

ON THE LINEAR LEAST SQUARES PROBLEM  
WITH A QUADRATIC CONSTRAINT

by

Walter Gander

STAN-CS-78-697  
November 1978

COMPUTER SCIENCE DEPARTMENT  
School of Humanities and Sciences  
STANFORD UNIVERSITY



Reprint May, 2019

On the Linear Least Squares Problem  
with a Quadratic Constraint

by

Walter Gander<sup>1</sup>

---

<sup>1</sup>Neu-Technikum Buchs, CH-9470 Buchs, Switzerland.  
Research supported in part by National Science Foundation Grant  
No. MCS78-17697.

# On the Linear Least Squares Problem

with a Quadratic Constraint

by

Walter Gander<sup>2</sup>

## Abstract

In this paper we present the theory and practical computational aspects of the linear least squares problem with a quadratic constraint. New theorems characterizing properties of the solutions are given and extended for the problem minimizing a general quadratic function subject to a quadratic constraint. For two important regularization methods we formulate dual equations which proved to be very useful for the application of smoothing of datas. The resulting algorithm is a numerically stable version of an algorithm proposed by Rutishauser. We show also how to choose a third order iteration method to solve the secular equation. However we are still far away from a foolproof machine independent algorithm.

---

<sup>2</sup>Computer Science Department, Stanford University, Stanford, Calif. 94305  
On leave from Neu-Technikum Buchs, Switzerland.  
The author has been supported by the Swiss National Science Foundation.  
Grant No. 5'521'330'61517.

To Maurice

## Acknowledgment

This work was done during my Stanford year 1977/78. I am greatly indebted to Prof. G. Golub for his valuable comments and suggestions, for his hospitality and for creating the extraordinary spirit of Serra House that attracts in one year most numerical analysts of the world. My stay would not have been possible without the support of Prof. P. Henrici one of my “godfathers” for the Swiss NSF. I would like to thank him for his continuous interest in my work and for his encouragement. I owe him much of my mathematical education. Also in Switzerland I have to thank Prof. J. Marti, my second “godfather” , who happened to work at a similar problem. I am indebted to Prof. Ch. van Loan, Cornell University, for his careful reading of the manuscript and for his various suggestions helping to improve my English. Thanks also to the “Forschungskommission” of ETHZ and the Swiss National Science Foundation who gave me the grant that enabled my stay at Stanford. Finally I have to thank my wife Heidi and my parents-in-law, A. and E. Wolf, for their non-mathematical assistance.

Stanford, Fall 1978  
Walter Gander

# Contents

<b>1</b>	<b>Basic Definitions and Remarks</b>	<b>9</b>
<b>2</b>	<b>The Least Squares Problem with a Quadratic Constraint</b>	<b>9</b>
2.1	Characterization of the Solution . . . . .	10
2.2	The Solutions of the Normal Equations . . . . .	13
2.3	The Solution for the Equality Constraint . . . . .	17
2.4	The Solution for the Inequality Constraint . . . . .	17
<b>3</b>	<b>The Relaxed Least Squares Problem</b>	<b>19</b>
3.1	Results from General Theory . . . . .	19
3.2	The Dual Normal Equations . . . . .	20
3.3	Eldén's Transformation . . . . .	22
3.4	Rutishauser's Relaxed and Doubly Relaxed Least Squares Problem . . . . .	23
<b>4</b>	<b>Minimum Norm Solution with Given Norm of the Residual</b>	<b>26</b>
4.1	The Dual Normal Equations . . . . .	26
4.2	Representation as Least Squares Problem . . . . .	27
<b>5</b>	<b>Computational Aspects</b>	<b>27</b>
5.1	Solution of a Relaxed Least Squares Problem with Band Matrix . . . . .	28
5.2	Solution of a Least Squares Problem with Two Band Matrices . . . . .	31
5.3	Bidiagonalization . . . . .	33
5.4	Computation of the Derivatives of the Length Function . . . . .	36
<b>6</b>	<b>One-point Iteration Methods to solve the Secular Equation</b>	<b>40</b>
6.1	Convergence Factors . . . . .	41
6.2	Third Order Iterative Methods . . . . .	43
6.3	The Convergence Factor for a Third Order Method . . . . .	46
6.4	Geometrical Interpretation of the Convergence Factors . . . . .	46
6.5	Solving the Secular Equation . . . . .	50
<b>7</b>	<b>Generalizations</b>	<b>53</b>
<b>8</b>	<b>Smoothing of Datas</b>	<b>55</b>

## 0. Introduction

In this paper we consider the linear least squares problem with a quadratic constraint. The matrices and vectors will be real and we use capital letters  $A, B, \dots$  for matrices and small letters  $\mathbf{a}, \mathbf{b}, \dots$  for vectors.  $\| \cdot \|$  will be used for the Euclidean vector-norm.

Given  $A, C, \mathbf{b}, \mathbf{d}$  and a number  $\alpha \geq 0$  we consider the problem to find  $\mathbf{x}$  such that

$$\text{subject to } \left. \begin{aligned} \|A\mathbf{x} - \mathbf{b}\| &= \min \\ \|C\mathbf{x} - \mathbf{d}\| &\leq \alpha. \end{aligned} \right\} \quad (1)$$

This problem is a generalization of the linear least squares problem with equality constraints

$$\text{subject to } \left. \begin{aligned} \|A\mathbf{x} - \mathbf{b}\| &= \min \\ C\mathbf{x} &= \mathbf{d} \end{aligned} \right\} \quad (2)$$

since for  $\alpha = 0$  every solution of (1) is also a solution of (2).

The motivation why to consider (1) rather than (2) is explained best by the following example. Let's assume we are given  $m$  values of a function  $f$

$$y_i = f(t_i), \quad i = 1, \dots, m$$

and we seek the coefficients of a polynomial

$$p(t) = \sum_{i=0}^{n-1} x_i t^i$$

that approximates  $f$ .

If we insist that  $p$  interpolate the given data, then we have to solve the system ( $m$  equations and  $n$  unknowns)

$$A\mathbf{x} = \mathbf{y} \quad (3)$$

with  $a_{ij} = t_i^j$ . If  $m > n$ , (3) may have no solution. If  $m < n$ , the solution is not unique. If  $m = n$ , there is a unique solution if  $t_i \neq t_j, i \neq j$ . However it is well known that  $A$  is ill conditioned which means that  $\mathbf{x}$  is difficult to compute accurately and that usually the norm of  $\mathbf{x}$  will be large. The polynomial  $p$  with large coefficients will be useless since cancellation will affect its evaluation. It may therefore be better not to interpolate but to approximate the data in the least squares sense. This leads to the unconstrained least squares problem

$$\|A\mathbf{x} - \mathbf{y}\| = \min. \quad (4)$$

This problem has for  $m \geq n$  and if  $A$  has full rank a unique solution. If  $A$  is rank deficient (4) has infinitely many solutions and therefore one usually looks for the solution with minimal norm

$$\text{subject to } \left. \begin{aligned} \min \|\mathbf{x}\| \\ \|A\mathbf{x} - \mathbf{y}\| &= \min. \end{aligned} \right\} \quad (5)$$

The solution of (5) is given explicitly by

$$\mathbf{x} = A^+ \mathbf{y}$$

where  $A^+$  is the pseudoinverse of  $A$ . However the solution of (4) [35] and (5) may also suffer having a large norm. The problem is then said to be ill posed and we consider two ways to regularize [42] the solution. We can prescribe a bound for  $\|\mathbf{x}\|$  and thus look for the solution of

$$\text{subject to } \left. \begin{aligned} \|A\mathbf{x} - \mathbf{y}\| &= \min \\ \|\mathbf{x}\| &\leq \alpha. \end{aligned} \right\} \quad (6)$$

Alternatively we may prefer that the deviation from the given data points be bounded. This leads to

$$\text{subject to } \left. \begin{aligned} \|\mathbf{x}\| &= \min \\ \|A\mathbf{x} - \mathbf{y}\| &\leq \alpha. \end{aligned} \right\} \quad (7)$$

The question how to choose  $\alpha$  is not simple to answer. It may be necessary to use some statistical tools to estimate  $\alpha$  [17].

We observe that (6) and (7) are least squares problems with a quadratic constraint. The degree of the polynomial ( $n - 1$ ) may be chosen independently of the number of given points. We may wish to compute a low degree polynomial ( $m \gg n$ ) or a polynomial with large degree  $m < n$  but with small coefficients. Problems (6) and (7) have a unique solution (if it exists) and we will show how to compute it efficiently.

Returning to our example, let us consider a partitioning of the data points in two sets. We may ask for a polynomial that interpolates the data of the first set exactly. This leads to a least squares problem with equality constraints:

$$\text{subject to } \left. \begin{array}{l} \|C\mathbf{x} - \mathbf{y}_2\| = \min \\ B\mathbf{x} = \mathbf{y}_1. \end{array} \right\} \quad (8)$$

where  $\mathbf{y}_1$  contains the function values of the first set. Instead of interpolating on the first set we could ask for a bound of the deviation and get again problem (1)

$$\text{subject to } \left. \begin{array}{l} \|C\mathbf{x} - \mathbf{y}_2\| = \min \\ \|B\mathbf{x} - \mathbf{y}_1\| \leq \alpha. \end{array} \right\} \quad (9)$$

In contrast to (6) and (7), (9) may not have a unique solution. As we shall see the solution is unique if and only if the nullspaces of  $C$  and  $B$  intersect trivially.



# 1 Basic Definitions and Remarks

The following notation will be used

Problem (P1):

$$\|A\mathbf{x} - \mathbf{b}\| = \min \tag{10}$$

subject to

$$\|C\mathbf{x} - \mathbf{d}\| \leq \alpha. \tag{11}$$

If we have  $\|C\mathbf{x} - \mathbf{d}\| = \alpha$  instead of (11) we will refer to the problem as (P1E). Two special cases of (P1) will be considered:

Problem (P2): like (P1) but  $C = I$ ,  $\mathbf{d} = 0$ .

Problem (P3): like (P1) but  $A = I$ ,  $\mathbf{b} = 0$ .

Finally (P2E) and (P3E) will be the corresponding problems with equality sign in the constraint.

The solution of (P1) is a stationary point of the Lagrange function (with the Lagrange multiplier  $\lambda$ )

$$L(\mathbf{x}, \lambda) = \|A\mathbf{x} - \mathbf{b}\|^2 + \lambda (\|C\mathbf{x} - \mathbf{d}\|^2 - \alpha^2)$$

and therefore a solution of  $\frac{\partial L}{\partial \mathbf{x}} = 0$  and  $\frac{\partial L}{\partial \lambda} = 0$ , which are the “normal equations”:

$$(A^\top A + \lambda C^\top C)\mathbf{x} = A^\top \mathbf{b} + \lambda C^\top \mathbf{d} \tag{12}$$

$$\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2. \tag{13}$$

If the matrix  $A^\top A + \lambda C^\top C$  is nonsingular, then we can define

$$f(\lambda) := \|C\mathbf{x}(\lambda) - \mathbf{d}\|^2 \tag{14}$$

where  $\mathbf{x}(\lambda)$  is the solution of (12). We will call  $f$  the “length function”. To determine a solution we have to solve the “secular equation”

$$f(\lambda) = \alpha^2. \tag{15}$$

Finally we observe that for  $\lambda > 0$  equations (12) and (13) are the normal equations of the least squares problem

$$\left\| \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x} - \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix} \right\| = \min.$$

A useful tool for the analysis of problem (P2) and (P3) is the singular value decomposition (SVD) [14] and its generalization (BSVD) [45] for (P1). These decompositions can also be used for the practical computations. However for the problems (P1E), (P2E) and (P3E) there are less expensive ways [8]. In some applications [33],  $A$  and  $C$  are band matrices and (P1) may be solved efficiently without transformations as we shall show.

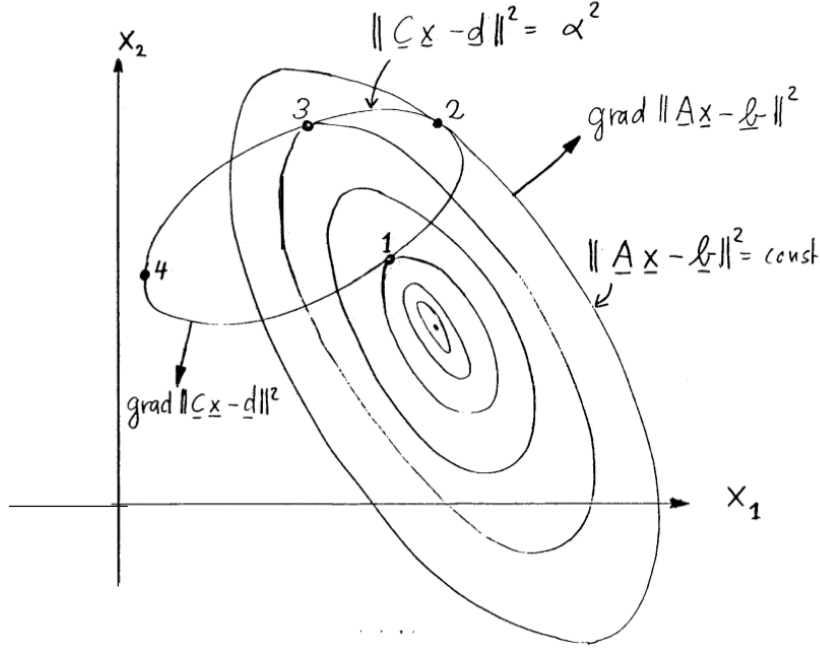
# 2 The Least Squares Problem with a Quadratic Constraint

Let  $A$  be an  $(m \times n)$  matrix,  $C$  a  $(p \times n)$  matrix,  $\mathbf{b}$  an  $m$  vector,  $\mathbf{d}$  a  $p$  vector, and  $\alpha$  a positive number.

We consider the problem to find an  $n$ -vector  $\mathbf{x}$  so that

$$\text{subject to } \left. \begin{aligned} \|A\mathbf{x} - \mathbf{b}\| &= \min \\ \|C\mathbf{x} - \mathbf{d}\| &= \alpha. \end{aligned} \right\} \tag{P1E}$$

For  $n = 2$  we can interpret this problem geometrically. The level lines of  $\|A\mathbf{x} - \mathbf{b}\|^2 = \text{const}$  are ellipses centered at  $A^+\mathbf{b}$ . The constraints  $\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2$  is also an ellipse.



We are looking for a point  $\mathbf{x}$  on the ellipse  $\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2$  which has the smallest value of  $\|A\mathbf{x} - \mathbf{b}\|^2$ . Clearly it is the point 1. At that point the gradients of the two functions

$$\|C\mathbf{x} - \mathbf{d}\|^2, \quad \|A\mathbf{x} - \mathbf{b}\|^2$$

are parallel, i.e., some  $\lambda$  exist so that

$$\text{grad } \|A\mathbf{x} - \mathbf{b}\|^2 = -\lambda \text{grad } \|C\mathbf{x} - \mathbf{d}\|^2.$$

If we rearrange this equation we have

$$(A^T A + \lambda C^T C)\mathbf{x} = A^T \mathbf{b} + \lambda C^T \mathbf{d}$$

which is equation (12) that we obtained using the technique of the Lagrange multipliers.

We observe that there are three other points for which the gradients of the two functions are parallel (2,3,4). However for these three points  $\lambda$  is negative since the gradients have the same directions.

We can think of problem (P1E) as blowing up a ballon centered at  $A^+ \mathbf{b}$  which has an ellipsoid shape until it touches the ellipsoid  $\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2$ .

If  $A^+ \mathbf{b}$  is outside of  $\|C\mathbf{x} - \mathbf{d}\|^2$  then at the point of touch the gradients will have opposite sign, i.e.  $\lambda > 0$ . However if  $A^+ \mathbf{b}$  is inside the gradients have the same sign and  $\lambda < 0$ .

For the problem with inequality constraints (P1), (P2) and (P3) we have only to consider the case that  $A^+ \mathbf{b}$  is outside of  $\|C\mathbf{x} - \mathbf{d}\|^2$ . If it is inside then  $\mathbf{x} = A^+ \mathbf{b}$  solves the problem.

## 2.1 Characterization of the Solution

The solution of (P1E) is among the solutions  $(\mathbf{x}, \lambda)$  of the normal equations (see Section 1):

$$\begin{aligned} (A^T A + \lambda C^T C)\mathbf{x} &= A^T \mathbf{b} + \lambda C^T \mathbf{d} \\ \|C\mathbf{x} - \mathbf{d}\|^2 &= \alpha^2. \end{aligned} \tag{16}$$

The following theorem compares two solutions of these equations:

**Theorem 1.** *If  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  are solutions of the normal equations (16), then*

$$\|A\mathbf{x}_2 - \mathbf{b}\|^2 - \|A\mathbf{x}_1 - \mathbf{b}\|^2 = \frac{\lambda_1 - \lambda_2}{2} \|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2. \tag{17}$$

*Proof.* Since  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  are solutions of (16) we have

$$A^\top A \mathbf{x}_1 - A^\top \mathbf{b} = -\lambda_1 C^\top C \mathbf{x}_1 + \lambda_1 C^\top \mathbf{d} \quad (18)$$

$$A^\top A \mathbf{x}_2 - A^\top \mathbf{b} = -\lambda_2 C^\top C \mathbf{x}_2 + \lambda_2 C^\top \mathbf{d} \quad (19)$$

$\mathbf{x}_2^\top$  (19) -  $\mathbf{x}_1^\top$  (18) gives

$$\|A \mathbf{x}_2\|^2 - \|A \mathbf{x}_1\|^2 - \mathbf{b}^\top A(\mathbf{x}_2 - \mathbf{x}_1) = \lambda_1(\|C \mathbf{x}_1\|^2 - \mathbf{d}^\top C \mathbf{x}_1) - \lambda_2(\|C \mathbf{x}_2\|^2 - \mathbf{d}^\top C \mathbf{x}_2) \quad (20)$$

$\mathbf{x}_2^\top$  (18) -  $\mathbf{x}_1^\top$  (19) gives

$$-\mathbf{b}^\top A(\mathbf{x}_2 - \mathbf{x}_1) = \lambda_1(-\mathbf{x}_2^\top C^\top C \mathbf{x}_1 + \mathbf{d}^\top C \mathbf{x}_2) - \lambda_2(-\mathbf{x}_1^\top C^\top C \mathbf{x}_2 + \mathbf{d}^\top C \mathbf{x}_1). \quad (21)$$

Observe that

$$\|A \mathbf{x}_2 - \mathbf{b}\|^2 - \|A \mathbf{x}_1 - \mathbf{b}\|^2 = \|A \mathbf{x}_2\|^2 - \|A \mathbf{x}_1\|^2 - 2 \mathbf{b}^\top A(\mathbf{x}_2 - \mathbf{x}_1).$$

So that if we add (20) - (21) we get

$$\begin{aligned} \|A \mathbf{x}_2 - \mathbf{b}\|^2 - \|A \mathbf{x}_1 - \mathbf{b}\|^2 &= \lambda_1 (\|C \mathbf{x}_1\|^2 - \mathbf{d}^\top C \mathbf{x}_1) - \mathbf{x}_1^\top C^\top C \mathbf{x}_2 + \mathbf{d}^\top C \mathbf{x}_2 \\ &\quad - \lambda_2 (\|C \mathbf{x}_2\|^2 - \mathbf{d}^\top C \mathbf{x}_2) - \mathbf{x}_1^\top C^\top C \mathbf{x}_2 + \mathbf{d}^\top C \mathbf{x}_1. \end{aligned} \quad (22)$$

Now we have

$$\begin{aligned} \|C \mathbf{x}_1 - \mathbf{d}\|^2 &= \|C \mathbf{x}_2 - \mathbf{d}\|^2 = \alpha^2 \\ \implies \|C \mathbf{x}_1\|^2 - 2 \mathbf{d}^\top C \mathbf{x}_1 + \|\mathbf{d}\|^2 &= \|C \mathbf{x}_2\|^2 - 2 \mathbf{d}^\top C \mathbf{x}_2 + \|\mathbf{d}\|^2 \\ \implies \|C \mathbf{x}_1\|^2 - \mathbf{d}^\top C \mathbf{x}_1 + \mathbf{d}^\top C \mathbf{x}_2 &= \|C \mathbf{x}_2\|^2 - \mathbf{d}^\top C \mathbf{x}_2 + \mathbf{d}^\top C \mathbf{x}_1 \end{aligned} \quad (23)$$

From (23) we conclude that the factors of  $\lambda_1$  and  $\lambda_2$  in (22) are the same. Therefore they also equal their arithmetic mean which is

$$\frac{1}{2} (\|C \mathbf{x}_1\|^2 - 2 \mathbf{x}_1^\top C^\top C \mathbf{x}_2 + \|C \mathbf{x}_2\|^2) = \frac{1}{2} \|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2.$$

□

**Corollary 1.** *The solution of (PIE) is the solution  $\mathbf{x}(\lambda)$  of the normal equations (16) with the largest  $\lambda$ .*

*Proof.* From (17) we have that if  $\lambda_1 > \lambda_2$ , then  $\|A \mathbf{x}_2 - \mathbf{b}\|^2 > \|A \mathbf{x}_1 - \mathbf{b}\|^2$ . □

The next theorem gives a result very similar to Theorem 1.

**Theorem 2.** *Assume  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  are solutions of the normal equations (16). Assume that  $|\lambda_1| + |\lambda_2| \neq 0$ . Then*

$$\|A \mathbf{x}_2 - \mathbf{b}\|^2 - \|A \mathbf{x}_1 - \mathbf{b}\|^2 = \frac{\lambda_2 - \lambda_1}{\lambda_2 + \lambda_1} \|A(\mathbf{x}_1 - \mathbf{x}_2)\|^2.$$

*Proof.* From  $(A^\top A + \lambda C^\top C) \mathbf{x} = A^\top \mathbf{b} + \lambda C^\top \mathbf{d}$  we have

$$\lambda_1 C^\top C \mathbf{x}_1 - \lambda_1 C^\top \mathbf{d} = -A^\top A \mathbf{x}_1 + A^\top \mathbf{b} \quad (24)$$

$$\lambda_2 C^\top C \mathbf{x}_2 - \lambda_2 C^\top \mathbf{d} = -A^\top A \mathbf{x}_2 + A^\top \mathbf{b} \quad (25)$$

$\lambda_1 \mathbf{x}_1^\top$  (25) -  $\lambda_2 \mathbf{x}_2^\top$  (24) gives

$$\lambda_1 \lambda_2 ((\mathbf{x}_2 - \mathbf{x}_1)^\top C^\top \mathbf{d}) = (\lambda_2 - \lambda_1) \mathbf{x}_1^\top A^\top A \mathbf{x}_2 + (\lambda_1 \mathbf{x}_1 - \lambda_2 \mathbf{x}_2)^\top A^\top \mathbf{b}. \quad (26)$$

$\lambda_1 \mathbf{x}_2^\top$  (25) -  $\lambda_2 \mathbf{x}_1^\top$  (24) gives

$$\begin{aligned} & \lambda_1 \lambda_2 (\|C\mathbf{x}_2\|^2 - \|C\mathbf{x}_1\|^2 + (\mathbf{x}_1 - \mathbf{x}_2)C^\top \mathbf{d}) \\ & = \lambda_2 \|A\mathbf{x}_1\|^2 - \lambda_1 \|A\mathbf{x}_2\|^2 + (\lambda_1 \mathbf{x}_2 - \lambda_2 \mathbf{x}_1)^\top A^\top \mathbf{b}. \end{aligned} \quad (27)$$

Observe that

$$0 = \|C\mathbf{x}_2 - \mathbf{d}\|^2 - \|C\mathbf{x}_1 - \mathbf{d}\|^2 = \|C\mathbf{x}_2\|^2 - \|C\mathbf{x}_1\|^2 + 2(\mathbf{x}_1 - \mathbf{x}_2)^\top C^\top \mathbf{d}.$$

So that if we subtract (27) - (26) we get

$$0 = \lambda_2 \|A\mathbf{x}_1\|^2 - \lambda_1 \|A\mathbf{x}_2\|^2 + (\lambda_1 \mathbf{x}_2 - \lambda_2 \mathbf{x}_1 - \lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2)^\top A^\top \mathbf{b} + (\lambda_1 - \lambda_2) \mathbf{x}_1^\top A^\top A \mathbf{x}_2,$$

or by rearranging

$$\begin{aligned} & \lambda_1 (\|A\mathbf{x}_2\|^2 - \mathbf{x}_2^\top A^\top \mathbf{b} + \mathbf{x}_1^\top A^\top \mathbf{b} - \mathbf{x}_1^\top A^\top A \mathbf{x}_2) \\ & = \lambda_2 (\|A\mathbf{x}_1\|^2 - \mathbf{x}_1^\top A^\top \mathbf{b} + \mathbf{x}_2^\top A^\top \mathbf{b} - \mathbf{x}_1^\top A^\top A \mathbf{x}_2). \end{aligned} \quad (28)$$

Now the ( ) on the left hand side of (28) is

$$\begin{aligned} & \frac{1}{2} \|A\mathbf{x}_2\|^2 + \frac{1}{2} (\|A\mathbf{x}_2\|^2 - 2\mathbf{x}_1^\top A^\top A \mathbf{x}_2 + \|A\mathbf{x}_1\|^2) \\ & - \frac{1}{2} \|A\mathbf{x}_1\|^2 + \mathbf{x}_1^\top A^\top \mathbf{b} - \mathbf{x}_2^\top A^\top \mathbf{b} + \frac{1}{2} \|\mathbf{b}\|^2 - \frac{1}{2} \|\mathbf{b}\|^2 \\ & = \frac{1}{2} (\|A\mathbf{x}_2 - \mathbf{b}\|^2 - \|A\mathbf{x}_1 - \mathbf{b}\|^2 + \|A(\mathbf{x}_2 - \mathbf{x}_1)\|^2). \end{aligned}$$

The right hand side of (28) simplifies analogously and by rearranging we obtain:

$$(\lambda_1 + \lambda_2) (\|A\mathbf{x}_2 - \mathbf{b}\|^2 - \|A\mathbf{x}_1 - \mathbf{b}\|^2) = (\lambda_2 - \lambda_1) \|A(\mathbf{x}_2 - \mathbf{x}_1)\|^2. \quad (29)$$

We now prove by contradiction that  $\lambda_1 + \lambda_2 \neq 0$ . Assume  $(\mathbf{x}_1, \lambda)$  and  $(\mathbf{x}_2, -\lambda)$  were solutions of the normal equations (16) with  $\lambda > 0$ . Then by (29) we would have

$$\|A(\mathbf{x}_2 - \mathbf{x}_1)\|^2 = 0 \implies A\mathbf{x}_2 = A\mathbf{x}_1 \implies \|A\mathbf{x}_2 - \mathbf{b}\|^2 = \|A\mathbf{x}_1 - \mathbf{b}\|^2.$$

But by Corrolary 1

$$\lambda_1 = \lambda > -\lambda = \lambda_2 \implies \|A\mathbf{x}_2 - \mathbf{b}\|^2 > \|A\mathbf{x}_1 - \mathbf{b}\|^2.$$

Therefore  $\lambda_1 + \lambda_2 \neq 0$  and we may divide in (29).  $\square$

If we combine both results from Theorem 1 and Theorem 2 we have

**Corrolary 2.** *Let  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  be two solutions of the normal equations (16). If  $|\lambda_1| + |\lambda_2| \neq 0$ , then*

$$-\frac{\lambda_1 + \lambda_2}{2} \|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2 = \|A(\mathbf{x}_1 - \mathbf{x}_2)\|^2. \quad (30)$$

From (30) we see immediately:

**Corrolary 3.** *The normal equations (16) have at most one solution  $(\mathbf{x}^*, \lambda^*)$  with  $\lambda^* > 0$ . For every other solution  $(\mathbf{x}_2, \lambda_2)$  we have  $\lambda < -\lambda^*$ .*

The next theorem gives conditions for a unique solution of (P1E).

