

ON THE LINEAR LEAST SQUARES PROBLEM
WITH A QUADRATIC CONSTRAINT

by

Walter Gander

STAN-CS-78-697
November 19 78

COMPUTER SCIENCE DEPARTMENT
School of Humanities and Sciences
STANFORD UNIVERSITY



On the Linear Least Squares Problem
with a Quadratic Constraint

On the Linear Least Squares Problem
with a Quadratic Constraint

by
Walter Gander ^{*/}

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 of 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 data. 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.

by

Walter Gander*

^{*/} Computer Science Department, Stanford University, Stanford, Calif. 94305.
On leave from Neu-Technikum Ruchs, Switzerland.
The author has been supported by the Swiss National Science Foundation,
Grant No. 5'521'330'61517.

* Neu-Technikum **Buchs, CH-9470 Buchs**, Switzerland.
Research supported in part by National Science Foundation Grant
No. **MG578-17697**.

Acknowledgement.

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.

To Maurice

Stanford, Fall 1978
Walter Gander

Table of Contents

0. Introduction 1

1. Basic Definitions and Remarks 5

2. The Least Square Problem with Quadratic Constraint 9

 2.1 Characterization of the Solution 15

 2.2 The Solutions of the Normal Equations 23

 2.3 The Solution for the Equality Constraint (PIE). 24

 2.4 The Solution for Inequality Constraint (PI) 27

3. The Relaxed Least Squares Problem 28

 3.1 Results from the General Theory 29

 3.2 The **Dual Normal Equations** 33

 3.3 **Elden's** Transformation 35

 3.4 Rutishauser's Relaxed and Doubly Relaxed Least Squares Problem 42

4. Minimum Norm Solution with Given Norm of the Residual 43

 4.1 **The Dual Normal Equations** 44

 4.2 Representation as Least Squares Problems 45

5. Computational Aspects 46

 5.1 Solution of a Relaxed Least Squares Problem with **Band Matrix** 53

 5.2 Solution of a Least Squares Problem with Two Band Matrices 56

5.3 Bidiagonalization 62

5.4 **Computation** of the Derivatives of the Length **Function** 71

6. One Point Iterative Methods to Solve the Secular Equation 72

 6.1 Convergence Factors 77

 6.2 Third Order Iterative Methods 82

 6.3 The Convergence Factors for a Third Order Method 84

 6.4 **Geometrical** Interpretation of the Convergence Factors 90

7. Solving the Secular Equation 95

8. Generalizations 99

 Smoothing of Datas 112

References

0. Introduction.

In this paper we consider least squares problems with a quadratic constraint. The matrices and vectors will be real and we use capital letters A, B, \dots for matrices and small letters a, b, \dots for vectors. $\| \cdot \|$ will be used for the euclidian vectornorm.

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

$$\left. \begin{aligned} \|Ax - b\| &= \min \\ \|Cx - d\| &\leq \alpha \end{aligned} \right\} \quad (0.1)$$

subject to

This problem is a generalization of the least square problem with equality constraints

$$\left. \begin{aligned} \|Ax - b\| &= \min \\ Cx &= d \end{aligned} \right\} \quad (0.2)$$

subject to

since for $\alpha = 0$ every solution of (0.1) is also a solution of (0.2).

The motivation why to consider (0.1) rather than (0.2) 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)

$$\underline{Ax} = \underline{y} \quad (0.3)$$

with $a_{ij} = t_i^j$. If $m > n$, (0.3) 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 x is difficult to compute accurately and that usually the norm of x will be large. The polynomial p with large coefficient will be useless since cancellation will affect its evaluation. It may therefore be better not to interpolate but to approximate the datas in the least squares sense. This leads to the unconstrained least squares problem

$$\|\underline{Ax} - \underline{y}\| = \min . \quad (0.4)$$

This problem has for $m \geq n$ and if A has full rank a unique solution. If A is rank deficient (0.4) has infinitely many solutions and therefore one usually looks for the solution with minimum norm

$$\left. \begin{array}{l} \min \|\underline{x}\| \\ \text{subject to } \|\underline{Ax} - \underline{y}\| = \min . \end{array} \right\} \quad (0.5)$$

The solution of (0.5) is given explicitly by

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

where A^+ is the pseudoinverse of A . However the solution of (0.4) [35] and (0.5) 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 $\|\underline{x}\|$ and thus look for the solution of

$$\left. \begin{array}{l} \|\underline{Ax} - \underline{y}\| = \min \\ \text{subject to } \|\underline{x}\| \leq \alpha . \end{array} \right\} \quad (0.6)$$

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

$$\left. \begin{array}{l} \|\underline{x}\| = \min \\ \text{subject to } \|\underline{Ax} - \underline{y}\| \leq \alpha . \end{array} \right\} \quad (0.7)$$

The question how to choose α is not simple to answer. It may be necessary to use some statistical tools to estimate α [17]. we observe that (0.6) and (0.7) are least squares problems with a quadratic constraint. The degree of the polynomial (n-1) may be chosen independently of the number of given data 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 (0.6) and (0.7) have a unique solution (if it exists) and we will show how to compute it efficiently.

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

$$\left. \begin{array}{l} \|\underline{Cx} - \underline{y}_2\| = \min \\ \text{subject to } \underline{Bx} = \underline{y}_1 \end{array} \right\} \quad (0.8)$$

where \underline{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 (0.1):

$$\left. \begin{array}{l} \|\underline{Cx} - \underline{y}_2\| = \min \\ \text{subject to } \|\underline{Bx} - \underline{y}_1\| \leq \alpha . \end{array} \right\} \quad (0.9)$$

In contrast to (0.6) and (0.7), (0.9) may not have a unique solution. As we shall see the solution is unique if and only if the nullspaces of \underline{C} and \underline{B} intersect trivially.

1. Basic Definitions and Remarks.

The following notation will be used:

Problem (P1):

$$\|\underline{Ax} - \underline{b}\| = \min \quad (1.1)$$

subject to

$$\|\underline{Cx} - \underline{d}\| \leq \alpha \quad (1.2)$$

If we have $\|\underline{Cx} - \underline{d}\| = \alpha$ instead of (1.2) we will refer to the problem as (P1E). Two special cases of (P1) will be considered:

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

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

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 λ)

$$L(\underline{x}, \lambda) = \|\underline{Ax} - \underline{b}\|^2 + \lambda \{ \|\underline{Cx} - \underline{d}\|^2 - \alpha^2 \}$$

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

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x} = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d} \quad (1.3)$$

$$\|\underline{Cx} - \underline{d}\|^2 = \alpha^2 \quad (1.4)$$

If the matrix $\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}$ is nonsingular, then we can define

$$f(\lambda) := \|\underline{C} x(\lambda) - \underline{d}\|^2 \quad (1.5)$$

where $x(\lambda)$ is the solution of (1.3). We will call f the "length

function". To determine a solution we have to solve the "secular equation"

$$f(h) = \alpha^2 \quad (1.6)$$

Finally we observe that for $\lambda > 0$ equations (1.3) and (1.4) are the normal equations of the least squares problem

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

A useful tool for the analysis of problems (P2) and (P3) is the singular value decomposition (SVD) [14] and its generalization (BSVD) [15] 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], \underline{A} and C are band matrices and (P1) may be solved efficiently without transformations as we shall show.

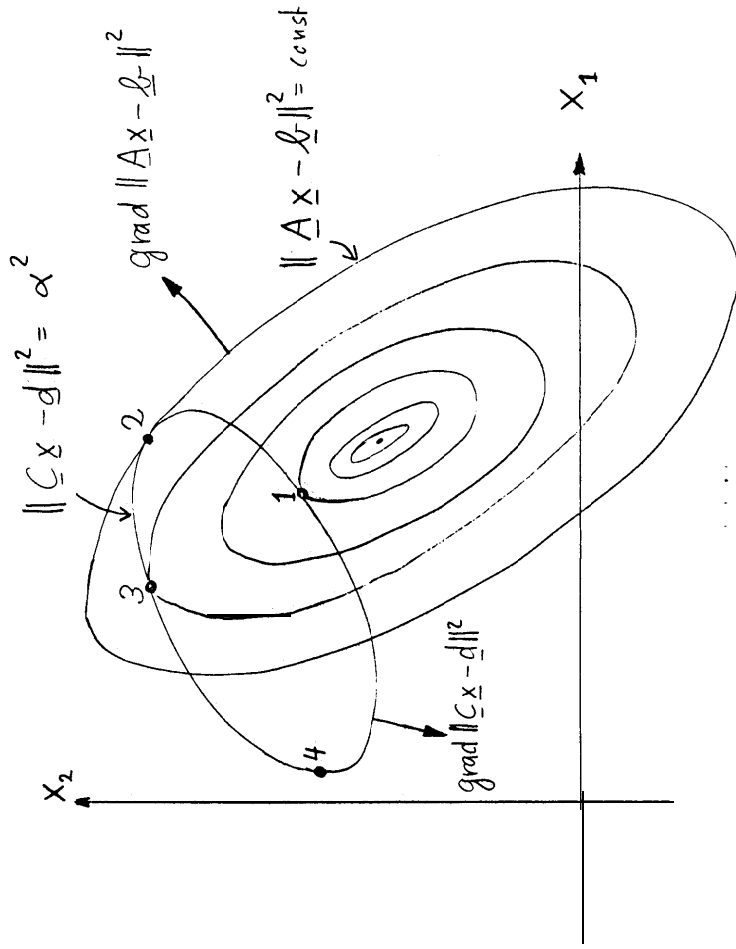
2. The Least Squares Problem with a Quadratic Constraint.

Let A be an $(m \times n)$ -matrix, C a $(p \times n)$ -matrix, b an m -vector, d a p -vector, and α a positive number.

We consider the problem to find an n -vector x so that

$$\left. \begin{aligned} \|\underline{Ax} - \underline{b}\| &= \min \\ \|\underline{Cx} - \underline{d}\| &= \alpha \end{aligned} \right\} \quad (PLE)$$

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



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

$$\left. \begin{aligned} \|\underline{Cx} - \underline{d}\|^2, \quad \|\underline{Ax} - \underline{b}\|^2 \\ (\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x} = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d} \\ \|\underline{Cx} - \underline{d}\|^2 = \alpha^2 \end{aligned} \right\} \quad (2.1)$$

are parallel, i.e., some λ exist so that

$$\text{grad } \|\underline{Ax} - \underline{b}\|^2 = -\lambda \text{grad } \|\underline{Cx} - \underline{d}\|^2.$$

If we rearrange this equation we have

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

which is equation (1.3) 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 λ is negative since the gradients have the same directions.

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

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

For the problems with inequality constraints (P1), (P2) and (P5) we have only to consider the case that $\underline{A}^+ \underline{b}$ is outside of $\|\underline{Cx} - \underline{d}\|^2$. If it is inside then $\underline{x} = \underline{A}^+ \underline{b}$ solves the problem.

2.1 Characterization of the Solution.

The solution of (PIE) is among the solutions (\underline{x}, λ) of the normal equations (see Section 1):

$$\left. \begin{aligned} (\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x} = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d} \\ \|\underline{Cx} - \underline{d}\|^2 = \alpha^2 \end{aligned} \right\} \quad (2.1)$$

The following theorem compares two solutions of these equations:

Theorem 2.1. If $(\underline{x}_1, \lambda_1)$ and $(\underline{x}_2, \lambda_2)$ are solutions of the normal equations (2.1), then

$$\|\underline{Ax}_2 - \underline{b}\|^2 - \|\underline{Ax}_1 - \underline{b}\|^2 = \frac{\lambda_1 - \lambda_2}{2} \|\underline{C}(\underline{x}_1 - \underline{x}_2)\|^2. \quad (2.2)$$

Proof. Since $(\underline{x}_1, \lambda_1)$, $(\underline{x}_2, \lambda_2)$ are solutions of (2.1) we have

$$\underline{A}^T \underline{Ax}_1 - \underline{A}^T \underline{b} = -\lambda_1 \underline{C}^T \underline{Cx}_1 + \lambda_1 \underline{C}^T \underline{d} \quad (2.3)$$

$$\underline{A}^T \underline{Ax}_2 - \underline{A}^T \underline{b} = -\lambda_2 \underline{C}^T \underline{Cx}_2 + \lambda_2 \underline{C}^T \underline{d} \quad (2.4)$$

\underline{x}_2^T (2.4) - \underline{x}_1^T (2.3) gives

$$\begin{aligned} \|\underline{Ax}_2\|^2 - \|\underline{Ax}_1\|^2 - \underline{b}^T \underline{A}(\underline{x}_2 - \underline{x}_1) \\ = \lambda_1 (\|\underline{Cx}_1\|^2 - \underline{d}^T \underline{Cx}_1) - \lambda_2 (\|\underline{Cx}_2\|^2 - \underline{d}^T \underline{Cx}_2) \end{aligned} \quad (2.5)$$

\underline{x}_2^T (2.3) - \underline{x}_1^T (2.4) gives

$$-\underline{b}^T (\underline{Ax}_2 - \underline{x}_1) = \lambda_1 (-\underline{x}_2^T \underline{C}^T \underline{Cx}_1 + \underline{d}^T \underline{Cx}_2) - \lambda_2 (-\underline{x}_1^T \underline{C}^T \underline{Cx}_2 + \underline{d}^T \underline{Cx}_1). \quad (2.6)$$

Observe that

$$\|Ax_2 - b\|^2 - \|Ax_1 - b\|^2 = \|Ax_2\|^2 - \|Ax_1\|^2 - 2b^T A(x_2 - x_1) \quad \cdot$$

So that if we add (2.5) - (2.6) we get

$$\begin{aligned} & \|Ax_2 - b\|^2 - \|Ax_1 - b\|^2 \\ &= \lambda_1 \{ \|Cx_1\|^2 - d^T Cx_1 - x_1^T C^T Cx_2 + d^T Cx_2 \} \\ & \quad - \lambda_2 \{ \|Cx_2\|^2 - d^T Cx_2 - x_1^T C^T Cx_1 + d^T Cx_1 \} \quad \cdot \end{aligned} \quad (2.7)$$

Now we have

$$\begin{aligned} & \|Cx_1 - d\|^2 = \|Cx_2 - d\|^2 = \alpha^2 \\ & \Rightarrow \|Cx_1\|^2 - 2d^T Cx_1 + \|d\|^2 = \|Cx_2\|^2 - 2d^T Cx_2 + \|d\|^2 \\ & \Rightarrow \|Cx_1\|^2 - d^T Cx_1 + d^T Cx_2 = \|Cx_2\|^2 - d^T Cx_2 + d^T Cx_1 \quad \cdot \end{aligned} \quad (2.8)$$

From (2.8) we conclude that the factors of λ_1 and λ_2 in (2.7) are the same. Therefore they also equal their arithmetic mean which is

$$\begin{aligned} & \frac{1}{2} \{ \|Cx_1\|^2 - 2x_1^T C^T Cx_2 + \|Cx_2\|^2 \} \\ &= \frac{1}{2} \|C(x_1 - x_2)\|^2 \quad \cdot \quad \square \end{aligned}$$

Corollary 2.1. The solution of (PIE) is the solution $x(\lambda)$ of the normal equations (2.1) with the largest λ .

Proof. From (2.2) we have that if $\lambda_1 > \lambda_2$, then

$$\|Ax_2 - b\|^2 > \|Ax_1 - b\|^2 \quad \cdot \quad \square$$

The next theorem gives a result very similar to Theorem 2.1.

