

Solving Rank-Deficient and Ill-posed Problems using UTV and QR Factorizations*

Leslie V. Foster
Department of Mathematics and Computer Science
San Jose State University
San Jose, CA 95192
foster@math.sjsu.edu

March 19, 2003

Abstract

The algorithm of Mathias and Stewart [A block QR algorithm and the singular value decomposition, *Linear Algebra and Its Applications*, 182:91-100, 1993] is examined as a tool for constructing regularized solutions to rank-deficient and ill-posed linear equations. The algorithm is based on a sequence of QR factorizations. If it is stopped after the first step it produces that same solution as the complete orthogonal decomposition used in LAPACK's xGELSY. However we show that for low-rank problems a careful implementation can lead to an order of magnitude improvement in speed over xGELSY as implemented in LAPACK. We prove, under assumptions similar to assumptions used by others, that if the numerical rank is chosen at a gap in the singular value spectrum and if the initial factorization is rank-revealing then, even if the algorithm is stopped after the first step, approximately half the time its solutions are closer to the desired solution than are the singular value decomposition (SVD) solutions. Conversely, the SVD will be closer approximately half the time and in this case overall the two algorithms are very similar in accuracy. We confirm this with numerical experiments. Although the algorithm works best for problems with a gap in the singular value spectrum, numerical experiments suggest that it may work well for problems with no gap.

1 Introduction

The solution to ill-posed or nearly rank-deficient linear equations is important in many applications [18]. To solve these systems some form of regularization is usually used. By regularization we mean the replacement of the original problem with a different, better posed problem. For example if the original problem is

$$\min \|b - Ax\| \tag{1}$$

where A is an $m \times n$ ill-conditioned matrix with $m \geq n$ and the norm is the Euclidean norm, it is often recommended to approximate A with an exactly rank-deficient matrix

*This research was supported in part by the Woodward bequest to the Department of Mathematics, San Jose State University. To appear in SIAM J. Matrix Anal. and Appl.

\widehat{A} and solve for the minimum norm solution to (1) with A replaced by \widehat{A} . To construct \widehat{A} it is useful to decompose A with a rank-revealing decomposition. The low-rank approximation \widehat{A} to A can be obtained by truncating such decompositions. The singular value decomposition (SVD) is a very good, but expensive, decomposition. We use the complete orthogonal or UTV decomposition $A = UTV^T$ with U orthogonal, V orthogonal and T triangular. Here the superscript T indicates transpose. Some of our results will apply to any UTV factorization and others to UTV factorizations produced by the algorithm of Mathias and Stewart [21]. This algorithm produces UTV factorizations by using a sequence of QR factorizations. We begin the algorithm with an initial UTV factorization of the form $A = UTV^T = QR\Pi^T$ where Q is an orthogonal matrix, R is upper triangular, and Π is the permutation matrix produced by the standard QR algorithm with pivoting [1, 4, 20]. (We also will discuss a variation where initially A^T is factored in this form). If the algorithm is stopped after the first step it produces the same solution as the complete orthogonal decomposition used in LAPACK's xGELSY. However we show that for low-rank problems a careful implementation can lead to an order of magnitude improvement in speed over the two routines, xGELSY and xGELSD, that LAPACK provides for solving rank-deficient problems. We prove, under assumptions similar to assumptions used by others about the true solution to (1) and the noise in b , that if the numerical rank is chosen at a gap in the singular value spectrum and if the initial factorization is rank-revealing [3, p. 22] then, even if the algorithm is stopped after the first step, approximately half the time its solutions are closer to the desired solution than are the singular value decomposition solutions. Conversely, the SVD will be closer approximately half the time and in this case overall the two algorithms are very similar in accuracy. We confirm this with numerical experiments. Although the algorithm works best for problems with a gap in the singular value spectrum, numerical experiments suggest that it may work well for problems with no gap.