**Theorem 3.** *The solution  $\mathbf{x}$  of (P1E) is unique (if it exists) if*

$$NS(A) \cap NS(C) = \{0\}$$

and

$$\lambda \neq -\mu_i$$

where  $(\mathbf{x}, \lambda)$  is a solution of the normal equations (16) and  $\mu_i$  is an eigenvalue of the generalized eigenvalue problem

$$\det(A^\top A - \mu C^\top C) = 0.$$

*Proof.* The proof is by contradiction. Assume  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  are solutions of the normal equations (16) which also solve (P1E).

If  $\lambda_1 \neq \lambda_2$  then we must have

$$\|A\mathbf{x}_1 - \mathbf{b}\|^2 = \|A\mathbf{x}_2 - \mathbf{b}\|^2 = \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2.$$

Theorem 1 implies

$$\|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2 = 0 \implies C(\mathbf{x}_1 - \mathbf{x}_2) = 0. \quad (31)$$

While Theorem2 gives

$$\|A(\mathbf{x}_1 - \mathbf{x}_2)\|^2 = 0 \implies A(\mathbf{x}_1 - \mathbf{x}_2) = 0. \quad (32)$$

Now if  $\mathbf{x}_1 \neq \mathbf{x}_2$  then (31) and (32) show that  $A$  and  $C$  have non-trivially intersecting nullspaces, which is a contradiction.

Therefore we must have  $\lambda_1 = \lambda_2 = \lambda$ . But in this case we then have

$$(A^\top A + \lambda C^\top C)(\mathbf{x}_1 - \mathbf{x}_2) = 0.$$

If  $\mathbf{x}_1 \neq \mathbf{x}_2$  then  $\lambda = -\mu_i$  which is also a contradiction. Therefore we must have  $\lambda_1 = \lambda_2$  and  $\mathbf{x}_1 = \mathbf{x}_2$ .  $\square$

$A$  and  $C$  have a trivial intersection of their nullspaces if and only if

$$\text{rank} \begin{pmatrix} A \\ C \end{pmatrix} = n.$$

Therefore a necessary condition for a unique solution is

$$m + p > n$$

which means that we must have “enough” equations to determine  $\mathbf{x}$ .

We shall assume in the following (till the end of the paper) that  $NS(A) \cap NS(C) = \{0\}$ . In Section2.2 we shall analyse the existence of solutions of the normal equations. We shall see that (P1E) has a solution if the constraint

$$\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2$$

is feasible, i.e., if

$$\min_{\mathbf{x}} \|C\mathbf{x} - \mathbf{d}\|^2 = \|(CC^+ - I)\mathbf{d}\|^2 < \alpha^2.$$

## 2.2 The Solutions of the Normal Equations

In this section we shall discuss the solutions of

$$(A^\top A + \lambda C^\top C)\mathbf{x} = A^\top \mathbf{b} + \lambda C^\top \mathbf{d} \quad (33)$$

$$\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2. \quad (34)$$

We note that we have assumed  $NS(A) \cap NS(C) = \{0\}$ .

If  $\lambda \neq -\mu_i$  where  $\mu_i \geq 0$  is an eigenvalue of the eigenvalue problem  $\det(A^\top A - \mu C^\top C) = 0$  then

$$\mathbf{x}(\lambda) = (A^\top A + \lambda C^\top C)^{-1} (A^\top \mathbf{b} + \lambda C^\top \mathbf{d}). \quad (35)$$

If we now choose  $\lambda$  so that the secular equation is satisfied:

$$f(\lambda) := \|C\mathbf{x}(\lambda) - \mathbf{d}\|^2 = \alpha^2 \quad (36)$$

then  $(\mathbf{x}(\lambda), \lambda)$  is a solution of (33), (34). We first discuss some properties of the length function  $f$ :

**Lemma 1.** *If  $f$  is defined by (36), (35), then*

1.  $f$  is a rational function defined on  $\mathbb{R} - \{-\mu_i \mid \det(A^\top A - \mu_i C^\top C) = 0\}$ ,

2.  $f(\lambda) \equiv \|\mathbf{d}\|^2 = \text{const}$  if  $A^\top \mathbf{b} = 0$  and  $C^\top \mathbf{d} = 0$ ,
3.  $f$  has at least one and at most  $n$  poles for some  $\lambda = -\mu_i$ ,
4.  $f(\lambda) > 0$ ,  $\lim_{\lambda \rightarrow \pm\infty} f(\lambda) = \|(CC^+ - I)\mathbf{d}\|^2$ .
5.  $f'(\lambda) < 0$  for  $0 < \lambda < \infty$ .

*Proof.* We use the generalized singular value decomposition BSVD [45] which gives the decomposition

$$\begin{aligned} U^\top AX &= D_A = \text{diag}(\alpha_1, \dots, \alpha_n), \alpha_i \geq 0, \\ V^\top CX &= D_C = \text{diag}(\gamma_1, \dots, \gamma_q), \gamma_i \geq 0, q = \min(n, p) \end{aligned} \quad (37)$$

where  $U$  ( $m \times m$ ) and  $V$  ( $p \times p$ ) are orthogonal, and  $X$  is ( $n \times n$ ) nonsingular, and  $D_A, D_C$  are diagonal matrices with  $\gamma_1 \geq \dots \geq \gamma_q$ . The decomposition exists only if  $m > n$  which we can assume since we can add in

$\|A\mathbf{x} - \mathbf{b}\|^2 = \left\| \begin{pmatrix} A \\ 0 \end{pmatrix} \mathbf{x} - \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix} \right\|^2$  zero rows in  $A$  and zero elements in  $\mathbf{b}$  without changing the solutions.

If  $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_r > \gamma_{r+1} = \dots = \gamma_q = 0$ , then the eigenvalue of  $\det(A^\top A - \mu C^\top C) = 0$  are given by

$$\mu_i = \frac{\alpha_i^2}{\gamma_i^2}, \quad i = 1, \dots, r.$$

Because we are assuming  $NS(A) \cap NS(C) = \{0\}$  we have

$$\alpha_i > 0, \quad \text{for } i = r+1, \dots, n. \quad [45]$$

By substituting (37) into (32) we have

$$X^{-\top} (D_A^\top D_A + \lambda D_C^\top D_C) X^{-1} \mathbf{x} = X^{-\top} D_A^\top U^\top \mathbf{b} + \lambda X^{-\top} D_C^\top V^\top \mathbf{d}$$

and putting

$$\mathbf{y} := X^{-\top} \mathbf{x}, \quad \mathbf{c} := U^\top \mathbf{b}, \quad \mathbf{e} := V^\top \mathbf{d}$$

we get

$$(D_A^\top D_A + \lambda D_C^\top D_C) \mathbf{y} = D_A^\top \mathbf{c} + \lambda C^\top \mathbf{e}. \quad (38)$$

Now

$$f(\lambda) = \|C\mathbf{x}(\lambda) - \mathbf{d}\|^2 = \|D_C \mathbf{y}(\lambda) - \mathbf{e}\|^2,$$

which is in components

$$\begin{aligned} f(\lambda) &= \sum_{i=1}^r \left( \gamma_i \frac{\alpha_i c_i + \lambda \gamma_i e_i}{\alpha_i^2 + \lambda \gamma_i^2} - e_i \right)^2 + \sum_{i=r+1}^p e_i^2 \\ f(\lambda) &= \sum_{i=1}^r \left( \alpha_i \frac{\gamma_i c_i - \alpha_i e_i}{\alpha_i^2 + \lambda \gamma_i^2} \right)^2 + \sum_{i=r+1}^p e_i^2. \end{aligned} \quad (39)$$

Obviously  $f(\lambda)$  is a rational function in  $\lambda$  with some poles. If  $\alpha_i(\gamma_i c_i - \alpha_i e_i) \neq 0$  for some  $1 \leq i \leq r$  then

$$\lambda = -\frac{\alpha_i^2}{\gamma_i^2} = -\mu_i$$

is a pole of  $f$ .

if  $A^\top \mathbf{b} = 0$  and  $C^\top \mathbf{d} = 0$  then from the definition of  $f$  (36),(35), we have

$$f(\lambda) = \|\mathbf{d}\|^2 = \text{const}.$$

From (36) we see that  $f(\lambda) > 0$ .

To prove the limit property in (4) we observe that the solution  $\mathbf{x}(\lambda)$  of (33) converges for  $\lambda \rightarrow \infty$  to a solution of  $C^\top C\mathbf{x} = C^\top \mathbf{d}$ . Though  $\mathbf{x}(\infty)$  may not be  $C^+ \mathbf{d}$  it has the same residual. Therefore:

$$\|C\mathbf{x}(\infty) - \mathbf{d}\|^2 = \|CC^+ \mathbf{d} - \mathbf{d}\|^2.$$

This property can of course also be seen from the representation (39).

To prove (5) we differentiate (36) and (33)

$$f'(\lambda) = 2\mathbf{x}'(\lambda)^\top C^\top (C\mathbf{x}(\lambda) - \mathbf{d}). \quad (40)$$

$$(A^\top A + \lambda C^\top C)\mathbf{x}'(\lambda) = -C^\top (C\mathbf{x}(\lambda) - \mathbf{d}). \quad (41)$$

Now since  $\lambda > 0$  we can perform the Cholesky decomposition

$$(A^\top A + \lambda C^\top C) = R_\lambda^\top R_\lambda. \quad (42)$$

Using (42) in (40) gives

$$f'(\lambda) = -2\|R_\lambda^{-\top} C^\top (C\mathbf{x}(\lambda) - \mathbf{d})\|^2 < 0. \quad (43)$$

□

Every solution  $\lambda$  of the secular equation (36) with  $\lambda \neq -\mu_i$  defines  $\mathbf{x}(\lambda)$  and  $(\mathbf{x}(\lambda), \lambda)$  is a solution of the normal equations (33),(34).

But there might be more solutions for some  $\lambda = -\mu_i$ . In this case the matrix of (33) is singular. If a solution should exist, the system must be consistent. This is especially true if  $A^\top \mathbf{b} = C^\top \mathbf{d} = 0$  then  $\mathbf{x}(\lambda)$  is an eigenvector to  $\lambda = -\mu_i$ . The next lemma describes the condition for consistency:

**Lemma 2.** *Let  $\mu_i$  be an eigenvalue of  $\det(A^\top A - \mu C^\top C) = 0$ . Then*

$$(A^\top A - \mu_i C^\top C)\mathbf{x} = A^\top \mathbf{b} - \mu_i C^\top \mathbf{d} \quad (44)$$

is a consistent system if one of the conditions holds:

(i)  $\lim_{\lambda \rightarrow -\mu_i} f(\lambda) = \lim_{\lambda \rightarrow -\mu_i} \|C\mathbf{x}(\lambda) - \mathbf{d}\|^2$  exists, i.e.,  $-\mu_i$  is not a pole of  $f$ .

(ii) Let  $J = \left\{ j \mid 1 \leq j \leq k, \frac{\alpha_j^2}{\gamma_j^2} = \mu_j \right\}$  then  $\alpha_j(c_j \gamma_j - \alpha_j e_j) = 0$  for  $j \in J$ , where all variables are defined in the proof of Lemma 1.

*Proof.* We prove first (ii) using the BSVD to transform (44).

$$\begin{aligned} A &= UD_A X^{-1}, & D_A &= \text{diag}(\alpha_1, \dots, \alpha_n) \\ C &= VD_C X^{-1}, & D_C &= \text{diag}(\gamma_1, \dots, \gamma_r, 0, \dots, 0). \end{aligned}$$

With  $\mathbf{y} := X^{-\top} \mathbf{x}$ ,  $\mathbf{c} := U^\top \mathbf{b}$ ,  $\mathbf{e} := V^\top \mathbf{d}$ ,  $(A^\top A + \lambda C^\top C)\mathbf{x} = A^\top \mathbf{b} + \lambda C^\top \mathbf{d}$  becomes

$$\left. \begin{aligned} (\alpha_k^2 + \lambda \gamma_k^2) y_k &= \alpha_k c_k + \lambda \gamma_k e_k, & k &= 1, \dots, r \\ \alpha_k^2 y_k &= \alpha_k c_k, & k &= r+1, \dots, n. \end{aligned} \right\} \quad (45)$$

Now observe that  $r = \text{rank of } C$  and because  $NS(A) \cap NS(C) = \{0\}$ ,  $\alpha_k \neq 0, k = r+1, \dots, n$ . If  $\lambda = -\mu_i = -\alpha_i^2/\gamma_i^2$ , then for every equation  $j$  of (45) with

$$\alpha_j^2 - \mu_j \gamma_j^2 = 0$$

i.e., for all  $j \in J = \left\{ j \mid 1 \leq j \leq k, \mu_i = \frac{\alpha_i^2}{\gamma_i^2} = \frac{\alpha_j^2}{\gamma_j^2} \right\}$ , the right hand side must be zero:

$$\alpha_j c_j - \frac{\alpha_i^2}{\gamma_i^2} \gamma_j e_j = \alpha_j c_j - \frac{\alpha_j^2}{\gamma_j^2} \gamma_j e_j = 0$$

$$\implies \alpha_j(\gamma_j c_j - \alpha_j e_j) = 0 \text{ for } j \in J$$

We note that the solution of (45) is given by

$$y_k = \begin{cases} (\alpha_k c_k - \mu_i \gamma_k e_k) / (\alpha_k^2 - \mu_i \gamma_k^2), & k \notin J \\ \text{arbitrary}, & k \in J \\ c_k / \alpha_k, & k = r+1, \dots, n. \end{cases} \quad (46)$$

Furthermore

$$\tilde{y}_k := \lim_{\lambda \rightarrow -\mu_i} y_k(\lambda) = \begin{cases} (\alpha_k c_k - \mu_i \gamma_k e_k) / (\alpha_k^2 - \mu_i \gamma_k^2), & k \notin J \\ e_k / \gamma_k, & k \in J \\ c_k / \alpha_k, & k = r+1, \dots, n. \end{cases} \quad (47)$$

From (47) we conclude that

$$\lim_{\lambda \rightarrow -\mu_i} \mathbf{x}(\lambda) = \lim_{\lambda \rightarrow -\mu_i} X \mathbf{y}(\lambda) = X \tilde{\mathbf{y}}$$

exists and therefore

$$\lim_{\lambda \rightarrow -\mu_i} \|C \mathbf{x}(\lambda) - \mathbf{d}\|^2$$

exists. This implies that  $f$  has no pole for  $\lambda = -\mu_i$ .

Conversely if  $f$  has not a pole for  $\lambda = -\mu_i$  then also

$$\lim_{\lambda \rightarrow -\mu_i} \mathbf{y}(\lambda)$$

must exist, which means that the system is consistent.  $\square$

We can now characterize the solutions of the normal equations precisely:

**Theorem 4.** *Let  $f$  be the length function defined by (36), (35). If  $f(\lambda) = \alpha^2$  and  $\det(A^\top A + \lambda C^\top C) \neq 0$  then there exists a unique  $\mathbf{x}(\lambda)$  so that  $(\mathbf{x}(\lambda), \lambda)$  solves the normal equations. If*

$$\det(A^\top A - \mu_i C^\top C) = 0 \quad \text{and} \quad \lim_{\lambda \rightarrow -\mu_i} f(\lambda) \leq \alpha^2 \quad (48)$$

*then there exists a  $\mathbf{x}(-\mu_i)$  so that  $(\mathbf{x}(-\mu_i), -\mu_i)$  solves the normal equations, but  $\mathbf{x}(-\mu_i)$  is only unique if  $\lim_{\lambda \rightarrow -\mu_i} f(\lambda) = \alpha^2$ .*

*Proof.* We have only to prove (48). We shall use the terminology defined on the proof of Lemma 2. For  $\lambda = -\mu_i$  the general solution  $\mathbf{y}$  of the transformed normal equations (45) is given by (46). We have to see if we can satisfy the equation

$$\|D_C \mathbf{y} - \mathbf{e}\|^2 = \alpha^2 \quad (49)$$

with this solution  $\mathbf{y}$ . In components (49) is

$$\sum_{k \notin J} (\gamma_k y_k - e_k)^2 + \sum_{k \in J} (\gamma_k y_k - e_k)^2 + \sum_{k=r+1}^n e_k^2 = \alpha^2. \quad (50)$$

We can choose the components of  $y_k$ ,  $k \in J$  arbitrarily. From (47) we see that

$$\lim_{\lambda \rightarrow -\mu_i} f(\lambda) = \lim_{\lambda \rightarrow -\mu_i} \|D_C \mathbf{y}(\lambda) - \mathbf{e}\|^2 = \|D_C \tilde{\mathbf{y}} - \mathbf{e}\|^2 \quad (51)$$

$$\lim_{\lambda \rightarrow -\mu_i} f(\lambda) = \sum_{k \notin J} (\gamma_k y_k - e_k)^2 + \sum_{k=r+1}^n e_k^2. \quad (52)$$

Therefore  $\|D_C \mathbf{y} - \mathbf{e}\|^2$  is minimized for  $\mathbf{y} = \tilde{\mathbf{y}}$ .

From (52) we see that if

$$\lim_{\lambda \rightarrow -\mu_i} f(\lambda) > \alpha^2$$



we cannot determine a solution  $\mathbf{y}(\lambda)$  that satisfies (49). If

$$\lim_{\lambda \rightarrow -\mu_i} f(\lambda) = \alpha^2$$

then the unique solution is  $\tilde{\mathbf{y}}(\lambda)$  (47). If

$$\lim_{\lambda \rightarrow -\mu_i} f(\lambda) < \alpha^2$$

then we can choose  $y_k, k \in J$  so that

$$\sum_{k \in J} (\gamma_k y_k - e_k)^2 = \alpha^2 - \lim_{\lambda \rightarrow -\mu_i} f(\lambda) \quad (53)$$

and  $\mathbf{y}(-\mu_i)$  is not unique.  $\square$

We remark that the components of  $\mathbf{y} : y_k, k \in J$  are (possibly) non-zero component of the eigenvector  $\mathbf{y}_{\mu_i}$  of

$$(D_A^\top D_A - \mu_i D_C^\top D_C) \mathbf{y}_{\mu_i} = 0.$$

Therefore a solution for the case where  $\lim_{\lambda \rightarrow -\mu_i} f(\lambda) < \alpha^2$  can also be written in the following way. Take any eigenvector  $\mathbf{y}_{\mu_i}$  belonging to  $\mu_i$  and determine  $\rho$  such that

$$\|D_C(\tilde{\mathbf{y}} + \rho \mathbf{y}_{\mu_i}) - \mathbf{d}\| = \alpha^2.$$

Then  $(\tilde{\mathbf{y}} + \rho \mathbf{y}_{\mu_i}, -\mu_i)$  is a solution of the normal equations.

### 2.3 The Solution for the Equality Constraint

We now consider problem (P1E)

$$\begin{aligned} & \|A\mathbf{x} - \mathbf{b}\|^2 = \min \\ \text{subject to} & \|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2. \end{aligned}$$

We continue to assume that the nullspaces of  $A$  and  $C$  intersect trivially. To find the solution we first have to compute the rightmost solution  $\lambda$  of the secular equation. If  $\lambda^* > 0$  then  $\mathbf{x}(\lambda^*)$  is the unique solution of (P1E). If  $\lambda^* < 0$  then we have to compute the smallest eigenvalue  $\mu_r$  of the generalized eigenvalue problem

$$\det(A^\top A - \mu C^\top C) = 0$$

and an eigenvector  $\mathbf{x}_r$ . If  $\lambda^* > \mu_r$  then again by Theorem 1  $\mathbf{x}(\lambda^*)$  is the unique solution. However if  $\lambda^* \leq \mu_r$  then a solution has the form

$$\mathbf{x}(\rho) = \lim_{\lambda \rightarrow -\mu_r} \mathbf{x}(\lambda) + \rho \mathbf{x}_r$$

where we have to determine  $\rho$  such that

$$\|C\mathbf{x}(\rho) - \mathbf{d}\|^2 = \alpha^2.$$

If  $A^\top \mathbf{b} = C^\top \mathbf{d} = 0$  then  $f(\lambda) = \|\mathbf{d}\|^2$  and  $\lambda^*$  does not exist. Then  $(\rho \mathbf{x}_r, -\mu_r)$  solves the problem if  $\|\mathbf{d}\|^2 \leq \alpha^2$  where  $\rho$  has to be chosen so that  $\|\rho C\mathbf{x}_r - \mathbf{d}\|^2 = \alpha^2$ .

### 2.4 The Solution for the Inequality Constraint

We are now ready to solve (P1):

$$\|A\mathbf{x} - \mathbf{b}\|^2 = \min \quad (54)$$

$$\text{subject to} \quad \|C\mathbf{x} - \mathbf{d}\|^2 \leq \alpha^2. \quad (55)$$

The weaker condition (55) simplifies the problem. Let  $M = \{\mathbf{x} \mid \|A\mathbf{x} - \mathbf{b}\| = \min\}$  denote the set of solutions of the unconstrained problem (54). If for some  $\mathbf{x} \in M$  we have  $\|C\mathbf{x} - \mathbf{d}\|^2 \leq \alpha^2$  then this  $\mathbf{x}$  is a solution.