Theorem 2.2. Assume (x_1, λ_1) and (x_2, λ_2) are solutions of the normal equations (2.1). Assume that $|\lambda_1| + |\lambda_2| \neq 0$. Then

$$\|Ax_2 - b\|^2 - \|Ax_1 - b\|^2 = \frac{\lambda_2^{-A}}{\lambda_2 + \lambda_1} \|A(x_1 - x_2)\|^2 \quad \cdot$$

Proof. From $(A^T A + \lambda C^T C)x = A^T b + \lambda C^T d$ we have

$$\lambda_1 C^T Cx_1 - \lambda_1 C^T d = -A^T Ax_1 + A^T b \quad (2.8)$$

$$\lambda_2 C^T Cx_2 - \lambda_2 C^T d = -A^T Ax_2 + A^T b \quad (2.9)$$

$\lambda_1 x_1^T (2.9) - \lambda_2 x_2^T (2.8)$ gives

$$\lambda_1 \lambda_2 \{ (x_2 - x_1)^T C^T d \} = (\lambda_2 - \lambda_1) x_1^T A^T Ax_2 + (\lambda_1 x_1 - \lambda_2 x_2)^T A^T b \quad \cdot \quad (2.10)$$

$\lambda_1 x_1^T (2.9) - \lambda_2 x_1^T (2.8)$ gives

$$\begin{aligned} & \lambda_1 \lambda_2 \{ \|Cx_2\|^2 - \|Cx_1\|^2 + (x_1 - x_2)^T C^T d \} \\ &= \lambda_2 \|Ax_1\|^2 - \lambda_1 \|Ax_2\|^2 + (\lambda_1 x_2 - \lambda_2 x_1)^T A^T b \quad \cdot \end{aligned} \quad (2.11)$$

Observe that

$$\begin{aligned} 0 &= \|Cx_2 - d\|^2 - \|Cx_1 - d\|^2 \\ &= \|Cx_1\|^2 - \|Cx_2\|^2 + 2(x_1 - x_2)^T C^T d \quad \cdot \end{aligned}$$

So that if we subtract (2.11) - (2.10) we get

But by Corollary 2.1,

$$\lambda_1 = A > -A = \lambda_2 \Rightarrow \|\underline{Ax}_2 - \underline{b}\|^2 > \|\underline{Ax}_1 - \underline{b}\|^2 .$$

Therefore $\lambda_1 + \lambda_2 \neq 0$ and we may divide in (2.13). \square

If we combine both results from Theorem 2.1 and Theorem 2.2 we have

Corollary 2.2. Let $(\underline{x}_1, \lambda_1)$ and $(\underline{x}_2, \lambda_2)$ be two solutions of the normal equations (2.1). If $|\lambda_1| + |\lambda_2| \neq 0$, then

$$-\frac{\lambda_1 + \lambda_2}{2} \|\underline{C}(\underline{x}_1 - \underline{x}_2)\|^2 = \|\underline{A}(\underline{x}_1 - \underline{x}_2)\|^2 . \quad (2.14)$$

From (2.14) we see immediately:

Corollary 2.3. The normal equations (2.1) have at most one solution $(\underline{x}^*, \lambda^*)$ with $A^* > 0$. For every other solution (\underline{x}, λ) we have

$$A < -\lambda^* .$$

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

Theorem 2.3. The solution \underline{x} of (P1E) is unique (if it exists) if

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

and

$$\lambda \neq -\mu_1$$

where (\underline{x}, λ) is a solution of the normal equation (2.1) and μ_1 is an eigenvalue of the generalized eigenvalue problem

$$\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0 .$$

.

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

or by rearranging

$$\begin{aligned} \lambda_1 \{ \|\underline{Ax}_2\|^2 - \underline{x}_2^T \underline{A}^T \underline{b} + \underline{x}_1^T \underline{A}^T \underline{b} - \underline{x}_1^T \underline{A}^T \underline{Ax}_2 \} \\ = \lambda_2 \{ \|\underline{Ax}_1\|^2 - \underline{x}_1^T \underline{A}^T \underline{b} + \underline{x}_2^T \underline{A}^T \underline{b} - \underline{x}_1^T \underline{A}^T \underline{Ax}_2 \} . \end{aligned} \quad (2.12)$$

Now the { } on the left hand side of (2.12) is

$$\begin{aligned} \frac{1}{2} \|\underline{Ax}_2\|^2 + \frac{1}{2} \{ \|\underline{Ax}_2\|^2 - 2\underline{x}_1^T \underline{A}^T \underline{Ax}_2 + \|\underline{Ax}_1\|^2 \} \\ - \frac{1}{2} \|\underline{Ax}_1\|^2 + \underline{x}_1^T \underline{A}^T \underline{b} - \underline{x}_2^T \underline{A}^T \underline{b} + \frac{1}{2} \|\underline{b}\|^2 - \frac{1}{2} \|\underline{b}\|^2 \\ = \frac{1}{2} \{ \|\underline{Ax}_2 - \underline{b}\|^2 - \|\underline{Ax}_1 - \underline{b}\|^2 + \|\underline{A}(\underline{x}_2 - \underline{x}_1)\|^2 \} . \end{aligned}$$

The right hand side of (2.12) simplifies analogously and by rearranging

we obtain:

$$(\lambda_1 + \lambda_2) \{ \|\underline{Ax}_2 - \underline{b}\|^2 - \|\underline{Ax}_1 - \underline{b}\|^2 \} = (\lambda_2 - \lambda_1) \|\underline{A}(\underline{x}_2 - \underline{x}_1)\|^2 . \quad (2.13)$$

We now prove by contradiction that $\lambda_1 + \lambda_2 \neq 0$. Assume $(\underline{x}_1, \lambda)$ and $(\underline{x}_2, -\lambda)$ were solutions of the normal equations (2.1) with $A > 0$. Then by (2.13) we would have

$$\begin{aligned} \|\underline{A}(\underline{x}_2 - \underline{x}_1)\|^2 &= 0 \\ \Rightarrow \underline{Ax}_2 &= \underline{Ax}_1 \\ \Rightarrow \|\underline{Ax}_2 - \underline{b}\|^2 &= \|\underline{Ax}_1 - \underline{b}\|^2 . \end{aligned}$$

Proof. The proof is by contradiction. Assume (x_1, λ_1) and (x_2, λ_2) are solutions of the normal equations (2.1) which also solve (PLE).

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

$$\|Ax_1 - b\|^2 = \|Ax_2 - b\|^2 = \min \|Ax - b\|^2.$$

Theorem 2.1 implies

$$\|C(x_1 - x_2)\|^2 = 0 \Rightarrow C(x_1 - x_2) = \underline{0}. \quad (2.15)$$

While Theorem 2.2 gives

$$\|A(x_1 - x_2)\|^2 = 0 \Rightarrow A(x_1 - x_2) = \underline{0}. \quad (2.16)$$

Now if $x_1 \neq x_2$ then (2.15) and (2.16) 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

$$(A^T A + \lambda C^T C)(x_1 - x_2) = \underline{0}.$$

If $x_1 \neq x_2$ then $\lambda = -\mu_i$ which is also a contradiction. Therefore we must have $\lambda_1 = \lambda_2$ and $x_1 = x_2$. \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 x .

We shall assume in the following (till the end of the paper) that $\text{NS}(A) \cap \text{NS}(C) = \{\underline{0}\}$. In Section 2.2 we shall analyze the existence of

solutions of the normal equations. We shall see that (PLE) has a solution if the constraint

$$\|Cx - d\| = \alpha^2$$

is feasible, i.e., if

$$\min_x \|Cx - d\|^2 = \|(C^+ C^+ - I)d\|^2 < \alpha^2.$$

2.2 The Solutions of the Normal Equations.

In this section we shall discuss the solutions of

$$(A^T A + \lambda C^T C)x = A^T b + A C^T d \quad (2.17a)$$

$$\|Cx - d\|^2 = \alpha^2. \quad (2.17b)$$

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

If $\lambda \neq -\mu_i$ where $\mu_i \geq 0$ is an eigenvalue of the eigenvalue

problem $\det(A^T A - \mu C^T C) = 0$ then

$$\underline{x}(\lambda) = (A^T A + \lambda C^T C)^{-1}(A^T b + \lambda C^T d). \quad (2.18)$$

If we now choose λ so that the secular equation is satisfied:

$$f(\lambda) := \|Cx(\lambda) - d\|^2 = \alpha^2 \quad (2.19)$$

then $(\underline{x}(\lambda), \lambda)$ is a solution of (2.17). We first discuss some properties of the length function f :

Lemma 2.1. If f is defined by (2.19), (2.18), then

- (1) f is a rational function defined on $\mathbb{R} - \{-\mu_i \mid \det(A^T A - \mu_i C^T C) = 0\}$ J
- (2) $f(\lambda) \equiv \|d\|^2 = \text{const}$ if $A^T b = \underline{0}$ and $C^T d = \underline{0}$,
- (3) f has at least one and at most n poles for same $\lambda = -\mu_i$,

$$(4) f(\lambda) > 0, \quad \lim_{\lambda \rightarrow +\infty} f(\lambda) = \|(\underline{C} \underline{C}^+ - \underline{I})\underline{d}\|^2.$$

$$(5) f(\lambda) < 0 \text{ for } 0 < \lambda < \infty.$$

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

$$\left. \begin{aligned} \underline{U}^T \underline{A} \underline{X} = \underline{D}_A &= \text{diag}(\alpha_1, \dots, \alpha_n), \alpha_i \geq 0, \\ \underline{V}^T \underline{C} \underline{X} = \underline{D}_C &= \text{diag}(\gamma_1, \dots, \gamma_q), \gamma_i \geq 0, \quad q = \min(n, p) \end{aligned} \right\} \quad (2.20)$$

where \underline{U} ($m \times m$) \underline{V} ($p \times p$) are orthogonal, \underline{X} ($n \times n$) is nonsingular, and \underline{D}_A , \underline{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

$$\|\underline{Ax} - \underline{b}\|^2 = \left\| \begin{pmatrix} \underline{A} \\ \underline{0} \end{pmatrix} \underline{x} - \begin{pmatrix} \underline{b} \\ \underline{0} \end{pmatrix} \right\|^2 \quad \text{zero rows in } \underline{A} \text{ and zero elements in } \underline{b}$$

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 $\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0$ are given by

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

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

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

By substituting (2.20) into (2.16) we have

$$\underline{X}^{-T} (\underline{D}_A^T \underline{D}_A + \lambda \underline{D}_C^T \underline{D}_C) \underline{X}^{-1} \underline{x} = \underline{X}^{-T} \underline{D}_A^T \underline{U} \underline{b} + \lambda \underline{X}^{-T} \underline{D}_C^T \underline{V} \underline{d}$$

and putting

$$\underline{y} := \underline{X}^{-1} \underline{x}, \quad \underline{c} := \underline{U} \underline{b}, \quad \underline{e} := \underline{V} \underline{d}$$

we get

$$(\underline{D}_A^T \underline{D}_A + \lambda \underline{D}_C^T \underline{D}_C) \underline{y} = \underline{D}_A^T \underline{c} + \lambda \underline{D}_C^T \underline{e}. \quad (2.21)$$

Now

$$f(\lambda) = \|\underline{C} \underline{x}(\lambda) - \underline{d}\|^2 = \|\underline{D}_C \underline{y}(\lambda) - \underline{e}\|^2$$

which is in components:

$$\begin{aligned} f(\lambda) &= \sum_{i=1}^r \left(\gamma_i \frac{\alpha_i c_i + \lambda \gamma_i e_i}{\alpha_i^2 + \lambda \gamma_i^2} - e_i \right)^2 + \sum_{i=r+1}^p e_i^2 \\ f(\lambda) &= \sum_{i=1}^r \left(\alpha_i \frac{\gamma_i c_i - \alpha_i e_i}{\alpha_i^2 + \lambda \gamma_i^2} \right)^2 + \sum_{i=r+1}^p e_i^2. \end{aligned} \quad (2.22)$$

Obviously $f(\lambda)$ is a rational function in λ with sane 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 $\underline{A}^T \underline{b} = \underline{0}$ and $\underline{C}^T \underline{d} = 0$ then from the definition of f (2.19), (2.18) we have

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

From (2.19) we see that $f(\lambda) > 0$.

To prove the limit property in (4) we observe that the solution $\underline{x}(\lambda)$ of (2.17a) converges for $\lambda \rightarrow \infty$ to a solution of $\underline{C}^T \underline{C} \underline{x} = \underline{C}^T \underline{d}$. Through $\underline{x}(\infty)$ may not be $\underline{C}^+ \underline{d}$ it has the same residual. Therefore:

$$\| \underline{C} \underline{x}(\infty) - \underline{d} \|^2 = \| \underline{C} \underline{c}^T \underline{d} - \underline{d} \|^2 .$$

This property can of course also be seen from the representation (2.22).

To prove (5) we differentiate (2.19) and (2.17a)

$$f'(h) = 2 \underline{x}'(\lambda)^T \underline{C}^T (\underline{C} \underline{x}(\lambda) - \underline{d}) \quad (2.23)$$

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x}'(\lambda) = -\underline{C}^T (\underline{C} \underline{x}(h) - \underline{d}) \quad (2.24)$$

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

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) = \underline{R}_\lambda^T \underline{R}_\lambda \quad (2.25)$$

Using (2.25) and (2.24) in (2.23) gives

$$f'(A) = -2 \| \underline{R}_\lambda^{-T} \underline{C}^T (\underline{C} \underline{x}(\lambda) - \underline{d}) \|^2 < \dots \quad (2.26)$$

Every solution λ of the secular equation (2.19) with $A \neq -\mu_i$ defines $\underline{x}(\lambda)$ (2.18) and $(\underline{x}(\lambda), \lambda)$ is a solution of the normal equations (2.17a,b).

But there might be more solutions for some $A = -\mu_i$. In this case the matrix of (2.17a) is singular. If a solution should exist, the system must be consistent. This is especially true if $\underline{A} \underline{b} = \underline{C}^T \underline{d} = 0$ then $\underline{x}(\lambda)$ is an eigenvector to $\lambda = -\mu_i$. The next lemma describes the condition for consistency:

Lemma 2.2. Let μ_i be an eigenvalue of $\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0$. Then

$$(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) \underline{x} = \underline{A}^T \underline{b} - \mu_i \underline{C}^T \underline{d} \quad (2.27)$$

is a consistent system if one of the conditions holds:

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

(ii) Let $J = \{j \mid 1 \leq j \leq k, \frac{\alpha_j^2}{\gamma_j} = \mu_i\}$, then $\alpha_j(\underline{c}_j^T \underline{y}_j - \alpha_j \underline{e}_j) = 0$

for $j \in J$, where all variables are defined in the proof of Lemma 2.1.

Proof. We prove first (ii) using the BSVD to transform (2.27).

$$\underline{A} = \underline{U} \underline{D} \underline{X}^{-1}, \quad \underline{D}_A = \text{diag}(\alpha_1, \dots, \alpha_n)$$

$$\underline{C} = \underline{V} \underline{D}_C \underline{X}^{-1}, \quad \underline{D}_C = \text{diag}(\gamma_1, \dots, \gamma_r, 0, \dots, 0) .$$

With $\underline{y} := \underline{X}^{-1} \underline{x}$, $\underline{e} = \underline{U}^T \underline{b}$, $\underline{e}_j = \underline{V}^T \underline{d}_j$, $(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x} = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d}$ becomes

$$\left. \begin{aligned} (\alpha_k^2 + \lambda \gamma_k^2) y_k &= \alpha_k c_k + \lambda \gamma_k e_k, \quad k = 1, \dots, r \\ \alpha_k^2 y_k &= \alpha_k c_k, \quad k = r+1, \dots, n . \end{aligned} \right\} \quad (2.28)$$

Now observe that $r = \text{rank of } \underline{C}$ and because $\text{NS}(\underline{A}) \cap \text{NS}(\underline{C}) = \{0\}$,

$\alpha_k \neq 0$, $k = r+1, \dots, n$. If $\lambda = -\mu_i = -\alpha_i^2 / \gamma_i^2$, then for every equation j of (2.28) with

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

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

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

$$\Rightarrow \alpha_j (\gamma_j e_j - \gamma_j e_j) = 0 \quad \text{for } j \in J .$$

We note that the solution of (2.28) is given by

$$\tilde{y}_k = \begin{cases} (\alpha_k^c \gamma_k e_k - \mu_1 \gamma_k e_k) / (\alpha_k^2 - \mu_1 \gamma_k^2), & k \neq J \\ \text{arbitrary}, & k \in J \\ c_k / \alpha_k, & k = r+1, \dots, n \end{cases} \quad (2.29)$$

Furthermore

$$\tilde{y}_k := \lim_{\lambda \rightarrow -\mu_1} y_k(\lambda) = \begin{cases} (\alpha_k^c \gamma_k e_k - \mu_1 \gamma_k e_k) / (\alpha_k^2 - \mu_1 \gamma_k^2), & k \neq J \\ e_k / \gamma_k, & k \in J \\ c_k / \alpha_k, & k = r+1, \dots, n \end{cases} \quad (2.30)$$

From (2.30) we conclude that

$$\lim_{A \rightarrow D} x(\lambda) = \lim_{A \rightarrow -\mu_1} X_{-} y(\lambda) = X \tilde{y}$$

exists and therefore

$$\lim_{\lambda \rightarrow -\mu_1} \|\underline{C} x(\lambda) - \underline{d}\|^2$$

exists. This implies that f has no pole for $A = -\mu_1$.

Conversely if f has not a pole for $\lambda = -\mu_1$ then also

$$\lim_{A \rightarrow -\mu_1} y(\lambda)$$

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

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

Theorem 2.4. Let f be the length function defined by (2.19), (2.18). If $f(A) = \alpha^2$ and $\det(A_{-}^T A + \lambda C_{-}^T C) \neq 0$ then there exists a unique $x(\lambda)$ so that $(x(\lambda), \lambda)$ solves the normal equations. If

$$\det(A_{-}^T A - \mu_1 C_{-}^T C) = 0 \quad \text{and} \quad \lim_{\lambda \rightarrow -\mu_1} f(A) \leq \alpha^2 \quad (2.31)$$

then there exist a $x(-\mu_1)$ so that $(x(-\mu_1), -\mu_1)$ solves the normal equations, but $x(-\mu_1)$ is only unique if $\lim_{A \rightarrow -\mu_1} f(A) = \alpha^2$.

