

# 8

## Total least squares

In atmospheric remote sensing, near real-time software processors frequently use approximations of the Jacobian matrix in order to speed up the calculation. If the forward model  $\mathbf{F}(\mathbf{x})$  depends on the state vector  $\mathbf{x}$  through some model parameters  $\mathbf{b}_k$ ,

$$\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{b}_1(\mathbf{x}), \dots, \mathbf{b}_N(\mathbf{x})),$$

then, an approximate expression of the Jacobian matrix

$$\mathbf{K} = \sum_{k=1}^N \frac{\partial \mathbf{F}}{\partial \mathbf{b}_k} \frac{\partial \mathbf{b}_k}{\partial \mathbf{x}},$$

can be obtained by assuming that some  $\mathbf{b}_k$  are insensitive to  $\mathbf{x}$ , i.e.,  $\partial \mathbf{b}_k / \partial \mathbf{x} = \mathbf{0}$ . For example, the limb radiance measured by a detector in the ultraviolet or visible spectral domains can be expressed as

$$I(\lambda, \mathbf{x}) = I_{\text{ss}}(\lambda, \mathbf{x}) + I_{\text{ms}}(\lambda, \mathbf{x}) = I_{\text{ss}}(\lambda, \mathbf{x}) [1 + c_{\text{ms}}(\lambda, \mathbf{x})], \quad (8.1)$$

where  $I_{\text{ss}}$  and  $I_{\text{ms}}$  are the single and multiple scattering terms,  $\lambda$  is the wavelength, and  $c_{\text{ms}}$  is a correction factor accounting for the multiple scattering contribution. As the computation of the derivative of  $c_{\text{ms}}$  is quite demanding, the Jacobian matrix calculation may involve only the derivative of  $I_{\text{ss}}$ . Similarly, in a line-by-line model, the absorption coefficient  $C_{\text{abs}m}$  of the gas molecule  $m$  is the product of the line strength  $S_{ml}$  and the normalized line shape function  $g_{ml}$  (cf. (1.12)),

$$C_{\text{abs}m}(\nu, T) = \sum_l S_{ml}(T) g_{ml}(\nu, T),$$

where  $\nu$  is the wavenumber,  $T$  is the temperature, and the summation is over all lines  $l$ . As the most important temperature dependence stems from the line strength, the derivative of the line shape function with respect to the temperature is sometimes ignored.

The total least squares (TLS) method is devoted to the solution of linear problems in which both the coefficient matrix and the data are subject to errors. The linear data model can be expressed as

$$\mathbf{y}^\delta = (\mathbf{K}_\Lambda - \Lambda) \mathbf{x} + \delta,$$

where the matrix  $\mathbf{K}_\Lambda$  is a perturbation of the exact (unknown) matrix  $\mathbf{K}$ ,  $\mathbf{K}_\Lambda = \mathbf{K} + \Lambda$ , and the data are affected by the instrumental noise  $\delta$ .

The TLS method was independently derived in several bodies of work by Golub and Van Loan (1980, 1996), and Van Huffel and Vanderwalle (1991). This literature has advanced the algorithmic and theoretical understanding of the method, as well as its application for computing stable solutions of linear systems of equations with highly ill-conditioned coefficient matrices. In this section we review the truncated and the regularized TLS methods for solving linear ill-posed problems, and reveal the similarity with the Tikhonov regularization. We then present a first attempt to extend the regularized TLS to nonlinear ill-posed problems.

## 8.1 Formulation

The linear model which encapsulates the uncertainties in the data vector and the coefficient matrix is of the form  $\mathbf{K}_\Lambda \mathbf{x} \approx \mathbf{y}^\delta$ . To sketch the TLS method, we introduce the augmented matrix  $\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}$  and consider the homogeneous system of equations

$$\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \mathbf{0}. \quad (8.2)$$

We then assume a singular value decomposition of the  $m \times (n+1)$  matrix,

$$\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} = \bar{\mathbf{U}} \bar{\Sigma} \bar{\mathbf{V}}^T, \quad (8.3)$$

and partition the matrices  $\bar{\mathbf{V}}$  and  $\bar{\Sigma}$  as follows:

$$\bar{\mathbf{V}} = [\bar{\mathbf{v}}_1, \dots, \bar{\mathbf{v}}_{n+1}] = \begin{bmatrix} \bar{\mathbf{V}}_{11} & \bar{\mathbf{v}}_{12} \\ \bar{\mathbf{v}}_{21}^T & \bar{v}_{22} \end{bmatrix}, \quad \bar{\mathbf{V}}_{11} \in \mathbb{R}^{n \times n}, \quad \bar{\mathbf{v}}_{12}, \bar{\mathbf{v}}_{21} \in \mathbb{R}^n, \quad (8.4)$$

and

$$\bar{\Sigma} = \begin{bmatrix} \bar{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \bar{\sigma}_{n+1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \bar{\Sigma}_1 = [\text{diag}(\bar{\sigma}_i)_{n \times n}],$$

respectively. If  $\bar{\sigma}_{n+1} \neq 0$ , then  $\text{rank}(\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}) = n+1$ , and the solution of the homogeneous system of equations (8.2) is the trivial solution. Thus, the last component of the solution vector is not  $-1$ , and to solve (8.2) it is necessary to reduce the rank of the augmented matrix from  $n+1$  to  $n$ . This can be achieved by approximating the rank- $(n+1)$  matrix  $\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}$  by a rank- $n$  matrix  $\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix}$ . As  $\text{rank}(\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix}) = n$ , we may assume that the last column vector of the matrix  $\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix}$  is a linear combination of the first  $n$  column vectors, i.e.,

$$\mathbf{y}_n = \sum_{i=1}^n x_i \mathbf{k}_i,$$

with  $\mathbf{K}_n = [\mathbf{k}_1, \dots, \mathbf{k}_n]$ , or equivalently that,

$$\mathbf{K}_n \mathbf{x} = \mathbf{y}_n,$$

with  $\mathbf{x} = [x_1, \dots, x_n]^T$ . The (matrix) approximation problem can be expressed as the constrained minimization problem

$$\begin{aligned} \min_{[\tilde{\mathbf{K}} \tilde{\mathbf{y}}] \in \mathbb{R}^{m \times (n+1)}} \quad & \left\| \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{K}} & \tilde{\mathbf{y}} \end{bmatrix} \right\|_F^2 \\ \text{subject to} \quad & \tilde{\mathbf{K}}\mathbf{x} = \tilde{\mathbf{y}}, \end{aligned} \tag{8.5}$$

where the Frobenius norm of the  $m \times n$  matrix  $\mathbf{A}$  is defined by

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n [\mathbf{A}]_{ij}^2}.$$

It should be pointed out that the ordinary least squares method minimizes the norm of the residual vector  $\mathbf{y}^\delta - \tilde{\mathbf{y}}$  under the assumption that  $\mathbf{K}_\Lambda = \tilde{\mathbf{K}}$ .

The solution to the minimization problem (8.5) is given by the Eckart–Young–Mirsky theorem (Golub and Van Loan, 1996): the matrix

$$\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix} = \sum_{i=1}^n \bar{\sigma}_i \bar{\mathbf{u}}_i \bar{\mathbf{v}}_i^T \tag{8.6}$$

is the closest rank- $n$  matrix to  $\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}$ , and we have

$$\left\| \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix} \right\|_F = \bar{\sigma}_{n+1} \bar{\mathbf{u}}_{n+1} \bar{\mathbf{v}}_{n+1}^T,$$

yielding

$$\left\| \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix} \right\|_F = \bar{\sigma}_{n+1}.$$

The homogeneous system of equations (8.2) is then replaced by a homogeneous system of equations involving the rank- $n$  matrix  $\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix}$ , that is,

$$\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \mathbf{0}. \tag{8.7}$$

Since (cf. (8.6))

$$\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix} \bar{\mathbf{v}}_{n+1} = \sum_{i=1}^n \bar{\sigma}_i (\bar{\mathbf{v}}_i^T \bar{\mathbf{v}}_{n+1}) \bar{\mathbf{u}}_i = \mathbf{0}, \tag{8.8}$$

we see that the vector  $a\bar{\mathbf{v}}_{n+1}$  is the general solution of the homogeneous system of equations (8.7) and that the scalar  $a$  is (uniquely) determined by imposing that the last component of the solution vector is  $-1$ . We obtain

$$\begin{bmatrix} \mathbf{x}_\Lambda^\delta \\ -1 \end{bmatrix} = -\frac{1}{[\bar{\mathbf{v}}_{n+1}]_{n+1}} \bar{\mathbf{v}}_{n+1}, \tag{8.9}$$

provided that  $[\bar{\mathbf{v}}_{n+1}]_{n+1} \neq 0$ . From (8.4), we find that the TLS solution can be expressed as

$$\mathbf{x}_\Lambda^\delta = -\frac{1}{\bar{v}_{22}} \bar{\mathbf{v}}_{12}. \tag{8.10}$$

Note that if  $\bar{\sigma}_{n+1}$  is a simple singular value, we have (cf. (8.8))  $\mathcal{N}(\begin{bmatrix} \mathbf{K}_n & \mathbf{y}_n \end{bmatrix}) = \text{span}\{\bar{\mathbf{v}}_{n+1}\}$ , and the TLS solution is unique.