If  $M \neq \{A^+\mathbf{b}\}$  then  $\mu_i = 0$  is an eigenvalue of  $\det(A^\top A - \mu C^\top C) = 0$ . Since 0 is never a pole of  $f$  ( the normal equations (44) are consistent for  $\lambda = 0$ ) we can define the number

$$\alpha_{\max}^2 := \lim_{\lambda \rightarrow 0} f(\lambda).$$

Let  $\alpha_{\min}^2 := \|(CC^+ - I)\mathbf{d}\|^2 = \lim_{\lambda \rightarrow \infty} f(\lambda)$ . Then (P1) has no solution if  $\alpha < \alpha_{\min}$ . If  $\alpha \geq \alpha_{\max}$  then

$$\mathbf{x}_0 = \lim_{\lambda \rightarrow 0} (A^\top A + \lambda C^\top C)^{-1} (A^\top \mathbf{b} + \lambda C^\top \mathbf{d})$$

is the solution. Observe that  $\mathbf{x}_0$  is in general not  $A^+\mathbf{b}$  [40]. Similarly if  $\alpha = \alpha_{\min}$  then we have to determine among the solutions  $\mathbf{x}$  that minimize  $\|C\mathbf{x} - \mathbf{d}\|^2$  the solution that minimizes  $\|A\mathbf{x} - \mathbf{b}\|^2$  which is

$$\mathbf{x}_\infty = \lim_{\mu \rightarrow 0} (\mu A^\top A + C^\top C)^{-1} (\mu A^\top \mathbf{b} + C^\top \mathbf{d}).$$

In general we will have

$$\alpha_{\min} < \alpha < \alpha_{\max}.$$

This means that we have to determine the unique positive solution  $\lambda^*$  of the secular equation

$$f(\lambda) = \alpha^2.$$

And  $\mathbf{x}(\lambda^*) = (A^\top A + \lambda^* C^\top C)^{-1} (A^\top \mathbf{b} + \lambda^* C^\top \mathbf{d})$  will solve (P1).

The extreme cases  $\alpha = \alpha_{\max}$  and  $\alpha = \alpha_{\min}$  correspond to the two problems

$$\left. \begin{array}{l} \min \|C\mathbf{x} - \mathbf{d}\| \\ \text{subject to } \|A\mathbf{x} - \mathbf{b}\| = \min \end{array} \right\} \text{ if } \alpha = \alpha_{\max}$$

$$\left. \begin{array}{l} \min \|A\mathbf{x} - \mathbf{b}\| \\ \text{subject to } \|C\mathbf{x} - \mathbf{d}\| = \min \end{array} \right\} \text{ if } \alpha = \alpha_{\min}.$$

Both are of course only nontrivial if  $A$  respectively  $C$  is rank deficient.

To compute the solution in the case  $\alpha_{\min} < \alpha < \alpha_{\max}$  we proceed iteratively. The following algorithm shows the basic idea.

(a) start with  $\lambda > 0$

(b) while not converged do

begin

solve the least squares problem

$$\begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x}_\lambda \approx \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix}$$

correct  $\lambda$  to solve  $f(\lambda) = \alpha^2$  where

$$f(\lambda) = \|C\mathbf{x}_\lambda - \mathbf{d}\|^2$$

end

At every step of the iteration we have to solve a least squares problem. If we use an orthogonalization procedure [48] it is more expensive than solving the normal equations by the Cholesky decomposition [38, 28, 33]. But it is well known [26] that it is numerically preferable to solve the least squares problem by orthogonal transformations.

We shall show in Section 5 how to reduce the amount of work using a structure of the matrices  $A$  and  $C$  or by using an orthogonal decomposition of them.

In the next two sections we shall consider the special cases  $A = I$  and  $C = I$ . It is interesting that we can formulate dual equations for these special cases. Furthermore Eldén [8] has showed how to reduce the general problem to a problem with  $C = I$ .

### 3 The Relaxed Least Squares Problem

We consider in this section the special case of a least squares problem with a quadratic constraint (P2E) where  $C = I$

$$\begin{aligned} & \|A\mathbf{x} - \mathbf{b}\| = \min \\ \text{subject to } & \|\mathbf{x}\| = \alpha \end{aligned} \quad (\text{P2E})$$

The problem with the inequality constraint  $\|\mathbf{x}\| \leq \alpha$  (P2) was called by Rutishauser [35] the relaxed least squares problem. Problem (P2E) can be interpreted to find the minimum of a quadratic form  $\|A\mathbf{x} - \mathbf{b}\|^2$  on the sphere  $\|\mathbf{x}\|^2 = \alpha^2$ , which is a special case of finding the stationary values of a quadratic form on a sphere. This problem has been analysed by Forsythe and Golub [10] in 1965. The two authors proved that if  $(\mathbf{x}_1, \lambda_1)$  and  $(\mathbf{x}_2, \lambda_2)$  are two solutions of the normal equations then

$$\lambda_1 > \lambda_2 \implies \|A\mathbf{x}_2 - \mathbf{b}\|^2 > \|A\mathbf{x}_1 - \mathbf{b}\|^2.$$

Their proof is rather complicated. Kahan [23] gave an elementary proof and stated the theorem

$$\|A\mathbf{x}_2 - \mathbf{b}\|^2 - \|A\mathbf{x}_1 - \mathbf{b}\|^2 = \frac{\lambda_1 - \lambda_2}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2 \quad (56)$$

which is Theorem 1 for  $C = I$ . Unfortunately this theorem was never published. In 1972 Spjøtvoll [39] gave a simpler proof than Forsythe-Golub but did not quite obtain Kahan's result. He showed that

$$\|A\mathbf{x}_2 - \mathbf{b}\|^2 - \|A\mathbf{x}_1 - \mathbf{b}\|^2 = (\lambda_1 - \lambda_2)(\alpha^2 - \mathbf{x}_1^\top \mathbf{x}_2)$$

but did not see that  $\alpha^2 = \frac{1}{2}(\|\mathbf{x}_1\|^2 + \|\mathbf{x}_2\|^2)$ .

#### 3.1 Results from General Theory

Considering the theory of Section 2 and putting  $C = I$  and  $\mathbf{d} = 0$  we get the normal equations

$$(A^\top A + \lambda I)\mathbf{x} = A^\top \mathbf{b} \quad (57)$$

$$\|\mathbf{x}\|^2 = \alpha^2. \quad (58)$$

Theorem 1 gives us Equation 56. This means that the solution of (P2E) is the solution  $(\mathbf{x}, \lambda)$  of 57,58 with largest  $\lambda$ . Since  $NS(A) \cap NS(C) = \{0\}$  we have from Theorem 3 that the solution of (P2E) is unique if  $\lambda \neq -\sigma_i^2$  (where  $\sigma_i$  is a singular value of  $A$ ).

The length function  $f$  can be written

$$f(\lambda) = \|(A^\top A + \lambda I)^{-1} A^\top \mathbf{b}\|^2. \quad (59)$$

Using the singular value decomposition of  $A$  [14]

$$A = U\Sigma V^\top$$

where

$$\begin{aligned} U, V & \text{ orthogonal} \\ \Sigma & = \text{diag}(\sigma_1, \dots, \sigma_p), p = \min(m, n) \end{aligned}$$

the right hand side of (59) becomes

$$f(\lambda) = \|(\Sigma^+ \Sigma + \lambda I)^{-1} \Sigma^\top \mathbf{c}\|^2 \quad \text{with } \mathbf{c} := U^\top \mathbf{b}.$$

In components this is

$$f(\lambda) = \sum_{i=1}^p \left( \frac{\sigma_i c_i}{\sigma_i^2 + \lambda} \right)^2. \quad (60)$$

If  $A^\top \mathbf{b} = 0$  then  $\mathbf{c} = 0$  and  $f(\lambda) \equiv 0$ . Theorem 4 and the fact that  $\lim_{\lambda \rightarrow \infty} f(\lambda) = 0$  implies that for every  $\alpha$  we have a solution of (P2E) which is not unique if  $\lambda = -\sigma_i^2$  and  $\lim_{\lambda \rightarrow -\sigma_i^2} f(\lambda) < \alpha^2$ .

This was already stated by Forsythe and Golub in 1965 [10].

For the problem with the inequality constraint  $\|\mathbf{x}\|^2 \leq \alpha^2$  (P2) we have either

$$\|A^+\mathbf{b}\|^2 \leq \alpha^2$$

then the solution is  $\mathbf{x} = A^+\mathbf{b}$  or if  $\|A^+\mathbf{b}\|^2 > \alpha^2$  then we have to determine the unique solution  $\lambda^*$  of the secular equation

$$f(\lambda) = \alpha^2 \quad \text{in } (0, \infty).$$

The solution is in that case  $\mathbf{x}(\lambda^*)$ . We note that we can reduce the length of  $\mathbf{x}$  arbitrarily by choosing  $\alpha$  small.

### 3.2 The Dual Normal Equations

**Theorem 5. (i)** *Let  $(\mathbf{x}, \lambda)$  be a solution of the primal normal equations*

$$(A^\top A + \lambda I)\mathbf{x} = A^\top \mathbf{b} \tag{61}$$

$$\|\mathbf{x}\|^2 = \alpha^2 \tag{62}$$

with  $\lambda \neq 0$ . Then  $(\mathbf{z}, \lambda)$  with

$$\mathbf{z} = \frac{1}{\lambda}(A\mathbf{x} - \mathbf{b}) \tag{63}$$

is a solution of the dual normal equations

$$(AA^\top + \lambda I)\mathbf{z} = -\mathbf{b} \tag{64}$$

$$\|A^\top \mathbf{z}\|^2 = \alpha^2 \tag{65}$$

**(ii)** *Let  $(\mathbf{z}, \lambda)$  be a solution of the dual normal equations (63), (64). Then  $(\mathbf{x}, \lambda)$  with*

$$\mathbf{x} = -A^\top \mathbf{z} \tag{66}$$

is a solution of the primal normal equations (61), (62).

*Proof. (i)*

$$\begin{aligned} (AA^\top + \lambda I)\frac{1}{\lambda}(A\mathbf{x} - \mathbf{b}) &= \frac{1}{\lambda}A \underbrace{(A^\top A + \lambda I)\mathbf{x}}_{A^\top \mathbf{b}} - \frac{1}{\lambda}AA^\top \mathbf{b} - \mathbf{b} \\ &= -\mathbf{b}. \end{aligned}$$

Furthermore

$$\|A^\top \mathbf{z}\|^2 = \|A^\top \frac{1}{\lambda}(A\mathbf{x} - \mathbf{b})\|^2 = \|\mathbf{x}\|^2 = \alpha^2.$$

**(ii)**

$$(A^\top A + \lambda I)(-A^\top \mathbf{z}) = -A^\top \underbrace{(AA^\top + \lambda I)\mathbf{z}}_{-\mathbf{b}}.$$

And

$$\| -A^\top \mathbf{z} \|^2 = \|\mathbf{x}\|^2 = \alpha^2.$$

□

Theorem 5 shows that we may simplify computations to solve (P2) or (P2E). Since  $A$  is an  $(m \times n)$  matrix one of the linear systems (61) or (63) will be smaller and it may be more economical to iterate with the smaller.

But Theorem 5 also gives theoretical equations.

**Corrolary 4.** *Ler  $A$  be an  $(m \times n)$  matrix and  $\lambda \neq -\sigma_i^2$  where  $\sigma_i$  is a singular value of  $A$ . Then*

$$(AA^\top + \lambda I)^{-1}A^\top = A^\top(AA^\top + \lambda I)^{-1} \quad (67)$$

$$(AA^\top + \lambda I)^{-1} = \frac{1}{\lambda} \left( I - A(A^\top A + \lambda I)^{-1}A^\top \right). \quad (68)$$

*Proof.* From Theorem 5 (67) follows by equating  $\mathbf{x}$ , and 68 by equating  $\mathbf{z}$  for the dual and primal equations.  $\square$

As an application we consider (68) for  $n = 1$  (i.e.,  $A$  is a vector) and  $\lambda = 1$ :

$$(I + \mathbf{a}\mathbf{a}^\top)^{-1} = I - \frac{1}{\mathbf{a}^\top \mathbf{a} + 1} \mathbf{a}\mathbf{a}^\top. \quad (69)$$

We observe that Corrolary 4 is a special case of the Sherman-Morrison-Woodbury formula:

Let  $A$  be  $(n \times n)$ ,  $U, V$   $(n \times p)$  matrices, then every solution  $\mathbf{x}$  of

$$(A + UV^\top)\mathbf{x} = \mathbf{b} \quad (70)$$

is also the solution of the augmented system

$$A\mathbf{x} + U\mathbf{y} = \mathbf{b} \quad (71)$$

$$V^\top \mathbf{x} - \mathbf{y} = 0. \quad (72)$$

Now assume that  $A$  is nonsingular. Then we have from (71)

$$\mathbf{x} = A^{-1}\mathbf{b} - A^{-1}U\mathbf{y}.$$

Introducing this in (72) we get

$$\begin{aligned} \mathbf{y} &= (I + V^\top A^{-1}U)^{-1}V^\top A^{-1}\mathbf{b} \\ \implies \mathbf{x} &= A^{-1}\mathbf{b} - A^{-1}U(I + V^\top A^{-1}U)^{-1}V^\top A^{-1}\mathbf{b}. \end{aligned}$$

But from (70) and (72) we have also

$$\begin{aligned} \mathbf{x} &= (A + UV^\top)^{-1}\mathbf{b} \\ \mathbf{y} &= V^\top (A + UV^\top)^{-1}\mathbf{b}. \end{aligned}$$

Equating both expressions for  $\mathbf{x}$  and  $\mathbf{y}$  we get the matrix equations

$$\begin{aligned} V^\top (A + UV^\top)^{-1} &= (I + V^\top A^{-1}U)^{-1}V^\top A^{-1} \\ (A + UV^\top)^{-1} &= A^{-1} (I - U(I + V^\top A^{-1}U)^{-1}V^\top A^{-1}). \end{aligned}$$

(Sherman-Morrison-Woodbury formulas)

We note that for  $A := \lambda I$  and  $U = V := A$  we get the equations (67) and (68).

We finally remark that if  $\lambda > 0$  the dual equations are the normal equations of the least squares problem

$$\begin{pmatrix} A^\top \\ \sqrt{\lambda}I \end{pmatrix} \mathbf{z} \approx \begin{pmatrix} 0 \\ -\frac{1}{\sqrt{\lambda}}\mathbf{b} \end{pmatrix}.$$

The dual equations have no solution if  $\lambda = 0$ . we see this clearly because of the factor  $-\frac{1}{\sqrt{\lambda}}$ , but also from (64) since it is clear that for  $\lambda = 0$  a solution  $\mathbf{z}$  of (64) may not exist for arbitrary  $\mathbf{b}$  and  $m > n$ .

### 3.3 Eldén's Transformation

Eldén [8] has shown how to transform a problem (P1) to a problem (P2) using orthogonal transformations. To illustrate his idea consider the least square representation of (P1)

$$\begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \approx \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix} \quad (73)$$

$$\|C\mathbf{x} - \mathbf{d}\| = \alpha^2 \quad (74)$$

We may assume that the rank of  $C = p \leq n$ . ( $C$  is a  $(p \times n)$  matrix.) If not then we can perform a QR decomposition with column pivoting

$$C = Q \begin{pmatrix} R \\ 0 \end{pmatrix} P^\top, \quad R = \begin{pmatrix} \text{trapezoid} \\ 0 \end{pmatrix}.$$

Observe that  $\|C\mathbf{x} - \mathbf{d}\| = \|RP^\top \mathbf{x} - Q^\top \mathbf{d}\|$ . Therefore we can replace  $C$  by  $R$ ,  $\mathbf{d}$  by  $Q^\top \mathbf{d}$ ,  $A$  by  $AP$ , and  $\mathbf{x}$  by  $\mathbf{y} = P^\top \mathbf{x}$  in (73), (74). The problem for  $\mathbf{y}$  has now a matrix  $C$  ( $p \times n$  with rank of  $C = p$ ). After it is solved we obtain  $\mathbf{x} = P\mathbf{y}$ . The same preprocessing we can apply to  $A$  and  $\mathbf{b}$  and so assume that  $\text{rank}(A) = m \leq n$ .

If  $p = n$  then we can make the change of variables

$$\mathbf{x}' := C\mathbf{x}, \quad A' := AC^{-1}$$

which transforms the problem to

$$\begin{pmatrix} A' \\ \sqrt{\lambda} I \end{pmatrix} \mathbf{x}' \approx \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix}$$

$$\|\mathbf{x}' - \mathbf{d}\|^2 = \alpha^2.$$

We assume now  $p < n$  and  $C$  has rank  $p$ . Then we make a QR decomposition of  $C^\top$ :

$$C^\top = (V_1, V_2) \begin{pmatrix} R \\ 0 \end{pmatrix} \begin{matrix} \} p \\ \} n - p \end{matrix} \quad (75)$$

with nonsingular triangular matrix  $R$ .

We change variables:

$$\mathbf{x} =: V_1 \mathbf{y}_1 + V_2 \mathbf{y}_2$$

and 73 becomes:

$$\begin{pmatrix} AV_1 & AV_2 \\ \sqrt{\lambda} R^\top & 0 \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} \approx \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} V_2^\top \mathbf{d} \end{pmatrix}. \quad (76)$$

Whereas 74 is now

$$\|R^\top V_1 \mathbf{y}_1 - \mathbf{d}\|^2 = \alpha^2. \quad (77)$$

Now we perform another QR decomposition

$$AV_2 = (Q_1, Q_2) \begin{pmatrix} U \\ 0 \end{pmatrix} \begin{matrix} \} p - n \end{matrix}. \quad (78)$$

Since we assumed that  $A$  has rank  $m$  and  $V_2$  is orthogonal,  $U$  is a nonsingular upper triangular matrix. Now (76) is the same problem as:

$$\begin{pmatrix} Q_1^\top AV_1, & U \\ Q_2^\top AV_1, & 0 \\ \sqrt{\lambda} R^\top, & 0 \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} \approx \begin{pmatrix} Q_1^\top \mathbf{b} \\ Q_2^\top \mathbf{b} \\ \sqrt{\lambda} V_2^\top \mathbf{d} \end{pmatrix}. \quad (79)$$

Now for every  $\mathbf{y}_1$  we can determine  $\mathbf{y}_2$  such that the first  $p - n$  equations are exactly satisfied. Therefore the problem splits as follows

$$\begin{pmatrix} Q_2^\top A V_1 \\ \sqrt{\lambda} R^\top \end{pmatrix} \mathbf{y}_1 \approx \begin{pmatrix} Q_2^\top \mathbf{b} \\ \sqrt{\lambda} V_2^\top \mathbf{d} \end{pmatrix} \quad (80)$$

with (77) as constraint. The final change of variables

$$\mathbf{y}'_1 := R^\top \mathbf{y}_1$$

gives the desired problem (P2).

### 3.4 Rutishauser's Relaxed and Doubly Relaxed Least Squares Problem

In [35] Rutishauser remarked that if we minimize

$$Q(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 \quad (81)$$

for an ill-conditioned matrix  $A$  (i.e.,  $\kappa = \|A\| \|A^+\| = \sigma_1/\sigma_r \gg 1$ , where  $\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0$  are the singular values of  $A$ ), that then the exact solution  $\hat{\mathbf{x}} = A^+ \mathbf{b}$  may not be satisfactory because  $\|\mathbf{x}\| \gg 1$  and evaluation of  $A\hat{\mathbf{x}}$  will be affected by cancellation.

Replacing (81) by

$$Q_\varepsilon(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \varepsilon^2 \|\mathbf{x}\|^2 \quad (82)$$

will give a shorter solution  $\mathbf{x}_\varepsilon$  which however does not minimize (81) anymore. This process of relaxing can be described as follows: The solution of (81) is the solution of the least squares problem

$$\mathbf{A}\mathbf{x} \approx \mathbf{b}. \quad (83)$$

The relaxed solution  $\mathbf{x}_\varepsilon$  of (82) is obtained by choosing some  $\varepsilon > 0$  and solving

$$\begin{pmatrix} A \\ \varepsilon I \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix}. \quad (84)$$

The connection to problem (P2) is as follows: Instead of presenting a bound for the length of  $\mathbf{x}$ :  $\|\mathbf{x}\|^2 \leq \alpha^2$  and solve the equation

$$f(\lambda) = \alpha^2$$

to determine  $\lambda > 0$  we simply set  $\lambda = \varepsilon^2 > 0$  and solve for  $\mathbf{x}$ .

Rutishauser [35] defines the double relaxed solution  $\mathbf{x}_{d_\varepsilon}$  to be the solution of

$$(A^\top A + \varepsilon^2 I + \varepsilon(A^\top A + \varepsilon^2 I)^{-1}) \mathbf{x} = A^\top \mathbf{b}. \quad (85)$$

Observe that  $\mathbf{x}_{d_\varepsilon}$  is obtained by relaxing the normal equations of the relaxed solution:

$$\begin{pmatrix} A^\top A + \varepsilon^2 I \\ \varepsilon^2 I \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} A^\top \mathbf{b} \\ 0 \end{pmatrix}. \quad (86)$$

If we put  $B := A^\top A + \varepsilon^2 I$  then the normal equations of (86) are

$$(B^\top B + \varepsilon^2 I) \mathbf{x} = B^\top A^\top \mathbf{b}. \quad (87)$$

But  $B$  is symmetric and non-singular, therefore we may multiply (87) from the left by  $B^{-1}$  and we get

$$(B + \varepsilon^2 B^{-1}) \mathbf{x} = A^\top \mathbf{b}$$

which is (85). From (85) we have

**Lemma 3.** *The double relaxed solution  $\mathbf{x}_{d_\varepsilon}$  minimizes*

$$Q_d(\mathbf{x}) = \|(A^\top A + \varepsilon^2 I) - A^\top \mathbf{b}\|^2 + \varepsilon^2 \|\mathbf{x}\|^2.$$

The aim of relaxing is to approximate the “ideal solution” [13] without explicitly computing the singular value decomposition. By the “ideal” solution we mean the following: Let

$$\begin{aligned} A &= U\Sigma V^\top, \quad U, V \text{ orthogonal} \\ \Sigma &= \text{diag}(\sigma_1, \dots, \sigma_n), \quad \sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0 \end{aligned}$$

be the singular value decomposition of  $A$ . If the data and the computation were exact then the shortest solution of

$$\|A\mathbf{x} - \mathbf{b}\|^2 = \min$$

would be  $\mathbf{x} = A^+\mathbf{b} = U\mathbf{y}$  where  $\mathbf{y}$  is defined by

$$y_i = \begin{cases} \frac{c_i}{\sigma_i} & \text{if } \sigma_i \neq 0 \quad (\mathbf{c} := U^\top \mathbf{b}) \\ 0 & \text{if } \sigma_i = 0 \end{cases}$$

In practical computation however it is not clear which  $\sigma_i$  are to be interpreted as zero [11], [26]. One is forced to make a rank decision, to prescribe a tolerance  $\tau$  and to compute

$$\tilde{y}_i = \begin{cases} \frac{c_i}{\sigma_i} & \text{if } \sigma_i \geq \tau \\ 0 & \text{if } \sigma_i < \tau \end{cases}. \quad (88)$$

We will call  $\tilde{y}_i$  the “ideal solution”.

Now we consider the function defined on  $[0, \infty)$

$$k(\sigma) := \begin{cases} 0 & \text{if } \sigma < \tau \\ \frac{1}{\sigma} & \text{if } \sigma \geq \tau \end{cases}. \quad (89)$$

The coefficient of  $\tilde{y}_i$  multiplying  $c_i$  is given by  $k(\sigma_i)$ . The length of  $\tilde{\mathbf{y}}$  depends on the choice of  $\tau$  and is

$$\|\tilde{\mathbf{y}}\|^2 = \sum_{\sigma_i \geq \tau} \left( \frac{c_i}{\sigma_i} \right)^2.$$

Let  $\gamma := \|\tilde{\mathbf{y}}\|/\|\mathbf{b}\|$ . We can ask the question: how to choose  $\varepsilon$  such that the relaxed and the doubly relaxed solution have the same or at least approximately the same length as the ideal solution. The answer is given in [13]. We have

$$\text{if } \varepsilon \geq \frac{1}{2\gamma} \implies \|\mathbf{x}_\varepsilon\| < \gamma\|\mathbf{b}\| = \|\tilde{\mathbf{y}}\|$$

and

$$\text{if } \varepsilon \geq \frac{1}{2\gamma^2} \implies \|\mathbf{x}_{d_\varepsilon}\| < \gamma\|\mathbf{b}\| = \|\tilde{\mathbf{y}}\|.$$

If we transform the relaxed problem (84) and the doubly relaxed problem (86) using the SVD we get

$$(\mathbf{y}_\varepsilon)_i = \frac{\sigma_i}{\sigma_i^2 + \varepsilon^2} c_i \quad (90)$$

$$(\mathbf{y}_{d_\varepsilon})_i = \frac{\sigma_i}{\sigma_i^2 + \varepsilon^2 + \frac{\varepsilon^2}{\sigma_i^2 + \varepsilon^2}} c_i. \quad (91)$$

We see that we can think of (90), (91) as approximations of (88). More precisely, we want to choose  $\varepsilon$  so that the two functions

$$k_\varepsilon := \frac{\sigma}{\sigma^2 + \varepsilon^2} \quad (92)$$



and

$$k_{d_\varepsilon} := \frac{\sigma}{\sigma^2 + \varepsilon^2 + \frac{\varepsilon^2}{\sigma^2 + \varepsilon^2}} \quad (93)$$

approximate the function  $k(\sigma)$  (89). Rutishauser [35] and Molinari [13] show that indeed  $k_{d_\varepsilon}$  is a better approximation than  $k_\varepsilon$ . More research could be done in this direction, e.g. relaxing the normal equations of the original problem

$$\begin{pmatrix} A^\top A \\ \varepsilon I \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} A^\top \mathbf{b} \\ 0 \end{pmatrix}$$

would yield the function

$$k_{1\varepsilon} := \frac{\sigma}{\sigma^2 + \frac{\varepsilon^2}{\sigma^2}}$$

which is somewhat between both discussed above. It is clear that those different relaxed solutions have to be computed without forming the normal equations. A program for  $\mathbf{x}_\varepsilon$  and  $\mathbf{x}_{d_\varepsilon}$  is given in [13].