Proof. We have only to prove (2.31). We use the terminology defined in the proof of Lemma 2.2. For $A = -\mu_1$ the general solution y of the transformed normal equations (2.28) is given by (2.29). We have to see if we can satisfy the equation

$$\|\underline{D} y - \underline{e}\| = \alpha^2. \quad (2.32)$$

with this solution y . In components (2.32) is

$$\sum_{k \neq J} (\gamma_k y_k - e_k)^2 + \sum_{k \in J} (\gamma_k y_k - e_k)^2 + \sum_{k=r+1}^n e_k^2 = \alpha^2. \quad (2.33)$$

We can choose the components of y_k , $k \in J$ arbitrarily. From (2.30) we see that

$$\begin{aligned} \lim_{\lambda \rightarrow -\mu_1} f(A) &= \lim_{\lambda \rightarrow -\mu_1} \|\underline{D} y(\lambda) - \underline{e}\|^2 = \|\underline{D} \tilde{y} - \underline{e}\|^2 \\ \lim_{\lambda \rightarrow -\mu_1} f(A) &= \sum_{k \neq J} (\gamma_k y_k - e_k)^2 + \sum_{k=r+1}^n e_k^2. \end{aligned} \quad (2.34)$$

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

From (2.34) we see that if

$$\lim_{A \rightarrow -\mu_1} f(A) > \alpha^2, \dots$$

we can not determine a solution $\underline{y}(\lambda)$ that satisfies (2.32). If

$$\lim_{h \rightarrow -\mu_1} f(A) = \alpha^2$$

then the unique solution is $\tilde{\underline{y}}$ (2.30). If

$$\lim_{A \rightarrow -\mu_1} f(A) < \alpha^2$$

then we can choose \underline{y}_k , $k \in J$ so that

$$\sum_{k \in J} (\underline{y}_k^T \underline{y}_k - e_k)^2 = \alpha^2 - \lim_{A \rightarrow -\mu_1} f(A) \quad (2.35)$$

and $\underline{y}(-\mu_1)$ is not unique. \square

We remark that the components of $\underline{y} : \underline{y}_k$, $k \in J$ are the (possibly) non-zero components of the eigenvector $\underline{y}_{-\mu_1}$ of

$$\begin{pmatrix} D^T & D \\ -A & -A \end{pmatrix} \underline{y}_{-\mu_1} - \underline{y}_{-\mu_1} \begin{pmatrix} D^T & D \\ -C & -C \end{pmatrix} \underline{y}_{-\mu_1} = \underline{0} .$$

Therefore a solution for the case where $\lim_{A \rightarrow -\mu_1} f(A) < \alpha^2$ can also be

written in the following way. Take any eigenvector $\underline{y}_{-\mu_1}$ belonging to μ_1

and determine ρ such that

$$\|D_C(\underline{y} + \rho \underline{y}_{-\mu_1}) - \underline{d}\| = \alpha^2 .$$

Then $(\underline{y} + \rho \underline{y}_{-\mu_1})$ is the solution of the normal equations. /

2.3 The Solution for the Equality Constraint.

We now consider problem (PLE)

$$\begin{aligned} \|\underline{Ax} - \underline{b}\|^2 &= \min \\ \text{subject to } \|\underline{Cx} - \underline{d}\| &= \alpha^2 . \end{aligned}$$

We continue to assume that the nullspaces of \underline{A} and \underline{C} intersect trivially.

To find the solution we first have to compute the rightmost solution \underline{A}

of the secular equation. If $\lambda^* > 0$ then $\underline{x}(\lambda^*)$ is the unique solution

of (PLE). If $\lambda^* < 0$ then we have to compute the smallest eigenvalue μ_r of the generalized eigenvalue problem

$$\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0$$

and an eigenvector \underline{x}_r . If $\lambda^* > -\mu_r$ then again by Theorem 2.1

$\underline{x}(\lambda^*)$ is the unique solution. However if $\lambda^* \leq -\mu_r$ then a solution has the form

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

where we have to determine ρ such that

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

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

2.4 The Solution for the Inequality Constraint.

We are now ready to solve (P1):

$$\|\underline{Ax} - \underline{b}\|^2 = \min \quad (2.36)$$

subject to

$$\|\underline{Cx} - \underline{d}\|^2 \leq \alpha^2 \quad (2.37)$$

The weaker condition (2.37) simplifies the problem. Let

$\mathbf{M} = \{\underline{x} \mid \|\underline{Ax} - \underline{b}\| = \min\}$ denote the set of solutions of the unconstrained problem (2.36). If for some $\mathbf{x} \in \mathbf{M}$ we have $\|\underline{Cx} - \underline{d}\|^2 \leq \alpha^2$ then this \mathbf{x} is a solution.

If $\mathbf{M} \neq \{\underline{A}^+ \underline{b}\}$ then $\mu_1 = 0$ is an eigenvalue of $\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0$. Since 0 is never a pole of f (the normal equations (2.27) are consistent for $\lambda = 0$) we can define the number

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

Let $\alpha_{\min}^2 := \|(C^+ - I)\underline{d}\|^2 = \lim_{\lambda \rightarrow \infty} f(\lambda)$. Then (P1) has no solution if

$\alpha < \alpha_{\min}$. If $a \geq \alpha_{\max}$ then

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

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

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

In general we will have

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

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

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

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

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

$$\left. \begin{array}{l} \min \|\underline{Cx} - \underline{d}\| \\ \|\underline{Ax} - \underline{b}\| = \min \end{array} \right\} \text{if } \alpha = \alpha_{\max}$$

$$\left. \begin{array}{l} \min \|\underline{Ax} - \underline{b}\|^2 \\ \|\underline{Cx} - \underline{d}\| = \min \end{array} \right\} \text{if } \alpha = \alpha_{\min}$$

Both are of course only nontrivial if \underline{A} respectively \underline{C} is rank deficient.

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

- (a) **start with** $\lambda > 0$
- (b) **while** not converged **do**
begin
 solve the least square problem

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

correct \underline{A} to solve $f(\lambda) = \alpha^2$ where

$$f(\lambda) = \|\underline{Cx}_{\lambda} - \underline{d}\|^2.$$

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 square 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 $\underline{A} = \underline{I}$ and $\underline{C} = \underline{I}$. 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 $\underline{C} = \underline{I}$.

3. The Relaxed Least Squares Problem.

We consider in this section the special case of a least squares problem with a quadratic constraint where $\underline{C} = \underline{I}$:

$$\begin{aligned} \|\underline{Ax} - \underline{b}\| &= \min \\ \|\underline{x}\| &= \alpha \end{aligned} \tag{P2E}$$

The problem with the inequality constraint $\|\underline{x}\| \leq \alpha$ (P2) was called by Rutishauser [35] the relaxed least squares problem. Problem (P2E) can be interpreted to find the minimum of a quadratic form $\|\underline{Ax} - \underline{b}\|^2$ on the sphere $\|\underline{x}\|^2 = \alpha^2$, which is a special case of finding the stationary values of a quadratic form on a sphere. This problem has been analyzed by Forsythe and Golub [10] in 1965. The two authors proved that if (λ_1, λ_1) and (λ_2, λ_2) are two solutions of the normal equations then

$$\lambda_1 > \lambda_2 \Rightarrow \|\underline{Ax}_2 - \underline{b}\|^2 > \|\underline{Ax}_1 - \underline{b}\|^2 .$$

Their proof is rather complicated. Kahan [25] gave an elementary proof and stated the theorem

$$\|\underline{Ax}_2 - \underline{b}\|^2 - \|\underline{Ax}_1 - \underline{b}\|^2 = \frac{\lambda_1 - \lambda_2}{2} \|\underline{x}_2 - \underline{x}_1\| \tag{3.1}$$

which is Theorem 2.1 for $\underline{C} = \underline{I}$. Unfortunately this theorem was never published. In 1972 Sjöqvoll [39] gave a simpler proof than Forsythe - Golub but did not quite obtain Kahan's result. He showed that

$$\|\underline{Ax}_2 - \underline{b}\|^2 - \|\underline{Ax}_1 - \underline{b}\|^2 = (\lambda_1 - \lambda_2) (\alpha^2 - \underline{x}_1^T \underline{x}_2)$$

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

3.1 Results from the General Theory.

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

$$(\underline{A}^T \underline{A} + \underline{A} \underline{I}) \underline{x} = \underline{A}^T \underline{b} \tag{3.2}$$

$$\|\underline{x}\|^2 = \alpha^2 . \tag{3.3}$$

Theorem 2.1 gives us equation (3.1). This means that the solution of

(P2E) is the solution (\underline{x}, λ) of (3.2), (3.3) with the largest λ .

Since $\text{NS}(\underline{A}) \cap \text{NS}(\underline{I}) = \{\underline{0}\}$ we have from Theorem 2.3 that the solution of (P2E) is unique if $\lambda \neq -\sigma_1^2$ (where σ_1 is a singular value of \underline{A}).

The length function f can be written

$$f(\lambda) = \|(\underline{A}^T \underline{A} + \underline{A} \underline{I})^{-1} \underline{A}^T \underline{b}\|^2 . \tag{3.4}$$

Using the singular value decomposition of \underline{A} [14]

$$\underline{A} = \underline{U} \underline{\Sigma} \underline{V}^T$$

where

$\underline{U}, \underline{V}$ orthogonal

$$\underline{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p) , \quad p = \min(m, n) .$$

The right hand side of (3.4) becomes

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

In components this is

$$f(\lambda) = \sum_{i=1}^p \left(\frac{\sigma_i c_i}{\sigma_i^2 + \lambda} \right)^2 . \tag{3.5}$$

If $\underline{A}^T \underline{b} = \underline{0}$ then $c = \underline{0}$ and $f(\lambda) \equiv 0$. Theorem 2.4 and the fact that $\lim_{\lambda \rightarrow \infty} f(\lambda) = 0$ implies that for every α 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 $\|\underline{x}\|^2 \leq \alpha^2$ (P2)

we have either

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

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

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

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

3.2 The Dual Normal Equations.

Theorem 3.1.

(i) Let (\underline{x}, λ) be a solution of the primal normal equation

$$(\underline{A}^T \underline{A} + \underline{A} \underline{I}) \underline{x} = \underline{A}^T \underline{b} \tag{3.6}$$

$$\|\underline{x}\|^2 = \alpha^2 \tag{3.7}$$

with $\lambda \neq 0$. Then (z, λ) with

$$z = \frac{1}{\lambda} (\underline{A} z - \underline{b}) \tag{3.8}$$

is a solution of the dual normal equations

...

$$(\underline{AA}^T + \lambda \underline{I})z = -b \tag{3.9}$$

$$\|\underline{A}^T z\|^2 = \alpha^2 \tag{3.10}$$

(ii) Let (z, λ) be a solution of the dual normal equations (3.8),

(3.9). Then (x, λ) with

$$x = -\underline{A}^T z \tag{3.11}$$

is a solution of the prime normal equations (3.6), (3.7).

Proof.

$$\begin{aligned} (i) \quad (\underline{AA}^T + \lambda \underline{I}) \frac{1}{\lambda} (\underline{Ax} - b) &= \frac{1}{\lambda} \underbrace{\underline{A}(\underline{AA}^T + \lambda \underline{I})z}_{\underline{A}^T b} - \frac{1}{\lambda} \underline{AA}^T b - b \\ &= -b \end{aligned}$$

Furthermore

$$\|\underline{A}^T z\|^2 = \|\underline{A}^T \frac{1}{\lambda} (\underline{Ax} - b)\|^2 = \|\underline{x}\|^2 = \alpha^2$$

$$(ii) \quad (\underline{A}^T \underline{A} + \lambda \underline{I})(-\underline{A}^T z) = -\underline{A}^T \underbrace{(\underline{AA}^T + \lambda \underline{I})z}_{-b}$$

And

$$\|-\underline{A}^T z\|^2 = \|\underline{x}\|^2 = \alpha^2 \quad \square$$

Theorem 3.1 shows that we may simplify computations to solve (P2) or (P2E). Since \underline{A} is an $(m \times n)$ matrix one of the linear systems (3.6) or (3.8) will be smaller and it may be more economical to iterate with the smaller.

But Theorem 3.1 also gives theoretical equations.

Corollary 3.1. Let A be an $(m \times n)$ matrix and $A \neq -\sigma_1^2$ where σ_1

is a singular value of \underline{A} . Then

$$(\underline{A}^T \underline{A} + \lambda \underline{I})^{-1} \underline{A}^T = \underline{A}^T (\underline{AA}^T + \lambda \underline{I})^{-1} \tag{3.12}$$

$$(\underline{AA}^T + \lambda \underline{I})^{-1} = \frac{1}{\lambda} (\underline{I} - \underline{A}(\underline{A}^T \underline{A} + \lambda \underline{I})^{-1} \underline{A}^T) \tag{3.13}$$

Proof. From Theorem 3.1 (3.12) follows by equating \underline{x} and (3.13) by equating \underline{z} for the dual and primal equations. \square

As an application we consider (3.13) for $n = 1$ (i.e., \underline{A} is a vector) and $\lambda = 1$:

$$(\underline{I} + \underline{aa}^T)^{-1} = \underline{I} - \frac{1}{(\underline{a}^T \underline{a} + 1)} \underline{aa}^T \tag{3.14}$$

We observe that Corollary 3.1 is a special case of the Sherman-Morrison-Woodbury formula:

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

$$(\underline{A} + \underline{U} \underline{V}^T) \underline{x} = \underline{b} \tag{3.15}$$

is also the solution of the augmented system

$$\underline{Ax} + \underline{Uy} = \underline{b} \tag{3.16a}$$

$$\underline{V}^T \underline{x} - \underline{y} = \underline{0} \tag{3.16b}$$

Now assume that \underline{A} is nonsingular. Then we have from (3.16a)

$$\underline{x} = \underline{A}^{-1} \underline{b} - \underline{A}^{-1} \underline{Uy}$$

Introducing this in (3.16b) we get

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)

$$\begin{pmatrix} \underline{A} \\ \sqrt{\lambda} \underline{C} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{b} \\ \sqrt{\lambda} \underline{d} \end{pmatrix} \quad (3.17)$$

$$\|\underline{C}\underline{x} - \underline{d}\|^2 = \alpha^2. \quad (3.18)$$

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

$$\underline{C} = \underline{Q} \begin{pmatrix} \underline{R} \\ \underline{0} \end{pmatrix} \underline{P}^T, \quad \underline{R} = \begin{pmatrix} \square & \\ & \square \end{pmatrix}.$$

Observe that $\|\underline{C}\underline{x} - \underline{d}\| = \|\underline{R} \underline{P}^T \underline{x} - \underline{Q}^T \underline{d}\|$. Therefore we can replace C by \underline{R} , \underline{d} by $\underline{Q}^T \underline{d}$, \underline{A} by $\underline{A}\underline{P}$, and \underline{x} by $\underline{y} = \underline{P}^T \underline{x}$ in (3.17), (3.18). The problem for \underline{y} has now a matrix \underline{C} ($p \times n$) with rank of $\underline{C} = p$.

After it is solved we obtain $\underline{x} = \underline{P}\underline{y}$. The same preprocessing we can apply to \underline{A} and \underline{b} and so assume that $\text{rank}(\underline{A}) = m \leq n$.

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

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

which transforms the problem to

$$\begin{pmatrix} \underline{A}' \\ \sqrt{\lambda} \underline{I} \end{pmatrix} \underline{x}' \approx \begin{pmatrix} \underline{b} \\ \sqrt{\lambda} \underline{d} \end{pmatrix} \quad \|\underline{x}' - \underline{d}\|^2 = \alpha^2.$$

$$\begin{aligned} \underline{y} &= (\underline{I} + \underline{V}^T \underline{A}^{-1} \underline{U})^{-1} \underline{V}^T \underline{A}^{-1} \underline{b} \\ \Rightarrow \underline{x} &= \underline{A}^{-1} \underline{b} - \underline{A}^{-1} \underline{U}(\underline{I} + \underline{V}^T \underline{A}^{-1} \underline{U})^{-1} \underline{V}^T \underline{A}^{-1} \underline{b}. \end{aligned}$$

But from (3.15) and (3.16b) we have also

$$\begin{aligned} \underline{x} &= (\underline{A} + \underline{U} \underline{V}^T)^{-1} \underline{b} \\ \underline{y} &= \underline{V}^T (\underline{A} + \underline{U} \underline{V}^T)^{-1} \underline{b}. \end{aligned}$$

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

$$\begin{aligned} \underline{V}^T (\underline{A} + \underline{U} \underline{V}^T)^{-1} &= (\underline{I} + \underline{V}^T \underline{A}^{-1} \underline{U})^{-1} \underline{V}^T \underline{A}^{-1} \\ (\underline{A} + \underline{U} \underline{V}^T)^{-1} &= \underline{A}^{-1} (\underline{I} - \underline{U} (\underline{I} + \underline{V}^T \underline{A}^{-1} \underline{U})^{-1} \underline{V}^T \underline{A}^{-1}). \end{aligned} \quad (\text{Shermann - Morrison - Woodbury formulas})$$

we note that for $\lambda := \lambda \underline{I}$ and $\underline{U} = \underline{V} := \underline{A}$ we get the equations (3.12) and (3.13).