## 8.2 Truncated total least squares

The truncated TLS method, which in general is devoted to numerically rank deficient problems, is also a suitable regularization method for discrete ill-posed problems. This technique is similar to the truncated SVD that treats small singular values of  $\mathbf{K}$  as zeros. In both methods, the redundant information in  $\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}$  and  $\mathbf{K}$ , respectively, associated to the small singular values, is discarded and the original ill-posed problem with a full rank matrix is replaced by a well-posed problem with a rank-deficient matrix. This approximation is achieved by means of the Eckart–Young–Mirsky theorem. For example, in the truncated SVD, the matrix  $\mathbf{K}$  with  $\text{rank}(\mathbf{K}) = n$  and singular value decomposition

$$\mathbf{K} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

is replaced by the matrix

$$\mathbf{K}_p = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^T,$$

with  $\text{rank}(\mathbf{K}_p) = p$ , and the regularized solution takes the form

$$\mathbf{x}_p^\delta = \arg \min_{\mathbf{x}} \|\mathbf{y}^\delta - \mathbf{K}_p \mathbf{x}\|^2 = \sum_{i=1}^p \frac{1}{\sigma_i} (\mathbf{u}_i^T \mathbf{y}^\delta) \mathbf{v}_i.$$

The major difference between the two methods lies in the way in which the approximation is performed: in the truncated SVD, the modification depends only on  $\mathbf{K}$ , while in the truncated TLS, the modification depends on both  $\mathbf{K}_\Lambda$  and  $\mathbf{y}^\delta$ . Thus, in the framework of the truncated TLS method we approximate the matrix  $\begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix}$  by the rank- $p$  matrix

$$\begin{bmatrix} \mathbf{K}_p & \mathbf{y}_p \end{bmatrix} = \sum_{i=1}^p \bar{\sigma}_i \bar{\mathbf{u}}_i \bar{\mathbf{v}}_i^T.$$

To determine the number  $p$  of large singular values or the truncation index, we may require a user-specified threshold or determine  $p$  adaptively. The null space of the approximation matrix is

$$\mathcal{N}(\begin{bmatrix} \mathbf{K}_p & \mathbf{y}_p \end{bmatrix}) = \text{span} \{ \bar{\mathbf{v}}_{p+1}, \dots, \bar{\mathbf{v}}_{n+1} \},$$

whence accounting for the partition

$$\bar{\mathbf{V}} = [\bar{\mathbf{v}}_1, \dots, \bar{\mathbf{v}}_{n+1}] = \begin{bmatrix} \bar{\mathbf{V}}_{11} & \bar{\mathbf{V}}_{12} \\ \bar{\mathbf{v}}_{21}^T & \bar{\mathbf{v}}_{22}^T \end{bmatrix}, \quad (8.11)$$

with  $\bar{\mathbf{V}}_{11} \in \mathbb{R}^{n \times p}$ ,  $\bar{\mathbf{V}}_{12} \in \mathbb{R}^{n \times (n-p+1)}$ , and

$$\begin{aligned} \bar{\mathbf{v}}_{21} &= \left[ [\bar{\mathbf{v}}_1]_{n+1}, \dots, [\bar{\mathbf{v}}_p]_{n+1} \right]^T \in \mathbb{R}^p, \\ \bar{\mathbf{v}}_{22} &= \left[ [\bar{\mathbf{v}}_{p+1}]_{n+1}, \dots, [\bar{\mathbf{v}}_{n+1}]_{n+1} \right]^T \in \mathbb{R}^{n-p+1}, \end{aligned}$$

we seek the solution as

$$\begin{bmatrix} \mathbf{x}_{\Lambda p}^\delta \\ -1 \end{bmatrix} = \sum_{i=p+1}^{n+1} a_i \bar{\mathbf{v}}_i = \begin{bmatrix} \bar{\mathbf{V}}_{12} \\ \bar{\mathbf{v}}_{22}^T \end{bmatrix} \mathbf{a}, \quad (8.12)$$

with  $\mathbf{a} = [a_{p+1}, \dots, a_{n+1}]^T \in \mathbb{R}^{n-p+1}$ . From the last equation we find that

$$\bar{\mathbf{v}}_{22}^T \mathbf{a} = -1,$$

or equivalently that

$$\sum_{i=p+1}^{n+1} a_i [\bar{\mathbf{v}}_i]_{n+1} = -1.$$

Since (cf. (8.12))

$$\left\| \begin{bmatrix} \mathbf{x}_{\Lambda p}^\delta \\ -1 \end{bmatrix} \right\|^2 = 1 + \|\mathbf{x}_{\Lambda p}^\delta\|^2 = \sum_{i=p+1}^{n+1} a_i^2, \quad (8.13)$$

we see that the minimum norm solution  $\mathbf{x}_{\Lambda p}^\delta$  requires a minimum value of  $\sum_{i=p+1}^{n+1} a_i^2$ . This can be obtained by solving the constrained minimization problem

$$\begin{aligned} \min_{a_i} \quad & \sum_{i=p+1}^{n+1} a_i^2 \\ \text{subject to} \quad & \sum_{i=p+1}^{n+1} a_i [\bar{\mathbf{v}}_i]_{n+1} = -1. \end{aligned}$$

In the framework of the Lagrange multiplier formalism, the first-order optimality conditions for the Lagrangian function

$$\mathcal{L}(\mathbf{a}, \lambda) = \frac{1}{2} \sum_{i=p+1}^{n+1} a_i^2 + \lambda \left( \sum_{i=p+1}^{n+1} a_i [\bar{\mathbf{v}}_i]_{n+1} + 1 \right),$$

yield

$$\begin{aligned} a_i + \lambda [\bar{\mathbf{v}}_i]_{n+1} &= 0, \quad i = p+1, \dots, n+1, \\ \sum_{i=p+1}^{n+1} a_i [\bar{\mathbf{v}}_i]_{n+1} &= -1, \end{aligned}$$

and we obtain

$$\mathbf{a} = -\frac{1}{\|\bar{\mathbf{v}}_{22}\|^2} \bar{\mathbf{v}}_{22}. \quad (8.14)$$

Hence, from (8.12) and (8.14), the minimum norm solution is given by

$$\mathbf{x}_{\Lambda p}^\delta = -\frac{1}{\|\bar{\mathbf{v}}_{22}\|^2} \bar{\mathbf{V}}_{12} \bar{\mathbf{v}}_{22}. \quad (8.15)$$

By (8.13), (8.14) and the Eckart–Young–Mirsky theorem, we have

$$\|\mathbf{x}_{\Lambda p}^\delta\|^2 = \frac{1}{\|\bar{\mathbf{v}}_{22}\|^2} - 1,$$

and

$$\|\mathbf{R}_{\Lambda p}^\delta\|_F^2 = \|[ \mathbf{K}_\Lambda \quad \mathbf{y}^\delta ] - [ \mathbf{K}_p \quad \mathbf{y}_p ]\|_F^2 = \bar{\sigma}_{p+1}^2 + \dots + \bar{\sigma}_{n+1}^2,$$

showing that the solution norm  $\|\mathbf{x}_{\Lambda p}^\delta\|$  increases monotonically with  $p$ , while the residual norm  $\|\mathbf{R}_{\Lambda p}^\delta\|_F$  decreases monotonically with  $p$ . These results recommend the discrepancy principle and the L-curve method for computing the truncation index.

In order to demonstrate the regularizing property of the truncated TLS method, we express  $\mathbf{x}_{\Lambda p}^\delta$  as the filtered sum

$$\mathbf{x}_{\Lambda p}^\delta = \sum_{i=1}^n f_i \frac{1}{\sigma_i} (\mathbf{u}_i^T \mathbf{y}^\delta) \mathbf{v}_i, \quad (8.16)$$

where  $(\sigma_i; \mathbf{v}_i, \mathbf{u}_i)$  is a singular system of  $\mathbf{K}_\Lambda$ . In Appendix I it is shown that if  $\text{rank}(\mathbf{K}_\Lambda) = n$  and  $\text{rank}([ \mathbf{K}_\Lambda \quad \mathbf{y}^\delta ]) = n + 1$ , and furthermore, if  $\mathbf{u}_i^T \mathbf{y}^\delta \neq 0$  for all  $i = 1, \dots, n$ , then the filter factors are given by

$$f_i = \frac{1}{\|\bar{\mathbf{v}}_{22}\|^2} \sum_{j=1}^p \frac{\sigma_j^2}{\bar{\sigma}_j^2 - \sigma_i^2} [\bar{\mathbf{v}}_j]_{n+1}^2, \quad (8.17)$$

and the estimates

$$1 < f_i \leq 1 + \left( \frac{\bar{\sigma}_{p+1}}{\sigma_i} \right)^2 + O\left( \frac{\bar{\sigma}_{p+1}^4}{\sigma_i^4} \right), \quad i = 1, \dots, p, \quad (8.18)$$

and

$$0 < f_i \leq \frac{1 - \|\bar{\mathbf{v}}_{22}\|^2}{\|\bar{\mathbf{v}}_{22}\|^2} \left( \frac{\sigma_i}{\bar{\sigma}_p} \right)^2 \left[ 1 + O\left( \frac{\sigma_i^2}{\bar{\sigma}_p^2} \right) \right], \quad i = p + 1, \dots, n \quad (8.19)$$

hold. From (8.18), (8.19) and the interlacing property of the singular values of  $[ \mathbf{K}_\Lambda \quad \mathbf{y}^\delta ]$  and  $\mathbf{K}_\Lambda$ ,

$$\bar{\sigma}_1 > \sigma_1 > \dots > \bar{\sigma}_p > \sigma_p > \bar{\sigma}_{p+1} > \sigma_{p+1} > \dots > \sigma_n > \bar{\sigma}_{n+1},$$

we see that for  $i \ll p$ ,  $(\bar{\sigma}_{p+1}/\sigma_i)^2 \ll 1$  and the filter factors are close to 1, while for  $i \gg p$ ,  $(\sigma_i/\bar{\sigma}_p)^2 \ll 1$  and the filter factors are very small. Thus, the filter factors of the truncated TLS method resemble the Tikhonov filter factors, and  $\mathbf{x}_{\Lambda p}^\delta$  is a filtered solution, with the truncation index  $p$  playing the role of the regularization parameter.