**Theorem 6.** Let  $A_\varepsilon = \begin{pmatrix} A \\ \varepsilon I \end{pmatrix}$  and  $A_\varepsilon^+ = (B_\varepsilon^+, C_\varepsilon)$ . ( $B_\varepsilon^+$  is a  $(n \times n)$  matrix). Then for  $\varepsilon$  sufficiently small we have

$$B_\varepsilon^+ = A^+ + \sum_{j=1}^{\infty} (-1)^j A^+ ((A^+)^\top A^+)^{2j} \varepsilon^{2j}.$$

*Proof.* Let  $A = U\Sigma V^\top$  with  $\Sigma = \begin{pmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ & & \sigma_r & \\ 0 & & & 0 \end{pmatrix}$  be the singular value decomposition, with

$\sigma_1 \geq \dots \geq \sigma_r > 0$ . Then

$$\begin{aligned} A_\varepsilon^+ &= (A^\top A + \varepsilon^2 I)^{-1} A_\varepsilon^\top = (V\Sigma^\top \underbrace{U^\top U}_{I_m} \Sigma V^\top + \varepsilon^2 I)^{-1} (V\Sigma^\top U^\top, I_\varepsilon) \\ &= V(\Sigma^\top \Sigma + \varepsilon^2 I)^{-1} (\Sigma^\top U^\top, \varepsilon V^\top) \\ &= V \left[ (\Sigma^\top \Sigma + \varepsilon^2 I)^{-1} \Sigma^\top \right] U^\top, \varepsilon (\Sigma^\top \Sigma + \varepsilon^2 I)^{-1} V^\top \\ &= (B_\varepsilon^+, C_\varepsilon) \end{aligned}$$

$$B_\varepsilon^+ = V(\Sigma^\top \Sigma + \varepsilon^2 I)^{-1} \Sigma^\top U^\top$$

$$B_\varepsilon^+ = V \begin{pmatrix} \frac{\sigma_1}{\sigma_1^2 + \varepsilon^2} & & & 0 \\ & \ddots & & \\ & & \frac{\sigma_r}{\sigma_r^2 + \varepsilon^2} & \\ 0 & & & 0 \end{pmatrix} U^\top,$$

now

$$\frac{\sigma_i}{\sigma_i^2 + \varepsilon^2} = \frac{1}{\sigma_i} \frac{1}{1 + (\varepsilon/\sigma_i)^2} = \frac{1}{\sigma_i} \sum_{j=0}^{\infty} (\varepsilon/\sigma_i)^{2j} \quad \text{for } |\varepsilon| < \sigma_r,$$

$$B_\varepsilon^+ = \sum_{j=0}^{\infty} V \underbrace{\begin{pmatrix} 1/\sigma_1^{2j+1} & & & 0 \\ & \ddots & & \\ & & 1/\sigma_r^{2j+1} & \\ 0 & & & 0 \end{pmatrix}}_{\Sigma^+ (\Sigma^{+\top} \Sigma^+)^{2j}} U^\top \varepsilon^{2j}$$

where

$$\Sigma^+ = \begin{pmatrix} & & m \\ 1/\sigma_1 & & 0 \\ & \ddots & \\ 0 & & 1/\sigma_r \end{pmatrix} n,$$

We observe

$$V\Sigma^+(\Sigma^{+\top}\Sigma^+)^{2j}U^\top = A^+(A^{+\top}A^+)^{2j}.$$

□

**Remark.** Theorem 6 shows that we can extrapolate  $\mathbf{x}_M = A^+\mathbf{b}$  from the relaxed solution  $\mathbf{x}_\varepsilon = B_\varepsilon^+\mathbf{b}$ . Using  $\varepsilon_i = \varepsilon_0/2^i$  we can apply the Romberg extrapolation since only even powers of  $\varepsilon$  occur. It has the advantage that the rank of  $A$  must not be determined.

## 4 Minimum Norm Solution with Given Norm of the Residual

In this section we discuss the special case of a least squares problem with a quadratic constraint where  $A = I$ :

$$\text{subject to } \left. \begin{array}{l} \|\mathbf{x}\| = \min \\ \|C\mathbf{x} - \mathbf{d}\| = \alpha. \end{array} \right\} \quad (\text{P3E})$$

The problem with inequality constraint  $\|C\mathbf{x} - \mathbf{d}\| \leq \alpha$  will be denoted by (P3). Geometrically (P3E) means that we are looking for a point on the ellipsoid  $\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2$  that is nearest to the origin. From the theory in Section 2 we have that the solution of (P3E) is a solution  $(\mathbf{x}, \lambda)$  of the normal equations

$$(I + \lambda C^\top C)\mathbf{x} = \lambda C^\top \mathbf{d} \quad (94)$$

$$\|C\mathbf{x} - \mathbf{d}\|^2 = \alpha^2 \quad (95)$$

with largest  $\lambda$ . Since  $NS(I) \cap NS(C) = \{0\}$  the solution is unique if  $\lambda \neq -\gamma_i^2$  where  $\gamma_i$  is a singular value of  $C$ . The solution exists only if

$$\alpha \geq \|(CC^+ - I)\mathbf{d}\| := \alpha_{\min}.$$

For problem (P3) the solution is  $\mathbf{x} = 0$  if  $\alpha \geq \|\mathbf{d}\|$ . The length function  $f$  is in this case

$$f(\lambda) = \|((\lambda C(I + \lambda C^\top C)^{-1}C^\top - I)\mathbf{d})\|^2. \quad (96)$$

By Corrolary 4, (68), this is equal to

$$f(\lambda) = \|(I + \lambda CC^\top)^{-1}\mathbf{d}\|^2 \quad (97)$$

which gives again the connection to the dual equations.

### 4.1 The Dual Normal Equations

We transform the normal equations (94) (95) as follows: We multiply (94) ffrom left by  $C$  giving

$$(I + \lambda CC^\top)C\mathbf{x} = \lambda CC^\top \mathbf{d}. \quad (98)$$

If we subtract from (98) on both sides  $\lambda CC^\top \mathbf{d} + \mathbf{d}$  we get

$$(I + \lambda CC^\top)(C\mathbf{x} - \mathbf{d}) = -\mathbf{d}.$$

Finally introducing the new variable  $\mathbf{z} = C\mathbf{x} - \mathbf{d}$  we get the theorem:

**Theorem 7. (i)** Let  $(\mathbf{x}, \lambda)$  be a solution of the primal normal equations

$$\begin{aligned} (I + \lambda C^\top C)\mathbf{x} &= \lambda C^\top \mathbf{d} \\ \|C\mathbf{x} - \mathbf{d}\|^2 &= \alpha^2 \end{aligned} \quad (99)$$

then  $(\mathbf{z}, \lambda)$  with  $\mathbf{z} := C\mathbf{x} - \mathbf{d}$  is a solution of the dual normal equations

$$\begin{aligned} (I + \lambda C C^\top)\mathbf{z} &= -\mathbf{d} \\ \|\mathbf{z}\|^2 &= \alpha^2. \end{aligned} \quad (100)$$

(ii) Let  $(\mathbf{z}, \lambda)$  be a solution of (100) then  $(\mathbf{x}, \lambda)$  with  $\mathbf{x} = -\lambda C^\top \mathbf{z}$  is a solution of (99).

*Proof.* (i)

$$\begin{aligned} (I + \lambda C C^\top)(C\mathbf{x} - \mathbf{d}) &= C(I + \lambda C^\top C)\mathbf{x} - \mathbf{d} - \lambda C C^\top \mathbf{d} \\ &= C(\lambda C^\top \mathbf{d}) - \mathbf{d} - \lambda C C^\top \mathbf{d} = -\mathbf{d} \end{aligned}$$

and

$$\alpha^2 = \|\mathbf{z}\|^2 = \|C\mathbf{x} - \mathbf{d}\|^2.$$

(ii)

$$\begin{aligned} (I + \lambda C^\top C)(-\lambda C^\top \mathbf{z}) &= -\lambda C^\top (I + \lambda C C^\top)\mathbf{z} = \lambda C^\top \mathbf{d}. \\ \|C\mathbf{x} - \mathbf{d}\|^2 &= \|C(-\lambda C^\top \mathbf{z}) - \mathbf{d}\|^2 = \|\mathbf{z}\|^2 = \alpha^2. \end{aligned}$$

□

Equating both expressions for  $\mathbf{x}$  and  $\mathbf{z}$  again gives us the identities of Corrolary 4. In contrast to the relaxed least squares problem, the dual equations exists for every  $\lambda$ .

## 4.2 Representation as Least Squares Problem

If we want to solve (P3) then the solution is  $\mathbf{x} = 0$  if  $\alpha > \|\mathbf{d}\|$  and exists only if  $\alpha > \alpha_{\min} = \|(CC^+ - I)\mathbf{d}\|$ . For

$$\alpha_{\min} < \alpha < \|\mathbf{d}\|,$$

we have to compute the solution iteratively solving the secular equation  $f(\lambda) = \alpha^2$  where

$$f(\lambda) = \|C\mathbf{x} - \mathbf{d}\|^2 = \|\mathbf{z}\|^2.$$

$\mathbf{x}$  and  $\mathbf{z}$  are obtained solving either

$$\begin{pmatrix} I \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} 0 \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix}$$

or the dual problem

$$\begin{pmatrix} I \\ \sqrt{\lambda} C^\top \end{pmatrix} \mathbf{z} \approx \begin{pmatrix} -\mathbf{d} \\ 0 \end{pmatrix}$$

## 5 Computational Aspects

To compute the solution of a least squares problem with a quadratic constraint we have to determine iteratively the largest solution of the secular equation. For the problems (P1E), (P2E) and (P3E) we can expect numerical difficulties since at every step of the iteration we have to solve a system of equations with the matrix

$$A^\top A + \lambda C^\top C \quad \text{where } \lambda < 0 \quad (101)$$

which is in general not positive definite. A good way to solve this system is to transform it to diagonal form using BSVD of van Loan [45]. If  $A = I$  or  $C = I$  we can use SVD to diagonalize

the system stably [48]. To avoid the forming of  $A^\top A + \lambda C^\top C$  we could consider the augmented  $(m+n+p) \times (m+n+p)$  system

$$\begin{pmatrix} -I & A & 0 \\ A^\top & 0 & C^\top \\ 0 & C & \frac{1}{\lambda}I \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ 0 \\ \mathbf{d} \end{pmatrix}. \quad (102)$$

Björk has shown in [1] that for the augmented system

$$\begin{pmatrix} I & A \\ A^\top & 0 \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix} \quad (103)$$

we can obtain the condition number  $2 \cdot \kappa(A)$  with appropriate scaling (instead of  $\kappa^2(A)$  for the normal equations  $A^\top A \mathbf{x} = A^\top \mathbf{b}$ ). However in (102) the condition number depends on the value of  $\lambda$  and may be large if  $\lambda$  is near  $-\mu_i$  where  $\mu_i$  is an eigenvalue of  $\det(A^\top A - \mu C^\top C) = 0$ .

The solution is numerically much better defined for problems (P1), P(2) and (P3) where  $\lambda > 0$ . We can formulate the normal equations as least squares problems and besides BSVD- and SVD-diagonalization, there are several different ways to compute stably and cheaply the solution, especially if we can use a structure of the matrix. Since we are solving the secular equation with some iterative method we may have to compute derivatives of the length function  $f$  (see Section 5.4). If the iterative method converges fast it may not be necessary to transform the problem at all (especially if  $A$  and  $C$  are sparse). In general, however, it appears to be best [8], [12] to bidiagonalize the matrix for problems (P2) and (P3) before iterating.

## 5.1 Solution of a Relaxed Least Squares Problem with Band Matrix

We consider (P2) with  $A$  respectively (P3) with  $C$  being a band matrix. For a given  $\lambda > 0$  we have to solve the least squares problem

$$\begin{pmatrix} A \\ \sqrt{\lambda}I \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix} \quad (104)$$

respectively

$$\begin{pmatrix} I \\ \sqrt{\lambda}C \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} 0 \\ \sqrt{\lambda}\mathbf{d} \end{pmatrix}. \quad (105)$$

If we interchange the equations in (105) we have in either case a least squares problem of the form

$$\begin{pmatrix} B \\ D \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{pmatrix} \quad (106)$$

where  $D$  is a diagonal matrix ( $I$  or  $\sqrt{\lambda}I$ ) and  $B$  is a band matrix ( $A$  or  $\sqrt{\lambda}C$ ).

Using Givens rotations we shall show how an orthogonal matrix  $G$  can be found such that

$$G \begin{pmatrix} B \\ D \end{pmatrix} = \begin{pmatrix} 0 \\ F \end{pmatrix} \quad (107)$$

where  $F$  is a band upper triangular matrix. If the same rotations are applied to the right hand side the solution can be found by backsubstitution using  $F$ .

To describe the transformation we assume that  $B$  is an  $(m \times n)$  matrix with

$$b_{ij} = 0 \quad \text{if } i - j > m_1 \text{ or } j - i > m_2.$$

This means that  $B$  has  $m_1$  diagonals below the main diagonal and  $m_2$  diagonals above, that contains the possible non-zero elements.

Let

$$\begin{aligned} k_{\min}(k) &:= \max\{k - m_2, 1\} \\ k_{\max}(k) &:= \min\{k + m_1, m\}. \end{aligned}$$

Then the  $k$ -th column of  $B$  has the only (possibly) non-zero elements

$$b_{ik}, \quad i = k_{\min}(k), \dots, k_{\max}(k).$$

In  $n$ -steps we now perform the transformation (107) annihilating in the  $k$ -th step the  $k$ -th column of  $B$  and producing the  $k$ -th row of  $F$ . Let  $F$  be at the beginning the diagonal matrix  $D$  and define by

$$\text{rot}(1, k, x)$$

the multiplication of  $\begin{pmatrix} B \\ F \end{pmatrix}$  from left by a Givens matrix that changes the  $i$ -th row of  $B$  and the  $k$ -th row of  $F$  annihilating the element  $x$ . Then the transformation (107) is described by

$$\begin{array}{l} \underline{\text{for}} \ k := 1 \ \underline{\text{step}} \ 1 \ \underline{\text{until}} \ n \ \underline{\text{do}} \\ \quad \underline{\text{for}} \ i := k_{\min}(k) \ \underline{\text{step}} \ 1 \ \underline{\text{until}} \ k_{\max}(k) \ \underline{\text{do}} \\ \quad \quad \text{rot}(i, k, b_{ik}) \end{array}$$

In the  $k$ -th step of this algorithm we annihilate  $m_1 + m_2 + 1$  elements (or less at the border of  $B$ ) using  $\sim 2(m_1 + m_2 + 1)^2$  multiplications. For  $i = k_{\min} = k - m_2$  the Givens rotation changes only 2 elements:  $b_{ik}$  and  $f_{kk}$ . Then for  $i = k_{\min} + 1$ , 4 elements are affected, etc., until for  $i = k_{\max} = k + m_1$  when  $b_{k+m_1, k}$  is rotated to zero, we have to change  $m_1 + m_2 + 1$  elements. The whole transformation requires therefore  $\sim 2n(m_1 + m_2 + 1)^2$  multiplications which is comparable to Gaussian elimination of the normal equations.

Especially if  $B$  is upper bidiagonal the first step to annihilate the first and the second column is done as follows:

$$\begin{array}{ccc} \begin{bmatrix} q_1 & e_1 & & & \\ & q_2 & e_2 & & \\ & & q_3 & e_3 & \\ & & & q_4 & e_4 \\ d_1 & & & & \\ & d_2 & & & \\ & & d_3 & & \\ & & & d_4 & \end{bmatrix} & \xrightarrow{\text{rot}(1, 1, q_1)} & \begin{bmatrix} 0 & e'_1 & & & \\ & q_2 & e_2 & & \\ & & q_3 & e_3 & \\ & & & q_4 & e_4 \\ u_1 & v_1 & & & \\ & d_2 & & & \\ & & d_3 & & \\ & & & d_4 & \end{bmatrix} \\ & & \swarrow \text{rot}(1, 2, e'_1) \\ \begin{bmatrix} 0 & 0 & & & \\ & q_2 & e_2 & & \\ & & q_3 & e_3 & \\ & & & q_4 & e_4 \\ u_1 & v_1 & & & \\ & d'_2 & & & \\ & & d_3 & & \\ & & & d_4 & \end{bmatrix} & \xrightarrow{\text{rot}(2, 2, q_2)} & \begin{bmatrix} 0 & 0 & & & \\ & 0 & e_2 & & \\ & & q_3 & e_3 & \\ & & & q_4 & e_4 \\ u_1 & v_1 & & & \\ & u_2 & v_2 & & \\ & & d_3 & & \\ & & & d_4 & \end{bmatrix} \end{array}$$

The following procedure performs the transformation

$$G \begin{pmatrix} B \\ D \end{pmatrix} = \begin{pmatrix} 0 \\ F \end{pmatrix}$$

where

$$B = \begin{pmatrix} q_1 & e_1 & & & 0 \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & e_{n-1} \\ 0 & & & & q_n \end{pmatrix} \quad F = \begin{pmatrix} u_1 & v_1 & & & 0 \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & v_{n-1} \\ 0 & & & & u_n \end{pmatrix}$$

and stores the Givens rotations in two arrays  $\text{co}[i]$ ,  $\text{si}[i]$ :

```

procedure updec(n,q,e,d,u,v,co,si);
value n; integer n;
array q,e,d,u,v,co,si;
begin integer i; real t,c,s,h;
  procedure rot(a,b);
  value a,b; real a,b;
  begin real t;
    if b = 0 then begin c := 0; s := 1 end
    else
      begin t := -a/b; c := 1/sqrt(1 + t ↑ 2); s := t * c end
    end
  for i := 1 step 1 until n do u[i] := d[i];
  for i := 1 step 1 until n do
  begin
    rot(q[i], u[i]);
    co[2 * i - 1] := c; si[2 * i - 1] := s;
    u[i] := -q[i] * s + u[i] * c;
    if i < n then
      begin
        v[i] := -e[i] * s; h := e[i] * c;
        rot(h, u[i + 1]);
        co[2 * i] := c; si[2 * i] := s;
        u[i + 1] := -h * s + u[i + 1] * c
      end if;
    end i;
  end updec;

```

The following procedure performs the transformation

$$G \begin{pmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{pmatrix}$$

and solves the bidiagonal system by backsubstitution.

```

procedure solve(n,u,v,co,si,c1,c2,y);
value n; integer n;
array u,v,co,si,c1,c2,y;
begin
  integer i; real t,c,s,h;
  for i := 1 step 1 until n do
  begin
    c := co[2 * i - 1]; s := si[2 * i - 1];
    h := c1[i] * c + c2[i] * s;
    c2[i] := -c1[i] * s + c2[i] * c;
    c1[i] := h;
    if i < n then
      begin
        c := co[2 * i]; s := si[2 * i];
        h := c1[i] * c + c2[i] * s;
        c2[i + 1] := -c1[i] * s + c2[i + 1] * c;
        c1[i] := h;
      end if;
    end i;
  y[n] := c2[n]/u[n];
  for i := n - 1 step -1 until 1 do

```

```

    y[i] := (c2[i] - v[i] * y[i + 1])/u[i];
  end solve

```

Alternatively we can avoid storing the rotations observing that if

$$\begin{pmatrix} B \\ D \end{pmatrix} \mathbf{y} \approx \begin{pmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{pmatrix}$$

and if

$$G \begin{pmatrix} B \\ D \end{pmatrix} = \begin{pmatrix} 0 \\ F \end{pmatrix}, \quad G \text{ orthogonal}$$

then the normal equations are

$$F^\top F \mathbf{y} = B^\top \mathbf{c}_1 + D \mathbf{c}_2.$$

Therefore we can also find  $\mathbf{y}$  solving

$$\begin{aligned} F^\top \mathbf{w} &= B^\top \mathbf{c}_1 + D \mathbf{c}_2 \\ F \mathbf{y} &= \mathbf{w}. \end{aligned}$$

The new procedure `updec` does not contain the variables `co` and `si`. Otherwise it is the same as before. The procedure `solve` however changes much:

```

procedure solve2(n,u,v,c,y);
value n; integer n;
array u,v,c,y;
comment solves F^T F y = c;
begin integer i;
  y[1] := c[1]/u[1];
  for i := 2 step 1 until n do
    y[i] := (c[i] - v[i - 1] * y[i - 1])/u[i];
    y[n] := y[n]/u[n];
    for i := n - 1 step -1 until 1 do
      y[i] := (y[i] - v[i] * y[i + 1])/u[i];
    end solve 2

```

Before calling `solve2` the right hand side  $\mathbf{c}$  has to be defined by

$$\mathbf{c} := B^\top \mathbf{c}_1 + D \mathbf{c}_2.$$

## 5.2 Solution of a Least Squares Problem with Two Band Matrices

The algorithm given in Section 5.1 can be generalized to solve a least squares problem of the type

$$\begin{pmatrix} A \\ C \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} \mathbf{b} \\ \mathbf{d} \end{pmatrix} \quad (108)$$

where now  $A$  and  $C$  are band matrices. We assume that  $C$  has a band width which is smaller or equal to the bandwidth of  $A$ . Using Givens rotations we transform

$$G \begin{pmatrix} A \\ C \end{pmatrix} = \begin{pmatrix} 0' \\ A' \end{pmatrix} \quad (109)$$

where  $A'$  has the same bandwidth as  $A$  and  $0'$  is zero up to some elements in the right bottom corner. Using the elements of  $0'$  and  $A'$  we then compute  $\mathbf{x}$  by back substitution.

In the  $k$ -th step of the algorithm we annihilate the  $k$ -th column of  $A$  and produce the  $k$ -th row of  $A'$ . We explain the rotations for the example where  $A$  has 4 and  $C$  has 3 diagonals.

$$\left[ \begin{array}{cccccccc} x & & & & & & & \\ x & x & & & & & & \\ x & x & x & & & & & \\ x & x & x & x & & & & \\ & x & x & x & x & & & \\ & & x & x & x & x & & \\ & & & x & x & x & & \\ & & & & x & x & & \\ x & x & x & & & & & \\ & x & x & x & & & & \\ & & x & x & x & & & \\ & & & x & x & x & & \end{array} \right] \xrightarrow{1} \left[ \begin{array}{cccccccc} 0 & \oplus & \oplus & & & & & \\ 0 & x & \oplus & & & & & \\ 0 & x & x & & & & & \\ 0 & x & x & x & & & & \\ & x & x & x & x & & & \\ & & x & x & x & x & & \\ & & & x & x & x & & \\ & & & & x & x & & \\ x & x & x & \oplus & & & & \\ & x & x & x & & & & \\ & & x & x & x & & & \\ & & & x & x & x & & \end{array} \right]$$

In step 1 we annihilate the elements of the first column of  $A$  using the first row of  $C$ . This produces the new non-zero elements  $\oplus$ . In a second “cleaning” step we zero the elements  $\oplus$  in the top part using the third row:

$$\left[ \begin{array}{ccc} 0 & \oplus & \oplus \\ 0 & x & \oplus \\ 0 & x & x \end{array} \right] \rightarrow \left[ \begin{array}{ccc} 0 & \oplus & 0 \\ 0 & x & 0 \\ 0 & x & x \end{array} \right] \rightarrow \left[ \begin{array}{ccc} 0 & 0 & 0 \\ 0 & x & 0 \\ 0 & x & x \end{array} \right]$$

After these two steps the remaining matrix is

$$\left[ \begin{array}{cccc} 0 & & & \\ 0 & x & & \\ 0 & x & x & \\ 0 & x & x & x \\ & x & x & x & x \\ & & x & x & x & x \\ & & & x & x & x \\ & & & & x & x \\ x & x & x & x & & \\ & x & x & x & & \\ & & x & x & x & \\ & & & x & x & x \end{array} \right]$$

Now we can use the second row of  $C$  to zero the second column of  $A$  etc. A special treatment is needed at the border of the matrix. Using the last row of  $C$  to zero out the  $n - 3$  column gives

$$\left[ \begin{array}{cccccc} 0 & & & & & \\ 0 & 0 & & & & \\ 0 & 0 & 0 & & & \\ 0 & 0 & 0 & 0 & x & x \\ & 0 & 0 & 0 & x & x \\ & & 0 & 0 & x & x \\ & & & 0 & x & x \\ x & x & x & x & & \\ & x & x & x & & \\ & & x & x & x & \\ & & & x & x & x \end{array} \right] \xrightarrow{2} \left[ \begin{array}{cccccc} 0 & & & & & \\ 0 & 0 & & & & \\ 0 & 0 & 0 & & & \\ 0 & 0 & 0 & 0 & 0 & 0 \\ & 0 & 0 & 0 & 0 & 0 \\ & & 0 & 0 & x & 0 \\ & & & 0 & x & x \\ x & x & x & x & & \\ & x & x & x & & \\ & & x & x & x & \\ & & & x & x & x \end{array} \right] \left. \begin{array}{l} \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \\ \vphantom{\left[ \right]} \end{array} \right\} \begin{array}{l} 0' \\ \\ \\ \\ A' \end{array}$$

In the cleaning step 2 we can zero all but 3 elements in the upper part of the matrix using the last row of  $A$ . Now the solution can be computed by backsubstitution: first  $x_6$  and  $x_7$  using  $0'$  and the other unknowns using  $A'$ . In Reinsch’s “Smoothing with Spline Functions” [33] the problem

$$(Q^T D^2 Q + \lambda T) \mathbf{z} = Q^T \mathbf{y} \quad (110)$$

subject to

$$\|DQ\mathbf{z}\|^2 = S^2 \quad (111)$$



occurs, where  $\lambda > 0$ ,  $D$  is an  $(n+1) \times (n+1)$  diagonal matrix,  $Q$  an  $(n+1) \times (n-1)$  tridiagonal matrix and  $T$  a positive definite tridiagonal  $(n-1) \times (n-1)$  matrix.  $S$  is a given constant and  $\mathbf{y}$  a given  $(n+1)$  vector. Reinsch solves the normal equation using the Cholesky decomposition of the penta-diagonal coefficient matrix. Using the above described algorithm we can solve (110) as a least squares problem without sacrificing the sparseness of the matrices. As preparation we compute the Cholesky decomposition of  $T$ :

$$T = B^\top B, \quad B \text{ upper bidiagonal}$$

and (110) then becomes

$$\begin{pmatrix} DQ \\ \sqrt{\lambda} B \end{pmatrix} \mathbf{z} \approx \begin{pmatrix} D^{-1} \mathbf{y} \\ 0 \end{pmatrix}. \quad (112)$$

The matrix of (112) has the form ( $n = 5$ )

$$\begin{bmatrix} x & & & & \\ x & x & & & \\ x & x & x & & \\ & x & x & x & \\ & & x & x & \\ x & x & & & \\ & x & x & & \\ & & x & x & \\ & & & x & \end{bmatrix} \xrightarrow{1} \begin{bmatrix} 0 & \oplus & & & \\ 0 & x & & & \\ 0 & x & x & & \\ & x & x & x & \\ & & x & x & \\ x & x & \oplus & & \\ & x & x & & \\ & & x & x & \\ & & & x & \end{bmatrix}$$

Zeroing the first column in step 1 leaves an element  $\oplus$  that has to be removed in the cleaning step.

