbf{z}^T C \mathbf{z} - 2\mathbf{b}^T \mathbf{z}$ . So in place of the original minimization we can solve the *Lagrange equations*

$$\begin{aligned} C\mathbf{z} &= \lambda\mathbf{z} + \mathbf{b}, \\ \mathbf{z}^T \mathbf{z} &= s^2, \\ \lambda &= \min. \end{aligned} \quad (13)$$

## 3. SOLVABILITY

We will now investigate the solvability of the Lagrange equations (13). Simultaneously this analysis will point out a first method to solve the problem.

### 3.1. Explicit Secular Equation

For our discussion, we need the eigenvalue decomposition

$$C = QDQ^T, \quad (14)$$

where

$$D = \text{diag}(\delta_1, \dots, \delta_n), \quad \delta_1 \leq \delta_2 \leq \dots \leq \delta_n$$

and

$$Q^T Q = I.$$

Thus the Lagrange equations (13) are transformed as follows:

$$QDQ^T \mathbf{z} = \lambda QQ^T \mathbf{z} + \mathbf{b}$$

$$\mathbf{z}^T \mathbf{z} = \mathbf{z}^T QQ^T \mathbf{z} = s^2.$$

With the definitions

$$\mathbf{u} := Q^T \mathbf{z}, \tag{15}$$

$$\mathbf{d} := Q^T \mathbf{b}, \tag{16}$$

this can be simplified to

$$D\mathbf{u} = \lambda \mathbf{u} + \mathbf{d},$$

$$\mathbf{u}^T \mathbf{u} = s^2, \tag{17}$$

$$\lambda = \min.$$

First let us suppose  $\lambda \in \lambda(D)$ . Then there exist diagonal elements  $\delta_i$  with  $\delta_i = \lambda$ . For the ensuing discussion the following index sets turn out to be useful:

$$I := \{i \mid \delta_i = \lambda\},$$

$$\bar{I} := \{i \mid \delta_i \neq \lambda\} = \{1, \dots, n\} \setminus I.$$

If there exists a corresponding  $\mathbf{u}$  for such a  $\lambda$ , it must hold that  $\forall i: \delta_i u_i = \lambda u_i + d_i$  with  $\mathbf{u}^T \mathbf{u} = s^2$ . Then for  $i \in I$ , it must be true that  $\lambda u_i = \lambda u_i + d_i$ , and this implies  $d_i = 0$ . And for  $i \in \bar{I}$ , it must hold that

$$u_i = \frac{d_i}{\delta_i - \lambda}.$$

The normalization condition  $\mathbf{u}^T \mathbf{u} = s^2$  can only be satisfied if

$$\sum_{i \in \bar{I}} u_i^2 = \sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \lambda} \right)^2 \leq s^2.$$

As a result we have the following three possibilities:

1. There exists no solution  $\mathbf{u}$  for a given  $\lambda \in \lambda(D)$  if  $\exists i \in I: d_i \neq 0$ , or

$$\sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \lambda} \right)^2 > s^2.$$

2. There exists a unique solution  $\mathbf{u}$  for a given  $\lambda \in \lambda(D)$  if  $\forall i \in I: d_i = 0$  and

$$\sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \lambda} \right)^2 = s^2.$$

Then  $\mathbf{u}$  can be calculated as

$$u_i = \begin{cases} \frac{d_i}{\delta_i - \lambda} & i \in \bar{I} \\ 0 & i \in I. \end{cases}$$

3. There exist several solutions  $\mathbf{u}$  for a given  $\lambda \in \lambda(D)$  if  $\forall i \in I: d_i = 0$ , and

$$\sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \lambda} \right)^2 < s^2.$$

The possible values for  $\mathbf{u}$  are given by

$$u_i = \frac{d_i}{\delta_i - \lambda}, \quad i \in \bar{I},$$

$$\sum_{i \in I} u_i^2 = s^2 - \sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \lambda} \right)^2 > 0.$$

Thus the set of all solutions  $\mathbf{u}$  constitutes a manifold of dimension  $|I| - 1$  because the  $u_i$  with  $i \in I$  can be chosen arbitrarily on the given hypersphere.

The second case is given by  $\lambda \notin \lambda(D)$ . Then the inverse  $(D - \lambda I)^{-1}$  exists and  $\mathbf{u}$  has the representation

$$\mathbf{u} = (D - \lambda I)^{-1} \mathbf{d}. \quad (18)$$

For  $\mathbf{u}$  to solve the normalization condition  $\mathbf{u}^T \mathbf{u} = s^2$  it must hold that

$$\mathbf{u}^T \mathbf{u} = \mathbf{d}^T (D - \lambda I)^{-2} \mathbf{d} = \sum_{i=1}^n \left( \frac{d_i}{\delta_i - \lambda} \right)^2 = s^2.$$

We define

$$f(\lambda) := \sum_{i=1}^n \left( \frac{d_i}{\delta_i - \lambda} \right)^2 - s^2 \quad (19)$$

as the so-called *explicit secular function* (see Figure 1). Thereby the Lagrange equations (13) have a unique solution  $\mathbf{u}$  for a given  $\lambda \notin \lambda(D)$  if and only if the *explicit secular equation*

$$f(\lambda) := \sum_{i=1}^n \left( \frac{d_i}{\delta_i - \lambda} \right)^2 - s^2 = 0 \quad (20)$$

is satisfied.