When the dimension of  $\mathbf{K}_\Lambda$  is not too large, the singular value decomposition of the augmented matrix  $[ \mathbf{K}_\Lambda \quad \mathbf{y}^\delta ]$  can be computed directly. For large-scale problems, this approach is computationally expensive and an iterative algorithm based on Lanczos bidiagonalization can be used instead (Fierro et al., 1997). The so-called Lanczos truncated TLS

algorithm uses the Lanczos bidiagonalization of the matrix  $\mathbf{K}_\Lambda$  to obtain, after  $p$  iteration steps, the factorization

$$\mathbf{K}_\Lambda \bar{\mathbf{V}}_p = \bar{\mathbf{U}}_{p+1} \mathbf{B}_p, \quad (8.20)$$

and projects the TLS problem onto the subspace spanned by  $\bar{\mathbf{U}}_{p+1} \in \mathbb{R}^{m \times (p+1)}$  and  $\bar{\mathbf{V}}_p \in \mathbb{R}^{n \times p}$ . The projection is a consequence of the assumption that for a sufficiently large  $p$ , all the large singular values of  $\mathbf{K}_\Lambda$ , which contribute to the regularized solution, have been captured. The projected TLS problem reads as

$$\begin{aligned} & \min_{[\tilde{\mathbf{K}}_p \tilde{\mathbf{y}}_p] \in \mathbb{R}^{m \times (n+1)}} \left\| \bar{\mathbf{U}}_{p+1}^T \left( \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{K}}_p & \tilde{\mathbf{y}}_p \end{bmatrix} \right) \begin{bmatrix} \bar{\mathbf{V}}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \right\|_{\mathbb{F}}^2 \\ & \text{subject to } \bar{\mathbf{U}}_{p+1}^T \tilde{\mathbf{K}}_p \bar{\mathbf{V}}_p \mathbf{z}_p = \bar{\mathbf{U}}_{p+1}^T \tilde{\mathbf{y}}_p, \end{aligned}$$

where we have set  $\mathbf{x} = \bar{\mathbf{V}}_p \mathbf{z}_p$  for some  $\mathbf{z}_p \in \mathbb{R}^p$ . Using the result (cf. (8.20) and (5.36))

$$\bar{\mathbf{U}}_{p+1}^T \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} \begin{bmatrix} \bar{\mathbf{V}}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{U}}_{p+1}^T \mathbf{K}_\Lambda \bar{\mathbf{V}}_p & \bar{\mathbf{U}}_{p+1}^T \mathbf{y}^\delta \end{bmatrix} = \begin{bmatrix} \mathbf{B}_p & \beta_1 \mathbf{e}_1^{(p+1)} \end{bmatrix},$$

the constrained minimization problem can be rewritten as

$$\begin{aligned} & \min_{[\tilde{\mathbf{B}}_p \tilde{\mathbf{e}}_p] \in \mathbb{R}^{(p+1) \times (p+1)}} \left\| \begin{bmatrix} \mathbf{B}_p & \beta_1 \mathbf{e}_1^{(p+1)} \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{B}}_p & \tilde{\mathbf{e}}_p \end{bmatrix} \right\|_{\mathbb{F}}^2 \\ & \text{subject to } \tilde{\mathbf{B}}_p \mathbf{z}_p = \tilde{\mathbf{e}}_p, \end{aligned} \quad (8.21)$$

where we have put  $\tilde{\mathbf{B}}_p = \bar{\mathbf{U}}_{p+1}^T \tilde{\mathbf{K}}_p \bar{\mathbf{V}}_p$  and  $\tilde{\mathbf{e}}_p = \bar{\mathbf{U}}_{p+1}^T \tilde{\mathbf{y}}_p$ . Thus, in each Lanczos step, we use the TLS algorithm for the small-scale problem (8.21) to compute a truncated TLS solution  $\mathbf{x}_{\Lambda p}^\delta$ . More precisely, assuming the singular value decomposition

$$\begin{bmatrix} \mathbf{B}_p & \beta_1 \mathbf{e}_1^{(p+1)} \end{bmatrix} = \bar{\bar{\mathbf{U}}} \bar{\bar{\Sigma}} \bar{\bar{\mathbf{V}}},$$

with

$$\bar{\bar{\mathbf{V}}} = \begin{bmatrix} \bar{\bar{\mathbf{V}}}_{11} & \bar{\bar{\mathbf{v}}}_{12} \\ \bar{\bar{\mathbf{v}}}_{21}^T & \bar{\bar{v}}_{22} \end{bmatrix}, \quad \bar{\bar{\mathbf{V}}}_{11} \in \mathbb{R}^{p \times p}, \quad \bar{\bar{\mathbf{v}}}_{12}, \bar{\bar{\mathbf{v}}}_{21} \in \mathbb{R}^p,$$

the TLS solution to (8.21) is (cf. (8.10))

$$\mathbf{z}_{\Lambda p}^\delta = -\frac{1}{\bar{\bar{v}}_{22}} \bar{\bar{\mathbf{v}}}_{12},$$

and the truncated TLS solution takes the form

$$\mathbf{x}_{\Lambda p}^\delta = \bar{\mathbf{V}}_p \mathbf{z}_{\Lambda p}^\delta = -\frac{1}{\bar{\bar{v}}_{22}} \bar{\mathbf{V}}_p \bar{\bar{\mathbf{v}}}_{12}.$$

In the Lanczos truncated TLS algorithm, the solution norm and the residual norm also possess monotonic behavior, i.e.,  $\|\mathbf{x}_{\Lambda p}^\delta\|$  is an increasing function of  $p$ , while  $\|\mathbf{r}_{\Lambda p}^\delta\|_{\mathbb{F}}$  is a decreasing function of  $p$  (Fierro et al., 1997).

Regularization parameter choice methods for truncated TLS are discrete methods. If explicit knowledge about the errors in  $\mathbf{K}_\Lambda$  and  $\mathbf{y}^\delta$  is available, the discrepancy principle can be used to compute the truncation index. When the errors in  $\mathbf{K}_\Lambda$  and  $\mathbf{y}^\delta$  are not available, error-free parameter choice methods can be employed. In this context, we mention that Sima and Van Huffel (2006) formulated the generalized cross-validation in the framework of the Lanczos truncated TLS, while the L-curve method has been applied by Fierro et al. (1997).

The truncated solution  $\mathbf{x}_{\Lambda p}^\delta$  is a filtered solution whose main contributions come from the first  $p$  singular vectors of  $\bar{\mathbf{K}}_\Lambda$  (Appendix I). Because these vectors are not always the best basis vectors for a regularized solution, we may implicitly include regularization in general form with  $\mathbf{L} \neq \mathbf{I}_n$ . This is done by transforming the problem involving  $\mathbf{K}_\Lambda$  and  $\mathbf{L}$  into a standard-form problem with the matrix  $\bar{\mathbf{K}}_\Lambda = \mathbf{K}_\Lambda \mathbf{L}^{-1}$ . Then, we apply the truncated TLS method to the standard-form problem to obtain a regularized solution  $\bar{\mathbf{x}}_{\Lambda p}^\delta$ , and finally, we transform  $\bar{\mathbf{x}}_{\Lambda p}^\delta$  back to the general-form setting by computing  $\mathbf{x}_{\Lambda p}^\delta = \mathbf{L}^{-1} \bar{\mathbf{x}}_{\Lambda p}^\delta$ . The conventional and the Lanczos versions of the truncated TLS method are outlined in Algorithms 13 and 14. It should be remarked that Algorithm 13 computes simultaneously the truncated SVD solution and the truncated TLS solution for a fixed value of the truncation index  $p$ .

### 8.3 Regularized total least squares for linear problems

Tikhonov regularization has been recast in the framework of the regularized TLS by Golub et al. (1999). To stress the differences and the similarities between the conventional Tikhonov regularization and the regularized TLS, we first note that Tikhonov regulariza-

---

**Algorithm 13.** Algorithm for computing the truncated SVD solution  $\mathbf{x}_p^\delta$  and the truncated TLS solution  $\mathbf{x}_{\Lambda p}^\delta$  for a fixed value of the truncation index  $p$ .

---

$\bar{\mathbf{K}}_\Lambda \leftarrow \mathbf{K}_\Lambda \mathbf{L}^{-1};$   
 {truncated SVD solution}  
 compute the SVD  $\bar{\mathbf{K}}_\Lambda = \mathbf{U} \Sigma \mathbf{V}^T;$   
 $\bar{\mathbf{x}}_p^\delta \leftarrow \sum_{i=1}^p (1/\sigma_i) (\mathbf{u}_i^T \mathbf{y}^\delta) \mathbf{v}_i;$   
 $\mathbf{x}_p^\delta \leftarrow \mathbf{L}^{-1} \bar{\mathbf{x}}_p^\delta;$   
 {truncated TLS solution}  
 compute the SVD  $[\bar{\mathbf{K}}_\Lambda \quad \mathbf{y}^\delta] = \bar{\mathbf{U}} \bar{\Sigma} \bar{\mathbf{V}}^T;$   
 partition  $\bar{\mathbf{V}} = \begin{bmatrix} \bar{\mathbf{V}}_{11} & \bar{\mathbf{V}}_{12} \\ \bar{\mathbf{v}}_{21}^T & \bar{\mathbf{v}}_{22}^T \end{bmatrix}$  with  $\bar{\mathbf{V}}_{11} \in \mathbb{R}^{n \times p};$   
 $\bar{\mathbf{x}}_{\Lambda p}^\delta \leftarrow - \left( 1 / \|\bar{\mathbf{v}}_{22}\|^2 \right) \bar{\mathbf{V}}_{12} \bar{\mathbf{v}}_{22};$   
 $\mathbf{x}_{\Lambda p}^\delta \leftarrow \mathbf{L}^{-1} \bar{\mathbf{x}}_{\Lambda p}^\delta;$