### 5.3 Bidiagonalization

If the matrices  $A$  and  $C$  are dense then for problems (P2), (P3) and (P1) (after performing Eldén's transformation, see Section 3.3) we have to solve a least squares problem of the form

$$\begin{pmatrix} A \\ \sqrt{\lambda} I \end{pmatrix} \mathbf{x} \approx \begin{pmatrix} \mathbf{b} \\ 0 \end{pmatrix} \quad (113)$$

for every new value of  $\lambda$ . If  $P$  and  $Q$  are orthogonal then an equivalent system to (113) is

$$\begin{aligned} \begin{pmatrix} P^\top & 0 \\ 0 & Q^\top \end{pmatrix} \begin{pmatrix} A \\ \sqrt{\lambda} I \end{pmatrix} Q Q^\top \mathbf{x} &\approx \begin{pmatrix} P^\top \mathbf{b} \\ 0 \end{pmatrix} \\ \implies \begin{pmatrix} P^\top A Q \\ \sqrt{\lambda} I \end{pmatrix} \mathbf{y} &\approx \begin{pmatrix} \mathbf{c} \\ 0 \end{pmatrix} \end{aligned} \quad (114)$$

with  $\mathbf{y} := Q^\top \mathbf{x}$  and  $\mathbf{c} := P^\top \mathbf{b}$ .

Now we can choose  $P$  and  $Q$  so that (114) is simpler to solve than (113). Moré [29] proposed to choose  $P$  and  $Q$  so that

$$B := P^\top A Q = \begin{pmatrix} 0 & \begin{array}{|c|} \hline \text{R} \\ \hline \end{array} \\ 0 & 0 \end{pmatrix}, \quad (115)$$

i.e., to perform the QR decomposition of  $A$  with column pivoting. However the solution of (114) in this case is still an  $n^3$  process. It has been pointed out by Moré (private communication) that in his application [29] only  $k \leq 2$  iterations are needed to solve the secular equation. Therefore it does not matter in this case if the solution of (114) is an  $n^3$  process. In general however we



```

t: for  $i := 1$  step 1 until  $n$  do
  begin
     $s := 0$ ;
    for  $k := i$  step 1 until  $m$  do  $s := s + a[k, i] \times y[k]$ ;
    for  $k := i$  step 1 until  $m$  do  $y[k] := y[k] - a[k, i] \times s$  ;
  end  $i$ ;
end

```

If we wish to compute  $\mathbf{y} = P\mathbf{x}$ , the only change we have to make in the above procedure is

```

t: for  $i := n$  step -1 until 1 do

```

The following procedure `bidia` bidiagonalizes  $A$ , i.e., computes the diagonal  $\mathbf{q}$  and the super-diagonal  $\mathbf{e}$  of  $B$  and stores the transformation vectors  $\mathbf{w}_j$  and  $\mathbf{v}_j$  in  $A$  (117):

```

procedure bidia(m,n,a,q,e)
value m,n; integer m,n;
array a,q,e;
begin
  integer i,j,k; real s, fak;
  for  $i := 1$  step 1 until  $n$  do
    begin
      comment transforms  $(a_{ii}, \dots, a_{mi})$  to  $(q_i, 0, \dots, 0)$ ;
       $s := 0$ ;
      for  $j := 1$  step 1 until  $m$  do  $s := s + a[j, i] \uparrow 2$ ;
      if  $s = 0$  then  $q[i] := 0$  else
        begin
           $s := \text{sqrt}(s)$ ;
           $q[i] := \text{if } a[i, i] > 0 \text{ then } -s \text{ else } s$ ;
           $fak := \text{sqrt}(s \times (s + \text{abs}(a[i, i])))$ ;
           $a[i, i] := a[i, i] - q[i]$ ;
          for  $k := 1$  step 1 until  $m$  do  $a[k, i] := a[k, i]/fak$ ;
          for  $j := i + 1$  step 1 until  $n$  do
            begin
               $s := 0$ ;
              for  $k := i$  step 1 until  $m$  do  $s := s + a[k, i] \times a[k, j]$ ;
              for  $k := i$  step 1 until  $m$  do  $a[k, j] := a[k, j] - a[k, i] \times s$ ;
            end  $j$ ;
            end  $s$ ;
          comment transform  $(a_{i,i+1}, \dots, a_{in})$  to  $(e_i, 0, \dots, 0)$ ;
          if  $i = n$  then goto ende;
           $s := 0$ ;
          for  $j := i + 1$  step 1 until  $n$  do  $s := s + a[j, i] \uparrow 2$ ;
          if  $s = 0$  then  $e[i] := 0$  else
            begin
               $s := \text{sqrt}(s)$ ;
               $e[i] := \text{if } a[i, i + 1] > 0 \text{ then } -s \text{ else } s$ ;
               $fak := \text{sqrt}(s \times (s + \text{abs}(a[i, i + 1])))$ ;
               $a[i, i + 1] := a[i, i + 1] - e[i]$ ;
              for  $k := i + 1$  step 1 until  $n$  do  $a[i, k] := a[i, k]/fak$ ;
              for  $j := i + 1$  step 1 until  $m$  do
                begin
                   $s := 0$ ;
                  for  $k := i + 1$  step 1 until  $n$  do  $s := s + a[j, k] \times a[i, k]$ ;
                  for  $k := i + 1$  step 1 until  $n$  do  $a[j, k] := a[j, k] - a[i, k] \times s$ ;
                end
              end
            end
          end
        end
      end
    end
  end

```

```

    end j;
  end s;
end i;
ende:
  end bidia;

```

The decomposition using `bidia` requires  $\sim 2(mn^2 - n^3/3)$  multiplications. If  $m > n$  it is possible to bidiagonalize  $A$  even cheaper [26], [5], [8]. The idea is to transform  $A$  first to an upper triangular matrix  $R$  and then bidiagonalize  $R$ . In this case the operation count is  $\sim mn^2 + n^3$ . An alternative approach to bidiagonalizing  $A$  has been suggested by Golub and Kahan [14]. This algorithm uses the matrix  $A$  not explicitly. Only the operator  $A\mathbf{x}$  is needed. Unfortunately it is not numerically stable but nevertheless it seems to be useful for sparse matrices [31].

## 5.4 Computation of the Derivatives of the Length Function

If we want to solve the secular equation

$$f(\lambda) = \alpha^2 \quad (118)$$

using some high order iteration method we have to compute derivatives of  $f$ . For the general problem (P1) we have

$$\begin{aligned} (A^\top A + \lambda C^\top C)\mathbf{x}(\lambda) &= A^\top \mathbf{b} + \lambda C^\top \mathbf{d} \\ f(\lambda) &= \|C\mathbf{x}(\lambda) - \mathbf{d}\|^2. \end{aligned} \quad (119)$$

By differentiating (119) we get

$$\begin{aligned} (A^\top A + \lambda C^\top C)\mathbf{x}' &= -C^\top (C\mathbf{x} - \mathbf{d}) \\ (A^\top A + \lambda C^\top C)\mathbf{x}'' &= -2C^\top C\mathbf{x}'. \end{aligned} \quad (120)$$

In general we have for  $k \geq 2$

$$(A^\top A + \lambda C^\top C)\mathbf{x}^{(k)} = -kC^\top C\mathbf{x}^{(k-1)}. \quad (121)$$

**Lemma 4.** *If  $B_\lambda := A^\top A + \lambda C^\top C$  and  $\mathbf{x}^{(k)}$  is the solution of (121) then*

$$\begin{aligned} (i) \quad \mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)} &= \begin{cases} \mathbf{x}^\top (A^\top \mathbf{b} + \lambda C^\top \mathbf{d}), & k = 0 \\ -\mathbf{x}'^\top C^\top (C\mathbf{x} - \mathbf{d}), & k = 1 \\ -k\mathbf{x}^{(k)\top} C^\top C\mathbf{x}^{(k-1)}, & k \geq 2 \end{cases} \\ (ii) \quad \mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)} &= \begin{cases} \frac{1}{2}\mathbf{x}^\top B_\lambda \mathbf{x}'' + \mathbf{x}'^\top C^\top \mathbf{d}, & k = 1 \\ \frac{k}{k+1}\mathbf{x}^{(k-1)\top} B_\lambda \mathbf{x}^{(k+1)}, & k \geq 2 \end{cases} \\ (iii) \quad \|C\mathbf{x}^{(k)}\|^2 &= \begin{cases} -\mathbf{x}^\top B_\lambda \mathbf{x}' + \mathbf{x}^\top C^\top \mathbf{d}, & k = 0 \\ -\frac{1}{k+1}\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)}, & k \geq 1 \end{cases} \\ (iv) \quad \|C\mathbf{x}^{(k)}\|^2 &= \begin{cases} \frac{1}{2}\mathbf{x}^{(2)\top} C^\top (C\mathbf{x} - \mathbf{d}), & k = 1 \\ \frac{k}{k+1}\mathbf{x}^{(k+1)\top} C^\top C\mathbf{x}^{(k-1)}, & k \geq 2 \end{cases} \\ (v) \quad \frac{d}{d\lambda}(\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)}) &= \begin{cases} -\|C\mathbf{x} - \mathbf{d}\|^2 + \|\mathbf{d}\|^2, & k = 0 \\ -(2k+1)\|C\mathbf{x}^{(k)}\|^2, & k \geq 1 \end{cases} \\ (vi) \quad \frac{d}{d\lambda}\|C\mathbf{x}^{(k)}\|^2 &= \begin{cases} -2\mathbf{x}'^\top B_\lambda \mathbf{x}' + 2\mathbf{x}'^\top C^\top \mathbf{d}, & k = 0 \\ -\frac{2}{k+1}\mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k+1)}, & k \geq 1. \end{cases} \end{aligned}$$

*Proof.* Equations (i) follow, by multiplying (119), (120), (121) by  $\mathbf{x}^{(k)\top}$  from the left.

Looking at two consecutive equations of (121):

$$B\mathbf{x}^{(k)} = -kC^\top C\mathbf{x}^{(k-1)} \quad (122)$$

$$B\mathbf{x}^{(k+1)} = -(k+1)C^\top C\mathbf{x}^{(k)}, \quad (123)$$

we obtain (ii) by multiplying the first (122) by  $\mathbf{x}^{(k)\top}$  and the second by  $\mathbf{x}^{(k-1)\top}$  and by eliminating the expression  $\mathbf{x}^{(k)\top} C^\top C\mathbf{x}^{(k-1)}$ . To prove (iii) we multiply the equation for  $\mathbf{x}^{(k+1)}$  of (121) by  $\mathbf{x}^{(k)}$  and (120) by  $\mathbf{x}^\top$ . Equation (iv) follows by multiplying (122) by  $\mathbf{x}^{(k+1)\top}$  and (123) by  $\mathbf{x}^{(k)\top}$  and subtracting both equations.

Finally observe that  $B'_\lambda = C^\top C$ , therefore

$$\begin{aligned} \frac{d}{d\lambda}(\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)}) &= \mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k)} + \mathbf{x}^{(k)\top} (B'_\lambda \mathbf{x}^{(k)} + B_\lambda \mathbf{x}^{(k+1)}) \\ &= 2\mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k)} + \|C\mathbf{x}\|^2. \end{aligned}$$

Now for  $k = 0$  using (iii) we have

$$\frac{d}{d\lambda}(\mathbf{x}^\top B_\lambda \mathbf{x}) = -\|C\mathbf{x}\|^2 + 2\mathbf{x}^\top C^\top \mathbf{d} = -\|C\mathbf{x} - \mathbf{d}\|^2 + \|\mathbf{d}\|^2.$$

For  $k \geq 1$  using (iii) we get

$$\frac{d}{d\lambda}(\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)}) = -(2k+1)\|C\mathbf{x}^{(k)}\|^2.$$

Using (i) we can write

$$\frac{d}{d\lambda} \|C\mathbf{x}^{(k)}\|^2 = -\frac{2}{k+1} \mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k+1)}.$$

□

**Corrolary 5.** *If  $\mathbf{x}$  is the solution of (119) then*

$$\left. \begin{aligned} \int f(\lambda) d\lambda &= \lambda \|\mathbf{d}\|^2 - \mathbf{x}^\top (A^\top A + \lambda C^\top C) \mathbf{x} + \text{const.} \\ &= -\lambda \mathbf{d}^\top (C\mathbf{x} - \mathbf{d}) - \mathbf{b}^\top A\mathbf{x} + \text{const.} \end{aligned} \right\} \quad (124)$$

*Proof.* From equation (v) of Lemma 4 we have

$$f(\lambda) = \|C\mathbf{x} - \mathbf{d}\|^2 = \|\mathbf{d}\|^2 - \frac{d}{d\lambda}(\mathbf{x}^\top B_\lambda \mathbf{x}).$$

By integrating and using (119) the result follows immediately. □

**Corrolary 6.** *If  $B_\lambda = (A^\top A + \lambda C^\top C)$  then for  $k \geq 1$ ,*

$$\begin{aligned} (i) \quad \frac{d}{d\lambda}(\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k)}) &= \frac{2k+1}{k+1} \mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)} \\ (ii) \quad \frac{d}{d\lambda}(\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)}) &= 2\mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k+1)}. \end{aligned}$$

**Corrolary 7.** *If  $\mathbf{x}^{(k)}$  is a solution of (121) then for  $k \geq 1$*

$$\begin{aligned} (i) \quad \frac{d}{d\lambda} \|C\mathbf{x}^{(k)}\|^2 &= 2\mathbf{x}^{(k+1)\top} C^\top C\mathbf{x}^{(k)} \\ (ii) \quad \frac{d}{d\lambda} \left( \mathbf{x}^{(k+1)\top} C^\top C\mathbf{x}^{(k)} \right) &= \frac{2k+3}{k+1} \|C\mathbf{x}^{(k+1)}\|^2. \end{aligned}$$

**Theorem 8.** Let  $\mathbf{x}$ ,  $\mathbf{x}'$  and  $\mathbf{x}^{(k)}$  be solutions (119), (120 and (121 and let  $B_\lambda = A^\top A + \lambda C^\top C$ . Then

$$f(\lambda) = \|C\mathbf{x} - \mathbf{d}\|^2 = -\mathbf{x}^\top B_\lambda \mathbf{x}' - \mathbf{d}^\top (C\mathbf{x} - \mathbf{d}) \quad (125)$$

$$f'(\lambda) = 2\mathbf{x}^\top C^\top (C\mathbf{x} - \mathbf{d}) = -2\mathbf{x}^\top B_\lambda \mathbf{x}' \quad (126)$$

and for  $k \geq 1$

$$f^{(2k)}(\lambda) = (k+1)\gamma_{2k}\|C\mathbf{x}^{(k)}\|^2 = -\gamma_{2k}\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)} \quad (127)$$

$$\begin{aligned} f^{(2k+1)}(\lambda) &= (k+1)\gamma_{2k+1}\mathbf{x}^{(k)\top} C^\top C\mathbf{x}^{(k+1)} \\ &= -\gamma_{2k+1}\mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k+1)} \end{aligned} \quad (128)$$

where

$$\left. \begin{aligned} \gamma_{2k} &= \frac{1 \cdot 3 \cdot 5 \cdots (2k+1)}{(k+1)!} 2^k \\ \gamma_{2k+1} &= \frac{1 \cdot 3 \cdot 5 \cdots (2k+1)}{(k+1)!} 2^{k+1} \end{aligned} \right\} \quad (129)$$

*Proof.* The proof is by induction. If

$$f(\lambda) = \|C\mathbf{x} - \mathbf{d}\|^2$$

then differentiating and by Lemma 4 (i) we have

$$f'(\lambda) = 2\mathbf{x}^\top C^\top (C\mathbf{x} - \mathbf{d}) = -2\mathbf{x}^\top B_\lambda \mathbf{x}'.$$

Differentiating again using (v) and (iii) of Lemma 4 we get

$$f''(\lambda) = 2 \cdot 3 \|C\mathbf{x}'\|^2 = -3\mathbf{x}^\top B_\lambda \mathbf{x}''$$

which is (127) for  $k = 1$ . Now we can use Corrolary 6 and 7 to compute the higher derivatives.

Assume

$$f^{(2k)}(\lambda) = (k+1)\gamma_{2k}\|C\mathbf{x}^{(k)}\|^2 = -\gamma_{2k}\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)}.$$

then

$$\begin{aligned} f^{(2k+1)}(\lambda) &= (k+1)\gamma_{2k} \cdot 2\mathbf{x}^{(k)\top} C^\top C\mathbf{x}^{(k+1)} \\ &= -\gamma_{2k} \cdot 2\mathbf{x}^{(k)\top} B_\lambda \mathbf{x}^{(k+1)} \end{aligned}$$

which is (128) for  $\gamma_{2k+1} = 2\gamma_{2k}$ . Again differentiating using Corrolary 6 and 7 yields

$$\begin{aligned} f^{(2k+2)}(\lambda) &= \gamma_{2k+1}(2k+3)\|C\mathbf{x}^{(k+1)}\|^2 \\ &= -\gamma_{2k+1} \frac{2k+3}{k+2} \mathbf{x}^{(k+1)\top} B_\lambda \mathbf{x}^{(k+2)} \end{aligned}$$

which is 127) for  $k+1$  and  $\gamma_{2k+2} = \gamma_{2k+1} \frac{2k+3}{k+2}$ . From this recursion for  $\gamma_i$  it is easy to verify (129).  $\square$

Theorem 8 shows that we can compute cheaply derivatives of  $f$ . To compute  $\mathbf{x}^{(k)}$  we have to solve a linear system with the same matrix  $B_\lambda$  as for  $\mathbf{x}$ . Therefore we can use a factorization of  $B_\lambda$ . If  $\lambda > 0$  then  $\mathbf{x}, \mathbf{x}', \dots, \mathbf{x}^{(k)}$  is a solution of the least squares problem

$$\begin{aligned} \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x} &\approx \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix} \\ \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x}' &\approx -\frac{1}{\sqrt{\lambda}} \begin{pmatrix} 0 \\ C\mathbf{x} - \mathbf{d} \end{pmatrix} \\ \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} \mathbf{x}^{(k)} &\approx -\frac{k}{\sqrt{\lambda}} \begin{pmatrix} 0 \\ C\mathbf{x}^{(k-1)} \end{pmatrix}, \quad k \geq 2. \end{aligned}$$

If we transform by another orthogonal matrix  $G$

$$G \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} = \begin{pmatrix} R_\lambda \\ 0 \end{pmatrix}, \quad R_\lambda \text{ upper triangular,}$$

then  $R^\top R$  is the Cholesky decomposition of  $B_\lambda = A^\top A + \lambda C^\top C$ . To solve

$$B_\lambda \mathbf{x}^{(k)} = -k C^\top C \mathbf{x}^{(k-1)}$$

we have two possibilities

(1) compute  $\mathbf{y}^{(k)}$  by forward substitution from

$$R_\lambda^\top \mathbf{y}^{(k)} = -k C^\top (C \mathbf{x}^{(k-1)});$$

(2) compute  $\mathbf{y}^{(k)}$  using  $G$

$$\begin{pmatrix} \mathbf{y}^{(k)} \\ \mathbf{h} \end{pmatrix} = G \begin{pmatrix} 0 \\ -\frac{k}{\sqrt{\lambda}} C \mathbf{x}^{(k-1)} \end{pmatrix}.$$

Then we obtain  $\mathbf{x}^{(k)}$  by backsubstitution in

$$R_\lambda \mathbf{x}^{(k)} = \mathbf{y}^{(k)}.$$

If we define  $\mathbf{z}^{(k)} := C \mathbf{x}^{(k)}$  then

$$f^{(2k-1)}(\lambda) = k \gamma_{2k-1} \mathbf{z}^{(k-1)\top} \mathbf{z}^{(k)}$$

or equivalently

$$f^{(2k-1)}(\lambda) = -\gamma_{2k-1} \|\mathbf{y}^{(k)}\|^2.$$

Similarly

$$f^{(2k)}(\lambda) = (k+1) \gamma_{2k} \|\mathbf{z}^{(k)}\|^2$$

or if  $\mathbf{y}^{(k+1)}$  has been computed

$$f^{(2k)}(\lambda) = -\gamma_{2k} \mathbf{y}^{(k)\top} \mathbf{y}^{(k+1)}.$$

In Section 6 we shall consider third order iteration methods to solve the secular equation for  $\lambda > 0$ . We need therefore the values of  $f$ ,  $f'$  and  $f''$ , which can be computed as follows:

1. Compute  $G$  and  $R_\lambda$  so that

$$G \begin{pmatrix} A \\ \sqrt{\lambda} C \end{pmatrix} = \begin{pmatrix} R_\lambda \\ 0 \end{pmatrix}$$

2. Compute  $\mathbf{y}$  from

$$R^\top \mathbf{y} = A^\top \mathbf{b} + \lambda C^\top \mathbf{d}$$

or using  $G$  by

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{h} \end{pmatrix} = G \begin{pmatrix} \mathbf{b} \\ \sqrt{\lambda} \mathbf{d} \end{pmatrix}.$$

3. Compute  $\mathbf{x}$  from  $R_\lambda \mathbf{x} = \mathbf{y}$  and form  $\mathbf{z} := C \mathbf{x} - \mathbf{d}$ .

4.  $f(\lambda) = \|\mathbf{z}\|^2$ .

5. Compute  $\mathbf{y}'$  by solving

$$R^\top \mathbf{y}' = -C^\top \mathbf{z}$$

or by using  $G$

$$\begin{pmatrix} \mathbf{y}' \\ \mathbf{h}' \end{pmatrix} := G \begin{pmatrix} 0 \\ -\frac{1}{\sqrt{\lambda}} \mathbf{z} \end{pmatrix}$$

6.  $f'(\lambda) := -2\|\mathbf{y}'\|^2$ .
7. Compute  $\mathbf{x}'$  from  $R_\lambda \mathbf{x}' = \mathbf{y}'$  and form  $\mathbf{z}' := C\mathbf{x}'$ .
8.  $f''(\lambda) = 6\|\mathbf{z}'\|^2$ .

If we do not want to store  $G$  then we have to compute  $\mathbf{y}'$  in 5 by forward substitution. However,  $\mathbf{y}$  in step 2 can be computed together with the decomposition in step 1 without forming  $G$  explicitly. Eldén gives in [8] similar recursions to compute the derivatives.