If  $d_i = 0$  for all  $i$ , the secular function (20) degenerated into  $f(\lambda) \equiv -s^2 < 0$  and therefore possesses no solutions. In this case the desired  $\lambda$  lies



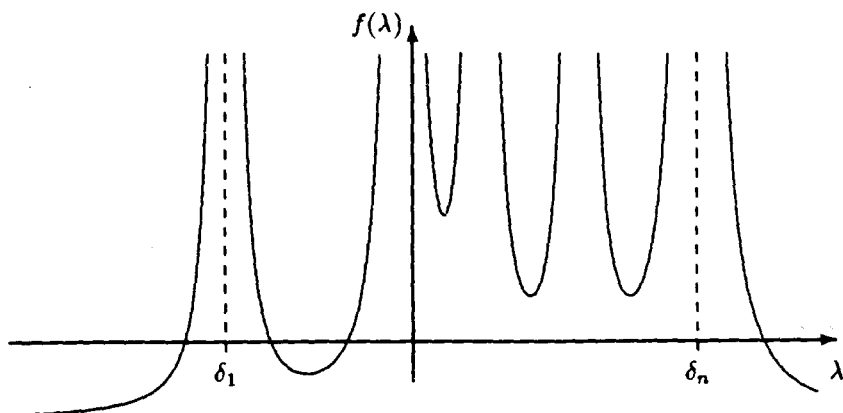


FIG. 1. Graph of the secular function.

in the spectrum  $\lambda(D)$ . Thus let  $k$  be the index of the first  $d_i \neq 0$ , i.e.  $d_k \neq 0$  and  $\forall i < k: d_i = 0$ . So we can write the secular function (19) as

$$f(\lambda) = \sum_{i=k}^n \left( \frac{d_i}{\delta_i - \lambda} \right)^2 - s^2$$

with  $\delta_k \neq 0$ . For  $\lambda$  increasing from  $-\infty$  to  $\delta_k$ ,  $f(\lambda)$  increases strictly, since the derivative

$$f'(\lambda) = \sum_{i=k}^n \frac{2d_i^2}{(\delta_i - \lambda)^3}$$

is positive for  $-\infty < \lambda < \delta_k$ . From the limits

$$\lim_{\lambda \rightarrow -\infty} f(\lambda) = -s^2,$$

$$\lim_{\lambda \rightarrow \delta_k^-} f(\lambda) = +\infty,$$

it immediately follows that for  $\lambda < \delta_k$  there exists exactly one solution.

Then the desired smallest  $\lambda$  can be located either inside or outside the spectrum  $\lambda(D)$ . With the help of the index set  $I := \{i \mid \delta_i = \delta_1\}$  we can distinguish two alternatives:

1. It holds that  $\forall i \in I: d_i = 0$ , and

$$\sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \delta_1} \right)^2 \leq s^2.$$

It follows that for  $\lambda = \delta_1$ , there exists a solution of the Lagrange equations (13).  $f(\lambda)$  possesses no more solutions for  $\lambda < \delta_1$ , since

$$f(\delta_1) = \sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \delta_1} \right)^2 - s^2 \leq 0.$$

Thus with  $\lambda = \delta_1$  we have found the smallest  $\lambda$ .

2. Or it holds that  $\exists i \in I: d_i \neq 0$ , or

$$\sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \delta_1} \right)^2 > s^2.$$

It follows that  $f(\lambda)$  has a singularity for  $\lambda = \delta_1$ , or else

$$f(\delta_1) = \sum_{i \in \bar{I}} \left( \frac{d_i}{\delta_i - \delta_1} \right)^2 - s^2 > 0.$$

Therefore  $f(\lambda)$  has exactly one solution for  $\lambda < \delta_1$ . This solution represents the desired smallest  $\lambda$ .

Hence, in both alternatives, the smallest  $\lambda$  always satisfies the condition  $\lambda \leq \delta_1$ .

### 3.2. *Implicit Secular Equation*

The above discussion on the location of the smallest  $\lambda$  that solves the Lagrange equations (13) can be carried out even without the calculation of the eigenvalue decomposition (14) of  $C$ . This is useful when we want to avoid this factorization numerically.

As indicated before, we know that the desired  $\lambda$  satisfies  $\lambda \leq \delta_1$ . So for  $\lambda = \delta_1$  the following cases can be distinguished:

1. The equation  $Cz = \delta_1 z + \mathbf{b}$  can be inconsistent, i.e.,

$$(C - \delta_1 I)(C - \delta_1 I)^+ \mathbf{b} \neq \mathbf{b}.$$

In this case we have  $\lambda < \delta_1$ , and the *implicit secular equation*

$$f(\lambda) = \mathbf{b}^T (C - \lambda I)^{-2} \mathbf{b} - s^2 = 0 \quad (21)$$

must be solved.