---



---

**Algorithm 14.** Lanczos truncated TLS algorithm with  $p_{\max} > 1$  iterations.
 

---

$\beta_1 \leftarrow \|\mathbf{y}^\delta\|$ ;  $\bar{\mathbf{u}} \leftarrow (1/\beta_1)\mathbf{y}^\delta$ ;  
 $\mathbf{q} \leftarrow \mathbf{L}^{-T}\mathbf{K}^T\bar{\mathbf{u}}$ ;  $\alpha_1 \leftarrow \|\mathbf{q}\|$ ;  $\bar{\mathbf{v}}_1 \leftarrow (1/\alpha_1)\mathbf{q}$ ;  
**for**  $p = 1, p_{\max}$  **do**  
    $\mathbf{p} \leftarrow \mathbf{KL}^{-1}\bar{\mathbf{v}}_p - \alpha_p\bar{\mathbf{u}}$ ;  $\beta_{p+1} \leftarrow \|\mathbf{p}\|$ ;  $\bar{\mathbf{u}} \leftarrow (1/\beta_{p+1})\mathbf{p}$ ;  
   **if**  $p > 1$  **then**  
     set  $\mathbf{A} = \begin{bmatrix} \mathbf{B}_p & \beta_1\mathbf{e}_1^{(p+1)} \end{bmatrix} = \begin{bmatrix} \alpha_1 & 0 & \dots & 0 & \beta_1 \\ \beta_2 & \alpha_2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \alpha_p & 0 \\ 0 & 0 & \dots & \beta_{p+1} & 0 \end{bmatrix}$ ;  
     compute the SVD  $\mathbf{A} = \bar{\bar{\mathbf{U}}}\bar{\bar{\Sigma}}\bar{\bar{\mathbf{V}}}$ ;  
     partition  $\bar{\bar{\mathbf{V}}} = \begin{bmatrix} \bar{\bar{\mathbf{V}}}_{11} & \bar{\bar{\mathbf{V}}}_{12} \\ \bar{\bar{\mathbf{v}}}_{21}^T & \bar{\bar{v}}_{22} \end{bmatrix}$  with  $\bar{\bar{\mathbf{V}}}_{11} \in \mathbb{R}^{p \times p}$ ;  
      $\bar{\mathbf{x}}_{\Lambda p}^\delta \leftarrow -(1/\bar{\bar{v}}_{22})\sum_{j=1}^p [\bar{\bar{\mathbf{V}}}_{12}]_j \bar{\mathbf{v}}_j$ ;  
      $\mathbf{x}_{\Lambda p}^\delta \leftarrow \mathbf{L}^{-1}\bar{\mathbf{x}}_{\Lambda p}^\delta$ ;  
   **end if**  
   **if**  $p < p_{\max}$  **then**  
      $\mathbf{q} \leftarrow \mathbf{L}^{-T}\mathbf{K}^T\bar{\mathbf{u}} - \beta_{p+1}\bar{\mathbf{v}}_p$ ;  $\alpha_{p+1} \leftarrow \|\mathbf{q}\|$ ;  $\bar{\mathbf{v}}_{p+1} \leftarrow (1/\alpha_{p+1})\mathbf{q}$ ;  
   **end if**  
**end for**

---

tion has an important equivalent formulation as

$$\begin{aligned} \min_{\mathbf{x}} \|\mathbf{y}^\delta - \mathbf{K}\mathbf{x}\|^2 & \quad (8.22) \\ \text{subject to } \|\mathbf{L}\mathbf{x}\| & \leq \varepsilon, \end{aligned}$$

where  $\varepsilon$  is a positive constant. The constrained least squares problem (8.22) can be solved by using the Lagrange multiplier formalism. Considering the Lagrangian function

$$\mathcal{L}(\mathbf{x}, \alpha) = \|\mathbf{y}^\delta - \mathbf{K}\mathbf{x}\|^2 + \alpha \left( \|\mathbf{L}\mathbf{x}\|^2 - \varepsilon^2 \right),$$

it can be shown that if  $\varepsilon \leq \|\mathbf{L}\mathbf{x}^\delta\|$ , where  $\mathbf{x}^\delta$  is the least squares solution of the equation  $\mathbf{K}\mathbf{x} = \mathbf{y}^\delta$ , then the solution  $\mathbf{x}_\varepsilon^\delta$  to (8.22) is identical to the Tikhonov solution  $\mathbf{x}_\alpha^\delta$ , with  $\alpha$  solving the equation

$$\|\mathbf{L}\mathbf{x}_\alpha^\delta\|^2 = \varepsilon^2. \quad (8.23)$$

To carry this idea over to the TLS setting, we add the bound  $\|\mathbf{L}\mathbf{x}\| \leq \varepsilon$  to the ordinary problem (8.5), in which case, the new problem statement becomes

$$\begin{aligned} \min_{[\tilde{\mathbf{K}} \tilde{\mathbf{y}}] \in \mathbb{R}^{m \times (n+1)}} \left\| \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{K}} & \tilde{\mathbf{y}} \end{bmatrix} \right\|_{\mathbb{F}}^2 & \quad (8.24) \\ \text{subject to } \tilde{\mathbf{K}}\mathbf{x} = \tilde{\mathbf{y}} \text{ and } \|\mathbf{L}\mathbf{x}\| & \leq \varepsilon. \end{aligned}$$

The corresponding Lagrangian function is

$$\mathcal{L}(\tilde{\mathbf{K}}, \mathbf{x}, \alpha) = \left\| \begin{bmatrix} \mathbf{K}_\Lambda & \mathbf{y}^\delta \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{K}} & \tilde{\mathbf{K}}\mathbf{x} \end{bmatrix} \right\|_{\mathbb{F}}^2 + \alpha \left( \|\mathbf{L}\mathbf{x}\|^2 - \varepsilon^2 \right),$$

and the Lagrange multiplier  $\alpha$  is non-zero if the inequality constraint is active. In fact, the solution  $\mathbf{x}_{\Lambda\varepsilon}^\delta$  to (8.24) is different from the TLS solution  $\mathbf{x}_\Lambda^\delta$ , whenever  $\varepsilon$  is less than  $\|\mathbf{L}\mathbf{x}_\Lambda^\delta\|$ .

To characterize  $\mathbf{x}_{\Lambda\varepsilon}^\delta$ , we set the partial derivatives of the Lagrangian function to zero. Differentiation with respect to the entries in  $\tilde{\mathbf{K}}$  yields

$$\tilde{\mathbf{K}} - \mathbf{K}_\Lambda - \mathbf{r}\mathbf{x}^T = \mathbf{0}, \quad (8.25)$$

with  $\mathbf{r} = \mathbf{y}^\delta - \tilde{\mathbf{K}}\mathbf{x}$ , while differentiation with respect to the entries in  $\mathbf{x}$  gives

$$-\tilde{\mathbf{K}}^T \mathbf{r} + \alpha \mathbf{L}^T \mathbf{L}\mathbf{x} = \mathbf{0}. \quad (8.26)$$

Setting the partial derivative with respect to  $\alpha$  to zero also yields

$$\|\mathbf{L}\mathbf{x}\|^2 = \varepsilon^2. \quad (8.27)$$

Making use of the expression of  $\mathbf{r}$ , we rearrange (8.26) as

$$\left( \tilde{\mathbf{K}}^T \tilde{\mathbf{K}} + \alpha \mathbf{L}^T \mathbf{L} \right) \mathbf{x} = \tilde{\mathbf{K}}^T \mathbf{y}^\delta. \quad (8.28)$$

Now, by (8.25) and (8.26), we have  $\mathbf{K}_\Lambda = \tilde{\mathbf{K}} - \mathbf{r}\mathbf{x}^T$  and  $\tilde{\mathbf{K}}^T \mathbf{r} = \alpha \mathbf{L}^T \mathbf{L}\mathbf{x}$ , respectively, and so, we obtain

$$\mathbf{K}_\Lambda^T \mathbf{K}_\Lambda = \tilde{\mathbf{K}}^T \tilde{\mathbf{K}} - \alpha \mathbf{x}\mathbf{x}^T \mathbf{L}^T \mathbf{L} + \|\mathbf{r}\|^2 \mathbf{x}\mathbf{x}^T - \alpha \mathbf{L}^T \mathbf{L}\mathbf{x}\mathbf{x}^T \quad (8.29)$$

and

$$\mathbf{K}_\Lambda^T \mathbf{y}^\delta = \tilde{\mathbf{K}}^T \mathbf{y}^\delta - (\mathbf{r}^T \mathbf{y}^\delta) \mathbf{x}. \quad (8.30)$$

Inserting (8.29) and (8.30) into (8.28), and using the identities (cf. (8.27))

$$\mathbf{x}\mathbf{x}^T \mathbf{L}^T \mathbf{L}\mathbf{x} = \varepsilon^2 \mathbf{x}, \quad \|\mathbf{r}\|^2 \mathbf{x}\mathbf{x}^T \mathbf{x} = \|\mathbf{r}\|^2 \|\mathbf{x}\|^2 \mathbf{x},$$

and

$$\mathbf{L}^T \mathbf{L}\mathbf{x}\mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|^2 \mathbf{L}^T \mathbf{L}\mathbf{x},$$

we arrive at

$$\left( \mathbf{K}_\Lambda^T \mathbf{K}_\Lambda + \alpha_{\mathbf{I}} \mathbf{I}_n + \alpha_{\mathbf{L}} \mathbf{L}^T \mathbf{L} \right) \mathbf{x} = \mathbf{K}_\Lambda^T \mathbf{y}^\delta, \quad (8.31)$$

with

$$\alpha_{\mathbf{I}} = \alpha \varepsilon^2 - \|\mathbf{r}\|^2 \|\mathbf{x}\|^2 - \mathbf{r}^T \mathbf{y}^\delta \quad (8.32)$$

and

$$\alpha_{\mathbf{L}} = \alpha \left( 1 + \|\mathbf{x}\|^2 \right). \quad (8.33)$$

The next step of our derivation is the elimination of the Lagrange multiplier  $\alpha$  in the expressions of  $\alpha_{\mathbf{I}}$  and  $\alpha_{\mathbf{L}}$ . First, we use the relation (cf. (8.25))

$$\mathbf{r} = \mathbf{y}^\delta - \tilde{\mathbf{K}}\mathbf{x} = \mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x} - \|\mathbf{x}\|^2 \mathbf{r},$$

to obtain

$$(1 + \|\mathbf{x}\|^2) \mathbf{r} = \mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}, \quad (8.34)$$

and further,

$$(1 + \|\mathbf{x}\|^2) \|\mathbf{r}\|^2 = \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2}. \quad (8.35)$$

On the other hand, scalar multiplying (8.26) by  $\mathbf{x}$  gives

$$\alpha = \frac{\mathbf{x}^T \tilde{\mathbf{K}}^T \mathbf{r}}{\|\mathbf{L}\mathbf{x}\|^2} = \frac{1}{\varepsilon^2} (\mathbf{r}^T \mathbf{y}^\delta - \|\mathbf{r}\|^2). \quad (8.36)$$