## 6 One-point Iteration Methods to solve the Secular Equation

In this chapter we discuss how to find the solution  $\lambda^* > 0$  of the secular equation

$$f(\lambda) = \alpha^2 \tag{130}$$

that we need to solve problem (P1), (P2) or (P3). The leght function  $f$  is a positive rational function for all  $\lambda$  and decreasing in  $(0, \infty)$ . If we start with  $\lambda^* > 0$  Newton's method will produce a strictly increasing sequence  $\{\lambda_n\}$  which converges globally to  $\lambda^*$ . However as it was observed in [34] if  $\alpha$  is small, convergence is slow and as a remedy Reinsch suggested solving the equation

$$\frac{1}{\sqrt{f}} - \frac{1}{\alpha} = 0 \tag{131}$$

instead of  $\sqrt{f} - \alpha = 0$  using Newton's method. Indeed convergence is much better in this case. There are several possible interpretations and explanations for this fact. For our purpose we compare the Newton step resulting for the three equations

$$\begin{aligned} g_1(\lambda) &:= f(\lambda) - \alpha^2 = 0 \\ g_2(\lambda) &:= \sqrt{f(\lambda)} - \alpha = 0 \\ g_3(\lambda) &:= \frac{1}{\sqrt{f}} - \frac{1}{\alpha} = 0. \end{aligned}$$

A short calculation yields the following Newton iteration functions for the three equations:

$$\lambda - \frac{f - \alpha^2}{f'} \quad \text{for } g_1, \tag{132}$$

$$\lambda - \frac{f - \alpha^2}{f'} \frac{2}{1 + \frac{\alpha}{\sqrt{f}}} \quad \text{for } g_2, \tag{133}$$

$$\lambda - \frac{f - \alpha^2}{f'} \frac{2\frac{\sqrt{f}}{\alpha}}{1 + \frac{\alpha}{\sqrt{f}}} \quad \text{for } g_3. \tag{134}$$

For  $\lambda = 0$  we typically have  $f(0) \gg \alpha^2$ . Since  $f' < 0$  for  $\lambda > 0$  the step in (133) will be about twice as big as in (132). But (134) will produce an even larger step, proportional to  $\frac{\sqrt{f}}{\alpha}$  the discrepancy of  $\sqrt{f}$  and  $\alpha$ . If  $f \approx \alpha^2$  then all three steps are of about the same size.

We can look at the two functions

$$\begin{aligned} h_1(\lambda) &= 2 \left/ \left( 1 + \frac{\alpha}{\sqrt{f}} \right) \right. \\ h_2(\lambda) &= h_1(\lambda) \frac{\sqrt{f}}{\alpha} \end{aligned}$$

as functions that help to accelerate the global convergence of (132) by preserving the order (all iterations are Newton sequences and of second order).



## 6.1 Convergence Factors

For the following discussion we change notation. Let  $f$  be a given function. We are looking for a number  $s$  such that  $f(s) = 0$ . We assume that  $f$  has sufficient continuous derivatives in a neighborhood of  $s$  and furthermore we assume that  $s$  is a simple zero of  $f$ . We consider one point iteration methods without memory [43]

$$\left\{ \begin{array}{l} x_0 \text{ arbitrary} \\ x_{n+1} = F(x_n), \quad n = 0, 1, \dots \end{array} \right\} \quad (135)$$

where

$$F(x) = x - \frac{f(x)}{f'(x)} G(x), \quad (136)$$

and  $G(x)$  is an appropriate chosen function which we will call the “convergence factor”. The idea is to choose  $G$  so that the global convergence using (136) is better than Newton’s iteration ( $G(x) \equiv 1$ ). As has been pointed out by Kahan [24], every sequence generated by an iteration (135) can be interpreted as a sequence obtained by applying Newton’s method to a certain equation

$$g(x) = 0. \quad (137)$$

Indeed if we put

$$x - \frac{g(x)}{g'(x)} = F(x)$$

a short calculation yields (with some  $c \neq 0$ )

$$g(x) = c \cdot \exp \left( \int \frac{dx}{x - F(x)} \right). \quad (138)$$

Solving (137) with (138) using Newton’s method yields the sequence (135). Some example may illustrate the point.

- (1) let  $x_{n+1} = 1 + \frac{x_n}{2}$ ,  $x_0 = 0$ . This sequence converges linearly to  $s = 2$ . Indeed using (138) we get

$$g(x) = \exp \left( \int dx / \left( \frac{x}{2} - 1 \right) \right) = \left( \frac{x}{2} - 1 \right)^2.$$

$g(x)$  has a double zero  $s = 2$  and therefore Newton’s iteration converges linearly.

- (2) Consider the Halley iteration formula  $x_{n+1} = F(x_n)$  with

$$F(x) = x - \frac{2f(x)f'(x)}{2f'(x)^2 - f''(x)f(x)}.$$

Using (138) we obtain

$$g(x) = \exp \left( \int \left( \frac{f(x)}{f'(x)} - \frac{f''(x)}{2f'(x)} \right) dx \right) = \frac{f(x)}{\sqrt{f'(x)}}.$$

Therefore using Halley’s method for  $f(x) = 0$  is applying Newton’s method to

$$g(x) = \frac{f(x)}{\sqrt{f'(x)}} = 0, \quad [3].$$

- (3) Consider the iteration

$$x_{n+1} = x_n - \frac{f(x_n) - \alpha}{f'(x_n)} \cdot \frac{f(x_n)}{\alpha} \quad (139)$$

to solve  $f(x) = \alpha$ . Using (138) we get

$$g(x) = 1 - \frac{\alpha}{f(x)}$$

or since  $g$  is only determined to a constant we may also divide by  $-\alpha$  and get

$$g(x) = \frac{1}{f(x)} - \frac{1}{\alpha}.$$

Therefore (139) is obtained by applying Newton's method to

$$\frac{1}{f(x)} - \frac{1}{\alpha} = 0$$

instead of  $f(x) - \alpha = 0$ .

For the class of iteration functions  $F(x)$  (136) we would like to consider however it will not be possible in general to evaluate the integral for  $g$  in (138) explicitly and thus provide a nice explanation of the iteration function.

The order of convergence [22] of an iteration formula

$$x_{n+1} = F(x_n)$$

is

$$\begin{aligned} &\text{one (or linear convergence), if } |F'(s)| < 1 \\ &\text{two (or quadratic convergence), if } F'(s) = 0 \\ &\text{three (or cubic convergence), if } F'(s) = F''(s) = 0 \\ &\vdots \\ &m, \text{ if } F'(s) = F''(s) = \dots = F^{(m-1)}(s) = 0. \end{aligned}$$

For  $G(x) \equiv 1$  the iteration with

$$F(x) = x - \frac{f(x)}{f'(x)}G(x) \tag{140}$$

is Newton's method and it is of second order for a zero of  $f$  with multiplicity one. Differentiating (140) gives

$$F'(x) = 1 - \left( \frac{f(x)}{f'(x)} \right)' G(x) - \frac{f(x)}{f'(x)} G'(x).$$

Since  $f(s) = 0$  we have  $F'(s) = 0$  if  $G(s) = 1$  and  $f(s) \cdot G'(s) = 0$ . Therefore we have

**Lemma 5.** *Let  $G$  be differentiable,  $G(s) = 1$ ,  $f(s) \cdot G'(s) = 0$ . Then the iteration*

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}G(x_n)$$

*is of second order for simple zeros  $s$  of  $f$ .*

Since  $s$  is unknown we have to choose

$$G(x) = H(f(x), f'(x), \dots).$$

### Examples

(4)  $G(x) = H(f) = 1 + f(x)$

$$\implies F(x) = x - \frac{f(x)}{f'(x)}(1 + f(x)). \tag{141}$$

(141) yields a second order iteration formula. Indeed using (138) we see that the iteration is obtained solving  $g(x) = 0$  with Newton where

$$g(x) = \exp\left(\int \frac{f'(x) dx}{f(x)(1 + f(x))}\right) = \frac{1}{1 + f(x)} - 1.$$

(141) is a special case ( $a = 1$ ) of

(5)  $H(f) = \frac{\alpha + f}{\alpha}$

for some  $\alpha \neq 0$ . Therefore the iteration

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \frac{\alpha + f(x_n)}{\alpha}$$

is the same as applying Newton's method to

$$g(x) = \frac{1}{f(x) + \alpha} - \frac{1}{\alpha} = 0, \quad \alpha \neq 0.$$

(6)

$$G(x) = \frac{1}{1 - \frac{1}{2} \frac{f(x)f''(x)}{f'(x)^2}}.$$

This choice is Halley's method. Again we have  $G(s) = 1$ .

## 6.2 Third Order Iterative Methods

We consider again the iteration

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} G(x_n). \quad (142)$$

In section 6.1 we saw that if  $G(s) = 1$  then (142) is quadratic convergent. Let  $u(x) := \frac{f(x)}{f'(x)}$ . Then

$$F(x) = x - u(x)G(x) \quad (143)$$

and we can ask for the conditions that the iteration (142) be cubically convergent. Differentiating (143) we get (dropping the argument)

$$\begin{aligned} F' &= 1 - u'G - uG' \\ F'' &= -u''G - 2u'G' - uG'' \\ u &= f/f' \\ u' &= 1 - \frac{ff''}{f'^2} \\ u'' &= -\frac{f''}{f'} + 2\frac{ff''^2}{f'^3} - \frac{ff'''}{f'^2}. \end{aligned}$$

Now we want to have

$$F'(s) = F''(s) = 0. \quad (144)$$

Since  $u(s) = 0$ ,  $u'(s) = 1$  and  $u''(s) = \frac{1}{2} \frac{f''(s)}{f'(s)}$ , this is the case if

$$G(s) = 1 \quad \text{and} \quad G'(s) = \frac{1}{2} \frac{f''(s)}{f'(s)}. \quad (145)$$

**Theorem 9.** Let  $H \in C^2[-a, a]$  for some  $a > 0$ . The iteration  $x_{n+1} = F(x_n)$  with

$$F(x) = x - \frac{f(x)}{f'(x)} \cdot H\left(\frac{f(x)f''(x)}{f'(x)^2}\right) \quad (146)$$

converges cubically to a simple zero of  $f$  if and only if

$$H(0) = 1 \quad \text{and} \quad H'(0) = \frac{1}{2}.$$

*Proof.* The function  $F$  in (146) is the special case of (143) with

$$G(x) = H\left(\frac{f(x)f''(x)}{f(x)^2}\right).$$

Let  $t(x) = \frac{f(x)f''(x)}{f(x)^2}$ . Then  $t(x) = 1 - u'(x)$  and therefore

$$t'(s) = -u''(s) = \frac{f''(s)}{f'(s)}.$$

Clearly  $G(s) = 1 \iff H(0) = 1$  and

$$\begin{aligned} G'(s) &= H'(t(s))t'(s) \\ &= H'(0) \cdot \frac{f''(s)}{f'(s)} = \frac{1}{2} \frac{f''(s)}{f'(s)} \\ \iff H'(0) &= \frac{1}{2}. \end{aligned}$$

□

Many well known third order iterative methods are special cases of Theorem 9. Let

$$t := \frac{f(x)f''(x)}{f'(x)^2} \tag{147}$$

(1) Euler's formula

$$H(t) = \frac{2}{1 + \sqrt{1 - 2t}} = 1 + \frac{1}{2}t + \frac{1}{2}t^2 + \frac{5}{8}t^3 + \dots$$

clearly  $H(0) = 1$  and  $H'(0) = \frac{1}{2}$ .

(2) Halley's formula.

$$H(t) = \frac{1}{1 - \frac{1}{2}t} = 1 + \frac{1}{2}t + \frac{1}{4}t^2 + \frac{1}{8}t^3 + \dots$$

(3) Quadratic inverse interpolation [43].

$$H(t) = 1 + \frac{1}{2}t.$$

(4) Ostrowski's square root iteration [30].

$$H(t) = \frac{1}{\sqrt{1 - t}} = 1 + \frac{1}{2}t + \frac{3}{8}t^2 + \frac{5}{16}t^3 + \dots$$

(5) Hansen-Patrick family [19].

$$H(t) = \frac{\alpha + 1}{\alpha + \sqrt{1 - (\alpha + 1)t}} = 1 + \frac{1}{2}t + \frac{\alpha + 3}{8}t^2 + \dots$$

(6) To solve the secular equation we will use

$$H(t) = e^{\frac{1}{2}t} = 1 + \frac{1}{2}t + \frac{1}{8}t^2 + \frac{1}{48}t^3 + \dots$$

This formula has the advantage that  $H(t) > 0$  even if the starting point is far from the solution  $s$ . For large  $t$ , the other formulas may not work (wrong sign in  $H$ , argument of the square root is negative).

The following lemma is useful for construction if third order iteration methods.

**Lemma 6.** Let  $H_1(t)$  and  $H_2(t)$  be two functions with

$$H_i(0) = a \quad \text{and} \quad H'_i(0) = b, \quad i = 1, 2.$$

The the three mean functions

$$\begin{aligned} A &= (H_1 + H_2)/2 \\ B &= \sqrt{H_1 H_2} \\ C &= 2 \left/ \left( \frac{1}{H_1} + \frac{1}{H_2} \right) \right. \end{aligned}$$

have the same property.

*Proof.* Obviously  $A(0) = B(0) = C(0) = a$ .

$$\begin{aligned} A' &= (H'_1 + H'_2)/2 = b \\ B' &= (H'_1 H_2 + H_1 H'_2)/(2\sqrt{H_1 H_2}) \\ B'(0) &= (ba + ab)/(2a) = b \\ C &= \frac{2H_1 H_2}{H_1 + H_2} \\ C' &= 2 \frac{(H_1 + H_2)(H'_1 H_2 + H_1 H'_2) - (H'_1 + H'_2)H_1 H_2}{(H_1 + H_2)^2} \\ C'(0) &= 2 \frac{2a(2ab) - 2ba^2}{4a^2} = b. \end{aligned}$$

□

Any of the three means of the examples (1) to (6) yields a new third order iterative method.

Notice that not every third order iterative method must have the form (147). For example, consider

$$G(x) = H(t(x)) + f^2(x),$$

with  $H(0) = 1$ ,  $H'(0) = 1/2$ , and  $t(x)$  defined by (147). Clearly this choice of the convergence factor  $G$  leads to a third order formula which has not the form of Theorem 9.

However it is possible to describe all the third order iteration methods. Let

$$F_k := x + \sum_{i=1}^{k-1} (-1)^i \frac{f^i}{i!} \left( \frac{1}{f'} \frac{d}{dx} \right)^{i-1} \frac{1}{f'}$$

denote the Schröder iteration function (see [7] or [21]). The iteration

$$x_{n+1} = F_k(x_n)$$

is of  $k$ -th order. Now every  $k$ -th order iteration function can be written as

$$F(x) = F_k(x) + f^k(x)\varphi(x)$$

where  $\varphi$  is an arbitrary function [7], [43]. Therefore for  $k = 2$  we obtain the most general  $G$  for a third order method

$$G(x) = 1 + \frac{1}{2}t(x) + f^2(x)\varphi(x)$$

which we can write

$$G(x) = H(t(x)) + f^2(x)\psi(x)$$

with arbitrary  $\psi$ .

### 6.3 The Convergence Factor for a Third Order Method

Assume  $x_{n+1} = x_n - K(x_n)$  is a third order iterative method for solving  $f(x) = 0$ . We clearly have

$$K(s) = K''(s) = 0, \quad K'(s) = 1. \quad (148)$$

Now consider the iteration  $x_{n+1} = F(x_n)$  with

$$F(x) = x - K(x)G(x). \quad (149)$$

We have

$$\left. \begin{aligned} F' &= 1 - K'G - KG' \\ F'' &= -K''G - 2K'G' - KG'' \end{aligned} \right\} \quad (150)$$

**Lemma 7.** *The iteration  $x_{n+1} = F(x_n)$  with  $F$  defined in (149) is also of third order if  $G \in C^2[a, b]$  with  $s \in (a, b)$  and*

$$G(s) = 1, G'(s) = 0. \quad (151)$$

*Proof.* If we use (148) and (151) in (150) then we have  $F'(s) = F''(s) = 0$ .  $\square$

If we choose  $G$  so that

$$G''(s) = -\frac{1}{3}K'''(s)$$

then we would have a fourth order iteration. However we think that a third order method with good global convergence is more useful than a local convergent fourth order formula.

The following functions are examples of possible convergence factors for a third order iteration:

(1)  $G(x) = (t^\beta(x) + t^{-\beta}(x))/2$

where  $t(x) = \frac{f(x) + \alpha}{\alpha} > 0$  and  $\alpha, \beta \in \mathbb{R}$ .

(2)  $G(x) = \cosh(f(x) \cdot r(x))$

where  $r \in C^2[a, b]$  with  $s \in (a, b)$ .

(3)  $G(x) = L(f^2(x) \cdot r(x))$

where  $L \in C^2[-d, d]$ ,  $d > 0$  and  $L(0) = 1$ ,  $r \in C^2[a, b]$  with  $s \in (a, b)$ .

We know now how to choose convergence factors to preserve or increase the order of an iteration formula. Our aim is however to improve global convergence. Let

$$F(x) = x - K(x)G(x)$$

be the iteration function. Then  $G$  has to be chosen so that

$$g(x) = \exp\left(\int \frac{dx}{K(x) \cdot G(x)}\right) \quad (152)$$

is nearly linear if  $x$  is far away from  $s$ . Since the iteration  $x_{n+1} = f(x_n)$  is the same as applying Newton's method to  $g(x) = 0$ , global convergence will be good if  $g$  is linear far away from the solution  $s$ . However, since (152) may be impossible to compute explicitly, this requirement seems not to be practical to determine  $G$ . We have to estimate  $G$  by other means. In Section 6.4 we shall interpret some  $G$  geometrically. Those functions can serve as models for others.

### 6.4 Geometrical Interpretation of the Convergence Factors

It is possible to derive iterative methods as follows:

(i) choose a "simple" function  $h$  so that

$$f^{(i)}(x) = h^{(i)}(x), \quad i = 0, 1, \dots, k;$$

- (ii) solve analytically  $h(z) = \alpha$  obtaining  $z = z(x)$ ;
- (iii) use the iteration  $x_{n+1} = z(x_n)$  to solve the equation  $f(x) = \alpha$

The function  $h$  should be simple and so that  $h(z) = \alpha$  can be solved analytically.  $h$  approximates  $f$  and some derivatives locally at one point  $x$  and we thus obtain a one point iteration formula without memory.

Instead of approximating  $f$  one can also find iteration methods by approximating the inverse function locally:

- (i) Choose a function  $h(y)$  so that

$$\left(f^{[-1]}(y)\right)^{(i)} = h^{(i)}(y), \quad i = 0, 1, \dots, k.$$

- (ii) Put  $y = f(x_n)$  and use the iteration formula

$$x_{n+1} = h(\alpha)$$

to solve  $f(x) = \alpha$ .

Since we do not know  $f^{[-1]}$ , the derivatives must be replaced using derivatives of  $f$  [21]:

$$\begin{aligned} f^{[-1]}(y)' &= \frac{1}{f'(f^{[-1]}(y))} \\ f^{[-1]}(y)'' &= -\frac{f''(f^{[-1]}(y))}{f'^3(f^{[-1]}(y))}. \end{aligned}$$

If  $h$  approximates  $f$  resp  $f^{[-1]}$  well we can expect a good global convergence. We give in the following some examples of methods derived by interpolation. The convergence factors of these examples may help to choose a method analytically.

- (1) Newton's method  $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$  is obtained approximating  $f$  or  $f^{[-1]}$  locally by a linear function  $h$  ( $h(x) = ax + b$  resp.  $h(y) = ay + b$ ).
- (2) We choose  $h(z) = \frac{a}{(z-b)^2}$  and determine  $a, b$  so that  $h^{(i)}(x) = f^{(i)}(x), i = 1, 2$ , giving (dropping the argument  $x$ ):

$$\left. \begin{aligned} a &= 4f^3/f'^2 \\ b &= -x - 2f/f'. \end{aligned} \right\} \quad (153)$$

Now  $h(z) = \alpha$  gives

$$z = -b \pm \sqrt{\frac{a}{\alpha}}$$

using (153) we get

$$z = x - \frac{f}{f'} 2 \left( \pm \sqrt{\frac{a}{\alpha}} - 1 \right).$$

If we use only the + sign we have the iteration

$$\left. \begin{aligned} x_{n+1} &= x_n - \frac{f(x_n) - \alpha}{f'(x_n)} G(x_n) \\ \text{with } G(x) &= 2\sqrt{\frac{f}{\alpha}} \left/ \left( 1 + \sqrt{\frac{f}{\alpha}} \right) \right. \end{aligned} \right\} \quad (154)$$

But (154) is Reinsch's proposal to solve  $\frac{1}{\sqrt{f}} - \frac{1}{\sqrt{\alpha}} = 0$  with Newton instead of  $f - \alpha = 0$ . We know from Section 3 that the length function  $f$  has the form

$$f(x) = \sum_{i=1}^n \frac{c_i^2}{(x - x_i)^2}.$$

Therefore it is reasonable to approximate  $f$  by  $h = a/(x - b)^2$ .

(3) Instead of solving  $f(x) - \alpha = 0$  with Newton's method we can also solve

$$g(x) = f(x)^{-1/\beta} - \alpha^{-1/\beta} \quad (155)$$

if  $f(x) > 0$  and  $\alpha > 0$ . Newton's iteration for (155) yields the same iteration function as if we approximate  $f$  and  $f'$  locally by

$$h(x) = \frac{a}{(x+b)^\beta}.$$

The resulting iteration formula is

$$x_{n+1} = x_n - \frac{f(x_n) - \alpha}{f'(x_n)} G(x_n)$$

with

$$G(x) = \frac{\beta f}{f - \alpha} \left( \sqrt[\beta]{\frac{f}{\alpha}} - 1 \right). \quad (156)$$

If we let  $\beta \rightarrow \infty$  in (156) we get

$$G(x) = \frac{f}{f - \alpha} \ln \left( \frac{f}{\alpha} \right). \quad (157)$$

This convergence factor is also obtained by solving  $\ln(f) - \ln(\alpha) = \ln(f/\alpha) = 0$  with Newton of by approximating locally  $f$  and  $f'$  by

$$h(x) = ae^{bx}.$$

The next examples use functions  $h$  that approximate  $f$ ,  $f'$  and  $f''$  locally

(4) We choose  $h(z) = a/(z+b) + c$  and determine  $a$ ,  $b$ ,  $c$  so that

$$h^{(i)}(x) = f^{(i)}(x), i = 0, 1, 2.$$

We obtain (dropping the argument  $x$ ):

$$\begin{aligned} a &= -4f'^3/f''^2 \\ b &= -x - 2f'/f'' \\ c &= f - 2f'^2/f''. \end{aligned}$$

Now solving  $h(z) = \alpha$  gives the iteration (with  $f_n^{(i)} = f^{(i)}(x_n)$ )

$$x_{n+1} = x_n - \frac{f_n - \alpha}{f'_n} \frac{1}{1 - \frac{1}{2} \frac{(f_n - \alpha) f''_n}{f_n'^2}} \quad (158)$$

which is Halley's formula to solve  $g(x) = f(x) - \alpha = 0$ . Recall (see Section 6.1) that the same iteration is obtained if we solve with Newton's method

$$g(x) = \frac{f(x) - \alpha}{\sqrt{f'(x)}}.$$

If we approximate the inverse function  $f^{[-1]}(y)$  locally by  $h(y) = a/(y+b) + c$  we get

$$\begin{aligned} a &= 4f'^3/f''^2 \\ b &= -f + 2f'^2/f'' \\ c &= x + 2f'/f''. \end{aligned}$$



Finally putting  $x_n = x$  and  $x_{n+1} = h(\alpha)$  yields again (158). Therefore Halley's formula is obtained by locally approximating  $f$  or  $f^{[-1]}$  by a hyperbola.

This rational approximation of  $f$  has the following property. Suppose we want to solve  $f(x) = \alpha$  using Halley's method. We obtain the same iteration (158 for

$$g_1(x) = f(x) - \alpha = 0$$

or

$$g_2(x) = \frac{1}{f(x)} - \frac{1}{\alpha} = 0, \quad (159)$$

i.e., for Halley's method the transformation (159) has no effect. But as we saw (159) yield another convergence factor for Newton's method.