2. Now we assume that the equation  $Cz = \delta_1 z + \mathbf{b}$  is consistent. The expression  $(C - \delta_1 I)^+ \mathbf{b}$  represents the solution with smallest norm. If  $\|(C - \delta_1 I)^+ \mathbf{b}\| > s$ , the normalization condition cannot be satisfied, and we have to solve again the secular equation (21).

3. If however  $\|(C - \delta_1 I)^+ \mathbf{b}\| = s$ , we have found the unique solution of the Lagrange equations (13).

4. Finally it can happen that  $\|(C - \delta_1 I)^+ \mathbf{b}\| < s$ . Let  $\xi^{(1)}, \dots, \xi^{(k)}$  denote the  $k$  orthonormal eigenvectors corresponding to the eigenvalue  $\delta_1$ . Since

$$(C - \delta_1 I)^+ \mathbf{b} \perp \mathcal{N}(C - \delta_1 I),$$

$$\xi^{(i)} \in \mathcal{N}(C - \delta_1 I),$$

every vector

$$\mathbf{z} = (C - \delta_1 I)^+ \mathbf{b} + c_1 \xi^{(1)} + \dots + c_k \xi^{(k)}$$

with

$$c_1^2 + \dots + c_k^2 = s^2 - \|(C - \delta_1 I)^+ \mathbf{b}\|^2$$

solves the Lagrange equations (13). Therefore the set of solutions constitutes a manifold of dimension  $k - 1$ .

### 3.3. Condition of the Secular Equation

The calculation of the smallest zero  $\lambda$  of the secular equation (20), (21) is a delicate procedure. Even small errors  $\Delta\lambda$  can result in large deviations  $\Delta\mathbf{x}$  of the solution  $\mathbf{x}$  and  $\Delta\min$  of the minimal value  $\min$ . To illustrate the point

we will approximately determine the deviation  $\Delta \mathbf{x}$  and  $\Delta \min$  if  $\Delta \lambda$  is given. We will assume that for small  $\Delta \lambda$  the quantities  $\Delta \mathbf{x}$  and  $\Delta \min$  are essentially linearly dependent on  $\Delta \lambda$ .

Starting with the smallest zero  $\lambda$  of the explicit secular equation (20),  $\mathbf{x}$  and  $\min$  are computed as follows:  $\mathbf{u} := (D - \lambda I)^{-1} \mathbf{d}$ ,  $\mathbf{z} := Q\mathbf{u}$ ,

$$\mathbf{x} := P \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix},$$

and

$$\min := \mathbf{x}^T A \mathbf{x} = \begin{bmatrix} \mathbf{y}^T & \mathbf{z}^T \end{bmatrix} \begin{bmatrix} B & \Gamma^T \\ \Gamma & C \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix}.$$

If instead of the theoretical zero  $\lambda$ , the value  $\lambda + \Delta \lambda$  is used, the result is as follows:

$$\begin{aligned} \mathbf{u} + \Delta \mathbf{u} &= [D - (\lambda + \Delta \lambda)I]^{-1} \mathbf{d} \\ &= (D - \lambda I)^{-1} \mathbf{d} + (D - \lambda I)^{-2} \mathbf{d} \Delta \lambda \\ &\quad + O(\Delta \lambda^2). \end{aligned}$$

From this it follows that  $\Delta \mathbf{u} \approx (D - \lambda I)^{-2} \mathbf{d} \Delta \lambda$ . Now, we have  $\mathbf{z} + \Delta \mathbf{z} = Q(\mathbf{u} + \Delta \mathbf{u})$ , so  $\Delta \mathbf{z} = Q \Delta \mathbf{u}$ , and

$$\mathbf{x} + \Delta \mathbf{x} = P \begin{bmatrix} \mathbf{y} \\ \mathbf{z} + \Delta \mathbf{z} \end{bmatrix}$$

implies

$$\Delta \mathbf{x} = P \begin{bmatrix} \mathbf{0} \\ \Delta \mathbf{z} \end{bmatrix}.$$

Finally,

$$\begin{aligned} \min + \Delta \min &= (\mathbf{x} + \Delta \mathbf{x})^T A (\mathbf{x} + \Delta \mathbf{x}) \\ &= \mathbf{x}^T A \mathbf{x} + 2\mathbf{x}^T A \Delta \mathbf{x} + \Delta \mathbf{x}^T A \Delta \mathbf{x}, \end{aligned}$$

so that

$$\Delta \min \approx 2\mathbf{x}^T A \Delta \mathbf{x}.$$

Now considering that

$$\begin{aligned}
 2\mathbf{x}^T A \Delta \mathbf{x} &= 2\mathbf{x}^T P P^T A P P^T \Delta \mathbf{x} \\
 &= 2 \begin{bmatrix} \mathbf{y}^T & \mathbf{z}^T \end{bmatrix} \begin{bmatrix} B & \Gamma^T \\ \Gamma & C \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \Delta \mathbf{z} \end{bmatrix} \\
 &= 2\mathbf{y}^T \Gamma^T \Delta \mathbf{z} + 2\mathbf{z}^T C \Delta \mathbf{z} \\
 &= 2(\mathbf{z}^T C - \mathbf{b}^T) \Delta \mathbf{z},
 \end{aligned}$$