Considering the parameter  $\alpha_{\mathbf{I}}$ , we insert (8.35) and (8.36) into (8.32), and find that

$$\alpha_{\mathbf{I}} = - \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2}. \quad (8.37)$$

Turning now to the parameter  $\alpha_{\mathbf{L}}$ , we use (8.33) and (8.36) to get

$$\alpha_{\mathbf{L}} = \alpha (1 + \|\mathbf{x}\|^2) = \frac{1}{\varepsilon^2} (\mathbf{r}^T \mathbf{y}^\delta - \|\mathbf{r}\|^2) (1 + \|\mathbf{x}\|^2). \quad (8.38)$$

Finally, a relationship connecting  $\alpha_{\mathbf{L}}$  and  $\alpha_{\mathbf{I}}$  can be derived as follows: by (8.35) and (8.37), we have  $\alpha_{\mathbf{I}} = -\|\mathbf{r}\|^2 (1 + \|\mathbf{x}\|^2)$ , whence using (8.34), (8.38) becomes

$$\alpha_{\mathbf{L}} = \frac{1}{\varepsilon^2} [\mathbf{y}^{\delta T} (\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}) + \alpha_{\mathbf{I}}]. \quad (8.39)$$

To evaluate the approximation error  $\|[\mathbf{K}_\Lambda \quad \mathbf{y}^\delta] - [\tilde{\mathbf{K}} \quad \tilde{\mathbf{y}}]\|_{\mathbf{F}}$ , we use the relation (cf. (8.25))

$$[\mathbf{K}_\Lambda \quad \mathbf{y}^\delta] - [\tilde{\mathbf{K}} \quad \tilde{\mathbf{K}}\mathbf{x}] = [\mathbf{K}_\Lambda - \tilde{\mathbf{K}} \quad \mathbf{r}] = [-\mathbf{r}\mathbf{x}^T \quad \mathbf{r}] = -\mathbf{r} \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix}^T,$$

together with (8.35) and (8.37), to obtain

$$\|[\mathbf{K}_\Lambda \quad \mathbf{y}^\delta] - [\tilde{\mathbf{K}} \quad \tilde{\mathbf{y}}]\|_{\mathbf{F}}^2 = (1 + \|\mathbf{x}\|^2) \|\mathbf{r}\|^2 = -\alpha_{\mathbf{I}}. \quad (8.40)$$

Collecting all results we conclude that  $\mathbf{x}_{\Lambda_\varepsilon}^\delta$  is the solution of equation (8.31) with  $\alpha_{\mathbf{I}}$  and  $\alpha_{\mathbf{L}}$  given by (8.37) and (8.39), respectively. The main features of the regularized TLS are presented below (Golub et al., 1999).

- (1) If the matrix  $\alpha_{\mathbf{I}}\mathbf{I}_n + \alpha_{\mathbf{L}}\mathbf{L}^T\mathbf{L}$  is positive definite, then the regularized TLS solution corresponds to the Tikhonov solution with the penalty term  $\alpha_{\mathbf{I}}\|\mathbf{x}\|^2 + \alpha_{\mathbf{L}}\|\mathbf{L}\mathbf{x}\|^2$ . If the matrix  $\alpha_{\mathbf{I}}\mathbf{I}_n + \alpha_{\mathbf{L}}\mathbf{L}^T\mathbf{L}$  is indefinite or negative definite, there is no equivalent interpretation.

- (2) For a given  $\varepsilon$ , there are several pairs of parameters  $\alpha_T$  and  $\alpha_L$  and thus several solutions  $\mathbf{x}_{\Lambda\varepsilon}^\delta$  that satisfy (8.31), (8.37) and (8.39). However, from (8.40), we see that only the solution with the smallest value of  $|\alpha_T|$  solves the constrained minimization problem (8.24).
- (3) If  $\varepsilon < \|\mathbf{L}\mathbf{x}_\Lambda^\delta\|$ , where  $\mathbf{x}_\Lambda^\delta$  is the TLS solution (8.10), the inequality constraint is binding, the Lagrange multiplier  $\alpha$  is positive and by (8.33), it follows that  $\alpha_L > 0$ . From (8.37) it is apparent that  $\alpha_T$  is always negative and thus adds some deregularization to the solution. The residual (8.40) is a monotonically decreasing function of  $\varepsilon$ , and so,  $\alpha_T$  is a monotonically increasing function of  $\varepsilon$ . If  $\varepsilon = \|\mathbf{L}\mathbf{x}_\Lambda^\delta\|$ , the Lagrange multiplier  $\alpha$  is zero and the regularized TLS solution  $\mathbf{x}_{\Lambda\varepsilon}^\delta$  coincides with the TLS solution  $\mathbf{x}_\Lambda^\delta$ ; for larger  $\varepsilon$ , the constraint is never again binding and so, the solution remains unchanged.

To compute the regularized TLS solution  $\mathbf{x}_{\Lambda\varepsilon}^\delta$  we have to solve a nonlinear problem, and several techniques have been proposed in the literature. In Golub et al. (1999),  $\alpha_L$  is considered as free parameter, a corresponding value is computed for  $\alpha_T$ , and the system of equations (8.31) is solved in an efficient way. The idea is to transform (8.31) into the augmented system of equations

$$\begin{bmatrix} \mathbf{I}_m & \mathbf{0} & \mathbf{K}_\Lambda \\ \mathbf{0} & \mathbf{I}_n & \sqrt{\alpha_L}\mathbf{L} \\ \mathbf{K}_\Lambda^T & \sqrt{\alpha_L}\mathbf{L}^T & -\alpha_T\mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{r} \\ \mathbf{s} \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{y}^\delta \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix},$$

to reduce  $\mathbf{K}_\Lambda$  to an  $n \times n$  bidiagonal form by means of orthogonal transformations, to apply Elden's algorithm to annihilate the matrix term containing the factor  $\sqrt{\alpha_L}$ , and finally, to use a symmetric perfect shuffle reordering to obtain a symmetric, tridiagonal, indefinite matrix of size  $2n \times 2n$  containing the parameter  $\alpha_T$  on the main diagonal.

In Guo and Renault (2002), a shifted inverse power method is used to obtain the eigenpair

$$\left( \lambda, \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} \right)$$

for the problem

$$\mathbf{B}(\mathbf{x}) \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x} \\ -1 \end{bmatrix}, \quad (8.41)$$

where

$$\mathbf{B}(\mathbf{x}) = \begin{bmatrix} \mathbf{K}_\Lambda^T \mathbf{K}_\Lambda + \alpha_L(\mathbf{x}) \mathbf{L}^T \mathbf{L} & \mathbf{K}_\Lambda^T \mathbf{y}^\delta \\ \mathbf{y}^{\delta T} \mathbf{K}_\Lambda & -\alpha_L(\mathbf{x}) \varepsilon^2 + \mathbf{y}^{\delta T} \mathbf{y}^\delta \end{bmatrix}$$

is an  $(n+1) \times (n+1)$  matrix,  $\lambda = -\alpha_T$ , and  $\alpha_L$  is given by (cf. (8.37) and (8.39))

$$\alpha_L(\mathbf{x}) = \frac{1}{\varepsilon^2} \left[ \mathbf{y}^{\delta T} (\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}) - \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2} \right]. \quad (8.42)$$

In Renault and Guo (2005), the solution of the eigenproblem (8.41) is considered together with the solution of a nonlinear equation which guarantees the bound  $\|\mathbf{L}\mathbf{x}\| = \varepsilon$ . To describe the main features of this algorithm, we consider the decomposition

$$\mathbf{B}(\alpha_L) = \mathbf{M} + \alpha_L \mathbf{N},$$

where

$$\mathbf{M} = \begin{bmatrix} \mathbf{K}_\Lambda^T \mathbf{K}_\Lambda & \mathbf{K}_\Lambda^T \mathbf{y}^\delta \\ \mathbf{y}^{\delta T} \mathbf{K}_\Lambda & \mathbf{y}^{\delta T} \mathbf{y}^\delta \end{bmatrix}, \quad \mathbf{N} = \begin{bmatrix} \mathbf{L}^T \mathbf{L} & \mathbf{0} \\ \mathbf{0} & -\varepsilon^2 \end{bmatrix},$$

and denote by

$$\left( \lambda_{\alpha_L}, \begin{bmatrix} \mathbf{x}_{\alpha_L} \\ -1 \end{bmatrix} \right)$$

the eigenpair corresponding to the smallest eigenvalue of  $\mathbf{B}(\alpha_L)$ . For a fixed  $\varepsilon$ , we introduce the function

$$g(\mathbf{x}) = \frac{\|\mathbf{L}\mathbf{x}\|^2 - \varepsilon^2}{1 + \|\mathbf{x}\|^2},$$

and compute  $\hat{\alpha}_L$  such that  $\mathbf{x}_{\hat{\alpha}_L}$  solves the equation

$$g(\mathbf{x}_{\alpha_L}) = 0; \quad (8.43)$$

$\mathbf{x}_{\hat{\alpha}_L}$  is then the regularized TLS solution of (8.24). To justify this algorithm, we assume that  $\mathbf{x}_{\hat{\alpha}_L}$  satisfies the eigensystem equation

$$\mathbf{B}(\hat{\alpha}_L) \begin{bmatrix} \mathbf{x}_{\hat{\alpha}_L} \\ -1 \end{bmatrix} = \lambda_{\hat{\alpha}_L} \begin{bmatrix} \mathbf{x}_{\hat{\alpha}_L} \\ -1 \end{bmatrix}, \quad (8.44)$$

and is also a solution of equation (8.43). The first block equation of the eigenvalue problem (8.44) gives (8.31) with  $\hat{\alpha}_T = -\lambda_{\hat{\alpha}_L}$ , while the second block equation yields (8.39). Multiplying the eigensystem equation by  $[\mathbf{x}_{\hat{\alpha}_L}^T, -1]$ , we find that

$$\lambda_{\hat{\alpha}_L} = \frac{1}{1 + \|\mathbf{x}_{\hat{\alpha}_L}\|^2} \left[ \|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}_{\hat{\alpha}_L}\|^2 + \hat{\alpha}_L \left( \|\mathbf{L}\mathbf{x}_{\hat{\alpha}_L}\|^2 - \varepsilon^2 \right) \right]. \quad (8.45)$$

Since  $g(\mathbf{x}_{\hat{\alpha}_L}) = 0$ , it follows that  $\|\mathbf{L}\mathbf{x}_{\hat{\alpha}_L}\|^2 = \varepsilon^2$ , and (8.45) becomes

$$\lambda_{\hat{\alpha}_L} = \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}_{\hat{\alpha}_L}\|^2}{1 + \|\mathbf{x}_{\hat{\alpha}_L}\|^2}; \quad (8.46)$$

thus  $\hat{\alpha}_T = -\lambda_{\hat{\alpha}_L}$  satisfies indeed (8.37). In summary,  $\mathbf{x}_{\hat{\alpha}_L}$  solves equation (8.31) with  $\hat{\alpha}_T$  as in (8.37) and  $\hat{\alpha}_L$  as in (8.39). Since  $\lambda_{\hat{\alpha}_L}$  is the smallest eigenvalue of  $\mathbf{B}$ , the present approach explicitly computes a solution with the smallest value of  $|\alpha_T|$ .