- (5) Euler's method is obtained by approximating  $f$ ,  $f'$ , and  $f''$  locally with

$$h(z) = az^2 + bz + c.$$

We get

$$\begin{aligned} a &= f'/2, & b &= f' - f''x \\ c &= f - f'x + \frac{f''}{2}x^2 \end{aligned}$$

and

$$x_{n+1} = x_n - \frac{f_n - \alpha}{f'_n} G(x_n)$$

with

$$G(x) = 2 \left/ \left( 1 + \sqrt{1 - 2 \frac{(f - \alpha)f''}{f'^2}} \right) \right.$$

- (6) Approximating  $f^{[-1]}$  by a parabola  $h(y) = ay^2 + by + c$ , yields

$$\begin{aligned} a &= -\frac{1}{2} \frac{f''}{f'^3} \\ b &= \frac{1}{f'} + \frac{f''}{f'^3} f \\ c &= x - \frac{f}{f'} - \frac{1}{2} \frac{f''}{f'^3} f^2. \end{aligned}$$

We obtain the iteration formula for  $f(x) - \alpha = 0$

$$x_{n+1} = x_n - \frac{f_n - \alpha}{f'_n} \left( 1 + \frac{1}{2} \frac{(f_n - \alpha)f''_n}{f'^2_n} \right)$$

- (7) Approximating  $f$  locally by  $h(z) = ae^{bz} + c$  yields

$$\begin{aligned} a &= \frac{f'^2}{f''} \exp\left(-\frac{f''}{f'}x\right) \\ b &= f''/f' \\ c &= f - f'^2/f''. \end{aligned}$$

Solving  $h(z) = \alpha$  gives the iteration

$$x_{n+1} = x_n - \frac{f'_n}{f''_n} \ln\left(1 - \frac{(f_n - \alpha)f''_n}{f'^2_n}\right)$$

which we can write

$$x_{n+1} = x_n - \frac{f_n - \alpha}{f'_n} G(x_n)$$

with

$$\begin{aligned} G(x) &= -\ln(1 - t(x))/t(x) \\ t(x) &= (f(x) - \alpha)f''(x)/f'^2(x). \end{aligned}$$

From Theorem 9 we know that all the methods (4)–(7) are cubically convergent.

## 6.5 Solving the Secular Equation

A very good and simple method to solve  $f(\lambda) = \alpha^2$  is Reinsch's proposal

$$\lambda_{n+1} = \lambda_n - \frac{f(\lambda_n) - \alpha}{f'(\lambda_n)} G(\lambda_n) \quad (160)$$

where

$$G(\lambda) = 2 \frac{\sqrt{f(\lambda)}}{\alpha} \left/ \left( 1 + \frac{\alpha}{\sqrt{f(\lambda)}} \right) \right.$$

Iteration 160) can also be written as

$$\lambda_{n+1} = \lambda_n - \frac{f(\lambda_n)}{f'(\lambda_n)} \cdot 2 \left( \frac{\sqrt{f(\lambda_n)}}{\alpha} - 1 \right)$$

Reinsch proved in [34] that the sequence  $\lambda_n$  obtained with this iteration converges starting with  $\lambda_1 = 0$  monotonically increasing to the solution. This property together with good global convergence makes this iteration very attractive.

If we use another method with starting value  $\lambda_1 = 0$ , e.g. Halley's iteration, global convergence may not be guaranteed. As an example, consider Problem (P3)

$$\begin{aligned} &\min \|\mathbf{x}\| \\ &\text{subject to} \\ &\|A\mathbf{x} - \mathbf{b}\| \leq \alpha \end{aligned}$$

where  $n = m = 6$ ,  $A$  Hilbert matrix,  $a_{ij} = 1/(i + j - 1)$ ,  $\mathbf{b} = [1, 0, 0, 0, 0, 0]^\top$ .

If we start with  $\lambda_0 = 0.1$  we get  $\alpha_1 = -0.7975$  instead of a value bigger than  $\lambda_0$ . The other third order methods involving a term

$$1 - \frac{1}{2}t \text{ or } 1 - t, \quad t = \frac{ff''}{f'^2}$$

may also fail if  $t > 1$ . The quadratic inverse interpolation

$$H(t) = 1 + \frac{1}{2}t \quad (161)$$

yields a correction with the right sign but numerical experiments show that if Halley converges, it converges globally much faster.

An improvement of (161) is to use

$$H(t) = \exp(t/2). \quad (162)$$

However global convergence has to be accelerated by a convergence factor  $H$ .

Therefore we propose to use

$$\lambda_{n+1} = \lambda_n - \frac{f(\lambda_n) - \alpha^2}{f'(\lambda_n)} \exp\left(\frac{1}{2} \frac{ff''}{f'^2}\right) \cdot H(\lambda_n) \quad (163)$$

with

$$H(\lambda) = \frac{1}{2} \left( \frac{\sqrt{f(\lambda)}}{\alpha} + \frac{\alpha}{\sqrt{f(\lambda)}} \right).$$

Since we cannot prove monotonicity or global convergence properties, it is necessary to take measures against too big steps. If the step is too big (usually at the beginning) we use Reinsch's iteration (160). For the above mentioned example with the Hilbert matrix we obtain using (163) and  $\lambda_0 = 0.1$  the solution of the secular equation  $\lambda \approx 7.7 \times 10^{11}$  in four iterations.

To illustrate the behaviour of the different methods we computed the only positive solution of the equation

$$f(\lambda) = 1$$

where

$$f(\lambda) = 0.6 + \sum_{i=1}^{20} \frac{2 + 0.8^i}{(\lambda + 0.8^i)^2}.$$

We used ALGOL W and double precision (ca. 16 decimal digits) on the IBM 360 of SLAC. The results are displayed in Table 1. Newton's method uses 21 steps to converge to machine precision. Reinsch's variant converges after 7 steps. The improved Halley's method needs only four steps.

**Table 1**

Example for the behaviour of different methods to solve  $g(\lambda) = f(\lambda) - \alpha = 0$ .

Second order methods  $\lambda_{n+1} = F(\lambda_n)$ :

Newton

Reinsch

Newton for  $\ln(f/\alpha) = 0$

$$F(\lambda) = \lambda - \frac{g}{g'}$$

$$\lambda - \frac{f}{f'} 2 \left( \frac{f}{\alpha} - 1 \right)$$

$$\lambda - \frac{f}{f'} \ln(f/\alpha)$$

$\alpha$	$\alpha$	$\alpha$
0.00783209855661240	3.20297843199116	0.0034187403647234
0.0209101058943323	7.50002221663786	0.530746065218083
0.0424438763172333	9.69455596729093	1.97350054359389
0.0774671705806476	10.2422619393792	4.58739428667985
0.13329409609385	10.2695361392028	7.47964284219632
0.223664677428177	10.2700022243362	9.5278543635909
0.365485460269776	10.2700022243362	10.2135024725941
0.567031564547462	10.2700022243362	10.2696683666710
0.928977169378511	10.2700022243362	10.2700022126256
1.44940702417730	10.2700022243362	10.2700022243362
2.22759607120317	10.2700022243362	10.2700022243362
3.36004538158582	10.2700022243362	10.2700022243362
4.92591522457886	10.2700022243362	10.2700022243362
6.87134069977040	10.2700022243362	10.2700022243362
8.79350313070678	10.2700022243362	10.2700022243362
9.92730829940186	10.2700022243362	10.2700022243362
10.2574259482621	10.2700022243362	10.2700022243362
10.2699795745767	10.2700022243362	10.2700022243362
10.2700022242627	10.2700022243362	10.2700022243362

Third order methods  $\lambda_{n+1} = \lambda_n - \frac{g}{g'} H(t)$  with  $t := \frac{gg''}{g'^2}$ :

The following two methods fail for starting values  $< 6.8$ :

Ostrowski

Euler

$H(t) = 1/(1-t)^{0.5}$	$2/(1+(1-2t)^{0.5})$
6.67134069977040	6.87134069977040
11.3578564654989	8.79350313070678
10.2615901313624	10.4062426841626
10.2700022251005	10.2699502114743
10.2700022243362	10.2700022243362
10.2700022243362	10.2700022243362

**Table 1 (cont.)**

Halley

Inverse interp.

$1/(1-t/2)$	$1+t/2$	$\exp(t/2)$
0.0494307764343190	0.0144232340488969	0.0181702328293991
0.253229265175140	0.0460967940921278	0.0830975341263060
0.986381372252434	0.113492541266632	0.169596042052953
3.21756490475624	0.253299632263121	0.413519282789038
7.65449586417966	0.536339890697594	0.951636742732932
10.2238730991803	1.0531852895540	2.08316964604166
10.2700019998135	2.14957041757370	4.28001192403909
10.2700022243362	4.03642343012226	7.65289089221572
10.2700022243362	6.91272434349204	10.0374630402632
10.2700022243362	9.6092062586452	10.2692135997469
10.2700022243362	10.2634961753614	10.2700022243362
10.2700022243362	10.2700022160612	10.2700022243362
10.2700022243362	10.2700022243362	10.2700022243362
10.2700022243362	10.2700022243362	10.2700022243362

Third order methods with convergence factor  $G = \left( \frac{\sqrt{f}}{\alpha} + \frac{\alpha}{\sqrt{f}} \right) / 2 :$

Halley

$H(t) = G/(1 - t/2)$	$G \exp(t/2)$
5.07846031436599	1.88673371095732
9.89047170050903	5.5133405779731
10.2699142055335	9.03682504447327
10.2700022243362	10.2700022243362
10.2700022243362	10.2700022243362
10.2700022243362	10.2700022243362

## 7 Generalizations

Let  $A$  be an  $(n \times n)$  matrix,  $\mathbf{b}$  a given  $n$  vector,  $C$  an  $m \times n$  matrix,  $\mathbf{d}$  a given  $m$  vector and  $\gamma, \alpha$  real numbers. We consider the problem

$$F(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{x} + \gamma = \min \quad (164)$$

subject to

$$\|C\mathbf{x} - \mathbf{d}\| = \alpha. \quad (165)$$

Using a Lagrange multiplier  $\lambda/2$  the solution of (164), (165) is a stationary point  $(\mathbf{x}, \lambda)$  of

$$L(\mathbf{x}, \lambda) = F(\mathbf{x}) + \frac{\lambda}{2} (\|C\mathbf{x} - \mathbf{d}\|^2 - \alpha^2). \quad (166)$$

$\frac{\partial L}{\partial \mathbf{x}} = 0$  and  $\frac{\partial L}{\partial \lambda} = 0$  gives

$$\left. \begin{aligned} (A + A^\top + \lambda C^\top C) \mathbf{x} &= -\mathbf{b} + \lambda C^\top \mathbf{d} \\ \|C\mathbf{x} - \mathbf{d}\|^2 &= \alpha^2. \end{aligned} \right\} \quad (167)$$

**Theorem 10.** Let  $(\mathbf{x}_i, \lambda_i)$ ,  $i = 1, 2$ , be two solutions of (167), then

$$F(\mathbf{x}_1) - F(\mathbf{x}_2) = \frac{\lambda_2 - \lambda_1}{4} \|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2. \quad (168)$$

*Proof.* This theorem is a generalization of Theorem 1 and the proof is very similar. We have

$$(A + A^\top + \lambda_1 C^\top C) \mathbf{x}_1 = -\mathbf{b} + \lambda_1 C^\top \mathbf{d} \quad (169)$$

$$(A + A^\top + \lambda_2 C^\top C) \mathbf{x}_2 = -\mathbf{b} + \lambda_2 C^\top \mathbf{d}. \quad (170)$$

$\mathbf{x}_1^\top$  (169) gives

$$2 \mathbf{x}_1^\top A \mathbf{x}_1 + \lambda_1 \|C \mathbf{x}_1\|^2 = -\mathbf{x}_1^\top \mathbf{b} + \lambda_1 \mathbf{x}_1^\top C^\top \mathbf{d},$$

or rearranged

$$2(\mathbf{x}_1^\top A \mathbf{x}_1 + \mathbf{x}_1^\top \mathbf{b}) = \mathbf{x}_1^\top \mathbf{b} - \lambda_1 (\|C \mathbf{x}_1\|^2 - \mathbf{x}_1^\top C^\top \mathbf{d}). \quad (171)$$

Similarly we have

$$2(\mathbf{x}_2^\top A \mathbf{x}_2 + \mathbf{x}_2^\top \mathbf{b}) = \mathbf{x}_2^\top \mathbf{b} - \lambda_2 (\|C \mathbf{x}_2\|^2 - \mathbf{x}_2^\top C^\top \mathbf{d}). \quad (172)$$

Now  $\mathbf{x}_2^\top$  (169)  $-\mathbf{x}_1^\top$  (170) gives

$$\begin{aligned} (\lambda_1 - \lambda_2) \mathbf{x}_2^\top C^\top C \mathbf{x}_1 &= -\mathbf{x}_2^\top \mathbf{b} + \mathbf{x}_1^\top \mathbf{b} + \lambda_1 \mathbf{x}_2^\top C^\top \mathbf{d} - \lambda_2 \mathbf{x}_1^\top C^\top \mathbf{d} \\ \implies (\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{b} &= (\lambda_1 - \lambda_2) \mathbf{x}_2^\top C^\top C \mathbf{x}_1 - \mathbf{d}^\top (\lambda_1 C \mathbf{x}_2 - \lambda_2 C \mathbf{x}_1). \end{aligned} \quad (173)$$

Subtracting (171) - (172) and replacing  $(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{b}$  by (173) we get

$$\begin{aligned} 2(F(\mathbf{x}_1) - F(\mathbf{x}_2)) &= -\lambda_1 \{ \|C \mathbf{x}_1\|^2 - \mathbf{x}_1^\top C^\top \mathbf{d} - \mathbf{x}_2^\top C^\top C \mathbf{x}_1 + \mathbf{x}_2^\top C^\top \mathbf{d} \} \\ &\quad + \lambda_2 \{ \|C \mathbf{x}_2\|^2 - \mathbf{x}_2^\top C^\top \mathbf{d} - \mathbf{x}_2^\top C^\top C \mathbf{x}_1 + \mathbf{x}_1^\top C^\top \mathbf{d} \}. \end{aligned}$$

Now like in the proof of Theorem 1 we know that both  $\{ \}$  are equal and therefore we replace them by their arithmetic mean which gives the desired result.  $\square$

**Corrolary 8.** *The solution of (164), (165) is the solution  $(\mathbf{x}, \lambda)$  of (167) with largest  $\lambda$ .*

*Proof.* From (168) we have

$$\lambda_2 < \lambda_1 \implies F(\mathbf{x}_2) > F(\mathbf{x}_1)$$

□

**Theorem 11.** *Let  $(\mathbf{x}_i, \lambda_i)$ ,  $i = 1, 2$ , be two solutions of (167), then*

$$(\lambda_1 + \lambda_2)(F(\mathbf{x}_1) - F(\mathbf{x}_2)) = (\lambda_1 - \lambda_2)(\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2). \quad (174)$$

*Proof.*

$$\lambda_1 C^\top C \mathbf{x}_1 - \lambda_1 C^\top \mathbf{d} = -(A + A^\top) \mathbf{x}_1 - \mathbf{b} \quad (175)$$

$$\lambda_2 C^\top C \mathbf{x}_2 - \lambda_2 C^\top \mathbf{d} = -(A + A^\top) \mathbf{x}_2 - \mathbf{b}. \quad (176)$$

$\lambda_1 \mathbf{x}_1^\top$  (176)  $-\lambda_2 \mathbf{x}_2^\top$  (175) gives

$$\lambda_1 \lambda_2 ((\mathbf{x}_2 - \mathbf{x}_1)^\top C^\top \mathbf{d}) = (\lambda_2 - \lambda_1) \mathbf{x}_1^\top (A + A^\top) \mathbf{x}_2 - (\lambda_1 \mathbf{x}_1 - \lambda_2 \mathbf{x}_2)^\top \mathbf{b}. \quad (177)$$

$\lambda_1 \mathbf{x}_2^\top$  (176)  $-\lambda_2 \mathbf{x}_1^\top$  (175) gives

$$\lambda_1 \lambda_2 (\|C \mathbf{x}_2\|^2 - \|C \mathbf{x}_1\|^2 + (\mathbf{x}_1 - \mathbf{x}_2)^\top C^\top \mathbf{d}) = 2\lambda_1 \mathbf{x}_2^\top A \mathbf{x}_2 + 2\lambda_2 \mathbf{x}_1^\top A \mathbf{x}_1 - (\lambda_1 \mathbf{x}_2 - \lambda_2 \mathbf{x}_1)^\top \mathbf{b}. \quad (178)$$

Observe that

$$0 = \|C \mathbf{x}_2 - \mathbf{d}\|^2 - \|C \mathbf{x}_1 - \mathbf{d}\|^2 = \|C \mathbf{x}_2\|^2 - \|C \mathbf{x}_1\|^2 + 2((\mathbf{x}_1 - \mathbf{x}_2)^\top C^\top \mathbf{d}).$$

So that if we subtract (178)–(177) we get

$$0 = 2\lambda_2 \mathbf{x}_1^\top A \mathbf{x}_1 - 2\lambda_1 \mathbf{x}_2^\top A \mathbf{x}_2 - (\lambda_1 \mathbf{x}_2 - \lambda_2 \mathbf{x}_1 - \lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2)^\top \mathbf{b} + (\lambda_1 - \lambda_2) \mathbf{x}_1^\top (A + A^\top) \mathbf{x}_2.$$

Rearranged we have

$$\begin{aligned} & \lambda_1 \{2\mathbf{x}_2^\top A \mathbf{x}_2 + \mathbf{x}_2^\top \mathbf{b} - \mathbf{x}_1^\top \mathbf{b} - \mathbf{x}_1^\top (A + A^\top) \mathbf{x}_2\} \\ & = \lambda_2 \{2\mathbf{x}_1^\top A \mathbf{x}_1 + \mathbf{x}_1^\top \mathbf{b} - \mathbf{x}_2^\top \mathbf{b} - \mathbf{x}_1^\top (A + A^\top) \mathbf{x}_2\}. \end{aligned} \quad (179)$$

Now observe that the left hand side of (179) is

$$\{ \quad \} = F(\mathbf{x}_2) - F(\mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2).$$

Similarly the right hand side simplifies. Therefore

$$\begin{aligned} & \lambda_1 \{F(\mathbf{x}_2) - F(\mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2)\} \\ & = \lambda_2 \{F(\mathbf{x}_1) - F(\mathbf{x}_2) + (\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2)\} \end{aligned}$$

or rearranged

$$(\lambda_1 + \lambda_2)(F(\mathbf{x}_1) - F(\mathbf{x}_2)) = (\lambda_1 - \lambda_2)(\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2).$$

□

**Corrolary 9.** *Let  $(\mathbf{x}_i, \lambda_i)$ ,  $i = 1, 2$ , be two solutions of (167) with  $\lambda_1 \neq \lambda_2$ . Then*

$$(\mathbf{x}_1 - \mathbf{x}_2)^\top A(\mathbf{x}_1 - \mathbf{x}_2) = -\frac{\lambda_1 + \lambda_2}{4} \|C(\mathbf{x}_1 - \mathbf{x}_2)\|^2. \quad (180)$$

*Proof.* We combine the results of Theorems 10 and 11 .

□

## 8 Smoothing of Datas

Rutishauser proposed in [36] a method to smooth datas. We present here a modified version. Let

$$d_i, \quad i = 1, \dots, n$$

be given datas of a smooth function. However let's assume the  $d_i$  are perturbed by measurement errors. We look for a new set of datas

$$x_i, \quad i = 1, \dots, n$$

that does not deviate too much from the  $d_i$  and that is smoother. If we assume that the  $d_i$  are equidistant then we may want to solve

$$\sum_{i=2}^{n-1} (x_{i+1} - 2x_i + x_{i-1})^2 = \min \quad (181)$$

subject to

$$\sum_{i=1}^n (x_i - d_i)^2 \leq n\delta^2. \quad (182)$$

Here  $\delta^2$  is a measure for the variance, the mean deviation we want to allow the new data  $x_i$  to differ from  $d_i$ :

$$\delta \approx \sqrt{\frac{\sum_{i=1}^n (x_i - d_i)^2}{n}}.$$

The larger  $\delta$  the more the  $x_i$  will be smoothed.

Introducing the  $(n-2) \times n$  tridiagonal matrix

$$A = \begin{pmatrix} 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & & 1 & -2 & 1 \end{pmatrix}$$

and the vectors  $\mathbf{d}$  and  $\mathbf{x}$ , equations (181) and (182) become

$$\left. \begin{array}{l} \|A\mathbf{x}\| = \min \\ \text{subject to} \\ \|\mathbf{x} - \mathbf{d}\| \leq \alpha := \sqrt{n} \cdot \delta \end{array} \right\} \quad (183)$$

and we have problem (P1) for  $\mathbf{b} = 0$  and  $C = I$ . The normal equations are

$$(A^T A + \lambda I)\mathbf{x} = \lambda \mathbf{d} \quad (184)$$

$$\|\mathbf{x} - \mathbf{d}\|^2 = \alpha^2. \quad (185)$$

Instead of solving (184), Rutishauser proposed to choose some “smoothing” parameter  $\gamma$  and to solve

$$(I + \gamma A^T A)\mathbf{x} = \mathbf{d} \quad (186)$$

which is (184) for  $\lambda = 1/\gamma$ . The condition number of the matrix in (186) is  $16\gamma$ . For large  $n$  one may have to choose  $\gamma$  large which leads to numerical problems [36]. However large  $\gamma$  or small  $\lambda$  may be meaningful: for  $\gamma \rightarrow \infty$  ( $\lambda \rightarrow 0$ ) the new values  $\mathbf{x}$  lie on a straight line, i.e., are values of a linear function. Of course one would like to have as limit the linear regression to the datas  $\mathbf{d}$ .

If we make a change of variables in (183)

$$\begin{aligned} \mathbf{w} &:= \mathbf{x} - \mathbf{d} \\ \mathbf{b} &:= -A\mathbf{d} \end{aligned} \quad (187)$$

then our problem becomes

$$\left. \begin{array}{l} \|A\mathbf{w} - \mathbf{b}\| = \min \\ \text{subject to} \\ \|\mathbf{w}\| \leq \alpha \end{array} \right\} \quad (188)$$

Problem (188) leads to a relaxed least squares problem

$$\left. \begin{array}{l} (A^\top A + \lambda I)\mathbf{w} = A^\top \mathbf{b} \\ \|\mathbf{w}\| = \alpha \end{array} \right\} \quad (189)$$

We can use the dual equations for (189) (see Section 3:

$$(AA^\top + \lambda I)\mathbf{z} = -\mathbf{b} = A\mathbf{d} \quad (190)$$

$$\|A^\top \mathbf{z}\| = \alpha \quad (191)$$

with

$$\mathbf{w} = -A^\top \mathbf{z}. \quad (192)$$

Observe that now  $\lambda \rightarrow 0$  causes no problems since  $AA^\top$  is not singular. Furthermore since  $\mathbf{b} = -A\mathbf{d}$  we can write (190) as

$$\begin{pmatrix} A^\top \\ \sqrt{\lambda}I \end{pmatrix} \mathbf{z} \approx \begin{pmatrix} \mathbf{d} \\ 0 \end{pmatrix} \quad (193)$$

$A$  is tridiagonal, therefore we shall use the algorithm described in Section 5.1 to solve for a given  $\lambda \geq 0$  the least squares problem (193).