we get for the deviations  $\Delta \mathbf{x}$  and  $\Delta \min$  the values

$$\begin{aligned}
 \Delta \mathbf{x} &= P \begin{bmatrix} \mathbf{0} \\ Q(D - \lambda I)^{-2} \mathbf{d} \end{bmatrix} \Delta \lambda \\
 &\equiv \kappa(\mathbf{x}) \Delta \lambda
 \end{aligned} \tag{22}$$

and

$$\begin{aligned}
 \Delta \min &= 2(\mathbf{z}^T C - \mathbf{b}^T) Q(D - \lambda I)^{-2} \mathbf{d} \Delta \lambda \\
 &\equiv \kappa(\min) \Delta \lambda.
 \end{aligned} \tag{23}$$

Here the quantities

$$\kappa(\mathbf{x}) := P \begin{bmatrix} \mathbf{0} \\ Q(D - \lambda I)^{-2} \mathbf{d} \end{bmatrix}, \tag{24}$$

$$\kappa(\min) := 2(\mathbf{z}^T C - \mathbf{b}^T) Q(D - \lambda I)^{-2} \mathbf{d} \tag{25}$$

are the *condition vector* of  $\mathbf{x}$  and the *condition number* of  $\min$ . In the actual computation  $\kappa(\mathbf{x})$  and  $\kappa(\min)$  can be calculated as well, and one can get an estimate of the numerical error.

The value of  $\|\kappa(\mathbf{x})\|$  is bounded as follows:

$$\begin{aligned}
 \|\kappa(\mathbf{x})\| &= \left\| P \begin{bmatrix} \mathbf{0} \\ Q(D - \lambda I)^{-2} \mathbf{d} \end{bmatrix} \right\| \\
 &\leq \|(D - \lambda I)^{-2}\| \|\mathbf{d}\| = \frac{\|\mathbf{d}\|}{(\delta_1 - \lambda)^2}.
 \end{aligned} \tag{26}$$

If  $\mathbf{d}$  happens to be an eigenvector to the eigenvalue  $\delta_1$  of  $D$ , then we have an equality. Similarly, we can bound the condition number of  $\min$ :

$$\begin{aligned} |\kappa(\min)| &= \left| 2\mathbf{x}^T A P \begin{bmatrix} \mathbf{0} \\ Q(D - \lambda I)^{-2} \mathbf{d} \end{bmatrix} \right| \\ &< 2\|A\| \|(D - \lambda I)^{-2}\| \|\mathbf{d}\| \\ &= 2\|A\| \frac{\|\mathbf{d}\|}{(\delta_1 - \lambda)^2}. \end{aligned} \quad (27)$$

However, we have a strict inequality now. We can prove this by contradiction. Without loss of generality, let us assume  $P = I$  and  $Q = I$ . That is, we could apply the transformations (4) and (14) beforehand and start right away with a matrix  $A$ , where the trailing  $n$ -by- $n$  submatrix is diagonal, and with an upper triangular matrix  $N$ . Suppose that we are given a problem where the bound on  $|\kappa(\min)|$  is attained. This essentially implies that

$$|\mathbf{x}^T A \Delta \mathbf{x}| = \|\mathbf{x}\| \|A\| \|\Delta \mathbf{x}\|.$$

This equality can only be satisfied if  $\Delta \mathbf{x}$  is an eigenvector to the largest eigenvalue (in absolute value) of  $A$ . Furthermore,  $\mathbf{x}$  and  $\Delta \mathbf{x}$  must be parallel, i.e.  $\mathbf{x}$  is a multiple of  $\Delta \mathbf{x}$ . As the first  $m$  elements of  $\Delta \mathbf{x}$  are zero, so are the corresponding elements of  $\mathbf{x}$ . Equation (6) implies  $\mathbf{y} = \mathbf{0}$ , and from (9) it follows that  $\mathbf{b} = \mathbf{0}$ . But from (16) we have  $\mathbf{d} = \mathbf{0}$ , so the desired smallest  $\lambda$  of the Lagrange equations (13) lies in the spectrum  $\lambda(C)$ , and we would not be solving the secular equation at all.

Thus it is obvious that we have to face large errors if the smallest zero  $\lambda$  of the secular equation (20), (21) is near the smallest eigenvalues  $\delta_1$  of  $C$ . Since the norm of the matrix  $A$  is normally bigger than 1, an inaccurately determined zero affects the minimal value  $\min$  more than the solution vector  $\mathbf{x}$ .

#### 4. ZERO FINDER

Now, we want to calculate the smallest zero of the secular equation (20) with  $d_k \neq 0$ . We will solve it by using an iterative method. Suppose we know

an approximation  $\lambda^{(i)}$ . Then we can approximate  $f(\lambda)$  with the replacement function

$$g(\lambda) = \frac{a}{(b - \lambda)^2} - s^2 \quad (28)$$

in such a way that

$$g(\lambda^{(i)}) = f(\lambda^{(i)}),$$

$$g'(\lambda^{(i)}) = f'(\lambda^{(i)}).$$

The zero of  $g(\lambda)$  will determine the next approximation  $\lambda^{(i+1)}$ .