For a practical implementation of the method of Renault and Guo we note the following results:

- (1) if  $\lambda_{n+1} > 0$  is the smallest eigenvalue of the matrix  $\mathbf{B}$  and  $\mathbf{v}_{n+1}$  is the corresponding eigenvector, then  $\lambda_{\alpha_L} = \lambda_{n+1}$  and

$$\begin{bmatrix} \mathbf{x}_{\alpha_L} \\ -1 \end{bmatrix} = -\frac{1}{[\mathbf{v}_{n+1}]_{n+1}} \mathbf{v}_{n+1};$$

- (2)  $g(\mathbf{x}_{\alpha_L})$  is a monotonically decreasing function of  $\alpha_L$ , and there exists only one solution  $\hat{\alpha}_L$  of the equation  $g(\mathbf{x}_{\alpha_L}) = 0$ .

Algorithm 15 computes the Tikhonov solution and the regularized TLS solution for a fixed value of the parameter  $\alpha$  corresponding to the method of Tikhonov regularization. Both solutions are related to each other through the constraint norms. The input parameter  $\alpha$  is used to determine the bound  $\varepsilon$  and to estimate a bisection interval for  $\alpha_L$ . The algorithm also computes the ‘equivalent’ regularization matrix defined as

$$\alpha \mathbf{L}_{\text{eq}}^T \mathbf{L}_{\text{eq}} = \hat{\alpha}_I \mathbf{I}_n + \hat{\alpha}_L \mathbf{L}^T \mathbf{L}. \quad (8.47)$$

This factorization is performed by using the Cholesky method with added multiple of identity, which takes into account that for large negative values of  $\hat{\alpha}_I$ , the matrix  $\hat{\alpha}_I \mathbf{I}_n + \hat{\alpha}_L \mathbf{L}^T \mathbf{L}$  may not be positive definite. Note that strategies based on modifying a Cholesky factorization or a symmetric indefinite factorization of a non-positive definite Hessian are standard approaches in the framework of Newton’s method (Nocedal and Wright, 2006).

---

**Algorithm 15.** Algorithm for computing the regularized TLS solution by solving the eigenvalue problem (8.41). The regularization parameter  $\alpha$  corresponds to the method of Tikhonov regularization. The algorithm computes the solution  $\mathbf{x}_{\hat{\alpha}_L}$ , the regularization parameters  $\hat{\alpha}_L$  and  $\hat{\alpha}_I$ , and the equivalent regularization matrix  $\mathbf{L}_{\text{eq}}$ .

---

```

compute the Tikhonov solution  $\mathbf{x}_\alpha^\delta$  for  $\alpha$ , i.e.,  $\mathbf{x}_\alpha^\delta = (\mathbf{K}_\Lambda^T \mathbf{K}_\Lambda + \alpha \mathbf{L}^T \mathbf{L})^{-1} \mathbf{K}_\Lambda^T \mathbf{y}^\delta$ ;
 $\varepsilon \leftarrow \|\mathbf{L} \mathbf{x}_\alpha^\delta\|$ ;
compute the matrices  $\mathbf{M}$  and  $\mathbf{N}$ ;
estimate a bisection interval  $[\alpha_{L \min}, \alpha_{L \max}]$  for  $\alpha_L$  around  $\alpha$ ;
solve  $g(\alpha_L) = 0$  in  $[\alpha_{L \min}, \alpha_{L \max}]$  using FuncEval ( $\alpha_L, \varepsilon, \mathbf{M}, \mathbf{N}; g, \mathbf{x}_{\alpha_L}, \alpha_I$ );
store the solution  $\hat{\alpha}_L$  and the corresponding  $\mathbf{x}_{\hat{\alpha}_L}$  and  $\hat{\alpha}_I$ ;
{regularization matrix using Cholesky factorization with added multiple of identity}
choose the tolerance  $\varepsilon_\alpha$ , e.g.,  $\varepsilon_\alpha = 0.001$ ;
 $\Delta \alpha \leftarrow \varepsilon_\alpha |\hat{\alpha}_I|$ ; stop  $\leftarrow$  false;
while stop = false do
    attempt to apply the Cholesky factorization to obtain  $\mathbf{L}_{\text{eq}}^T \mathbf{L}_{\text{eq}} = \hat{\alpha}_I \mathbf{I}_n + \hat{\alpha}_L \mathbf{L}^T \mathbf{L}$ ;
    if factorization is successful then
        stop  $\leftarrow$  true;
    else
         $\hat{\alpha}_I \leftarrow \hat{\alpha}_I + \Delta \alpha$ ;
    end if
end while
 $\mathbf{L}_{\text{eq}} \leftarrow (1/\sqrt{\alpha}) \mathbf{L}_{\text{eq}}$ .

{for given  $\alpha_L$ , the routine computes  $g(\alpha_L)$ ,  $\mathbf{x}_{\alpha_L}$  and  $\alpha_I$ }
routine FuncEval ( $\alpha_L, \varepsilon, \mathbf{M}, \mathbf{N}; g, \mathbf{x}_{\alpha_L}, \alpha_I$ )
 $\mathbf{B} \leftarrow \mathbf{M} + \alpha_L \mathbf{N}$ ;
compute the smallest eigenvalue  $\lambda_{n+1}$  of  $\mathbf{B}$  and the eigenvector  $\mathbf{v}_{n+1}$ ;
compute  $\mathbf{x}_{\alpha_L}$  as  $\begin{bmatrix} \mathbf{x}_{\alpha_L} \\ -1 \end{bmatrix} = -(1/[\mathbf{v}_{n+1}]_{n+1}) \mathbf{v}_{n+1}$ ;
 $\alpha_I \leftarrow -\lambda_{n+1}$ ;
 $g \leftarrow (\|\mathbf{L} \mathbf{x}_{\alpha_L}\|^2 - \varepsilon^2) / (1 + \|\mathbf{x}_{\alpha_L}\|^2)$ .

```

---

In Sima et al. (2003), the objective function is the so-called orthogonal distance, and the constrained minimization problem takes the form (cf. (8.37) and (8.40))

$$\begin{aligned} \min_{\mathbf{x}} \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2} \\ \text{subject to } \|\mathbf{L}\mathbf{x}\| \leq \varepsilon. \end{aligned}$$

The first-order optimality conditions for the Lagrangian function

$$\mathcal{L}(\mathbf{x}, \lambda) = \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{1 + \|\mathbf{x}\|^2} + \lambda \left( \|\mathbf{L}\mathbf{x}\|^2 - \varepsilon^2 \right),$$

yield

$$\mathbf{D}(\mathbf{x}) \mathbf{x} + \lambda \mathbf{L}^T \mathbf{L} \mathbf{x} = \mathbf{d}(\mathbf{x}), \quad \|\mathbf{L}\mathbf{x}\|^2 = \varepsilon^2, \quad (8.48)$$

with

$$\mathbf{D}(\mathbf{x}) = \frac{\mathbf{K}_\Lambda^T \mathbf{K}_\Lambda}{1 + \|\mathbf{x}\|^2} - \frac{\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2}{(1 + \|\mathbf{x}\|^2)^2} \mathbf{I}_n, \quad \mathbf{d}(\mathbf{x}) = \frac{\mathbf{K}_\Lambda^T \mathbf{y}^\delta}{1 + \|\mathbf{x}\|^2}.$$

The problem (8.48) is first transformed into the standard form and then solved iteratively by using a fixed point iteration method. Assuming that  $\mathbf{L}$  is square and nonsingular, the transformation to the standard form gives

$$(\mathbf{W} + \lambda \mathbf{I}_n) \bar{\mathbf{x}} = \mathbf{h}, \quad \|\bar{\mathbf{x}}\|^2 = \varepsilon^2, \quad (8.49)$$

with  $\bar{\mathbf{x}} = \mathbf{L}\mathbf{x}$ ,  $\mathbf{W} = \mathbf{L}^{-T} \mathbf{D} \mathbf{L}^{-1}$  and  $\mathbf{h} = \mathbf{L}^{-T} \mathbf{d}$ . Note that since  $\mathbf{D}$  is a symmetric matrix,  $\mathbf{W}$  is also a symmetric matrix. Let us now consider the problem

$$(\mathbf{W} + \lambda \mathbf{I}_n)^2 \mathbf{u} = \mathbf{h}, \quad \mathbf{h}^T \mathbf{u} = \varepsilon^2 \quad (8.50)$$

for  $\mathbf{u} \in \mathbb{R}^n$ . Setting

$$\bar{\mathbf{x}} = (\mathbf{W} + \lambda \mathbf{I}_n) \mathbf{u},$$

and taking into account that, due to the symmetry of  $\mathbf{W} + \lambda \mathbf{I}_n$ , there holds

$$\varepsilon^2 = \mathbf{h}^T \mathbf{u} = \mathbf{u}^T (\mathbf{W} + \lambda \mathbf{I}_n)^2 \mathbf{u} = \|\bar{\mathbf{x}}\|^2,$$

we see that the problems (8.49) and (8.50) are equivalent. Further, using the identity

$$\mathbf{h} = \frac{1}{\varepsilon^2} (\mathbf{h}^T \mathbf{u}) \mathbf{h} = \frac{1}{\varepsilon^2} \mathbf{h} \mathbf{h}^T \mathbf{u},$$

we deduce that (8.50) can be transformed into the quadratic eigenvalue problem

$$\left( \lambda^2 \mathbf{I}_n + 2\lambda \mathbf{W} + \mathbf{W}^2 - \frac{1}{\varepsilon^2} \mathbf{h} \mathbf{h}^T \right) \mathbf{u} = \mathbf{0}. \quad (8.51)$$