The paper is organized as follows. Following this introduction, in Section 2 we discuss UTV factorizations in general. Section 3 discusses the algorithm in [21]. Section 4 focuses on perturbation errors and Section 5 on regularization errors. Section 6 describes implementation of the algorithm and numerical experiments. Section 7 has conclusions.

2 UTV factorizations

Consider any UTV factorization of A , $A = UTV^T$. Let k be the rank of the low-rank approximation to A . It is useful to partition the factorization as follows. If T is lower triangular ($T = L$) we partition UTV^T as

$$A = UTV^T = ULV^T = \begin{pmatrix} \widehat{U} & U_0 \end{pmatrix} \begin{pmatrix} \widehat{L} & 0 \\ H & E \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \widehat{V} & V_0 \end{pmatrix}^T. \quad (2)$$

If T is upper triangular ($T = R$) we partition UTV^T as

$$A = UTV^T = URV^T = \begin{pmatrix} \widehat{U} & U_0 \end{pmatrix} \begin{pmatrix} \widehat{R} & F \\ 0 & G \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \widehat{V} & V_0 \end{pmatrix}^T. \quad (3)$$

In these equations \widehat{U} is $m \times k$, U_0 is $m \times (m - k)$, \widehat{V} is $n \times k$, V_0 is $n \times (n - k)$, \widehat{L} is $k \times k$, H is $(n - k) \times k$, E is $(n - k) \times (n - k)$, \widehat{R} is $k \times k$, F is $k \times (n - k)$, and G is

$(n - k) \times (n - k)$. In equations (2) and (3) U_0 corresponds to the last two block rows in the block triangular matrices. If we do not need to distinguish whether T is lower or upper triangular we will let \widehat{T} represent either \widehat{L} or \widehat{R} . In each case we consider two low-rank approximations to A . If T is either lower or upper triangular we will call $\widehat{U}\widehat{T}\widehat{V}^T$ the corner low-rank approximation to A . If T is lower triangular we call $U[\widehat{L}^T \ H^T \ 0]^T \widehat{V}^T$ the block-column low-rank approximation to A . Similarly if T is upper triangular we call $\widehat{U}[\widehat{R} \ F]V^T$ the block-row low-rank approximation to A .

We will also partition the SVD of A in a similar manner to (2) and (3).

$$A = U_S D V_S^T = \begin{pmatrix} \widehat{U}_S & U_{S_0} \end{pmatrix} \begin{pmatrix} \widehat{D} & 0 \\ 0 & D_0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \widehat{V}_S & V_{S_0} \end{pmatrix}^T.$$

$\widehat{A}_S = \widehat{U}_S \widehat{D} \widehat{V}_S^T$ is the rank k approximation produced by the SVD. We will use $s_1 \geq s_2 \geq \dots \geq s_n$ to indicate the singular values of A . We will also use $\sigma_k(A)$, $1 \leq k \leq n$, to indicate the k^{th} singular value of a matrix A . Note that the SVD produces the best rank k approximation to A in the sense that $\|A - \widetilde{A}\|$ is minimized over all rank k matrices \widetilde{A} when $\widetilde{A} = \widehat{A}_S$ [3, p. 12].

When solving equation (1) we will consider the regularized solution $x_T = \widehat{A}_T^+ b$ where the superscript $+$ indicates pseudoinverse and \widehat{A}_T is either a corner or block-row/column low-rank approximation to A corresponding to a UTV factorization of A . It will be clear from the context whether x_T refers to a corner or block-row/column low-rank approximation. We call x_T the truncated UTV solution to (1). We assume in the rest of this paper that \widehat{T} and \widehat{D} are nonsingular. In this case the corner low-rank solution has a simple form $x_T = \widehat{V}\widehat{T}^{-1}\widehat{U}^T b$. The truncated SVD (TSVD) approximate solution to (1) is $x_S = \widehat{V}_S \widehat{D}^{-1} \widehat{U}_S^T b$.