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

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

The dual equations have no solution if $\lambda = 0$. We see this clearly because of the factor $-\frac{1}{\sqrt{\lambda}}$, but also from (3.9) since it is clear that for $\lambda = 0$ a solution \underline{z} of (3.9) may not exist for arbitrary \underline{b} and $m > n$.

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

$$C^T = \begin{pmatrix} V_1 & V_2 \\ 0 & 0 \end{pmatrix} \begin{matrix} R \\ 0 \\ 0 \end{matrix} \quad \begin{matrix} \text{p} \\ \text{n-p} \\ \text{n-p} \end{matrix} \quad (3.19)$$

with nonsingular triangular matrix R .

We change variables:

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

and (3.17) becomes:

$$\begin{pmatrix} AV_1 & AV_2 \\ \sqrt{\lambda} R^T & 0 \end{pmatrix} \begin{pmatrix} \underline{y}_1 \\ \underline{y}_2 \end{pmatrix} \approx \begin{pmatrix} \underline{b} \\ \sqrt{\lambda} V_2^T \underline{d} \end{pmatrix}. \quad (3.21)$$

Whereas (3.18) is now

$$\|R^T V_1 \underline{y}_1 - \underline{d}\|^2 = \alpha^2. \quad (3.22)$$

Now we perform another QR decomposition

$$\underline{A} \underline{V}_2 = (\underline{Q}_1, \underline{Q}_2) \begin{pmatrix} U \\ 0 \end{pmatrix} \quad \begin{matrix} \text{p-n} \\ \text{p-n} \end{matrix} \quad (3.23)$$

Since we assumed that A has rank m and V_2 is orthogonal, U is a non-singular upper triangular matrix. Now (3.21) is the same problem

$$\begin{pmatrix} \underline{Q}_1^T A V_1 & \underline{U} \\ \underline{Q}_2^T A V_1 & 0 \\ \sqrt{\lambda} R^T & 0 \end{pmatrix} \begin{pmatrix} \underline{y}_1 \\ \underline{y}_2 \end{pmatrix} \approx \begin{pmatrix} \underline{Q}_1^T \underline{b} \\ \underline{Q}_2^T \underline{b} \\ \sqrt{\lambda} V_2^T \underline{d} \end{pmatrix}. \quad (3.25)$$

as:

$$\underline{Q}_e(\underline{x}) = \|\underline{Ax} - \underline{b}\|^2 + \epsilon^2 \|\underline{x}\|^2 \quad (3.27)$$

will give a shorter solution \underline{x}_{-e} which however does not minimize (3.27) any more. The process of relaxing can be described as follows: The solution of (3.26) is the solution of the least squares problem

$$\underline{Ax} \approx \underline{b}. \quad (3.28)$$

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

$$\begin{pmatrix} \underline{Q}_2^T A V_1 \\ \sqrt{\lambda} R^T \end{pmatrix} \underline{y}_1 \approx \begin{pmatrix} \underline{Q}_2^T \underline{b} \\ \sqrt{\lambda} V_2^T \underline{d} \end{pmatrix} \quad (3.25)$$

with (3.22) as constraint. The final change of variables

$$\underline{y}_1 := R^T \underline{y}_1$$

gives the desired problem (P2).

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

In [35] Rutishauser remarked that if we minimize

$$Q(\underline{x}) = \|\underline{Ax} - \underline{b}\|^2 \quad (3.26)$$

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

Replacing (3.26) by

$$\underline{Q}_e(\underline{x}) = \|\underline{Ax} - \underline{b}\|^2 + \epsilon^2 \|\underline{x}\|^2 \quad (3.27)$$

will give a shorter solution \underline{x}_{-e} which however does not minimize (3.27) any more. The process of relaxing can be described as follows: The solution of (3.26) is the solution of the least squares problem

$$\underline{Ax} \approx \underline{b}. \quad (3.28)$$

The relaxed solution \underline{x}_ϵ of (3.27) is obtained by choosing some $\epsilon > 0$ and solving

$$\begin{pmatrix} \underline{A} \\ \epsilon \underline{I} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{b} \\ \underline{0} \end{pmatrix} . \quad (3.29)$$

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

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

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

Rutishauser [35] defines the double relaxed solution $\underline{x}_{d\epsilon}$ to be the solution of

$$\begin{pmatrix} \underline{A}^T \underline{A} + \epsilon^2 \underline{I} + \epsilon (\underline{A}^T \underline{A} + \epsilon^2 \underline{I})^{-1} \end{pmatrix} \underline{x} = \underline{A}^T \underline{b} . \quad (3.30)$$

Observe that $\underline{x}_{d\epsilon}$ is obtained by relaxing the normal equations of the relaxed solution:

$$\begin{pmatrix} \underline{A}^T \underline{A} + \epsilon^2 \underline{I} \\ \epsilon \underline{I} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{A}^T \underline{b} \\ \underline{0} \end{pmatrix} . \quad (3.31)$$

If we put $\underline{B} := \underline{A}^T \underline{A} + \epsilon^2 \underline{I}$ then the normal equations of (3.31) are

$$\begin{pmatrix} \underline{B} \underline{B} + \epsilon^2 \underline{I} \end{pmatrix} \underline{x} = \underline{B}^T \underline{A}^T \underline{b} . \quad (3.32)$$

But \underline{B} is symmetric and non-singular, therefore we may multiply 6.32 from left by \underline{B}^{-1} and we get

$$(\underline{B} + \epsilon^2 \underline{B}^{-1}) \underline{x} = \underline{A}^T \underline{b}$$

which is (3.30). From (3.30) we have

Lemma 3.1. The double relaxed solution $\underline{x}_{d\epsilon}$ minimizes

$$Q_d(\underline{x}) = \|(\underline{A} \underline{A} + \epsilon \underline{I}) \underline{x} - \underline{A}^T \underline{b}\|^2 + \epsilon \|\underline{x}\|^2 .$$

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

$$\underline{A} = \underline{u} \underline{\Sigma} \underline{v}^T , \underline{U} , \underline{V} \text{ orthogonal}$$

$$\underline{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n) , \sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0$$

be the singular value decomposition of \underline{A} . If the **datas** and the computation were exact then the shortest solution of

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

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

$$\underline{Y}_i = \begin{cases} \frac{c_i}{\sigma_i} & \text{if } \sigma_i \neq 0 \\ 0 & \text{if } \sigma_i = 0 . \end{cases} \quad (c := \underline{U}^T \underline{b})$$

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

$$\tilde{\underline{Y}}_i = \begin{cases} \frac{c_i}{\sigma_i} & \text{if } \sigma_i \geq \tau \\ 0 & \text{if } \sigma_i < \tau \end{cases} . \quad (3.33)$$

We will call $\tilde{\underline{Y}}_i$ the "ideal solution".

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

$$k(u) := \begin{cases} 0 & \text{if } \sigma < \tau \\ \frac{1}{\sigma} & \text{if } \sigma > \tau \end{cases} \quad (3.34)$$

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

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

Let $\gamma := \|\tilde{y}\| / \|b\|$. We can ask the question: how to choose ϵ 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

$$\begin{aligned} \text{if } \epsilon \geq \frac{1}{2\gamma} &\Rightarrow \|x_\epsilon\| < \gamma \|b\| = \|\tilde{y}\| \\ \text{if } \epsilon \geq \frac{1}{2\gamma^2} &\Rightarrow \|x_{d\epsilon}\| < \gamma \|b\| = \|\tilde{y}\|. \end{aligned}$$

and

If we transform the relaxed problem (3.29) and the doubly relaxed problem (3.31) using SVD we get

$$(y_\epsilon)_i = \frac{\sigma_i}{\sigma_i^2 + \epsilon} c_i \quad (3.35)$$

$$(y_{d\epsilon})_i = \frac{\sigma_i}{\sigma_i^2 + \epsilon + \frac{\epsilon^2}{\sigma_i^2}} c_i. \quad (3.36)$$

We see that we can think of (3.35), (3.36) as approximations of (3.33). More precisely, we want to choose ϵ so that the two functions

$$k_\epsilon(\sigma) := \frac{\sigma}{\sigma^2 + \epsilon} \quad (3.37)$$

and

$$k_{d\epsilon}(\sigma) := \frac{\sigma}{\sigma^2 + \epsilon + \frac{\epsilon^2}{\sigma^2}} \quad (3.38)$$

approximate the function $k(\sigma)$ (3.34). Rutishauser [35] and Molinari [13] show that indeed $k_{d\epsilon}$ is a better approximation than k_ϵ . More research could be done in this direction, e.g. relaxing the normal equations of the original problem

$$\begin{pmatrix} A^T A \\ \epsilon I \end{pmatrix} x \approx \begin{pmatrix} A^T b \\ 0 \end{pmatrix}$$

would yield the function

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

which is somehow between both discussed above. It is clear that those different relaxed solutions have to be computed without forming the normal equations. A program for x_ϵ and $x_{d\epsilon}$ is given in [13].

Theorem 3.2. Let $A_\epsilon = \begin{pmatrix} A \\ \epsilon I \end{pmatrix}$ and $A_\epsilon^+ = \begin{pmatrix} B^+ \\ C^+ \end{pmatrix}$. (B^+ is an $n \times n$

matrix.) Then for ϵ sufficiently small we have

$$B_\epsilon^+ = A^+ + \sum_{j=1}^{\infty} (-1)^j A^+ (A^+)^{2j} \epsilon^{2j},$$

Proof. Let $A = U \Sigma V^T$ 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

$$\begin{aligned} \underline{A}_{-\epsilon}^+ &= (\underline{A}^T \underline{A} + \epsilon^2 \underline{I})^{-1} \underline{A}^T = \underbrace{(V \Sigma^T U^T U \Sigma V^T + \epsilon^2 \underline{I})^{-1}}_{\underline{I}_m} (V \Sigma^T U^T, \underline{I}_{-\epsilon}) \\ &= V(\Sigma^T \Sigma + \epsilon^2 \underline{I})^{-1} (\Sigma^T U^T, \epsilon V^T) \\ &= V((\Sigma^T \Sigma + \epsilon^2 \underline{I})^{-1} \Sigma^T) U^T, \epsilon (\Sigma^T \Sigma + \epsilon^2 \underline{I})^{-1} V^T \\ &= [\underline{B}_{-\epsilon}^+, \underline{C}_{-\epsilon}^+] \end{aligned}$$

$$\underline{B}_{-\epsilon}^+ = V(\Sigma^T \Sigma + \epsilon^2 \underline{I})^{-1} \Sigma^T U^T$$

$$\underline{B}_{-\epsilon}^+ = \underline{V} \begin{array}{c|c} \frac{\sigma_1}{a_1^2 + \epsilon^2} & 0 \\ \cdot & \cdot \\ \frac{\sigma_k}{\sigma_r^2 + \epsilon^2} & \cdot \\ \cdot & \cdot \\ 0 & 0 \end{array} \underline{U}^T$$

$$\text{now } \frac{a_1}{a_1^2 + \epsilon^2} = \frac{1}{\sigma_1} \frac{1}{1 + (\epsilon/\sigma_1)^2} = \frac{1}{\sigma_1} \sum_{j=0}^{\infty} (\epsilon/\sigma_1)^{2j} \text{ for } |\epsilon| < \sigma_r,$$

$$\underline{B}_{-\epsilon}^+ = \sum_{j=0}^{\infty} \underline{V} \underbrace{\begin{bmatrix} 1/\sigma_1^{2j+1} & & & 0 \\ & \ddots & & \\ & & 1/\sigma_r^{2j+1} & \\ 0 & & & 0 \end{bmatrix}}_{\Sigma^+ (\Sigma^T \Sigma)^+ \Sigma^j} \underline{U}^T$$

$$\text{where } \underline{\Sigma}^+ = \begin{bmatrix} 1/\sigma_1 & & & 0 \\ & \ddots & & \\ & & \cdot & \\ 0 & & & 1/\sigma_r \end{bmatrix}_n,$$

we observe

$$V \Sigma^+ (\Sigma^T \Sigma)^+ \Sigma^j U^T = \underline{A}^+ (\underline{A}^+ \underline{A})^+ \Sigma^j \quad \cdot \quad 0$$

Remark. Theorem 3.2 shows that we can extrapolate $\underline{x}_M = \underline{A}^+ \underline{b}$ from the relaxed solution $\underline{x}_{\epsilon} = \underline{B}_{-\epsilon}^+ \underline{b}$. Using $\epsilon_1 = \epsilon_0/2^1$ we can apply the Romberg extrapolation since only even powers of ϵ occur. It has the advantage that the rank of \underline{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 $A = I$:

$$\left. \begin{aligned} \|x\| &= \min & (P3E) \\ \|\underline{Cx} - \underline{d}\| &= \alpha . \end{aligned} \right\}$$

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

$$(I + \lambda C^T C) \underline{x} = \lambda C^T \underline{d} \quad (4.1)$$

$$\|\underline{Cx} - \underline{d}\|^2 = \alpha^2 \quad ; \quad (4.2)$$

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

$$\alpha \geq \|(C^T C - I)\underline{d}\| := \alpha_{\min} .$$

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

$$f(h) = \|(\lambda C(I + \lambda C^T C)^{-1} C^T - I)\underline{d}\|^2 . \quad (4.3)$$

By Corollary 3.1, (3.13), this is equal to

$$f(h) = \|(\underline{I} + \lambda \underline{CC}^T)^{-1} \underline{d}\|^2 \quad (4.5)$$

which gives again the connection to the dual equations.

4.1 The Dual Normal Equations.

We transform the normal equations (4.1), (4.2) as follows: We multiply (4.1) from left by C giving

$$(I + \lambda \underline{CC}^T) \underline{Cx} = \lambda \underline{CC}^T \underline{d} . \quad (4.6)$$

If we subtract from (4.6) on both sides $\lambda \underline{CC}^T \underline{d} + \underline{d}$ we get

$$(I + \lambda \underline{CC}^T)(\underline{Cx} - \underline{d}) = -\underline{d} .$$

Finally introducing the new variable $\underline{z} = \underline{Cx} - \underline{d}$ we get the theorem:

Theorem 4.1.

(i) Let (\underline{x}, λ) be a solution of the primal normal equations

$$\begin{aligned} (\underline{I} + \lambda \underline{C}^T \underline{C}) \underline{x} &= \lambda \underline{C}^T \underline{d} \\ \|\underline{Cx} - \underline{d}\|^2 &= \alpha^2 \end{aligned} \quad (4.7)$$

then (\underline{z}, λ) with $\underline{z} := \underline{Cx} - \underline{d}$ is a solution of the dual normal equations

$$\begin{aligned} (\underline{I} + \lambda \underline{CC}^T) \underline{z} &= -\underline{d} \\ \|\underline{z}\|^2 &= \alpha^2 . \end{aligned} \quad (4.8)$$

(ii) Let (\underline{z}, λ) be a solution of (4.8), then (\underline{x}, λ) with $\underline{x} = -h \underline{C}^T \underline{z}$ is a solution of (4.7).

Proof.

$$\begin{aligned} (i) \quad (\underline{I} + \lambda \underline{CC}^T)(\underline{Cx} - \underline{d}) &= \underline{C}(\underline{I} + \lambda \underline{C}^T \underline{C}) \underline{x} - \underline{d} - \lambda \underline{CC}^T \underline{d} \\ &= \underline{C}(\lambda \underline{C}^T \underline{d}) - \underline{d} - \lambda \underline{CC}^T \underline{d} = -\underline{d} . \end{aligned}$$

and

$$\alpha^2 = \|\underline{z}\|^2 = \|\underline{c}\underline{x} - \underline{d}\|^2.$$

$$(ii) \quad (\underline{I} + h \underline{c}^T \underline{c})(-\lambda \underline{c}^T \underline{z}) = -h \underline{c}^T(\underline{x} + h \underline{c} \underline{c}^T \underline{z}) = h \underline{c}^T \underline{b}.$$

$$\|\underline{c}\underline{x} - \underline{d}\|^2 = \|\underline{c}(-\lambda \underline{c}^T \underline{z}) - \underline{d}\|^2 = \|\underline{z}\|^2 = \alpha^2. \quad \square$$

Equating both expressions for x and z again gives us the identities of Corollary 3.1.

in contrast to the relaxed least squares problem, the dual equations exists for every λ .