This quadratic eigenvalue problem is solved in order to find the largest eigenvalue  $\lambda$  and the corresponding eigenvector  $\mathbf{u}$  scaled so that  $\mathbf{h}^T \mathbf{u} = \varepsilon^2$ . As all matrices in (8.51) are

real and symmetric, the quadratic eigenvalues are real and come in complex conjugate pairs. Moreover, the special form of the quadratic eigenvalue problem (8.51) implies that the rightmost (largest real) eigenvalue is real and positive. The solution of the original problem is then recovered by first computing  $\bar{\mathbf{x}} = (\mathbf{W} + \lambda \mathbf{I}_n) \mathbf{u}$  and then  $\mathbf{x} = \mathbf{L}^{-1} \bar{\mathbf{x}}$ .

---

**Algorithm 16.** Algorithm for computing the regularized TLS solution by solving the quadratic eigenvalue problem (8.51). The regularization parameter  $\alpha$  corresponds to the method of Tikhonov regularization. The algorithm computes the solution  $\mathbf{x}$ , the regularization parameters  $\alpha_L$  and  $\alpha_I$ , and the equivalent regularization matrix  $\mathbf{L}_{\text{eq}}$ .

---

choose the tolerances  $\varepsilon_1$  and  $\varepsilon_x$  for the convergence test;

compute the Tikhonov solution  $\mathbf{x}_\alpha^\delta$  for  $\alpha$ , i.e.,  $\mathbf{x}_\alpha^\delta = (\mathbf{K}_\Lambda^T \mathbf{K}_\Lambda + \alpha \mathbf{L}^T \mathbf{L})^{-1} \mathbf{K}_\Lambda^T \mathbf{y}^\delta$ ;

$\varepsilon \leftarrow \|\mathbf{L} \mathbf{x}_\alpha^\delta\|$ ;

$\bar{\mathbf{K}}_\Lambda \leftarrow \mathbf{K}_\Lambda \mathbf{L}^{-1}$ ;

*stop*  $\leftarrow$  false;  $k \leftarrow 0$ ;  $\mathbf{x} \leftarrow \mathbf{x}_\alpha^\delta$ ; {starting vector}

**while** *stop* = false **do**

$r \leftarrow \|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2 / (1 + \|\mathbf{x}\|^2)$ ;  $c \leftarrow 1 / (1 + \|\mathbf{x}\|^2)$ ;

$\mathbf{W} \leftarrow c \bar{\mathbf{K}}_\Lambda^T \bar{\mathbf{K}}_\Lambda - r c \mathbf{L}^{-T} \mathbf{L}^{-1}$ ;  $\mathbf{h} \leftarrow c \bar{\mathbf{K}}_\Lambda^T \mathbf{y}^\delta$ ;

set  $\mathbf{A} = \begin{bmatrix} -2\mathbf{W} & -\mathbf{W}^2 + \varepsilon^{-2} \mathbf{h} \mathbf{h}^T \\ \mathbf{I}_n & \mathbf{0} \end{bmatrix}$ ;

compute the largest eigenvalue  $\lambda$  and the corresponding eigenvector  $\begin{bmatrix} \mathbf{v} \\ \mathbf{u} \end{bmatrix}$  of  $\mathbf{A}$ ;

$\mathbf{u} \leftarrow (\varepsilon^2 / \mathbf{h}^T \mathbf{u}) \mathbf{u}$ ; {scale  $\mathbf{u}$ }

$\mathbf{W} \leftarrow \mathbf{W} + \lambda \mathbf{I}_n$ ;

$\mathbf{x} \leftarrow \mathbf{L}^{-1} \mathbf{W} \mathbf{u}$ ;

{convergence test}

**if**  $k > 0$  **and**  $|\lambda - \lambda_{\text{prv}}| \leq \varepsilon_1 \lambda$  **and**  $\|\mathbf{x} - \mathbf{x}_{\text{prv}}\| \leq \varepsilon_x \|\mathbf{x}\|$  **then**

*stop*  $\leftarrow$  true;

**else**

$\lambda_{\text{prv}} \leftarrow \lambda$ ;  $\mathbf{x}_{\text{prv}} \leftarrow \mathbf{x}$ ;

$k \leftarrow k + 1$ ;

**end if**

**end while**

$\alpha_L \leftarrow \lambda (1 + \|\mathbf{x}\|^2)$ ;  $\alpha_I \leftarrow -\|\mathbf{y}^\delta - \mathbf{K}_\Lambda \mathbf{x}\|^2 / (1 + \|\mathbf{x}\|^2)$ ;

compute  $\mathbf{L}_{\text{eq}}$  as in Algorithm 15

---

The quadratic eigenvalue problem (8.51) is equivalent to the linear eigenvalue problem

$$\begin{bmatrix} -2\mathbf{W} & -\mathbf{W}^2 + \frac{1}{\varepsilon^2} \mathbf{h} \mathbf{h}^T \\ \mathbf{I}_n & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{u} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{v} \\ \mathbf{u} \end{bmatrix},$$

and this can be solved by using for example, the routine DGEEV from the LAPACK library (Anderson et al., 1995), or the routine DNAUPD from the ARPACK library (Maschhoff and Sorensen, 1996). The DNAUPD routine is more efficient because it calculates only the largest eigenvalue and the corresponding eigenvector by using Arnoldi's method (Arnoldi,



1951). The Algorithm 16 generates a sequence  $\{(\lambda_k, \mathbf{x}_k)\}$  by solving the quadratic eigenvalue problem (8.51) at each iteration step  $k$ . From the analysis of Sima et al. (2003) we infer the following results:

- (1)  $\mathbf{x}_k$  should correspond to the largest eigenvalue  $\lambda_k > 0$  since only then the algorithm converges;
- (2) the orthogonal distance decreases at each iteration step;
- (3) any limit point of the sequence  $\{(\lambda_k, \mathbf{x}_k)\}$  solves equation (8.48).

The last result suggests that instead of requiring the convergence of the sequence  $\{(\lambda_k, \mathbf{x}_k)\}$  we may check if equation (8.48) is satisfied within a prescribed tolerance at each iteration step.

#### 8.4 Regularized total least squares for nonlinear problems

As stated in Chapter 6, the solution of a nonlinear ill-posed problem by means of Tikhonov regularization is equivalent to the solution of a sequence of ill-posed linearizations of the forward model about the current iterate. Essentially, at the iteration step  $k$ , we solve the linearized equation

$$\mathbf{K}_{\alpha k} \Delta \mathbf{x} = \mathbf{y}_k^\delta, \quad (8.52)$$

with  $\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}_a$ ,  $\mathbf{K}_{\alpha k} = \mathbf{K}(\mathbf{x}_{\alpha k}^\delta)$ , and

$$\mathbf{y}_k^\delta = \mathbf{y}^\delta - \mathbf{F}(\mathbf{x}_{\alpha k}^\delta) + \mathbf{K}_{\alpha k}(\mathbf{x}_{\alpha k}^\delta - \mathbf{x}_a),$$

via Tikhonov regularization with the penalty term  $\|\mathbf{L}\Delta \mathbf{x}\|^2$  and the regularization parameter  $\alpha$ . If  $\Delta \mathbf{x}_{\alpha k}^\delta$  is the minimizer of the Tikhonov function

$$\mathcal{F}_{1\alpha k}(\Delta \mathbf{x}) = \|\mathbf{y}_k^\delta - \mathbf{K}_{\alpha k} \Delta \mathbf{x}\|^2 + \alpha \|\mathbf{L}\Delta \mathbf{x}\|^2, \quad (8.53)$$

the new iterate is given by  $\mathbf{x}_{\alpha k+1}^\delta = \mathbf{x}_a + \Delta \mathbf{x}_{\alpha k}^\delta$ , and the constraint norm can be readily computed as

$$\varepsilon = \|\mathbf{L}\Delta \mathbf{x}_{\alpha k}^\delta\|. \quad (8.54)$$

In the framework of the regularized TLS, we assume that  $\mathbf{K}_{\alpha k}$  is contaminated by errors, and instead of minimizing (8.53) we solve the problem

$$\min_{[\tilde{\mathbf{K}} \tilde{\mathbf{y}}] \in \mathbb{R}^{m \times (n+1)}} \left\| \begin{bmatrix} \mathbf{K}_{\alpha k} & \mathbf{y}_k^\delta \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{K}} & \tilde{\mathbf{y}} \end{bmatrix} \right\|_{\mathbf{F}}^2 \quad (8.55)$$

subject to  $\tilde{\mathbf{K}}\Delta \mathbf{x} = \tilde{\mathbf{y}}$  and  $\|\mathbf{L}\Delta \mathbf{x}\| \leq \varepsilon$ ,

with  $\varepsilon$  being given by (8.54). The free parameter of the method is the Tikhonov regularization parameter  $\alpha$ , and the Algorithms 15 and 16 can be used to compute both the Tikhonov solution and the regularized TLS solution. Although the numerical implementation of the regularized TLS is very similar to that of Tikhonov regularization, the use of a step-length procedure is problematic. In principle it can be applied for the objective function

$$\mathcal{F}_\alpha(\mathbf{x}) = \frac{1}{2} \|\mathbf{f}_\alpha(\mathbf{x})\|^2, \quad \mathbf{f}_\alpha(\mathbf{x}) = \begin{bmatrix} \mathbf{F}(\mathbf{x}) - \mathbf{y}^\delta \\ \sqrt{\alpha} \mathbf{L}_{\text{eq}}(\mathbf{x} - \mathbf{x}_a) \end{bmatrix}, \quad (8.56)$$

but solving (8.55) is not equivalent to minimizing (8.56) at the iteration step  $k$  because  $\mathbf{L}_{eq}$  may not be the exact Cholesky factor of  $\hat{\alpha}_T \mathbf{I}_n + \hat{\alpha}_L \mathbf{L}^T \mathbf{L}$  (cf. (8.47)).