To evaluate the accuracy of x_T we will assume that there is an underlying noiseless solution x_0 to equation (1) such that $Ax_0 = b_0$ and that in (1) $b = b_0 + \delta b$ where δb is a noise vector in the right hand side b . We will prove theorems and carry out numerical experiments that evaluate x_T based on the value of $\|x_T - x_0\|$ and will compare $\|x_T - x_0\|$ with $\|x_S - x_0\|$. We might note that other authors [6, 8, 10] have focused on bounding $\|x_T - x_S\|$. In many cases the goal in solving (1) is to recover an underlying solution x_0 that is different from x_S [18, 22]. In these cases comparison of $\|x_T - x_0\|$ with $\|x_S - x_0\|$ is of interest.

Suppose that C is some regularization operator so that $x = Cb$ is the regularized solution to (1). If x_0 is the underlying noiseless solution then

$$x - x_0 = Cb - x_0 = (CA - I)x_0 + C(\delta b) \quad \text{and} \quad (4)$$

$$\|x - x_0\| \leq \|(CA - I)x_0\| + \|C(\delta b)\|. \quad (5)$$

The two terms on the right are called, respectively, the regularization error and the perturbation error. In the case that C corresponds to a corner low-rank solution calculated using a truncated UTV factorization, where T is lower triangular, we have a sharper result than (5):

$$\|x - x_0\|^2 = \|(CA - I)x_0\|^2 + \|C(\delta b)\|^2. \quad (6)$$

This result follows since, if $C = \widehat{V}\widehat{L}^{-1}\widehat{U}^T$, then $C^T(CA - I) = 0$ follows easily.

Our first theorem relates $\|x_T - x_0\|$ and $\|x_S - x_0\|$.

Theorem 1 *Define*

$$\tilde{U} = U^T U_S = \begin{pmatrix} \hat{U}^T \hat{U}_S & \hat{U}^T U_{S0} \\ U_0^T \hat{U}_S & U_0^T U_{S0} \end{pmatrix} = \begin{pmatrix} \tilde{U}_{11} & \tilde{U}_{12} \\ \tilde{U}_{21} & \tilde{U}_{22} \end{pmatrix}, \quad (7)$$

$$\tilde{V} = V^T V_S = \begin{pmatrix} \hat{V}^T \hat{V}_S & \hat{V}^T V_{S0} \\ V_0^T \hat{V}_S & V_0^T V_{S0} \end{pmatrix} = \begin{pmatrix} \tilde{V}_{11} & \tilde{V}_{12} \\ \tilde{V}_{21} & \tilde{V}_{22} \end{pmatrix}, \quad (8)$$

$$M = \begin{pmatrix} -\hat{D}^{-1} \tilde{V}_{21}^T \tilde{V}_{21} \hat{D}^{-1} & \tilde{U}_{11}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12} \\ \tilde{U}_{12}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{11} & \tilde{U}_{12}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12} \end{pmatrix} \quad (9)$$

and

$$N = \begin{pmatrix} \tilde{V}_{21}^T \tilde{V}_{21} & \tilde{V}_{21}^T \tilde{V}_{22} \\ \tilde{V}_{22}^T \tilde{V}_{21} & -\tilde{V}_{12}^T \tilde{V}_{12} \end{pmatrix} \quad (10)$$

Also let $\tilde{\delta}b = U_S^T \delta b$ and $\tilde{x}_0 = V_S^T x_0$ and let x_T be the corner low-rank solution to (1) calculated from a truncated UTV factorization with T lower triangular. Then

$$\|x_T - x_0\|^2 = \|x_S - x_0\|^2 + \tilde{\delta}b^T M \tilde{\delta}b + \tilde{x}_0^T N \tilde{x}_0. \quad (11)$$

Proof. First note that

$$T = U^T U_S D V_S^T V = \tilde{U} D \tilde{V}^T = \begin{pmatrix} \tilde{U}_{11} & \tilde{U}_{12} \\ \tilde{U}_{21} & \tilde{U}_{22} \end{pmatrix} \begin{pmatrix} \hat{D} & 0 \\ 0 & D_0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \tilde{V}_{11} & \tilde{V}_{12} \\ \tilde{V}_{21} & \tilde{V}_{22} \end{pmatrix}^T. \quad (12)$$