A short calculation yields the values

$$a = 4 \frac{(f(\lambda^{(i)}) + s^2)^3}{f'^2(\lambda^{(i)})},$$

$$b = \lambda^{(i)} + 2 \frac{f(\lambda^{(i)}) + s^2}{f'(\lambda^{(i)})},$$

and for the zero  $\lambda^{(i+1)} = b - \sqrt{a}/s$  of the replacement function  $g(\lambda)$  we get

$$\lambda^{(i+1)} = \lambda^{(i)} - 2 \frac{f(\lambda^{(i)}) + s^2}{f'(\lambda^{(i)})} \left( \frac{\sqrt{f(\lambda^{(i)}) + s^2}}{s} - 1 \right).$$

It can be shown that this iteration process will yield a strictly decreasing sequence of approximations  $\lambda^{(i)}$ . The reader is referred to [9, 12].

#### 4.1. Initial Value

Now, in order to start the iteration we need to construct an initial value. For a first guess the reduced secular equation

$$\left( \frac{d_k}{\delta_k - \lambda} \right)^2 - s^2 = 0$$

is useful. This leads to the initial value

$$\lambda^{(0)} = \delta_k - \frac{|d_k|}{s}.$$

For this  $\lambda^{(0)}$  it is obvious that  $\lambda^{(0)} < \delta_k$  with  $f(\lambda^{(0)}) \geq 0$ .

#### 4.2. Stopping Criterion

In theory, the iteration process yields a strictly decreasing sequence of approximations  $\lambda^{(i)}$ , but this property does not persist in finite arithmetic. It is therefore reasonable to terminate the iteration when the strict monotonicity is lost, namely if  $\lambda^{(i+1)} \geq \lambda^{(i)}$ . This method has the advantage that it is machine-independent and that it does not need any knowledge of machine accuracy.

#### 4.3. Implicit Secular Equation

If we do not want to compute the eigenvalue decomposition (14) of  $C$ , we have to consider the evaluation of the *implicit secular function*

$$f(\lambda) = \mathbf{b}^T(C - \lambda I)^{-2}\mathbf{b} - s^2 \quad (29)$$

and its derivative

$$f'(\lambda) = 2\mathbf{b}^T(C - \lambda I)^{-3}\mathbf{b}.$$

With the definitions

$$\mathbf{u} := (C - \lambda I)^{-1}\mathbf{b},$$

$$\mathbf{u}' := (C - \lambda I)^{-1}\mathbf{u},$$

these values can be expressed as

$$f(\lambda) = \mathbf{u}^T\mathbf{u} - s^2,$$

$$f'(\lambda) = 2\mathbf{u}^T\mathbf{u}'.$$

Therefore each iteration step requires the solution of two linear systems with the matrix  $C - \lambda I$ .



Again with explicit secular equation, the quantity

$$\lambda^{(0)} = \delta_k - \frac{|d_k|}{s}$$

yields an initial value with  $f(\lambda^{(0)}) \geq 0$  and  $\lambda^{(0)} < \delta_k$ . If the column vector  $\mathbf{q}_k$  of  $Q$  were known, we could compute the value  $d_k$ , since (16) implies  $d_k = \mathbf{q}_k^T \mathbf{b}$ . If, however, we have an eigenvector  $\xi$  with  $\|\xi\| = 1$  corresponding to the smallest eigenvalue  $\delta_1$  of  $C$ , it is straightforward to use the quantity

$$\lambda^{(0)} := \delta_1 - \frac{|\xi^T \mathbf{b}|}{s}$$

as an initial value. Then it can be shown [12] that  $f(\lambda^{(0)}) \geq 0$  and  $\lambda^{(0)} < \delta_1$  holds when  $\xi^T \mathbf{b} \neq 0$ .

## 5. QUADRATIC EIGENVALUE PROBLEM

The two previously mentioned methods have the property that they reduce the problem to finding the solution of a one-dimensional secular equation. Considering the problem from another point of view, the Lagrange equations (13) can be reduced to a quadratic eigenvalue problem. For the derivation let us assume  $\lambda \notin \lambda(C)$ . In this case  $\mathbf{z}$  can be written as

$$\mathbf{z} = (C - \lambda I)^{-1} \mathbf{b}.$$

Taking into account the normalization condition for  $\mathbf{z}$ , we get the secular function

$$f(\lambda) = \mathbf{b}^T (C - \lambda I)^{-2} \mathbf{b} - s^2,$$

of which the zeros are to be computed. The task looks different if we make the definition

$$\gamma := (C - \lambda I)^{-2} \mathbf{b},$$

so that

$$(C - \lambda I)^2 \gamma = \mathbf{b}.$$

Instead of the secular equations, we have to solve the system

$$\mathbf{b}^T \gamma - s^2 = 0, \quad (30)$$

$$(C - \lambda I)^2 \gamma = \mathbf{b}. \quad (31)$$

The first condition (30) can also be formulated as

$$1 = \frac{1}{s^2} \mathbf{b}^T \gamma.$$

Using this factor 1 as a coefficient of  $\mathbf{b}$  in (31), we get the *quadratic eigenvalue problem*

$$(C - \lambda I)^2 \gamma = \frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma. \quad (32)$$

Note that the restriction  $\lambda \notin \lambda(C)$  is no longer necessary. Of course, we must face the fact that the set of solutions for  $\lambda$  has been extended by these manipulations, for two equations cannot be formulated as a single one without consequences. Subsequently we will compare the solutions of the quadratic eigenvalue problem (32) with those of the Lagrange equations (13).