4.2 Representation as Least Squares Problem.

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

$$\alpha_{\min} < \alpha < \|\underline{d}\|.$$

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

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

\underline{x} and \underline{z} are obtained solving either

$$\begin{pmatrix} \underline{I} \\ \sqrt{\lambda} \underline{c} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{0} \\ \sqrt{\lambda} \underline{d} \end{pmatrix}$$

or the problem

$$\begin{pmatrix} \underline{I} \\ \sqrt{\lambda} \underline{c}^T \end{pmatrix} \underline{z} \approx \begin{pmatrix} -\underline{d} \\ \underline{0} \end{pmatrix}$$

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 problem (PIE), (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

$$\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C} \quad \text{where } h < 0 \quad (5.1)$$

which is in general not positive definite. A good way to solve this system is to transform it to diagonal form using BSVD of von Loan [45]. If

$$\underline{A} = \underline{I} \quad \text{or} \quad \underline{C} = \underline{I} \quad \text{we can use SVD to diagonalize the system stably [48].}$$

To avoid the forming of $\underline{A}^T \underline{A}$ and $\underline{C}^T \underline{C}$ we could consider the augmented $(m+n+p) \times (m+n+p)$ system

$$\begin{pmatrix} -1 & \underline{A} & \underline{0} \\ \underline{A}^T & \underline{0} & \underline{C}^T \\ 0 & \underline{C} & \frac{1}{\lambda} \underline{I} \end{pmatrix} \begin{pmatrix} \underline{y} \\ \underline{x} \\ \underline{z} \end{pmatrix} = \begin{pmatrix} \underline{b} \\ \underline{0} \\ \underline{d} \end{pmatrix} \quad (5.2)$$

Björk has shown [1] that for the augmented system

$$\begin{pmatrix} \underline{I} & \underline{A} \\ \underline{A}^T & \underline{0} \end{pmatrix} \begin{pmatrix} \underline{x} \\ \underline{z} \end{pmatrix} = \begin{pmatrix} \underline{b} \\ \underline{0} \end{pmatrix} \quad (5.3)$$

We can obtain the condition number $2 \kappa(\underline{A})$ with appropriate scaling (instead of $\kappa^2(\underline{A})$ for the normal equations $\underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b}$). However in (5.2) the condition number depends on the value of λ and may be large if h is near $-\mu_1$ where μ_1 is an eigenvalue of $\det(\underline{A}^T \underline{A} - \mu \underline{C}^T \underline{C}) = 0$.

The solution is numerically much better defined for problem (P1), (P2) 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 \underline{A} and \underline{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 \underline{A} respectively (P3) with \underline{C} being a band matrix. For a given $\lambda > 0$ we have to solve the least squares problems

$$\begin{pmatrix} \underline{A} \\ \sqrt{\lambda} \underline{I} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{b} \\ \underline{0} \end{pmatrix} \quad (5.4)$$

respectively

$$\begin{pmatrix} \underline{I} \\ \sqrt{\lambda} \underline{C} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{0} \\ \sqrt{\lambda} \underline{d} \end{pmatrix} \quad (5.5)$$

If we interchange the equations in (5.5) we have in either case a least squares problem of the form

$$\begin{pmatrix} \underline{B} \\ \underline{D} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{c}_1 \\ \underline{c}_2 \end{pmatrix} \quad (5.6)$$

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

Using Givens-rotations we shall show how an orthogonal matrix \underline{G} can be found such that

$$\underline{G} \begin{pmatrix} \underline{B} \\ \underline{D} \end{pmatrix} = \begin{pmatrix} \underline{0} \\ \underline{F} \end{pmatrix} \quad (5.7)$$

where \underline{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 \underline{F} .

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

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

This means that \underline{B} has m_1 diagonals below the main diagonal and m_2 diagonals above, that contains the possibly non-zero elements.

Let

$$k_{\min}(k) := \max\{k-m_2, 1\}$$

$$k_{\max}(k) := \min\{k+m_1, m\}.$$

Then the k -th column of \underline{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 (5.7) annihilating in the k -th step the k -th column of \underline{B} and producing the k -th row of \underline{F} . Let \underline{F} be at the beginning the diagonal matrix \underline{D} and defined by

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

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

The following procedure performs the transformation

$$G \begin{pmatrix} c_1 \\ \vdots \\ c_l \end{pmatrix}$$

and solves the bidiagonal system by backsubstitution.

```

procedure solve (n,u,v,c0,si,cl,c2,y);
  value n; integer n;
  array c0, si, cl, c2, y;
  begin
    integer i; real t,c,s,h;
    for i := 1 step 1 until n do
      begin
        c := c0[2*i-1]; s := si[2*i-1];
        h := cl[i]*c + c2[i]*s;
        c2[i] := -cl[i]*s + c2[i]*c;
        cl[i] := h;
        if i < n then
          begin
            c := c0[2*i]; s := si[2*i];
            h := cl[i]*c + c2[i]*s;
            c2[i+1] := -cl[i]*s + c2[i+1]*c;
            cl[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;
  end

```

and stores the Givens rotations in two arrays `co[i]`, `si[i]`

```

procedure updec (n,q,e,d,u,v,c0,si);
  value n; integer n;
  array q,e,d,u,v,c0,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*t); 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]);
        c0[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]);
            c0[2*i] := c; si[2*i] := s;
            u[i+1] := -h*s + u[i+1]*c;
          end if;
        end i;
      end updec;
  end

```

Alternatively we can avoid storing the rotations observing that if

$$\begin{pmatrix} \underline{B} \\ \underline{D} \end{pmatrix} \underline{y} \approx \begin{pmatrix} \underline{c1} \\ \underline{c2} \end{pmatrix}$$

and if

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

then the normal equations are

$$\underline{F}^T \underline{F} \underline{y} = \underline{F}^T \underline{c1} + \underline{D} \underline{c2}$$

Therefore we can also find \underline{y} solving

$$\underline{F}^T \underline{w} = \underline{B}^T \underline{c1} + \underline{D} \underline{c2}$$

$$\underline{F} \underline{y} = \underline{w}$$

The new procedure updec does not contain the variable $\underline{c0}$ and \underline{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 solve2;

```

Before calling solve2 the right hand side c has to be defined by

$$\underline{c} := \underline{B}^T \underline{c1} + \underline{D} \underline{c2}$$

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 square problem of the type

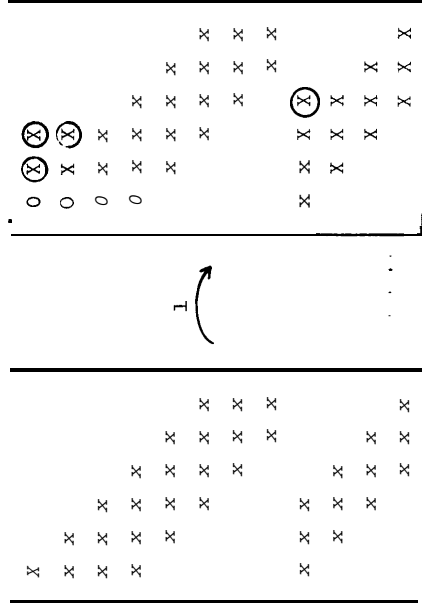
$$\begin{pmatrix} \underline{A} \\ \underline{C} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{b} \\ \underline{d} \end{pmatrix} \quad (5.8)$$

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

$$\underline{G} \begin{pmatrix} \underline{A} \\ \underline{C} \end{pmatrix} = \begin{pmatrix} \underline{O}' \\ \underline{A}' \end{pmatrix} \quad (5.9)$$

where \underline{A}' has the same bandwidth as \underline{A} and \underline{O}' is zero up to some elements in the right bottom corner. Using the elements of \underline{O}' and \underline{A}' we then compute \underline{x} by backsubstitution.

In the k -th step of the algorithm we annihilate the k -th column of \underline{A} and produce the k -th row of \underline{A}' . We explain the rotation: for the example where \underline{A} has 4 and \underline{C} has 3 diagonals.



$T = \underline{P}^T \underline{B}$, \underline{B} upper bidiagonal

and (5.10) then becomes

$$\begin{pmatrix} \underline{D} & \underline{Q} \\ \sqrt{\lambda} & \underline{B} \end{pmatrix} \underline{z} \approx \begin{pmatrix} \underline{D}^{-1} \underline{y} \\ \underline{0} \end{pmatrix} \quad (5.12)$$

The matrix of (5.12) has the form ($n = 5$)

$$\begin{array}{c} \left[\begin{array}{cccccc} X & & & & & \\ x & x & & & & \\ x & x & x & & & \\ x & x & x & x & & \\ & x & x & x & x & \\ & x & x & x & x & X \\ & & & & & X \end{array} \right] \begin{array}{c} \underline{0} \\ \underline{1} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \end{array} \end{array} \approx \begin{array}{c} \left[\begin{array}{cccccc} \underline{0} & \underline{X} & & & & \\ 0 & X & & & & \\ 0 & X & X & & & \\ & X & X & X & & \\ & & X & X & X & \\ & & & X & X & X \\ & & & & X & X \\ & & & & & X \end{array} \right] \begin{array}{c} \underline{c} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \end{array} \end{array}$$

Zeroing the first column in step 1 leaves one element \underline{X} that has to be removed in the cleaning step.

5.3 Bidiagonalization.

If the matrices \underline{A} and \underline{C} are dense then for problems (P2), (P3) and (P1) (after performing Eiden's transformation, see Section 3.3) we have to solve a least square problem of the form

$$\begin{pmatrix} \underline{A} \\ \sqrt{\lambda} & \underline{I} \end{pmatrix} \underline{x} \approx \begin{pmatrix} \underline{b} \\ \underline{0} \end{pmatrix} \quad (5.13)$$

for every new value of λ . If \underline{P} and \underline{Q} are orthogonal then an equivalent system to (5.13) is

$$\begin{pmatrix} \underline{P}^T & \underline{0} \\ \underline{0} & \underline{Q}^T \end{pmatrix} \begin{pmatrix} \underline{A} \\ \sqrt{\lambda} & \underline{I} \end{pmatrix} \underline{Q} \underline{Q}^T \underline{x} \approx \begin{pmatrix} \underline{P}^T \underline{b} \\ \underline{0} \end{pmatrix} \quad (5.14)$$

$$\Rightarrow \begin{pmatrix} \underline{P}^T \underline{A} \underline{Q} \\ \sqrt{\lambda} & \underline{I} \end{pmatrix} \underline{y} \approx \begin{pmatrix} \underline{c} \\ \underline{0} \end{pmatrix}$$

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

Now we can choose \underline{P} and \underline{Q} so that (5.14) is simpler to solve than (5.13). More [29] proposed to choose \underline{P} and \underline{Q} so that

$$\underline{B} := \underline{P}^T \underline{A} \underline{Q} = \begin{pmatrix} \underline{R} & \\ & \underline{0} & \underline{0} \end{pmatrix}, \quad (5.15)$$

i.e., to perform the QR decomposition of \underline{A} with column pivoting. However the solution of (5.14) in this case is still an n^3 process.

It has been pointed out by More/ (private communication) that in his application [29] only $l \leq 2$ iterations are needed to solve the secular equation. Therefore it does not matter in this case if the solution of (5.14) is an n^3 process. In general however we will have to perform much more iterations (see Section 6). With about double the amount of work to perform the decomposition (5.15) we can bidiagonalize \underline{A} :

$$\underline{B} := \underline{P}^T \underline{A} \underline{Q} = \begin{pmatrix} \underline{R} & \\ & \underline{0} \end{pmatrix} \quad (5.16)$$

as proposed by [8] and [12]. Then to solve (5.14) we can use the algorithm of Section 5.1 that requires only $\sim n$ multiplications.

If we wish to compute $y = P \cdot x$, 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** q and the superdiagonal e of B and stores the transformation

vectors w_{-j} and v_{-j} in \underline{A} (5.17):