For  $n = 8$ , e.g. the remaining matrix before and after the third step is:

$$\left[ \begin{array}{cccccccc} 0 & & & & & & & \\ 0 & 0 & & & & & & \\ 0 & 0 & r & & & & & \\ & 0 & s & t & & & & \\ & & 1 & -2 & 1 & & & \\ & & & 1 & -2 & 1 & & \\ & & & & 1 & -2 & 1 & \\ & & & & & 1 & & \\ x & x & x & & & & & \\ & x & x & x & & & & \\ & & \sqrt{\lambda} & & & & & \\ & & & \sqrt{\lambda} & & & & \\ & & & & \sqrt{\lambda} & & & \\ & & & & & \sqrt{\lambda} & & \\ & & & & & & \sqrt{\lambda} & \end{array} \right] \xrightarrow{\text{step 3}} \left[ \begin{array}{cccccccc} 0 & & & & & & & \\ 0 & 0 & & & & & & \\ 0 & 0 & 0 & & & & & \\ & 0 & 0 & r & & & & \\ & & 0 & s & t & & & \\ & & & 1 & -2 & 1 & & \\ & & & & 1 & -2 & 1 & \\ & & & & & 1 & & \\ x & x & x & & & & & \\ & x & x & x & & & & \\ & & d_{13} & d_{23} & d_{33} & & & \\ & & & \sqrt{\lambda} & & & & \\ & & & & \sqrt{\lambda} & & & \\ & & & & & \sqrt{\lambda} & & \\ & & & & & & \sqrt{\lambda} & \end{array} \right]$$

The following algorithm transforms

$$U^\top \begin{pmatrix} A^\top \\ \sqrt{\lambda}I \end{pmatrix} = \begin{pmatrix} 0 \\ R \end{pmatrix} \quad \text{and} \quad U^\top \begin{pmatrix} \mathbf{d} \\ 0 \end{pmatrix} = \begin{pmatrix} \mathbf{c} \\ \mathbf{y} \end{pmatrix}$$

where  $U$  is orthogonal and  $R$  upper triangular and tridiagonal:

$$R = \begin{bmatrix} d_{11} & d_{21} & d_{31} & & & & & \\ & d_{12} & d_{22} & d_{32} & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & \ddots & \ddots & d_{3_{n-4}} & & \\ & & & & \ddots & d_{2_{n-3}} & & \\ & & & & & d_{1_{n-2}} & & \end{bmatrix}.$$

The matrix  $A$  is not stored, we work only with three simple variables  $r$ ,  $s$ , and  $t$ . We use the notation

$$\begin{pmatrix} 0 & a'_2 & a'_3 \\ b'_1 & b'_2 & b'_3 \end{pmatrix} := G \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} \quad (194)$$



where  $G$  is a Givens matrix

$$G = \begin{pmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{pmatrix}$$

and  $\alpha$  has been chosen so that the zero in the left hand side of (194) appears.

```

begin
  comment zeroing columns 1 to  $n - 4$  ;
   $r := t := 1; s := -2; c_1 := d_1; c_2 := d_2;$ 
  for  $i := 1$  step 1 until  $n - 4$  do
    begin
       $\begin{pmatrix} 0 & c_i \\ d1_i & y_i \end{pmatrix} := G_1 \begin{pmatrix} r & c_i \\ \sqrt{\lambda} & 0 \end{pmatrix};$ 

       $\begin{pmatrix} 0 & r & c_{i+1} \\ d1_i & d2_i & y_i \end{pmatrix} := G_2 \begin{pmatrix} s & t & c_{i+1} \\ d1_i & 0 & y_i \end{pmatrix};$ 

       $\begin{pmatrix} 0 & s & t & c_{i+2} \\ d1_i & d2_i & d3_i & y_i \end{pmatrix} := G_3 \begin{pmatrix} 1 & -2 & 1 & d_{i+2} \\ d1_i & d2_i & 0 & y_i \end{pmatrix};$ 
    end;
  comment column  $n - 3$ ;
   $\begin{pmatrix} 0 & c_{n-3} \\ d1_{n-3} & y_{n-3} \end{pmatrix} := G_1 \begin{pmatrix} r & c_{n-3} \\ \sqrt{\lambda} & 0 \end{pmatrix};$ 

   $\begin{pmatrix} 0 & r & c_{n-2} \\ d1_{n-3} & d2_{n-3} & y_{n-3} \end{pmatrix} := G_2 \begin{pmatrix} s & t & c_{n-2} \\ d1_{n-3} & 0 & y_{n-3} \end{pmatrix};$ 

   $\begin{pmatrix} 0 & s & c_{n-1} \\ d1_{n-3} & d2_{n-3} & y_{n-3} \end{pmatrix} := G_3 \begin{pmatrix} 1 & -2 & d_{n-1} \\ d1_{n-3} & d2_{n-3} & y_{n-3} \end{pmatrix};$ 
  comment column  $n - 2$ ;
   $\begin{pmatrix} 0 & c_{n-2} \\ d1_{n-2} & y_{n-2} \end{pmatrix} := G_1 \begin{pmatrix} r & c_{n-2} \\ \sqrt{\lambda} & 0 \end{pmatrix};$ 

   $\begin{pmatrix} 0 & c_{n-1} \\ d1_{n-2} & y_{n-2} \end{pmatrix} := G_2 \begin{pmatrix} s & c_{n-1} \\ d1_{n-2} & y_{n-2} \end{pmatrix};$ 

   $\begin{pmatrix} 0 & c_n \\ d1_{n-2} & y_{n-2} \end{pmatrix} := G_3 \begin{pmatrix} 1 & d_n \\ d1_{n-2} & y_{n-2} \end{pmatrix};$ 
end.

```

To solve the problem we need the derivatives of  $f(\lambda)$  defined by

$$(AA^\top + \lambda I)\mathbf{z}(\lambda) = \mathbf{Ad}$$

$$f(\lambda) = \|A^\top \mathbf{z}(\lambda)\|^2.$$

We cannot apply the results of Section 5.4 since  $f$  is not of the same form. We have

$$(AA^\top + \lambda I)\mathbf{z}^{(k)} = -\mathbf{z}^{(k-1)}, \quad k \geq 1.$$

We shall use Reinsch's proposal and therefore we need only  $f'$ :

$$\begin{aligned} f'(\lambda) &= 2\mathbf{z}'^\top AA^\top \mathbf{z} \\ &= -2\mathbf{z}'^\top (\mathbf{z} + \lambda \mathbf{z}'). \end{aligned}$$

An ALGOL W procedure `smooth(n, δ, d, x)` is given at the end of this section. The following example has been computed with this procedure using single precision (ca. 6 decimal digits) on the SLAC computer (IBM 360).

### Example

We choose  $n = 30$  and

$$d_i = \sqrt{i} + 0.2 \sin(i), \quad i = 1, \dots, 30.$$

For this example we have  $f(0) = 1.825814$ , therefore if

$$\delta > \sqrt{\frac{f(0)}{n}} = 0.246699$$

the smoothed value  $x_i$  are the same as obtained by linear regression.

On the other hand if  $\delta = 10^{-4}$  we have  $\|\mathbf{x} - \mathbf{d}\| = 5.4 \times 10^{-4}$  and the  $x_i$  differ only in the fourth decimal from the  $d_i$ . If we interpret the  $d_i$  as perturbed value of  $\sqrt{i}$  then because the mean of  $|0.2 \sin(i)|$  is approximately

$$0.2 \frac{1}{\pi} \int_0^\pi \sin(x) dx = \frac{0.4}{\pi} \approx 1.3$$

we expect the best smoothing for  $\delta \approx 1.3$  which indeed can be seen clearly in Table 2.

**Table 2:** Smoothing of  $d_i = \sqrt{i} + 0.2 \sin(i)$ ,  $i = 1, \dots, 30$ .

$\delta =$	0.2467	0.2466	0.2	0.17	0.15	0.13	0.12
1.722253	1.722027	1.603430	1.507798	1.414468	1.261987	1.243047	
1.861272	1.851084	1.763119	1.654705	1.609509	1.496164	1.483359	
2.000287	2.000137	1.922686	1.861354	1.804054	1.727399	1.717025	
2.139301	2.133186	2.081965	2.037419	1.997584	1.955855	1.947420	
2.278316	2.275234	2.240736	2.212497	2.189486	2.182792	2.181764	
2.417335	2.417207	2.393729	2.386040	2.378851	2.406053	2.418491	
2.556347	2.556335	2.555611	2.557375	2.564402	2.619123	2.643021	
2.695338	2.695371	2.711038	2.725835	2.745208	2.815103	2.841743	
2.834349	2.834425	2.864746	2.890921	2.920292	2.992062	3.001115	
2.973354	2.973435	3.016541	3.052362	3.089548	3.154735	3.155210	
3.112376	3.112478	3.166302	3.210031	3.253120	3.310709	3.301394	
3.251370	3.251450	3.313931	3.363794	3.411095	3.464372	3.457992	
3.390381	3.390501	3.459267	3.513468	3.563265	3.613997	3.616666	
3.529382	3.529527	3.602212	3.658435	3.709314	3.754458	3.762660	
3.668403	3.668556	3.742742	3.799443	3.849200	3.883019	3.886796	
3.807439	3.807585	3.880661	3.936624	3.983325	4.002785	3.995223	
3.946456	3.946573	4.016701	4.069475	4.112409	4.120646	4.105030	
4.085469	4.085518	4.150387	4.198711	4.237081	4.241594	4.228400	
4.224489	4.224607	4.282032	4.324522	4.357614	4.364625	4.362045	
4.363507	4.363605	4.411738	4.447031	4.473988	4.484037	4.490573	
4.502520	4.502613	4.539614	4.566405	4.586239	4.594892	4.600982	
4.641560	4.641610	4.665621	4.682461	4.694700	4.697631	4.694865	
4.780583	4.780605	4.790689	4.797174	4.800344	4.797594	4.787060	
4.919610	4.919535	4.914412	4.909512	4.903647	4.899208	4.891846	
5.058637	5.058389	5.037205	5.020307	5.005067	5.003985	5.000236	
5.197662	5.197525	5.175234	5.129760	5.104605	5.103989	5.123779	
5.336608	5.336575	5.280628	5.238020	5.202195	5.192799	5.216838	
5.475712	5.475664	5.401550	5.345339	5.298068	5.268023	5.281142	
5.614742	5.614557	5.522194	5.452085	5.392826	5.333412	5.324587	
5.753768	5.753551	5.642737	5.558625	5.487171	5.395162	5.368006	
$\ \mathbf{x} - \mathbf{d}\ $	1.351226	1.350682	1.095428	0.9311236	0.8215933	0.7120381	0.6572665
$[\sum(\Delta^2 x_i)^2]^{1/2}$	$6.67 \cdot 10^{-5}$	$1.11 \cdot 10^{-4}$	0.002449	0.0150982	0.021434	0.0445737	0.0796041
$[\sum(x_i - \sqrt{i})^2]^{1/2}$	1.205420	1.204210	0.9014646	0.682826	0.5094146	0.306330	0.3031789
# of iteration	0	1	5	5	7	7	6
$\lambda = f^{-1}(n\delta^2)$	0	$3.83 \cdot 10^{-7}$	$2.79 \cdot 10^{-4}$	$7.60 \cdot 10^{-4}$	$2.01 \cdot 10^{-3}$	$3.15 \cdot 10^{-2}$	$8.89 \cdot 10^{-2}$

Table 2 (cont.)

$\delta =$	0.1	0.07	0.05	0.01	0.001	0.0001	$\sigma(4)$
1.238009	1.232636	1.223015	1.183728	1.169964	1.168462	1.168294	
1.483301	1.497992	1.514537	1.571436	1.595238	1.595784	1.596072	
1.710183	1.713635	1.723948	1.755995	1.760237	1.760276	1.760274	
1.930025	1.906127	1.891484	1.858504	1.849538	1.848726	1.848639	
2.167420	2.136563	2.11879	2.056374	2.045358	2.044389	2.044283	
2.425476	2.422603	2.415583	2.397161	2.393924	2.393637	2.393606	
2.674776	2.710691	2.731615	2.769294	2.776435	2.777078	2.777148	
2.877497	2.924973	2.956028	3.014403	3.025218	3.026191	3.026290	
3.022841	3.040714	3.053176	3.077481	3.081974	3.082378	3.082423	
3.139317	3.111163	3.092867	3.050061	3.054070	3.053531	3.053473	
3.271142	3.221653	3.188797	3.128712	3.117723	3.116734	3.116826	
3.439886	3.413122	3.395536	3.363274	3.357374	3.356843	3.356786	
3.626307	3.645697	3.659373	3.664522	3.689124	3.689530	3.689534	
3.789226	3.836155	3.868443	3.927829	3.936693	3.939671	3.939778	
3.904166	3.935042	3.956223	3.995197	4.002328	4.002970	4.003040	
3.986406	3.972642	3.963206	3.945898	3.942735	3.942450	3.942419	
4.077334	4.031558	4.000144	3.942432	3.931878	3.930929	3.930825	
4.206725	4.171137	4.146597	4.101509	4.093266	4.092824	4.092843	
4.355673	4.373278	4.378106	4.387066	4.388709	4.388857	4.388873	
4.515091	4.559080	4.588960	4.643706	4.653724	4.654625	4.654724	
4.622007	4.662054	4.689597	4.739809	4.748988	4.749015	4.749005	
4.690324	4.690964	4.689130	4.688752	4.688654	4.688645	4.688648	
4.757719	4.714632	4.687192	4.636802	4.627513	4.626677	4.626586	
4.861431	4.812959	4.782627	4.728794	4.716855	4.717960	4.717865	
5.004077	4.990499	4.983374	4.975096	4.973674	4.973543	4.973529	
5.150356	5.120525	5.200182	5.242207	5.250694	5.251448	5.251531	
5.258895	5.303250	5.329032	5.375997	5.386373	5.387322	5.387426	
5.307034	5.333555	5.342944	5.345037	5.345425	5.345654	5.345683	
5.318588	5.311031	5.302584	5.268337	5.254272	5.252623	5.252437	
5.319951	5.284638	5.272645	5.273405	5.278864	5.279541	5.279616	
$\ x - d\ $	0.547712e	0.383404e	0.2735594	0.054773	5.47 10 <sup>-3</sup>	5.47 10 <sup>-4</sup>	0
$[\sum(\Delta^2 x_i)^2]^{1/2}$	0.1683889	0.3160816	0.4206818	0.6471357	0.7032083	0.7097257	0.7103844
$[\sum(x_i - \sqrt{i})^2]^{1/2}$	0.3583993	0.4791697	0.5678519	0.7469465	0.7844216	0.7879522	0.7883406
# of iteration	6	5	5	5	4	3	
$\lambda = f^{-1}(n\delta^2)$	0.2640591	0.7694743	1.500484	13.16168	153.4504	1562.654	

```

PROCEDURE SMOOTH(INTEGER VALUE N;
                 REAL VALUE DELTA ;
                 REAL ARRAY D,X(*) ) ;
BEGIN
  REAL G,F0,FA,F1,F2,LAMB,LAMBN,WLAMB,R,S,T,ALPHA,ALPHA2,CO,SI,H ;
  REAL ARRAY C,AZ,AZS(1::N) ;
  REAL ARRAY D1,Y,Z,ZS(1::N-2) ;
  REAL ARRAY D2(1::N-3) ; REAL ARRAY D3(1::N-4) ;
  INTEGER I ;
  PROCEDURE ATZ(REAL ARRAY Z,Y(*) ) ;
  BEGIN INTEGER I ;
    Y(1) := Z(1) ;
    Y(2) := -2*Z(1) + Z(2) ;
    FOR I := 3 STEP 1 UNTIL N-2 DO
      Y(I) := Z(I-2) -2*Z(I-1) + Z(I) ;
    Y(N-1) := Z(N-3) -2*Z(N-2) ;
    Y(N) := Z(N-2) ;
  END ATZ ;
  PROCEDURE BACK(INTEGER VALUE N; REAL ARRAY A,B,C,D,X(*) ) ;
  BEGIN INTEGER I ;
    X(N) := D(N)/A(N) ;
    X(N-1) := (D(N-1) - X(N)*B(N-1))/A(N-1) ;
    FOR I := N-2 STEP -1 UNTIL 1 DO
      X(I) := (D(I) - X(I+1)*B(I) - X(I+2)*C(I))/A(I) ;
    END BACK ;
  PROCEDURE VORW(INTEGER VALUE N; REAL ARRAY A,B,C,D,X(*) ) ;
  BEGIN INTEGER I ;
    X(1) := D(1)/A(1) ;
    X(2) := (D(2) - X(1)*B(1))/A(2) ;
    FOR I := 3 STEP 1 UNTIL N DO
      X(I) := (D(I) - X(I-1)*B(I-1) - X(I-2)*C(I-2))/A(I) ;
    END VORW ;
  PROCEDURE ROT(REAL VALUE A; REAL VALUE B) ;
  BEGIN REAL T ;
    IF B = 0 THEN

```

```

        BEGIN CO := 0 ; SI := 1 END
ELSE
    BEGIN
        T := -A/B ; CO := 1/SQRT(1 + T**2) ;
        SI := T*CO ;
        END ;
END ROT ;
REAL PROCEDURE INPROD(INTEGER VALUE N; REAL ARRAY X,Y(*) ) ;
BEGIN INTEGER I ; REAL S ;
    S := 0 ;
    FOR I := 1 STEP 1 UNTIL N DO S := S + X(I)*Y(I) ;
    S
END INPROD ;
ALPHA := SQRT(N)*DELTA ; ALPHA2 := ALPHA**2 ;
LAMB := 0;
IT:WLAMB := SQRT(LAMB) ;
COMMENT COLUMNS 1 TO N-4
R := T := 1 ; S := -2 ; C(1) := D(1); C(2) := D(2) ;
FOR I := 1 STEP 1 UNTIL N-4 DO
BEGIN
    ROT(R,WLAMB) ;
    D1(I) := -SI*R + CO*WLAMB ;
    Y(I) := -SI*C(I) ; C(I) := CO*C(I) ;
    ROT(S,D1(I)) ;
    D1(I) := -SI*S + CO*D1(I) ;
    R := CO*T ; D2(I) := -SI*T ;
    H := CO*C(I+1)+SI*Y(I); Y(I):=-SI*C(I+1)+CO*Y(I) ;
    C(I+1) := H ;
    ROT(1,D1(I)) ;
    D1(I) := -SI + CO*D1(I) ;
    S := -2*CO + SI*D2(I) ; D2(I) := 2*SI + CO*D2(I) ;
    T := CO ; D3(I) := -SI ;
    H := CO*D2(I+2) + SI*Y(I) ;
    Y(I) := -SI*D(I+2) + CO*Y(I) ; C(I+2) := H ;
END ;
FOR I := N-3,N-2 DO
BEGIN
    COMMENT COLUMNS N-3 AND N-2
    ROT(R,WLAMB) ;
    D1(I) := -SI*R + CO*WLAMB ;
    Y(I) := -SI*C(I) ; C(I) := CO*C(I) ;
    ROT(S,D1(I)) ;
    D1(I) := -SI*S + CO*D1(I) ;
    IF I = N-3 THEN
    BEGIN
        R := CO*T ; D2(I) := -SI*T ;
    END
    H := CO*C(I+1) + SI*Y(I) ;
    Y(I) := -SI*C(I+1) + CO*Y(I) ; C(I+1) := H ;
    ROT(1,D1(I)) ;
    D1(I) := -SI + CO*D1(I) ;
    IF I = N-3 THEN
    BEGIN
        S := -2*CO + SI*D2(I) ; D2(I) := 2*SI + CO*D2(I) ;
    END
    C(I+2) := CO*D(I+2) + SI*Y(I) ;
    Y(I) := -SI*D(I+2) + CO*Y(I) ;
END I ;
IF LAMB = 0 THEN FA:= 2*INPROD(N,C,C) ;

```

```

BACK(N-2,D1,D2,D3,Y,Z) ;
ATZ(Z,AZ) ;
FO := INPROD(N,AZ,AZ) ;
VORW(N-2,D1,D2,D3,Z,ZS) ;
COMMENT ZS IS -Z'
BACK(N-2,D1,D2,D3,ZS,ZS) ;
ATZ(ZS,AZS) ;
F1 := -2*INPROD(N,AZS,AZ) ;
LAMB := LAMB -FO/F1*2*(SQRT(FO)/ALPHA-1) ;
IF (FA <= FO) OR (FO < ALPHA2) OR (LAMB <= LAMB)
THEN GOTO FIN ;
FA := FO; LAMB := LAMB ;
GOTO IT ;
FIN:
FOR I := 1 STEP 1 UNTIL N DO
X(I) := D(I) - AZ(I);
END SMOOTH ;

```

## References

- [1] Bjork, A., "Iterative Refinement of Linear Least Squares Solutions I", BIT 7 (1967), 257-278.
- [2] Bjork, A., "Solving Linear Least Squares Problems by Gram-Schmidt Orthogonalization", BIT 7 (1967), 1-21.
- [3] Brown, G.H. Jr., "On Halley's Variation of Newton's Method", American Mathematical Monthly 84, 9 (1977).
- [4] Bunch, J. R. and Rose, D. J., Sparse Matrix Computations, Academic Press, (1976).
- [5] Chan, T. F. C., "On Computing the Singular Value Decomposition", Stanford Computer Science Department Report STAN-CS-77-588, (1977).
- [6] Davis, M. and Dawson, B., "On Global Convergence of Halley's Iteration Formula", Numer. Math. 24 (1975).
- [7] Ehrmann, H., "Konstruktion und Durchfuerung von Iterationsverfahren hoeherer Ordnung", Arch. Rational Mech. Anal. 4 (1959).
- [8] Elden, L., "Algorithms for the Regularization of Ill-Conditioned Least Squares Problems", BIT 17 (1977), 134-145.
- [9] Elden, L., "Numerical Analysis of Regularization and Constrained Least Squares Methods", Thesis No. 21, Linköping University (1977).
- [10] Forsythe, G. E. and Golub, G. H. , "On the Stat. Values of a Second Degree Polynomial on the Unit Sphere", J. Soc. Indust. Appl. Math. 13 (1965).
- [11] Forsythe, G. E., Malcom, M. A. and Moler, C. B., Computer Methods for Mathematical Computations, Prentice Hall, 1977.
- [12] Gander, W., "How to apply the Dual Problem to Solve a Certain Constrained Minimum Norm Problem", SIAM Spring Meeting Madison 1978.
- [13] Gander, W., Molinari, L. and Svecowa, H., Numerische Prozeduren . . . , ISNM 33, Birkhaeuser Verlag 1977.
- [14] Golub, G. H. and Kahan, W., "Calculating the Singular Values and Pseudo-Inverse of a Matrix", SIAM J. Numer. Anal. 2, 2 (1965).

- [15] Golub, G. H., Klema, V. and Stewart, G. W., "Rank Degeneracy and Least Squares Problems", Stanford Computer Science Department Report STAN-CS-76-559 (1976)
- [16] Golub, G. H., "Some Modified Matrix Eigenvalue Problems", SIAM Review 15, 2 (1973).
- [17] Golub, G. H., Heath, M. and Wahba, G., "Generalised Cross-Validation", Stanford Computer Science Department Report STAN-CS-77-622 (1977)
- [18] Golub, G. H. and Luk, F. T., "Singular Value Decomposition: Applications and Computations", Internal Report Serra House Stanford Computer Science Department (1978).
- [19] Hansen, E. and Patrick, M., "A Family of Root Finding Methods", Numer. Math. 27 (1977), 257-269.
- [20] Hanson, R. J. and Phillips, J. L. , "An Adaptive Numerical Method for Fredholm Integral Equations", Numer. Math. 24 (1975), 291-307.
- [21] Henrici, P., Applied and Computational Complex Analysis, Wiley, 1974.
- [22] Henrici, P., Elements of Numerical Analysis, Wiley, 1964
- [23] Kahan, W., Manuscript in Box 5 G of the G.Forsythe archive, Stanford Main Library.
- [24] Kahan, W., "Why Use Tangents When secants Will Do", Gatlinburgh Conference 1977, Asilomar.
- [25] Kalman, R. E., "Algebraic Aspects of the Generalised Inverse . . .", Generalised Inverses and Applications, Academic Press, (1976).
- [26] Lawson, Ch. L. and Hanson, R. J., Solving Least Squares Problems, Prentice Hall, 1974.
- [27] Marquart, D. W., "An Algorithm for Least-Squares Estimation of Nonlinear Parameters", SIAM 11, 2 (1963).
- [28] Marti, J. T., "Minimum Norm Solutions of Fredholm Integral Equations of the First Kind", Report 78-01 Sem. f. angew. Math. ETHZ, (1978).
- [29] Moré, J. J., "The levenberg-Marquart Algorithm: Implementation and Theory", Dundee Conference on Numerical Analysis (1977).
- [30] Ostrowski, A. M., Solutions of Equations and System of Equations, Academic Press, (1973).
- [31] Paige, C. C., "Bidiagonalization of Matrices and Solution of Linear Equations", SIAM J. Numer. Anal. 11, 1 (1974).
- [32] Peters, G. and Wilkinson, J. H., " $Ax = \lambda Bx$  and the Generalized Eigenproblem", SIAM J. Numer. Anal. 7, 4 (1970).
- [33] Reinsch, Chr. H., "Smoothing by Spline Functions", Numer. Math. 10 (1967), 177-183.
- [34] Reinsch, Chr. H., "Smoothing by Spline Functions II", Numer. Math. 16 (1971), 451-454.
- [35] Rutishauser, H., "Once Again: The Least Square Problem", Linear Algebra and Its Applications 1 (1968), 479-488.
- [36] Rutishauser, H., "Vorlesungen ueber numerische Mathematik, hrsg. von Martin Gutknecht", Band I, II, Birkhaeuser Verlag, 1976.
- [37] Schek, H.-J. and Eggenesperger, R., "Least-Squares - Loesungen und Daempfung bei unterbestimmten Gleichungssystemen", Computing 19, (1977).
- [38] Schwarz, H. R., Rutishauser, H. and Stiefel, E., Matrizen-Numerik, Teubner Verlag, 1968.

- [39] Spjøtvoll, E., “A Note on a Theorem of Forsythe and Golub”, *SIAM J. Appl. Math.* 23,3 (1972).
- [40] Stewart, G. W., “On the Continuity of the Generalized Inverse”, *SIAM J. Appl. Math.* 17, 1 (1969).
- [41] Stewart, G. W., “On the Numerical Properties of an Iteration for Computing the Generalised Inverse”, CNA 12, Austin, Texas (1971),
- [42] Tikhonov, A. N., *Solution of Ill Posed Problems*, Wiley, 1977.
- [43] Traub, J. F., *Iterative Methods for Solution of Equations*, Prentice Hall, 1964.
- [44] Van Loan, Ch., F., “Lectures in Least Squares”, Technical Report TR 76-279, Cornell University (1976).
- [45] Van Loan, Ch. F., “Generalizing the Singular Value Decomposition”, *SIAM J. Numer. Anal.* 13, 1 (1976).
- [46] Varah, J. M., “On the Numerical Solution of Ill-Conditioned Linear Systems”, *SIAM J. Num. Anal.* 10 (1973).
- [47] Varah, J. M., “A Practical Examination of Num. Methods for Ill Posed Problems”, Technical Report 76-08, University of British Columbia, 1976.
- [48] Wilkinson, J. H. and Reinsch, C., *Linear Algebra*, Springer, 1971.