### 5.1. Solvability

We show the following. Assume  $\lambda$  and  $\mathbf{z}$  fulfill the Lagrange equations

$$C\mathbf{z} = \lambda\mathbf{z} + \mathbf{b},$$

$$\mathbf{z}^T \mathbf{z} = s^2.$$

**THEOREM 5.1.** *The quadratic eigenvalue problem*

$$(C - \lambda I)^2 \gamma = \frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma$$

*has a solution for this  $\lambda$ .*

*Proof.* In our proof, we have to distinguish whether  $\lambda$  lies in the spectrum  $\lambda(C)$  or not.

*Case 1:*  $\lambda \in \lambda(C)$ . Let  $\gamma$  be an eigenvector of  $C$  with the eigenvalue  $\lambda$ . Then  $\mathbf{b} = (C - \lambda I)\mathbf{z}$  implies  $\mathbf{b}^T \gamma = \mathbf{z}^T (C - \lambda I) \gamma = 0$ . From this it follows that

$$(C - \lambda I)^2 \gamma = \mathbf{0},$$

$$\frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma = \mathbf{0},$$

and hence  $\gamma$  satisfies the eigenvalue equation.

*Case 2:*  $\lambda \notin \lambda(C)$ . With the definition

$$\gamma := (C - \lambda I)^{-1} \mathbf{z}$$

it follows that

$$(C - \lambda I)^2 \gamma = (C - \lambda I) \mathbf{z} = \mathbf{b},$$

$$\frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma = \frac{1}{s^2} \mathbf{b} \mathbf{b}^T (C - \lambda I)^{-1} \mathbf{z}$$

$$= \frac{1}{s^2} \mathbf{b} \mathbf{z}^T \mathbf{z} = \mathbf{b},$$

and again  $\gamma$  satisfies the eigenvalue equation.

Therefore, we can construct a solution  $\gamma$  of the quadratic eigenvalue problem (32) in both cases. ■

Conversely, we can assume that  $\lambda$  and  $\gamma$  fulfill the quadratic eigenvalue equation

$$(C - \lambda I)^2 \gamma = \frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma, \quad \lambda \notin \lambda(C).$$

**THEOREM 5.2.**  $\lambda$  and  $\mathbf{z} := (C - \lambda I)^{-1} \mathbf{b}$  fulfill the Lagrange equations

$$C\mathbf{z} = \lambda \mathbf{z} + \mathbf{b},$$

$$\mathbf{z}^T \mathbf{z} = s^2.$$

*Proof.* Obviously the first equation is satisfied. Multiplying the quadratic eigenvalue equation (32) by  $(C - \lambda I)^{-2}$ , we get

$$\gamma = \frac{1}{s^2} (\mathbf{b}^T \gamma) (C - \lambda I)^{-2} \mathbf{b} \neq 0.$$

This implies  $\mathbf{b}^T \gamma \neq 0$  and

$$(C - \lambda I)^{-2} \mathbf{b} = \frac{s^2}{\mathbf{b}^T \gamma} \gamma.$$

The square of the norm of  $\mathbf{z}$  becomes

$$\mathbf{z}^T \mathbf{z} = \mathbf{b}^T (C - \lambda I)^{-2} \mathbf{b} = \mathbf{b}^T \frac{s^2}{\mathbf{b}^T \gamma} \gamma = s^2 \frac{\mathbf{b}^T \gamma}{\mathbf{b}^T \gamma} = s^2,$$

with which the second equation is satisfied too. ■

Finally, we assume that  $\lambda$  and  $\gamma$  fulfill the quadratic eigenvalue equation

$$(C - \lambda I)^2 \gamma = \frac{1}{s^2} \mathbf{b} \mathbf{b}^T \gamma, \quad \lambda \in \lambda(C).$$

We define  $\mathbf{u} := (C - \lambda I)^+ \mathbf{b}$ .

**THEOREM 5.3.** *For the solvability of the Lagrange equations*

$$C\mathbf{z} = \lambda \mathbf{z} + \mathbf{b},$$

$$\mathbf{z}^T \mathbf{z} = s^2$$

*we can make the following distinction: If the Lagrange equations are inconsistent, i.e.  $(C - \lambda I)\mathbf{u} \neq \mathbf{b}$ , or if  $\mathbf{u}^T \mathbf{u} > s^2$ , then there is no solution for*

this  $\lambda$ . On the other hand, if the equations are consistent, we have a unique solution  $\mathbf{z} = \mathbf{u}$  for  $\mathbf{u}^T \mathbf{u} = s^2$ , and several solutions for  $\mathbf{u}^T \mathbf{u} < s^2$ .