```
procedure bidia(m,n,a,g,e);
integer,m,n; _____ m,n;
array;
begin
  integer i,j,k; real s, fak;
  for i := 1 step 1 until n do
    begin
      comment transform (ai1, ..., ami) to (ei, 0, ..., 0);
      s := 0;
      for j := i step 1 until m do s := s + a[j,i]2;
      if s = 0 then q[i] := 0 else
        begin
          s := sqrt(s);
          q[i] := if a[i,i] > 0 then -s else s;
          fak := sqrt(s x (s + abs(a[i,i])));
          a[i,i] := a[i,i] - q[i];
          for k := i step 1 until m do a[k,i] := a[k,i]/fak;
          for j := i+1 step 1 until n do
            /
            ende;
          end bidia;
        end;
      comment transform (ai,i+1}, ..., ain) to (ei, 0, ..., 0);
      if i = n then goto ende;
      s := 0;
      for j := i+1 step 1 until n do s := s + a[i,j]2;
      if s = 0 then e[i] := 0 else
        begin
          s := sqrt(s);
          e[i] := if a[i,i+1] > 0 then -s else s;
          fak := sqrt(s x (s + 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] x a[i,k];
              for k := i+1 step 1 until n do a[j,k] := a[j,k] - a[i,k] x s;
            end;
          end s;
          end i;
        end;
      ende;
    end bidia;
  end;

```

The decomposition using `_bigdia` requires $\sim 2(m n^2 - n^3/3)$ multiplications.

If $m > n$ it is possible to bidiagonalize A even cheaper [26], [5 I, [8].

The idea is to transform \underline{A} first to an upper triangular matrix R and then bidiagonalize R . In this case the operation count is $\sim m n^2 + n^3$.

An alternative approach to bidiagonalizing \underline{A} has been suggested by Golub and Kahan [14]. This algorithm uses the matrix \underline{A} not explicitly. Only the operator $\underline{A} \underline{x}$ 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

$$f(\lambda) = \alpha^2 \tag{5.18}$$

using some high order iteration method we have to compute the derivatives

of f . For the general problem (Pl) we have

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x}(\lambda) = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d} \tag{5.19}$$

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

By differentiating (5.19) we get

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x}' = -\underline{C}^T (\underline{C} \underline{x} - \underline{d}) \tag{5.20}$$

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x}'' = -2 \underline{C}^T \underline{C} \underline{x}' .$$

In general we have for $k \geq 2$

$$(\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x}^{(k)} = -k \underline{C}^T \underline{C} \underline{x}^{(k-1)} . \tag{5.21}$$

Lemma 5.1. If $\underline{B}_\lambda := \underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}$ and $\underline{x}^{(k)}$ is the solution of (5.21)

then

$$(i) \quad \underline{x}^{(k)T} \underline{B}_\lambda \underline{x}^{(k)} = \begin{cases} \underline{x}^T (\underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d}) , & k = 0 \\ - \underline{x}'^T \underline{C}^T (\underline{C} \underline{x} - \underline{d}) , & k = 1 \\ - k \underline{x}^{(k)T} \underline{C}^T \underline{C} \underline{x}^{(k-1)} , & k \geq 2 \end{cases}$$

$$(ii) \quad \underline{x}^{(k)T} \underline{B}_{\lambda-\underline{x}^{(k-1)}} \underline{x}^{(k)} = \begin{cases} \frac{1}{2} \underline{x}^T \underline{B}_\lambda \underline{x}'' + \underline{x}^T \underline{C}^T \underline{d} , & k = 1 \\ \frac{k}{k+1} \underline{x}^{(k-1)T} \underline{B}_{\lambda-\underline{x}^{(k-1)}} \underline{x}^{(k+1)} , & k \geq 2 \end{cases}$$

$$(iii) \quad \|\underline{C} \underline{x}^{(k)}\|^2 = \begin{cases} - \underline{x}^T \underline{B}_\lambda \underline{x}' + \underline{x}^T \underline{C}^T \underline{d} , & k = 0 \\ - \frac{1}{k+1} \underline{x}^{(k)T} \underline{B}_{\lambda-\underline{x}^{(k+1)}} \underline{x}^{(k+1)} , & k > 1 \end{cases}$$

$$(iv) \quad \|\underline{C} \underline{x}^{(k)}\| = \begin{cases} \frac{1}{2} \underline{x}^T (\underline{C}^T \underline{C} \underline{x} - \underline{d}) , & k = 1 \\ \frac{k}{k+1} \underline{x}^{(k+1)T} \underline{C}^T \underline{C} \underline{x}^{(k-1)} , & k > 2 \end{cases}$$

$$(v) \quad \frac{d}{d\lambda} (\underline{x}^{(k)T} \underline{B}_\lambda \underline{x}^{(k)}) = \begin{cases} - \|\underline{C} \underline{x} - \underline{d}\|^2 + \|\underline{d}\|^2 , & k = 0 \\ - (2k+1) \|\underline{C} \underline{x}^{(k)}\|^2 , & k \geq 1 \end{cases}$$

$$(vi) \quad \frac{d}{d\lambda} \|\underline{C} \underline{x}^{(k)}\|^2 = \begin{cases} - 2 \underline{x}^T \underline{B}_\lambda \underline{x}' + 2 \underline{x}^T \underline{C}^T \underline{d} , & k = 0 \\ - \frac{k}{k+1} \underline{x}^{(k+1)T} \underline{B}_{\lambda-\underline{x}^{(k+1)}} \underline{x}^{(k+1)} , & k > 1 \end{cases}$$

Proof. Equations (i) follow, by multiplying (5.19), (5.20), (5.21) by $\underline{x}^{(k)T}$ from the left.

Looking at two consecutive equations of (5.21):

$$\underline{B}_\lambda \underline{x}^{(k)} = -k \underline{C}^T \underline{C} \underline{x}^{(k-1)} \quad (5.22)$$

$$\underline{B}_\lambda \underline{x}^{(k+1)} = -(k+1) \underline{C}^T \underline{C} \underline{x}^{(k)}. \quad (5.23)$$

We obtain (ii) by multiplying the first (5.22) by $\underline{x}^{(k)T}$ and the second by $\underline{x}^{(k-1)T}$ and by eliminating the expression $\underline{x}^{(k)T} \underline{C}^T \underline{C} \underline{x}^{(k-1)}$. To prove (iii) we multiply the equation for $\underline{x}^{(k+1)}$ of (5.21) by $\underline{x}^{(k)}$ and (5.20) by \underline{x}^T . Equation (iv) follows by multiplying (5.22) by $\underline{x}^{(k+1)T}$ and (5.23) by $\underline{x}^{(k)}$ and subtracting both equations.

Finally observe that $\underline{B}'_\lambda = \underline{C}^T \underline{C}$, therefore

$$\begin{aligned} \frac{d}{d\lambda} (\underline{x}^{(k)T} \underline{B}'_\lambda \underline{x}^{(k)}) &= \underline{x}^{(k+1)T} \underline{B}'_\lambda \underline{x}^{(k)} + \underline{x}^{(k)T} \{ \underline{B}'_\lambda \underline{x}^{(k)} + \underline{B}'_\lambda \underline{x}^{(k+1)} \} \\ &= 2 \underline{x}^{(k+1)T} \underline{B}'_\lambda \underline{x}^{(k)} + \|\underline{C}\underline{x}^{(k)}\|^2. \end{aligned}$$

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

$$\frac{d}{d\lambda} (\underline{x}^T \underline{B}'_\lambda \underline{x}) = -\|\underline{C}\underline{x}\|^2 + 2 \underline{x}^T \underline{C}^T \underline{d} = -\|\underline{C}\underline{x} - \underline{d}\|^2 + \|\underline{d}\|^2.$$

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

$$\frac{d}{d\lambda} (\underline{x}^{(k)T} \underline{B}'_\lambda \underline{x}^{(k)}) = -(2k+1) \|\underline{C}\underline{x}^{(k)}\|^2.$$

Using (i) we can write

$$\frac{d}{d\lambda} \|\underline{C}\underline{x}^{(k)}\|^2 = -\frac{2}{k+1} \underline{x}^{(k+1)T} \underline{B}'_\lambda \underline{x}^{(k+1)}, \quad \square$$

Corollary 5.1. If \underline{x} is the solution of (5.19) then

$$\left. \begin{aligned} \int f(\lambda) d\lambda &= \lambda \|\underline{d}\|^2 - \underline{x}^T (\underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}) \underline{x} + \text{const.} \\ &= -\underline{A} \underline{d}^T (\underline{C}\underline{x} - \underline{d}) - \underline{b}^T \underline{A}\underline{x} + \text{const.} \end{aligned} \right\} \quad (5.24)$$

Proof. From equation (v) of Lemma 5.1 we have

$$f(\lambda) = \|\underline{C}\underline{x} - \underline{d}\|^2 = \|\underline{d}\|^2 - \frac{d}{d\lambda} (\underline{x}^T \underline{B}'_\lambda \underline{x}).$$

By integrating and using (5.19) the result follows immediately. \square

Corollary 5.2. If $\underline{B}'_\lambda = \underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}$ then for $k > 1$,

$$(i) \quad \frac{d}{d\lambda} (\underline{x}^{(k)T} \underline{B}'_\lambda \underline{x}^{(k)}) = \frac{2k+1}{k+1} \underline{x}^{(k)T} \underline{B}'_\lambda \underline{x}^{(k+1)}$$

$$(ii) \quad \frac{d}{d\lambda} (\underline{x}^{(k)T} \underline{B}'_\lambda \underline{x}^{(k+1)}) = 2 \underline{x}^{(k+1)T} \underline{B}'_\lambda \underline{x}^{(k+1)}$$

Corollary 5.3. If $\underline{x}^{(k)}$ is a solution of (5.21), then for $k \geq 1$

$$(i) \quad \frac{d}{d\lambda} \|\underline{C}\underline{x}^{(k)}\|^2 = 2 \underline{x}^{(k+1)T} \underline{C}^T \underline{C} \underline{x}^{(k)}$$

$$(ii) \quad \frac{d}{d\lambda} (\underline{x}^{(k+1)T} \underline{C}^T \underline{C} \underline{x}^{(k)}) = \frac{2k+3}{k+1} \|\underline{C}\underline{x}^{(k+1)}\|^2.$$

Theorem 5.1. Let \underline{x} , \underline{x}' , and $\underline{x}^{(k)}$ be solutions of (5.19), (5.20), and (5.21) and let $\underline{B}'_\lambda = \underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}$. Then

$$f(A) = \|\underline{C}\underline{x} - \underline{d}\|^2 = -\underline{x}^T \underline{B}'_\lambda \underline{x}' - \underline{d}^T (\underline{C}\underline{x} - \underline{d}) \quad (5.25)$$

$$f'(A) = 2 \underline{x}^T \underline{C}^T (\underline{C}\underline{x} - \underline{d}) = -2 \underline{x}^T \underline{B}'_\lambda \underline{x}' \quad (5.26)$$

and for $k \geq 1$

$$f^{(2k)}(\lambda) = (k+1)\gamma_{2k} \|\underline{c}_x^{(k)}\|^2 = -\gamma_{2k} \underline{x}^{(k)T} \underline{B}_\lambda \underline{x}^{(k+1)} \quad (5.27)$$

$$\begin{aligned} f^{(2k+1)}(\lambda) &= (k+1)\gamma_{2k+1} \underline{x}^{(k)T} \underline{c}^T \underline{c} \underline{x}^{(k+1)} \\ &= -\gamma_{2k+1} \underline{x}^{(k+1)T} \underline{B}_\lambda \underline{x}^{(k+1)} \end{aligned} \quad (5.28)$$

where

$$\begin{aligned} \gamma_{2k} &= \frac{1 \cdot 3 \cdot 5 \cdots (2k+1)}{(k+1)!} 2^k \\ \gamma_{2k+1} &= \frac{1 \cdot 3 \cdot 5 \cdots (2k+1)}{(k+1)!} 2^{k+1} \end{aligned} \quad (5.29)$$

Proof. The proof is by induction. If

$$f(\lambda) = \|\underline{c}_x - \underline{d}\|^2$$

then by differentiating and by Lemma 5.1 (i) we have

$$f'(A) = 2 \underline{x}^T \underline{c}^T (\underline{c}_x - \underline{d}) = -2 \underline{x}^T \underline{B}_\lambda \underline{x}'$$

Differentiating again using (v) and (iii) of Lemma 5.1 we get

$$f''(A) = 2 \cdot 3 \|\underline{c}_x'\|^2 = -3 \underline{x}^T \underline{B}_\lambda \underline{x}''$$

which is (5.27) for $k = 1$. Now we can use Corollary 5.2 and 5.3 to compute the higher derivatives.

Assume

$$f^{(2k)}(\lambda) = (k+1)\gamma_{2k} \|\underline{c}_x^{(k)}\|^2 = -\gamma_{2k} \underline{x}^{(k)T} \underline{B}_\lambda \underline{x}^{(k+1)}$$

Then

$$\begin{aligned} f^{(2k+1)}(\lambda) &= (k+1)\gamma_{2k+1} \cdot 2 \underline{x}^{(k)T} \underline{c}^T \underline{c} \underline{x}^{(k+1)} \\ &= -\gamma_{2k+1} \cdot 2 \underline{x}^{(k+1)T} \underline{B}_\lambda \underline{x}^{(k+1)} \end{aligned}$$

which is (5.28) for $\gamma_{2k+1} = 2\gamma_{2k}$. Again differentiating using Corollary 5.2 and 5.3 yields

$$\begin{aligned} f^{(2k+2)}(\lambda) &= \gamma_{2k+1} (2k+3) \|\underline{c}_x^{(k+1)}\|^2 \\ &= -\gamma_{2k+1} \frac{2k+3}{k+2} \underline{x}^{(k+1)T} \underline{B}_\lambda \underline{x}^{(k+2)} \end{aligned}$$

which is (5.27) for $k+1$ and $\lambda_{2k+2} = \gamma_{2k+1} \frac{2k+3}{k+2}$. From this recursion for γ_i it is easy to verify (5.29). \square

Theorem 5.1 shows that we can compute cheaply derivatives of f . To compute $x^{(k)}$ we have to solve a linear system with the same matrix \underline{B}_λ as for x . Therefore we can use a factorization of \underline{B}_λ , If $A > 0$ then $\underline{x}, \underline{x}', \dots, \underline{x}^{(k)}$ is a solution of the least squares problem

$$\begin{pmatrix} A \\ \sqrt{\lambda} \underline{c} \end{pmatrix} \underline{x} \approx \begin{pmatrix} b \\ \sqrt{\lambda} \underline{d} \end{pmatrix}$$

$$\begin{pmatrix} A \\ \sqrt{\lambda} \underline{c} \end{pmatrix} \underline{x}' \approx -\frac{1}{\sqrt{\lambda}} \begin{pmatrix} 0 \\ \underline{c}_x - \underline{d} \end{pmatrix}$$

$$\begin{pmatrix} A \\ \sqrt{\lambda} \underline{c} \end{pmatrix} \underline{x}^{(k)} \approx -\frac{k}{\sqrt{\lambda}} \begin{pmatrix} 0 \\ \underline{c}_x^{(k-1)} \end{pmatrix}, \quad k > 2$$

If we transform by an orthogonal matrix \underline{G}

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

then $\underline{R}_\lambda^T \underline{R}_\lambda$ is the Cholesky decomposition of $\underline{R}_\lambda = \underline{A}^T \underline{A} + \lambda \underline{C}^T \underline{C}$. To solve

$$\underline{R}_\lambda \underline{y}^{(k)} = -\underline{k} \underline{C}^T \underline{C} \underline{x}^{(k-1)}$$

we have two possibilities

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

$$\underline{R}_\lambda^T \underline{y}^{(k)} = -\underline{k} \underline{C}^T (\underline{C} \underline{x}^{(k-1)});$$

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

$$\begin{pmatrix} \underline{y}^{(k)} \\ \underline{k} \end{pmatrix} = \underline{G} \begin{pmatrix} \underline{0} \\ -\underline{k} \underline{C} \underline{x}^{(k-1)} \end{pmatrix}.$$

Then we obtain $\underline{x}^{(k)}$ by back substitution in

$$\underline{R}_\lambda \underline{x}^{(k)} = \underline{y}^{(k)}.$$

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

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

or equivalently

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

Similarly

...

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

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

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

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

1. Compute \underline{G} and \underline{R}_λ so that

$$\underline{G} \begin{pmatrix} \underline{A} \\ \sqrt{\lambda} \underline{C} \end{pmatrix} = \begin{pmatrix} \underline{R}_\lambda \\ \underline{0} \end{pmatrix}.$$

2. Compute \underline{y} from

$$\underline{R}_\lambda^T \underline{y} = \underline{A}^T \underline{b} + \lambda \underline{C}^T \underline{d}$$

or using \underline{G} by

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

3. Compute \underline{x} from $\underline{R}_\lambda \underline{x} = \underline{y}$ and form $\underline{z} := \underline{C} \underline{x} - \underline{d}$.

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

5. Compute \underline{y}' by solving

$$\underline{R}_\lambda^T \underline{y}' = -\underline{C}^T \underline{z}$$

or by using \underline{G}

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

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

If we do not want to store \underline{G} then we have to compute \underline{y}' in step 5 by forward substitution. However, \underline{y} in step 2 can be computed together with the decomposition in step 1 without forming \underline{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 $A^* > 0$ of the secular equation

$$f(\lambda) = \alpha^2 \tag{6.1}$$

that we need to solve problem (P1), (P2) or (P3). The length function f is a positive rational function for all A and decreasing in $(0, \infty)$.

If we start with $A = 0$ Newton's method will produce a strictly increasing sequence $\{\lambda_n\}$ which converges globally to A^* . However as it was observed in [34] if α is small, convergence is slow and as a remedy Reinsch suggested solving the equation

$$\frac{1}{\sqrt{f}} - \frac{1}{\alpha} = 0 \tag{6.2}$$

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

$$\begin{aligned} \mathcal{E}_1(\lambda) &:= f(\lambda) - \alpha^2 = 0 \\ \mathcal{E}_2(\lambda) &:= \sqrt{f(\lambda)} - \alpha = 0 \\ \mathcal{E}_3(\lambda) &:= \frac{1}{\sqrt{f(\lambda)}} - \frac{1}{\alpha} = 0 \end{aligned}$$

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

$$\lambda = \frac{f - \alpha^2}{f'} \quad \text{for } \epsilon_1, \quad (6.3)$$

$$\lambda = \frac{f - \alpha^2}{f'} \frac{2}{1 + \frac{\alpha}{\sqrt{f}}} \quad \text{for } \epsilon_2, \quad (6.4)$$

$$\lambda = \frac{f - \alpha^2}{f'} \frac{2}{1 + \frac{\alpha}{\sqrt{f}}} \frac{2}{1 + \frac{\alpha}{\sqrt{f}}} \quad \text{for } \epsilon_3. \quad (6.5)$$

For $\lambda = 0$ we typically have $f(0) \gg \alpha^2$. Since $f' < 0$ for $A > 0$ the step in (6.4) will be about twice as big as in (6.3). But (6.5) will produce an even larger step, proportional to $\frac{\sqrt{f}}{\alpha}$ the discrepancy of \sqrt{f} and α . If $f \approx \alpha^2$ then all the three steps are of about the same size.

We can look at the two functions

$$h_1(\lambda) = 2 / \left(1 + \frac{\alpha}{\sqrt{f(\lambda)}} \right)$$

$$h_2(\lambda) = h_1(\lambda) \frac{\sqrt{f}}{\alpha}$$

as functions that help to accelerate the global convergence of (6.3) 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]

$$\left. \begin{array}{l} x_0 \text{ arbitrary} \\ x_{n+1} = F(x_n) \end{array} \right\} \quad n = 0, 1, \dots \quad (6.6)$$

where

$$F(x) = x - \frac{f(x)}{f'(x)}. \quad (6.7)$$

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 (6.7) is better than Newton's iteration ($G(x) \equiv 1$). As has been pointed out by Kahan [24], every sequence generated by an iteration (6.6) can be interpreted as a sequence obtained by applying Newton's method to a certain equation

$$g(x) = 0. \quad (6.8)$$

Indeed if we put

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

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

$$g(x) = c \cdot \exp \left(\int \frac{dx}{x - F(x)} \right). \quad (6.9)$$

Solving (6.8) with (6.9) using Newton's method yields the sequence (6.6). Some examples may illustrate the point.

(1) Let $x_{n+1} = 1 + \frac{x}{2}$, $x_0 = 0$. This sequence converges

linearly to $s = 2$. Indeed using (6.9) 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

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

using (6.9) we obtain

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

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

$$x_{n+1} = x_n - \frac{f(x_n) - \alpha}{f'(x_n)} \cdot \frac{f(x_n)}{\alpha} \quad (6.10)$$

to solve $f(x) = \alpha$. Using (6.9) we get

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

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

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

Therefore (6.10) 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 $F(x)$ (6.7) we would like to consider however if it will not be possible in general to evaluate the integral for g in (6.9) 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

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

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

$$F(x) = x - \frac{f(x)}{f'(x)} \cdot G(x) \quad (6.11)$$

is Newton's method and it is of second order for a zero of f with multiplicity one. Differentiation (6.11) 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 6.1. 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:

$$(4) \quad G(x) = H(f) = 1 + f(x) \Rightarrow F(x) = x - \frac{f(x)}{f'(x)} (1 + f(x)) \quad (6.13)$$

(6.13) yields a second order iteration formula. Indeed using (6.9)

we see that the iteration is obtained solving $g(x) = 0$ with Newton

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

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

$$(5) \quad H(f) = \frac{\alpha + f}{\alpha} \quad (6.14)$$

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

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

$$(6) \quad G(x) = 1 / \left(1 - \frac{1}{2} \frac{f(x)f''(x)}{(f'(x))^2}\right)$$

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

6.2 Third Order Iterative Methods.

We consider again the iteration

$$x_{n+1} = \tilde{F}(x_n) = x_n \frac{f(x_n)}{f'(x_n)} \cdot G(x_n) \quad (6.15)$$

In Section 6.1 we saw that if $G(s) = 1$ then (6.15) is quadratic convergent. Let $u(x) := \frac{f(x)}{f'(x)}$. Then

$$F(x) = x - u(x), \dots \quad (6.16)$$

and we can ask for the conditions that the iteration (6.15) be cubically convergent. Differentiating (6.16) we get (dropping the argument)

$$\begin{aligned} F' &= 1 - u'G - uG' \\ F'' &= -u''G - 2u'G' - uG'' \\ u &= f/f' \\ u' &= 1 - \frac{ff''}{f'^2} \\ u'' &= -\frac{f'''}{f'} + 2 \frac{ff'''}{f'^3} - \frac{f f''^2}{f'^4} \end{aligned}$$

Now we want to have

$$F'(s) = F''(s) = 0. \tag{6.17}$$

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

$$G(s) = 1 \text{ and } G'(s) = \frac{1}{2} \frac{f''(s)}{f'(s)}. \tag{6.18}$$

Theorem 6.1. Let $H \in C^2[-a, a]$ for some $a > 0$. The iteration $x_{n+1} = F(x_n)$ with

$$F(x) = x - \frac{f(x)}{f'(x)} \cdot H\left(\frac{f(x)f''(x)}{(f'(x))^2}\right) \tag{6.19}$$

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

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

Proof. The function F in (6.19) is the special case of (6.16) 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

$$t'(s) = -u''(s) = \frac{f''(s)^2}{f'(s)^3}.$$

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

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

$$\Leftrightarrow H'(0) = \frac{1}{2}. \quad \square$$

Many well known third order iterative methods are special cases of Theorem 6.1. Let

$$t := \frac{f(x)f''(x)}{(f'(x))^2}. \tag{6.20}$$

(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 $H'(0) = 1/2$.

(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 [50].

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

(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

$$H(t) = e^{\frac{1}{4}t} = 1 + \frac{1}{2}t + \frac{1}{8}t^2 + \frac{1}{48}t^3 + \dots$$

This formula has the advantage that $H(t) > 0$ even if the starting point is far away 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 of third order iteration methods.

Lemma 6.2. Let $H_1(t)$ and $H_2(t)$ be two functions with

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

Then the three mean functions

$$A = (H_1 + H_2)/2$$

$$B = \sqrt{H_1 H_2}$$

$$C = 2 / \left(\frac{1}{H_1} + \frac{1}{H_2} \right)$$

have the same property.

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

$$A' = (H_1' + H_2')/2 = b.$$

$$B' = (H_1' H_2 + H_1 H_2') / (2 \sqrt{H_1 H_2})$$

$$B'(0) = (b \cdot a + a \cdot b) / (2a) = b.$$

$$C = \frac{2 H_1 H_2}{H_1 + H_2}$$

$$C' = 2 \frac{(H_1 + H_2)(H_1' H_2 + H_1 H_2') - (H_1' + H_2') H_1 H_2}{(H_1 + H_2)^2}$$

$$C'(0) = 2 \frac{2a(2ab) - 2ba^2}{4a^2} = b. \quad \square$$

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 (6.19). For example, consider

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

with $H(0) = 1$, $H'(0) = 1/2$, and $t(x)$ defined by (6.20). Clearly

this choice of the convergence factor G leads to a third order formula which has not the form of Theorem 6.1.

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

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

denote the Schröder iteration functions (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

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

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

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

which we can write

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

with arbitrary ψ .

6.3 The Convergence Factor for a Third Order Method.

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

$$K(s) = K'(s) = 0, \quad K'(s) = 1. \tag{6.21}$$

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

$$F(x) = x - K(x) \cdot G(x). \tag{6.22}$$

We have

$$\left. \begin{aligned} F' &= 1 - K'G - KG' \\ F'' &= -K''G - 2K'G' - KG'' \end{aligned} \right\} \tag{6.23}$$

Lemma 6.3. The iteration $x_{n+1} = F(x_n)$ with F defined in (6.22) is also of third order if $G \in C^2[a, b]$ with $s \in (a, b)$ and

$$G(s) = 1, \quad G'(s) = 0. \tag{6.24}$$

Proof. If we use (6.21) and (6.24) in (6.23) then we have