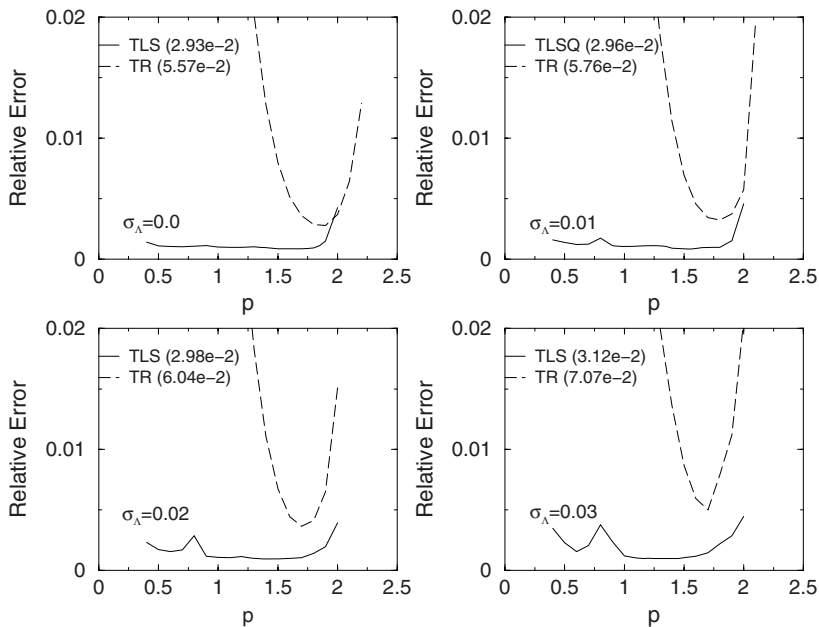
In our numerical analysis, we consider the  $O_3$  retrieval test problem and compute the Jacobian matrix  $\mathbf{K}_{ss}$  by assuming only the single scattering contribution (cf. (8.1)). Furthermore, at each iteration step, we perturb this matrix as

$$[\mathbf{K}_{k\alpha}]_{ij} = [\mathbf{K}_{ssk\alpha}]_{ij} + \sigma_\Lambda \varepsilon_{ij} [\mathbf{K}_{ssk\alpha}]_{ij},$$

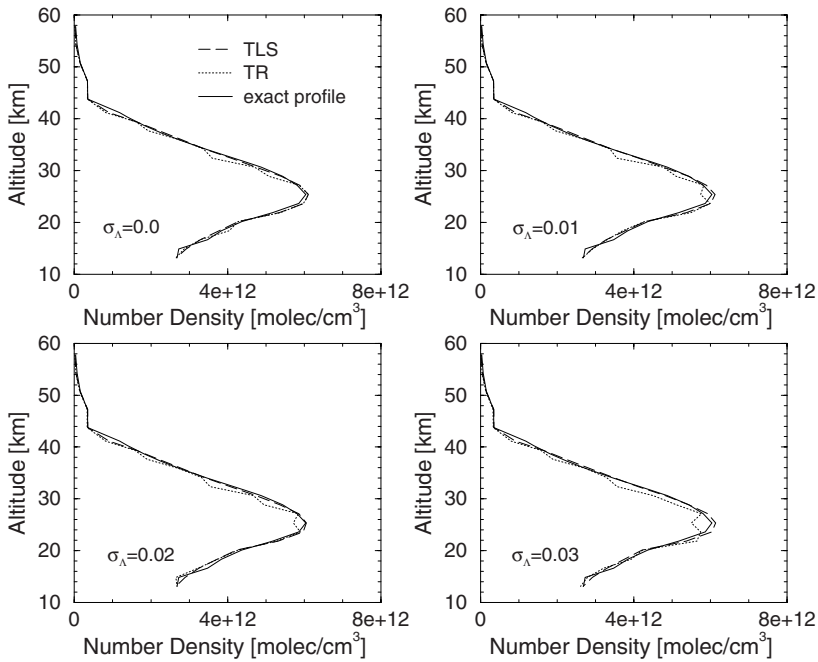
where the elements  $\varepsilon_{ij}$  are from a normal distribution with zero mean and unit variance. Figure 8.1 shows the relative errors in the Tikhonov and the regularized TLS solutions for four values of the standard deviation  $\sigma_\Lambda$ , namely 0, 0.01, 0.02 and 0.03. In all situations, the minimum solution error for the regularized TLS is clearly smaller than that for Tikhonov regularization. Even in the case  $\sigma_\Lambda = 0$  there is a solution improvement due to the approximate Jacobian calculation. The plots also show that the minima of the TLS errors are flat and this situation is beneficial for the inversion process.

In Figure 8.2 we plot the Tikhonov and the regularized TLS solutions, corresponding to the minimizers of the error curves in Figure 8.1. In fact, the improvement of the TLS error as compared to the Tikhonov error is due to the additional term  $\alpha_T \mathbf{I}_n$  in Eq. (8.31).

From the point of view of their accuracy, the regularized TLS algorithms solving the eigenvalue problem (8.41) and the quadratic eigenvalue problem (8.51) are completely



**Fig. 8.1.** Relative errors in the Tikhonov and the regularized TLS solutions as a function of the exponent  $p$ , where  $\alpha = \sigma^p$  and  $\sigma$  is the noise standard deviation. The results correspond to the  $O_3$  retrieval test problem and are computed with the regularized TLS algorithm solving the quadratic eigenvalue problem (8.51). The numbers in parentheses indicate the minimum values of the relative solution error.



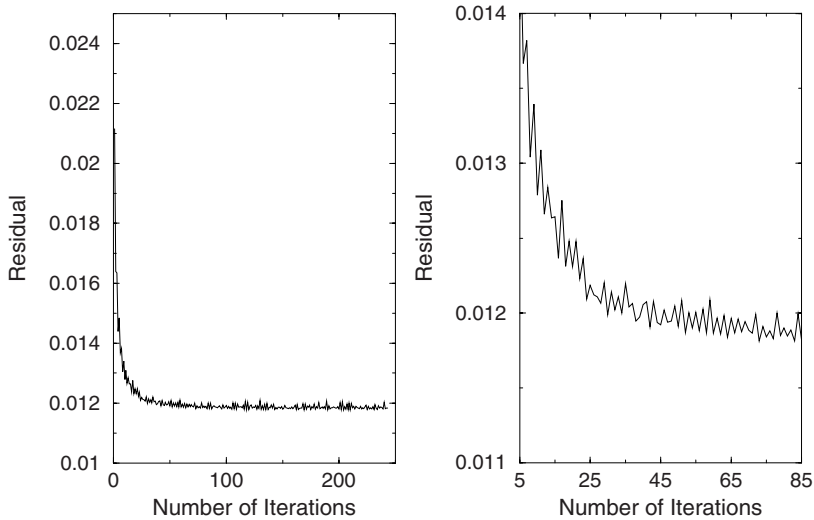
**Fig. 8.2.** Tikhonov (TR) and regularized TLS solutions corresponding to the minimizers of the error curves in Figure 8.1.

equivalent. However, the computation time of the algorithm based on a quadratic eigenvalue problem is on average 6 times smaller (Table 8.1). The main drawback of the regularized TLS is the extraordinarily large number of iteration steps (and so, computation time) as compared to Tikhonov regularization. The decrease of the solution error by a factor of 4–5 is accompanied by an increase of the computation time by a factor of 7–8.

The large number of iteration steps is also a consequence of the fact that we do not use a step-length procedure to guarantee a monotonic decrease of the residual norm (Figure 8.3). A step-length algorithm stops the iterative process too early (because the search direction is not a descent direction for the Tikhonov function), and as a result, the solu-

**Table 8.1.** Computation time in min:ss format. The numbers in parentheses indicate the number of iteration steps for Tikhonov regularization (TR) and the regularized TLS algorithms solving the eigenvalue problem (8.41) (TLS-EP) and the quadratic eigenvalue problem (8.51) (TLS-QEP).

Method	Standard deviation $\sigma_\Lambda$			
	0	0.01	0.02	0.03
TR	0:14 (4)	0:15 (6)	0:18 (8)	0:24 (16)
TLS-QEP	1:24 (108)	1:37 (124)	2:23 (202)	2:58 (243)
TLS-EP	8:01 (108)	9:57 (124)	13:13 (202)	19:17 (243)



**Fig. 8.3.** History of the residual norm in the case  $\sigma_{\Lambda} = 0.03$ . In the left panel the curves are plotted for all iteration steps, while in the right panel, the  $y$ -axis is zoomed out.

tion errors are not sufficiently small. For example, in the case  $\sigma_{\Lambda} = 0.03$ , the regularized TLS with a step-length algorithm terminates after 19 iteration steps with a solution error of  $1.56 \cdot 10^{-2}$ , and if the step-length algorithm is turned off, it terminates after 243 iteration steps with a solution error of  $9.77 \cdot 10^{-4}$ .

The design of an efficient regularized TLS algorithm for nonlinear problems is far from being complete. The selection of an optimal value of the regularization parameter by an a posteriori method will dramatically increase the computational effort, while the use of a variable regularization parameter computed for example, by using the L-curve method, is also problematic. In our numerical simulations, the L-curve either does not have a distinctive L-shape, or it predicts values of the regularization parameter that are too small.

The regularized TLS has been applied to atmospheric trace gas profile retrievals by Koner and Drummond (2008). In this work, the regularized TLS algorithm solving the quadratic eigenvalue problem (8.51) is used for the automatic determination of the regularization strength.