*Proof.* If  $(C - \lambda I)\mathbf{u} = (C - \lambda I)(C - \lambda I)^+ \mathbf{b} \neq \mathbf{b}$ , then  $\mathbf{b} \notin \mathcal{R}(C - \lambda I)$ , and the Lagrange equations have no solution for this  $\lambda$ . Let  $\mathbf{b} \in \mathcal{R}(C - \lambda I)$  in the following. Then  $(C - \lambda I)\mathbf{u} = (C - \lambda I)(C - \lambda I)^+ \mathbf{b} = \mathbf{b}$ , and the first Lagrange equation is satisfied.

The vector  $\mathbf{u}$  denotes the solution with smallest norm of the equation  $C\mathbf{u} = \lambda\mathbf{u} + \mathbf{b}$ . Therefore, if  $\|\mathbf{u}\| > s$ , the normalization constraint  $\mathbf{z}^T \mathbf{z} = s^2$  cannot be satisfied, and we have no solution of the Lagrange equations for this  $\lambda$ . In the case of  $\|\mathbf{u}\| = s$ , the vector  $\mathbf{u}$  denotes the unique solution.

Finally, let us assume that  $\|\mathbf{u}\| < s$ . The vectors  $\xi^{(1)}, \dots, \xi^{(k)}$  are chosen as  $k$  orthonormal eigenvectors corresponding to the eigenvalue  $\lambda$ . As

$$\mathbf{u} \perp \mathcal{N}(C - \lambda I),$$

$$\xi^{(i)} \in \mathcal{N}(C - \lambda I),$$

all the vectors

$$\mathbf{z} = \mathbf{u} + c_1 \xi^{(1)} + \dots + c_k \xi^{(k)}$$

with

$$c_1^2 + \dots + c_k^2 = s^2 - \|\mathbf{u}\|^2$$

are solutions of the Lagrange equations. Therefore, the set of solutions constitutes a manifold of dimension  $k - 1$ . ■

Discussing the solvability of the Lagrange equations (13), we derived the result that the solution with the smallest  $\lambda$  must satisfy  $\lambda \leq \delta_1$ , where  $\delta_1$  denotes the smallest eigenvalue of  $C$ . Let  $\lambda$  be hereafter the smallest eigenvalue of the quadratic eigenvalue equation (32). With the aforementioned theorems we can make the following distinction:

1. It can hold that  $\lambda < \delta_1$ . Then  $\lambda$  lies outside the spectrum  $\lambda(C)$ , and  $\mathbf{z} := (C - \lambda I)^{-1} \mathbf{b}$  fulfills the Lagrange equations (13). The solution is unique.
2. Now let  $\lambda = \delta_1$ . The vector  $\mathbf{u} := (C - \lambda I)^+ \mathbf{b}$  must satisfy the equation  $C\mathbf{u} = \lambda\mathbf{u} + \mathbf{b}$ . If  $\mathbf{u}^T \mathbf{u} = s^2$ , then  $\mathbf{z} := \mathbf{u}$  is the unique solution of the Lagrange equations (13).

3. Finally, if  $\mathbf{u}^T \mathbf{u} < s^2$ , we must find an eigenvector  $\xi$  to the eigenvalue  $\lambda$  of  $C$  with  $\xi^T \xi = s^2 - \mathbf{u}^T \mathbf{u}$ . Then,  $\mathbf{z} := \mathbf{u} + \xi$  represents one of the many solutions of the Lagrange equations (13).

### 5.2. Solving the Quadratic Eigenvalue Problem

The quadratic eigenvalue problem (32) can be reduced to an ordinary eigenvalue problem by properly chosen transformations. With the definition

$$\boldsymbol{\eta} := (C - \lambda I)\boldsymbol{\gamma} \quad (33)$$

the following equations can be established:

$$C\boldsymbol{\gamma} - \boldsymbol{\eta} = \lambda\boldsymbol{\gamma},$$

$$C\boldsymbol{\eta} - \frac{1}{s^2}\mathbf{b}\mathbf{b}^T\boldsymbol{\gamma} = \lambda\boldsymbol{\eta}.$$

In matrix terms this leads to

$$\begin{bmatrix} C & -I \\ -\frac{1}{s^2}\mathbf{b}\mathbf{b}^T & C \end{bmatrix} \begin{bmatrix} \boldsymbol{\gamma} \\ \boldsymbol{\eta} \end{bmatrix} = \lambda \begin{bmatrix} \boldsymbol{\gamma} \\ \boldsymbol{\eta} \end{bmatrix} \quad (34)$$

Thus we have transformed the original quadratic eigenvalue problem into an equivalent linear one that can be solved with traditional methods.