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

If we choose G so that

$$G''(s) = -\frac{1}{2} K''(s)$$

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:

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

where $t(x) = \frac{f(x) + \alpha}{\alpha} > 0$

and $\alpha, \beta \in \mathbb{R}$.

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

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

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

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

$$g(x) = \exp\left(\int \frac{dx}{K(x) \cdot G(x)}\right) \tag{6.25}$$

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 (6.25) 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 $f(x) = a$.

The function h should be simple 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

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

(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 $f^{[i]}$, the derivatives must be replaced using derivatives of f [21]:

$$(f^{[i]}(y))' = \frac{1}{f'(f^{[i]}(y))}$$

$$(f^{[i]}(y))'' = -\frac{f''(f^{[i]}(y))}{f'^2(f^{[i]}(y))}.$$

If h approximates f resp $f^{[i]}$ 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 $f^{[i]}$ 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):

$$\left. \begin{aligned} a &= 4 f^3 / f'^2 \\ b &= -x - 2 f / f' \end{aligned} \right\} \quad (6.26)$$

Now $h(z) = \alpha$ gives

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

using (6.26) we get

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

If we use only the + sign we have the iteration

$$\left. \begin{aligned} x_{n+1} &= x_n - \frac{f(x_n) - \alpha}{f'(x_n)} = G(x_n) \\ \text{with } G(x) &= 2 \sqrt{\frac{f}{\alpha}} / \left(1 + \sqrt{\frac{\alpha}{f}} \right) \end{aligned} \right\} \quad (6.27)$$

But (6.27) is Reinsch's proposal to solve $\frac{1}{\sqrt{f}} = \frac{1}{\sqrt{\alpha}}$ with Newton

instead of $f - a = 0$. We know from Section 3 that the length function f

has the form

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

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

$$g(x) = (f(x))^{-1/\beta} - \alpha^{-1/\beta} \quad (6.28)$$

if $f(x) > 0$ and $a > 0$. Newton's iteration for (6.28) yields the same iteration function as if we approximate f and f' locally by

$$h(x) = \frac{a}{(x+b)^{-1/\beta}}$$

The resulting iteration formula is

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

with

$$G(x) = \frac{\beta f}{f - \alpha} \left(\beta \sqrt{\frac{f}{\alpha}} - 1 \right) \quad (6.29)$$

If we let $\beta \rightarrow \infty$ in (6.29) we get

$$G(x) = \frac{f}{f - \alpha} \ln \left(\frac{f}{\alpha} \right) \quad (6.30)$$

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

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

The next examples use functions 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):

$$a = -4 f'^3 / f''^2$$

$$b = -x - 2 f' / f''$$

$$c = f - 2 f'^2 / f''$$

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

$$x_{n+1} = x_n - \frac{f_n - \alpha}{f'_n} \frac{1}{1 - \frac{1}{2} \frac{(f_n - \alpha) f''_n}{f'^2_n}} \quad (6.31)$$

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

$$g(x) = \frac{f(x) - \alpha}{\sqrt{f'(x)}} = 0$$

If we approximate the inverse function $f^{[-1]}(y)$ locally by

$h(y) = a/(y+b) + c$ we get

$$a = -4 f'^3 / f''^2$$

$$b = -f + 2 f'^2 / f''$$

$$c = x + 2 f' / f''$$

(6) Approximating $f^{[-1]}$ by a parabola

$$h(y) = ay^2 + by + c,$$

yields

$$a = -\frac{1}{2} \frac{f''}{f'^2}$$

$$b = \frac{1}{f'} + \frac{f''}{f'^2} f$$

$$c = x - \frac{f}{f'} - \frac{1}{2} \frac{f''}{f'^2} f^2.$$

We obtain the iteration formula for $f(x)$ - a = 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

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

yields

$$a = \frac{f'^2}{f''} \exp - \frac{f''}{f'} x$$

$$b = f''/f'$$

$$c = f - f'^2/f''.$$

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

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

Finally putting $x_n = x$ and $x_{n+1} = h(a)$ yields again (6.31).

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 (6.31) for

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

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

(6.32)

i.e., for Halley's method the transformation (6.32) has no effect. But as we saw (6.32) 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

$$a = f'/2, \quad b = f' - f''x$$

$$c = f - f'x + \frac{f''}{2} x^2$$

and

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

with

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

with

$$G(x) = -\ln(1 - t(x))/t(x)$$

$$t(x) = (f(x) - a) f''(x)/f'(x)^2$$

From Theorem 6.1 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 $f(h) = \alpha^2$ is Reinsch's

proposal

$$\lambda_{n+1} = \lambda_n - \frac{f(\lambda_n) - \alpha^2}{f'(\lambda_n)} G(\lambda_n) \quad (6.33)$$

where

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

Iteration (6.33) 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 λ_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)

$$\min \|x\|$$

subject to

$$\|Ax - b\| \leq \alpha$$

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

If we start with $\lambda_0 = 0.1$ we get $\lambda_1 = -0.7973$ instead of a value bigger than λ_0 . The other third order methods involving a term

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

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

$$H(t) = 1 + \frac{1}{t} t \quad (6.34)$$

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

An improvement of (6.34) is to use

$$H(t) = \exp(t/t) \quad (6.35)$$

However global convergence has to be accelerated by a convergence factor h .

Therefore we propose to use

$$\lambda_{n+1} = \lambda_n - \frac{f'' - \alpha^2}{f''} \exp\left(\frac{1}{\lambda} \frac{f''}{f'}\right) \cdot H(\lambda_n) \quad (6.36)$$

with

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

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 (6.33). For the

above mentioned example with the **Hilbert** matrix we obtain using (6.36) and $\lambda_0 = 0.1$ the solution of the secular equation $\lambda \approx 7.7 \times 10^{-11}$ in four iterations.

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

$$f(\lambda) = 1$$

where

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

We used ALGOL W and double precision (ca. 16 decimal digits) on the **IBM 360 of SIAC**. The results are displayed in Table 1. Newton's method uses 21 steps to converge to machine precision. **Reinsch's** variant converges after seven **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 $\ln(f/\alpha) = 0$
$F(\lambda) = \lambda - \frac{g}{g'}$	$A - \frac{f}{f'}, 2(\frac{f}{\alpha} - 1)$	$\lambda - \frac{f}{f'}, \ln(\frac{f}{\alpha})$
0.00783209855661240	5.20297243199116	0
0.0209101056943383	7.50002221683766	0.0634187403647834
0.0424438763172333	9.6945598729093	0.53074605218063
0.077467170580476	10.2422619393798	1.97352054359389
0.133629409609385	10.2699361392028	4.58739426667985
0.223664677428177	10.2700022239607	7.47964284219632
0.365485460269776	10.2700022243362	9.52785436035909
0.567031564547462	10.2700022243362	10.2135824785941
0.926977169378511	10.2700022243362	10.26966838666710
1.44940702417750	10.2700022243362	10.2700022126256
2.22759607120317	10.2700022243362	10.2700022243362
3.36604539153582	10.2700022243362	10.2700022243362
4.92591522457886	10.2700022243362	10.2700022243362
6.87134065977040	10.2700022243362	10.2700022243362
8.79350315070676	10.2700022243362	10.2700022243362
9.97230229940186	10.2700022243362	10.2700022243362
10.25742259482621	10.2700022243362	10.2700022243362
10.2669795745767	10.2700022243362	10.2700022243362
10.2700022242627	10.2700022243362	10.2700022243362

Third order Methods $\lambda_{n+1} = \lambda_n - \frac{g}{g'}$, $H(t)$ with $t := \frac{g}{g'}$

The following two methods fail for starting values < 6.8 :

Ostrowski	Euler
$H(t) = 1/(1-t)^{0.5}$	$2/(1 + (1-2t)^{0.5})$
6.87134069977040	6.87134069977040
11.3578564654989	8.79350313070076
10.2615901313624	10.4062426841626
10.2700022241005	10.2699542114743
10.2700022243362	10.2700022243362
10.2700022243362	10.2700022243362

Table 1 (cont.)

Halley	inverse interp.	1 + t/2	exp(t/2)
0.0494307764343190	0.0144232540488059	0.0161702328293991	0
0.253229265175140	0.0460967950941273	0.30975341263060	0
0.986381372252434	0.11349254126652	0.159596942052953	0
3.21756490475624	0.25329961220121	0.413519282769038	0
7.85449566417966	0.536339290697594	0.951636742732932	0
10.2238730991803	1.05310528695540	2.06516964524166	0
10.2700019998135	2.14957041757370	4.26001192403995	0
10.2700022243362	4.03642343012226	7.6528939221572	0
10.2700022243362	6.91272434349204	10.0374630402332	0
10.2700022243362	9.60092062566452	10.2698435997489	0
10.2700022243362	10.2634951753514	10.2700022243362	0
10.2700022243362	10.27000221606612	10.2700022243362	0
10.2700022243362	10.2700022243362	10.2700022243362	0
10.2700022243362	10.2700022243362	10.2700022243362	0

Third order methods with convergence factor $G = \left(\frac{\sqrt{F}}{\alpha} + \sqrt{F} \right) / 2$:

Halley	G / (1 - t/2)	G exp(t/2)
5.07846031438509	1.86673871095732	0
9.89047170050903	5.513033409770731	0
10.2699142055339	9.0900256437327	0
10.2700022543362	10.299326597900	0
10.2700022243362	10.27000221410736	0
10.2700022243362	10.2700022243362	0

7. Generalizations.

Let A be an $(n \times n)$ matrix, b a given n vector, \underline{c} an $m \times n$ matrix, d a given m vector and γ, α real numbers. We consider the problem

$$F(\underline{x}) = \underline{x}^T A \underline{x} + \underline{b}^T \underline{x} + \gamma = \min \tag{7.1}$$

subject to

$$\| \underline{C} \underline{x} - \underline{d} \| = a. \tag{7.2}$$

Using a Lagrange multiplier $\lambda/2$ the solution of (7.1), (7.2) is a stationary point (\underline{x}, λ) of

$$L(\underline{x}, \lambda) = F(\underline{x}) + \frac{\lambda}{2} \{ \| \underline{C} \underline{x} - \underline{d} \|^2 - \alpha^2 \}. \tag{7.3}$$

$\frac{\partial L}{\partial \underline{x}} = 0$ and $\frac{\partial L}{\partial \lambda} = 0$ gives

$$\left. \begin{aligned} (\underline{A} + \underline{A}^T + \lambda \underline{C}^T \underline{C}) \underline{x} &= -\underline{b} + \lambda \underline{C}^T \underline{d} \\ \| \underline{C} \underline{x} - \underline{d} \|^2 &= \alpha^2. \end{aligned} \right\} \tag{7.4}$$

Theorem 7.1. Let $(\underline{x}_i, \lambda_i)$, $i = 1, 2$, be two solutions of (7.4), then

$$F(\underline{x}_1) - F(\underline{x}_2) = \frac{\lambda_2 - \lambda_1}{4} \| \underline{C}(\underline{x}_1 - \underline{x}_2) \|^2. \tag{7.5}$$

Proof. This theorem is a generalization of Theorem 2.1 and the proof is very similar. We have

$$(\underline{A} + \underline{A}^T + \lambda_1 \underline{C}^T \underline{C}) \underline{x}_1 = -\underline{b} + \lambda_1 \underline{C}^T \underline{d} \tag{7.6}$$

$$(\underline{A} + \underline{A}^T + \lambda_2 \underline{C}^T \underline{C}) \underline{x}_2 = -\underline{b} + \lambda_2 \underline{C}^T \underline{d}. \tag{7.7}$$

\underline{x}_1^T (7.6) gives

Theorem 7.2. Let (x_i, λ_i) , $i = 1, 2$, be two solutions of (7.4), then

$$(\lambda_1 + \lambda_2) (F(x_1) - F(x_2)) = (\lambda_2 - \lambda_1) (x_1 - x_2)^T A (x_1 - x_2). \quad (7.11)$$

Proof.

$$\lambda_1 C^T C x_1 - \lambda_1 C^T d = -(A + \frac{\lambda_1}{A}) x_1 - b \quad (7.12)$$

$$\lambda_2 C^T C x_2 - \lambda_2 C^T d = -(A + \frac{\lambda_2}{A}) x_2 - b. \quad (7.13)$$

$\lambda_1 x_1^T (7.6) - \lambda_2 x_2^T (7.12)$ gives

$$\lambda_1 \lambda_2 \{ (x_2 - x_1)^T C^T d \} = (\lambda_2 - \lambda_1) x_1^T (A + \frac{\lambda_1}{A}) x_2 - (\lambda_1 x_1 - \lambda_2 x_2)^T b. \quad (7.14)$$

$\lambda_1 x_1^T (7.13) - \lambda_2 x_2^T (7.12)$ gives

$$\begin{aligned} \lambda_1 \lambda_2 \{ \|C x_2\|^2 - \|C x_1\|^2 + (x_1 - x_2)^T C^T d \} \\ = -2\lambda_1 x_2^T A x_2 + 2\lambda_2 x_1^T A x_1 - (\lambda_1 x_2 - \lambda_2 x_1)^T b. \end{aligned} \quad (7.15)$$

Observe that

$$0 = \|C x_2 - d\|^2 - \|C x_1 - d\|^2 = \|C x_2\|^2 - \|C x_1\|^2 + 2(x_1 - x_2)^T C^T d.$$

So that if we subtract (7.15) - (7.14) we get

$$\begin{aligned} 0 = 2\lambda_2 x_1^T A x_1 - 2\lambda_1 x_2^T A x_2 - (\lambda_1 x_2 - \lambda_2 x_1 - \lambda_1 x_1 + \lambda_2 x_2)^T b \\ + (\lambda_1 - \lambda_2) x_1^T (A + \frac{\lambda_1}{A}) x_2. \end{aligned}$$

$$2 x_1^T A x_1 + \lambda_1 \|C x_1\|^2 = -x_1^T b + \lambda_1 x_1^T C^T d,$$

or rearranged

$$2(x_1^T A x_1 + x_1^T b) = x_1^T b - \lambda_1 (\|C x_1\|^2 - x_1^T C^T d). \quad (7.8)$$

Similarly we have

$$2(x_2^T A x_2 + x_2^T b) = x_2^T b - \lambda_2 (\|C x_2\|^2 - x_2^T C^T d). \quad (7.9)$$

Now $x_2^T (7.6) - x_1^T (7.7)$ gives

$$\begin{aligned} (\lambda_1 - \lambda_2) x_2^T C^T C x_1 &= -x_2^T b + x_1^T b + \lambda_1 x_2^T C^T d - \lambda_2 x_1^T C^T d \\ \Rightarrow (x_1 - x_2)^T b &= (\lambda_1 - \lambda_2) x_2^T C^T C x_1 - d^T (\lambda_1 C x_2 - \lambda_2 C x_1). \end{aligned} \quad (7.10)$$

Subtracting (7.8) - (7.9) and replacing $(x_1 - x_2)^T b$ by (7.10) we get

$$\begin{aligned} 2(F(x_1) - F(x_2)) &= \\ &- \lambda_1 (\|C x_1\|^2 - x_1^T C^T d - x_2^T C^T C x_1 + x_2^T C^T d) \\ &+ \lambda_2 (\|C x_2\|^2 - x_2^T C^T d - x_2^T C^T C x_1 + x_1^T C^T d). \end{aligned}$$

Now like in the proof of Theorem 2.1 we know that both $\{ \}$ are equal and therefore we replace them by their arithmetic mean which gives the desired result. \square

Corollary 7.1. The solution of (7.1), (7.2) is the solution (x, λ) of

(7.4) with largest A .

Proof. From (7.5) we have

$$\lambda_2 < A_1 \Rightarrow F(x_2) > F(x_1). \quad \square$$

Rearranged we have

$$\begin{aligned} & \lambda_1 \{ 2 \underline{x}_2^T \underline{A} \underline{x}_2 + \underline{x}_2^T \underline{b} - \underline{x}_1^T \underline{b} - \underline{x}_1^T (\underline{A} + \underline{A}^T) \underline{x}_2 \} \\ & = \lambda_2 \{ 2 \underline{x}_1^T \underline{A} \underline{x}_1 + \underline{x}_1^T \underline{b} - \underline{x}_2^T \underline{b} - \underline{x}_1^T (\underline{A} + \underline{A}^T) \underline{x}_2 \} \end{aligned} \quad (7.16)$$

Now observe that the left hand side of (7.16) is

$$\{ \} = F(\underline{x}_2) - F(\underline{x}_1) + (\underline{x}_1 - \underline{x}_2)^T \underline{A} (\underline{x}_1 - \underline{x}_2) \cdot$$

Similarly the right hand side simplifies. Therefore

$$\begin{aligned} & \lambda_1 \{ F(\underline{x}_2) - F(\underline{x}_1) + (\underline{x}_1 - \underline{x}_2)^T \underline{A} (\underline{x}_1 - \underline{x}_2) \} \\ & = \lambda_2 \{ F(\underline{x}_1) - F(\underline{x}_2) + (\underline{x}_1 - \underline{x}_2)^T \underline{A} (\underline{x}_1 - \underline{x}_2) \} \end{aligned}$$

or rearranged

$$(\lambda_1 + \lambda_2) (F(\underline{x}_1) - F(\underline{x}_2)) = (\lambda_1 - \lambda_2) (\underline{x}_1 - \underline{x}_2)^T \underline{A} (\underline{x}_1 - \underline{x}_2) \cdot \square$$

Corollary 7.2. Let $(\underline{x}_1, \lambda_1)$, $i = 1, 2$, be two solutions of (7.4) with $\lambda_1 \neq \lambda_2$. Then

$$(\underline{x}_1 - \underline{x}_2)^T \underline{A} (\underline{x}_1 - \underline{x}_2) = -\frac{\lambda_1 + \lambda_2}{4} \|C(\underline{x}_1 - \underline{x}_2)\|^2 \cdot \quad (7.17)$$

Proof. We combine the results of Theorems 7.1 and 7.2. \square ;

8. Smoothing of Datas.

Rutishauser proposed in [31] 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 d_i and that is smoother. If we

assume that the d_i are equidistant then we might want to solve

$$\sum_{i=2}^{n-1} (x_{i+1} - 2x_i + x_{i-1})^2 = \min \quad (8.1)$$

subject to

$$\sum_{i=1}^n (x_i - d_i)^2 \leq n \varepsilon^2 \quad (8.2)$$

Here ε^2 is a measure for the variance, the mean deviation we want to allow the new data x_i to differ from d_i :

$$\varepsilon \approx \sqrt{\frac{\sum_{i=1}^n (x_i - d_i)^2}{n}}$$

The larger ε the more the x_i will be smoothed.

Introducing the $(n-2) \times n$ tridiagonal matrix

$$\underline{A} = \begin{pmatrix} 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & \ddots & \ddots & \ddots & \\ & & & & 1 & -2 & 1 \\ & & & & & 1 & -2 & 1 \end{pmatrix}$$

and the vectors \underline{d} and \underline{x} , equations (8.1) and (8.2) become

$$\left. \begin{aligned} \|\underline{Ax}\| &= \min \\ \text{subject to} \quad \|\underline{x} - \underline{d}\| &\leq \alpha := \sqrt{\sigma} \cdot \delta \end{aligned} \right\} \quad (8.3)$$

and we have problem (P1) for $\underline{b} = \underline{0}$ and $\underline{C} = \underline{I}_n$. The normal equations are

$$(\underline{A}^T \underline{A} + \lambda \underline{I}_n) \underline{x} = \lambda \underline{d} \quad (8.4)$$

$$\|\underline{x} - \underline{d}\|^2 = \alpha^2 \quad (8.5)$$

Instead of solving (8.4), Rutishauser proposed to choose some "smoothing" parameter γ and to solve

$$(\underline{I}_n + \gamma \underline{A}^T \underline{A}) \underline{x} = \underline{d} \quad (8.6)$$

which is (8.4) for $\lambda = 1/\gamma$. The condition number of the matrix in (8.6) is 16γ . For large n one may have to choose γ large which leads to numerical problems [36]. However large γ or small λ may be meaningful: for $\gamma \rightarrow \infty$ ($\lambda \rightarrow 0$) the new values \underline{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 data \underline{d} .

If we make the change of variables in (8.3)

$$\begin{aligned} \underline{w} &:= \underline{x} - \underline{d} \\ \underline{b} &:= -\underline{d} \end{aligned} \quad (8.7)$$

then our problem becomes

$$\left. \begin{aligned} \|\underline{Aw} - \underline{b}\| &= \min \\ \text{subject to} \quad \|\underline{w}\| &\leq \alpha \end{aligned} \right\} \quad (8.8)$$

Problem (8.8) leads to a relaxed least squares problem

$$\left. \begin{aligned} (\underline{A}^T \underline{A} + \lambda \underline{I}_n) \underline{w} &= \underline{A}^T \underline{b} \\ \|\underline{w}\| &= \alpha \end{aligned} \right\} \quad (8.9)$$

We can use the dual equations for (8.9) (see Section 3):

$$(\underline{A} \underline{A}^T + \lambda \underline{I}_n) \underline{z} = -\underline{b} = -\underline{Ad} \quad (8.10)$$

$$\|\underline{A}^T \underline{z}\| = \alpha \quad (8.11)$$

with

$$\underline{w} = -\underline{A}^T \underline{z} \quad (8.12)$$

Observe that now $\lambda \rightarrow 0$ causes no problems since $\underline{A} \underline{A}^T$ is not singular. Furthermore since $\underline{b} = -\underline{Ad}$ we can write (8.10) as

$$\begin{pmatrix} \underline{A}^T \\ \sqrt{\lambda} \underline{I}_n \end{pmatrix} \underline{z} \approx \begin{pmatrix} \underline{d} \\ \underline{0} \end{pmatrix} \quad (8.13)$$

\underline{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 (8.13). For $n = 8$, e.g. the remaining matrix before and after the third step is:

$$\begin{array}{c}
 \begin{array}{cccc}
 0 & 0 & & \\
 0 & 0 & r & \\
 0 & s & t & \\
 1 & -2 & 1 & \\
 1 & -2 & 1 & \\
 1 & -2 & 1 & \\
 & & & 1
 \end{array}
 & \xrightarrow{\text{step 3}} &
 \begin{array}{cccc}
 0 & 0 & & \\
 0 & 0 & 0 & \\
 0 & 0 & 0 & r \\
 0 & s & t & \\
 1 & -2 & 1 & \\
 1 & -2 & 1 & \\
 1 & -2 & 1 & \\
 & & & 1
 \end{array}
 \end{array}$$

The matrix A is not stored, we work only with three simple variables r , s , and t . We use the notation

$$\begin{pmatrix} 0 & a'_2 & a'_3 \\ b'_1 & b'_2 & b'_3 \end{pmatrix} := G \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} \tag{8.14}$$

where G is a Givens matrix

$$G = \begin{pmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{pmatrix}$$

and α has been chosen so that the zero in the left hand side of (8.14) 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

$$\begin{pmatrix} 0 & c_i \\ d_{1_i} & y_i \end{pmatrix} := G_1 \begin{pmatrix} r & c_i \\ \sqrt{\lambda} & 0 \end{pmatrix};$$

$$\begin{pmatrix} 0 & r & c_{i+1} \\ d_{1_i} & d_{2_i} & y_i \end{pmatrix} := G_2 \begin{pmatrix} s & t & c_{i+1} \\ d_{1_i} & 0 & y_i \end{pmatrix};$$

$$\begin{pmatrix} 0 & s & t & c_{i+2} \\ d_{1_i} & d_{2_i} & d_{3_i} & y_i \end{pmatrix} := G_3 \begin{pmatrix} 1 & -2 & 1 & d_{i+2} \\ d_{1_i} & d_{2_i} & 0 & y_i \end{pmatrix};$$

end;

The following algorithm transforms

$$U^T \begin{pmatrix} A \\ \sqrt{\lambda} I \end{pmatrix} = \begin{pmatrix} 0 \\ R \end{pmatrix} \quad \text{and} \quad U^T \begin{pmatrix} d \\ 0 \end{pmatrix} = \begin{pmatrix} c \\ y \end{pmatrix}$$

where U is orthogonal and R upper triangular and tridiagonal:

$$R = \begin{array}{cccc}
 d_{1_1} & d_{2_1} & d_{3_1} & \\
 d_{1_2} & d_{2_2} & d_{3_2} & \dots \\
 & & \dots & d_{3_{n-4}} \\
 & & & d_{2_{n-3}} \\
 & & & d_{1_{n-2}}
 \end{array}$$

We cannot apply all the results of Section 5.4 since f is not of the same form. We have

$$(AA^T + A I)z(k) = -z(k-1), \quad k \geq 1.$$

We shall use Reinsch's proposal and therefore we need only f' :

$$\begin{aligned} f'(A) &= 2 z^T A^T A z \\ &= -2 z^T (z + \lambda z'). \end{aligned}$$

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 $\|x - d\| = 5.4 \times 10^{-4}$ 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

...

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

$$(AA^T + \lambda I)z(\lambda) = A d$$

$$f(\lambda) = \|A^T z(\lambda)\|^2.$$

$$0.2 \frac{1}{\pi} \int_0^{\pi} \sin(x) dx \quad \frac{0.4}{\pi} \approx 1.3$$

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	1.861084	1.763119	1.684705	1.609509	1.496164	1.483359
2.	000287	2.000137	1.922686	1.861354	1.804054	1.727389	1.717025
2.	139301	2.139186	2.081965	2.031419	1.997584	1.955655	1.947423
2.	270316	2.278234	2.240736	2.212497	2.189486	2.182792	2.181764
2.	417336	2.417287	2.398729	2.386040	2.378851	2.406053	2.418491
2.	556347	2.556335	2.555611	2.557375	2.564482	2.619121	2.643821
2.	695338	2.695371	2.711038	2.725835	2.745208	2.815103	2.841743
2.	034349	2.834425	2.864746	2.890921	2.920292	2.992062	3.003115
2.	973354	2.973438	3.016541	3.052362	3.089548	3.154735	3.155210
3.	112376	3.112476	3.166302	3.210031	3.253120	3.310709	3.301294
3.	251370	3.251490	3.313931	3.363794	3.411095	3.464372	3.457992
3.	390381	3.390501	3.459267	3.513468	3.563265	3.613997	3.616866
3.	529382	3.529527	3.602212	3.658838	3.703311	3.754458	3.762860
3.	668403	3.668556	3.742742	3.799843	3.849200	3.883019	3.886796
3.	307439	3.807585	3.880861	3.936624	3.983325	4.002785	3.995223
3.	946456	3.946573	4.016701	4.069475	4.112409	4.120646	4.105030
4.	085469	4.085618	4.150337	4.196711	4.237081	4.241594	4.228400
4.	224488	4.224601	4.282032	4.324522	4.357614	4.364625	4.362045
4.	363507	4.363605	4.411738	4.447031	4.473988	4.404037	4.490513
4.	502520	4.502613	4.539614	4.566405	4.586239	4.594892	4.600982
4.	641560	4.641610	4.665841	4.682961	4.694780	4.697631	4.694865
4.	780583	4.780605	4.790689	4.797174	4.800344	4.797594	4.787060
4.	919610	4.919595	4.914412	4.90512	4.903647	4.899828	4.891946
5.	058637	5.058589	5.037205	5.020307	5.005067	5.003985	5.001230
5.	137662	5.137586	5.159234	5.129760	5.104605	5.103989	5.123779
5.	336688	5.336576	5.280628	5.238020	5.202195	5.192799	5.216838
5.	475712	5.475564	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	5.3601306
$\ \underline{x} - \underline{d} \ $	1.351226	1.350682	1.095428	0.9311236	0.8215933	0.7120381	0.6572663
$\left[\sum (\Delta^2 x_i)^2 \right]^{1/2}$	$6.67 \cdot 10^{-5}$	$1.11 \cdot 10^{-4}$	0.008449	0.0150982	0.021434	0.0445737	0.0796041
$\left[\sum (x_i - \sqrt{i})^2 \right]^{1/2}$	1.205420	1.204810	0.9014646	0.6828826	0.5084146	0.306330	0.3031789
# of iterations	0	1	5	5	7		
$\lambda = f^{-1}(n\delta^2)$	0	$3.83 \cdot 10^{-7}$	$2.79 \cdot 10^{-4}$	$7.60 \cdot 10^{-4}$	$2.01 \cdot 10^{-3}$	$3.15 \cdot 10^{-2}$	$8.89 \cdot 10^{-2}$