## 6. NUMERICAL RESULTS

Finally we present some results of the numerical experiments. All calculations were carried out on a VAX 8600 with floating point accelerator (single precision, 32-bit reals) under VMS 4.6. The IMSL-Library Edition 10.0/Version 1.0 served for the basic computations. The input data were produced by a random number generator. As a reference solution we employed an

TABLE 1  
 $n = 10, \|\kappa(\mathbf{x})\| = 0.104, \kappa(\min) = -0.282$

Algorithm	$\lambda$	min	Norm	Error	Time (sec)
Reference solution	-4.220885	-0.677439	1.000000		
Explicit secular equation	-4.220886	-0.677439	1.000000	$2.0 \times 10^{-7}$	0.05
Implicit secular equation	-4.220886	-0.677439	1.000000	$6.5 \times 10^{-8}$	0.10
Quadratic eigenvalue problem	-4.220890	-0.677438	1.000000	$5.4 \times 10^{-7}$	0.23

implementation in double precision. Tables 1 to 3 summarize the results. The columns have the following meaning:

$\lambda$ :	Smallest zero of the secular equation
min:	Minimum
Norm:	Norm of the solution $\mathbf{x}$
Error:	Norm of the difference $\mathbf{x} - \mathbf{x}_{\text{reference solution}}$
Time:	CPU time

Furthermore the caption contains the dimension  $n$  of the problem together with the condition numbers (24), (25) of the solutions.

TABLE 2  
 $n = 45, \|\kappa(\mathbf{x})\| = 930.8, \kappa(\min) = -13131$

Algorithm	$\lambda$	min	Norm	Error	Time (sec)
Reference solution	-7.054650	-7.054494	1.000000		
Explicit secular equation	-7.054657	-7.052025	0.999825	$1.8 \times 10^{-4}$	1.16
Implicit secular equation	-7.054650	-7.060933	1.000456	$4.6 \times 10^{-4}$	2.60
Quadratic eigenvalue problem	-7.055462	-2.290894	0.569900	$4.3 \times 10^{-1}$	10.4

TABLE 3  
 $n = 100$ ,  $\|\kappa(x)\| = 6.019$ ,  $\kappa(\min) = -125.2$

Algorithm	$\lambda$	min	Norm	Error	Time (sec)
Reference solution	-11.48939	-11.07225	1.000000		
Explicit secular equation	-11.48939	-11.07225	1.000000	$5.3 \times 10^{-6}$	9.90
Implicit secular equation	-11.48939	-11.07230	1.000002	$2.2 \times 10^{-6}$	24.5
Quadratic eigenvalue problem	-11.48976	-11.02595	0.997983	$2.2 \times 10^{-3}$	104.0

Based on the numerical calculations, the three solutions can be judged as follows:

*Explicit secular equation.* The zero of the secular equation (20) is determined to machine precision. The accuracy that can be expected from the condition numbers (24), (25) is achieved.

*Implicit secular equation.* This method achieves the same accuracy as the first one. However, the calculation of an eigenvalue decomposition (14) is replaced by the determination of a generalized inverse, which is in no way cheaper than the former operation.

*Quadratic eigenvalue problem.* It turns out that the smallest eigenvalue of the general matrix (34) can be calculated only very inexactly. We suppose that the transformation into a quadratic eigenvalue problem (32) impairs the condition of the problem. With large condition numbers (24), (25) all decimal places can be incorrect.

*We wish to thank Professor A. M. Lesk, who stimulated the research described in this paper through a personal communication.*

## REFERENCES

- 1 N. R. Draper, "Ridge analysis" of response surfaces, *Technometrics* 5:469-479 (1963).
- 2 G. E. Forsythe and G. H. Golub, On the stationary values of a second-degree polynomial on the unit sphere, *SIAM J. Appl. Math.* 13:1050-1068 (1963).
- 3 W. Gander, Least squares with a quadratic constraint, *Numer. Math.* 36:291-307 (1981).
- 4 G. H. Golub, Some modified matrix eigenvalue problems, *SIAM Rev.* 15:318-334 (1973).



- 5 G. H. Golub, M. Heath, and G. Wahba, Generalized cross-validation as a method for choosing a good ridge parameter, *Technometrics* 21:215–223 (1979).
- 6 G. H. Golub and C. F. Van Loan, *Matrix Computations*, Johns Hopkins U.P., Baltimore, 1983.
- 7 J. J. Moré, The Levenberg-Marquardt algorithm: Implementation and theory, in *Proceedings of the Biennial Conference Held at Dundee* (A. Dold and B. Eckmann, Eds.), Springer-Verlag, 1978, pp. 105–116.
- 8 J. J. Moré and D. C. Sorensen, Computing a trust region step, *SIAM J. Sci. Statist. Comput.* 4:553–572 (1983).
- 9 Chr. H. Reinsch, Smoothing by spline functions, II, *Numer. Math.* 16:451–454 (1971).
- 10 E. Spjøtvoll, A note on a theorem of Forsythe and Golub, *SIAM J. Appl. Math.* 23:307–311 (1972).
- 11 E. Spjøtvoll, Multiple comparison of regression functions, *Ann. Math. Statist.* 43:1076–1088 (1972).
- 12 U. von Matt, A Constrained Eigenvalue Problem, Diploma Thesis, Abteilung für Informatik, ETH Zürich, 1988.

*Received 30 September 1988; final manuscript accepted 6 October 1988*