```

PROCEDURE SMOOTH(INTEGER VALUE N;
REAL VALUE DELTA;
REAL ARRAY D,X(*));
BEGIN
REAL C,F0,FA,F1,F2,LAMB,LAMBH,WLAMB,R,S,T,ALPHA,ALPHA2,CO,SI,H;
REAL ARRAY C,AZ,AZS(1:N);
REAL ARRAY D1,X,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,E,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)*E(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 := SORT(N)*DELTA; ALPHA2 := ALPHA**2;
LAMB := 0;

```

Table 2 (cont.)

8 = $\| \bar{x} - \hat{d} \|$
 $\| z(x_T^{-1}, z_T^{-1}) \|$
 $\| z(x_T^{-1}, z_T^{-1}) \|$
 $\lambda = \tau^{-1} (m^2)$
of iteration

```

:
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*D(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) ;
  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 ;
END ;
FOR I := N-3,N-2 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) ;
VORH(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 - F0/F1**2*(SQRT(F0)/ALPHA - 1) ;
IF (FA <= F0) OR (F0 < ALPHA2) OR (LAMB <= LAHB)
THEN GOTO FIN ;
FA := F0 ; LAMB := LAMB ;
GOTO IT ;
FIN:
FOR I := 1 STEP 1 UNTIL N DO
  X(I) := D(I) - AZ(I) ;
END SMOOTH ;

```

References

- [1] Björk, A., "Iterative Refinement of Linear Least Squares Solutions I," BIT 7 (1967), 257-278.
- [2] Björk, 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] Bunb, 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] Davies, M. and Dawson, B., "On Global Convergence of Halley's Iteration Formula," Numer. Math. 24 (1975).
- [7] Ehrmann, H., "Konstruktion und Durchführung von Iterationsverfahren hoererer 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., Malcolm, M. A. and Moler, C. B., Computer Methods for Mathematical Computations, Prentice Hall, 147.
- [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 Svecova, 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., "Generalized 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 University 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," Gatlinburg Conference 1977, Asilomar.
- [25] Kalman, R. E., "Algebraic Aspects of the Generalized Inverse . . ." Generalized 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] Ostrzewski, A. M., Solution of Equations and Systems 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 = λBx 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 M. Gutknecht," Band I, II, Birkhaeuser Verlag, 1976.
- [37] Schek, H.-J. and Eggenberger, 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 (142).
- [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 Generalized 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. "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.