The perturbation error term for the SVD solution is $\|\hat{A}_S^+ \delta b\|$ where $\hat{A}_S^+ = \hat{V}_S \hat{D}^{-1} \hat{U}_S^T$ and for the UTV solution it is $\|\hat{A}_T^+ \delta b\|$ where $\hat{A}_T^+ = \hat{V} \hat{T}^{-1} \hat{U}^T$. Now $\|\hat{A}_S^+ \delta b\|^2 = \|\hat{D}^{-1} \hat{U}_S^T U_S \tilde{\delta}b\|^2 = \|\hat{D}^{-1} (I \ 0) \tilde{\delta}b\|^2$. Note that since \tilde{V} is orthogonal $I = \tilde{V}_{11}^T \tilde{V}_{11} + \tilde{V}_{21}^T \tilde{V}_{21}$ and therefore $\hat{D}^{-2} = \hat{D}^{-1} \tilde{V}_{11}^T \tilde{V}_{11} \hat{D}^{-1} + \hat{D}^{-1} \tilde{V}_{21}^T \tilde{V}_{21} \hat{D}^{-1}$. Rewriting (12) as $T\tilde{V} = \tilde{U}D$ and since T is lower triangular it follows that $\tilde{V}_{11} \hat{D}^{-1} = \hat{T}^{-1} \tilde{U}_{11}$. We may conclude that

$$\|\hat{A}_S^+ \delta b\|^2 = \tilde{\delta}b^T \begin{pmatrix} \tilde{U}_{11}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{11} + \hat{D}^{-1} \tilde{V}_{21}^T \tilde{V}_{21} \hat{D}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \tilde{\delta}b.$$

Next note that $\|\hat{A}_T^+ \delta b\|^2 = \|(\hat{V} \hat{T}^{-1} \hat{U}^T) U_S \tilde{\delta}b\|^2 = \|\hat{T}^{-1} (\tilde{U}_{11} \ \tilde{U}_{12}) \tilde{\delta}b\|^2$. Therefore

$$\|\hat{A}_T^+ \delta b\|^2 = \tilde{\delta}b^T \begin{pmatrix} \tilde{U}_{11}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{11} & \tilde{U}_{11}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12} \\ \tilde{U}_{12}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{11} & \tilde{U}_{12}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12} \end{pmatrix} \tilde{\delta}b.$$

It now follows that

$$\|\hat{A}_T^+ \delta b\|^2 = \|\hat{A}_S^+ \delta b\|^2 + \tilde{\delta}b^T M \tilde{\delta}b. \quad (13)$$

Using $\hat{A}_S^+ = \hat{V}_S \hat{D}^{-1} \hat{U}_S^T$ it follows that the regularization error term for the truncated SVD satisfies $\|(\hat{A}_S^+ A - I) x_0\|^2 = \|(\hat{A}_S^+ A - I) V_S \tilde{x}_0\|^2 = \|(\hat{V}_S \hat{V}_S^T - I) V_S \tilde{x}_0\|^2 = \|(0 \ I) \tilde{x}_0\|^2$. Also $\|(\hat{A}_T^+ A - I) x_0\|^2 = \|(\hat{V} \hat{V}^T - I) V_S \tilde{x}_0\|^2 = \|V_0 V_0^T V_S \tilde{x}_0\|^2 = \|V_0^T (\hat{V}_S \ V_{S0}) \tilde{x}_0\|^2 = \|(\tilde{V}_{21} \ \tilde{V}_{22}) \tilde{x}_0\|^2$. Using this result and $I = \tilde{V}_{22}^T \tilde{V}_{22} + \tilde{V}_{12}^T \tilde{V}_{12}$ (since \tilde{V} is orthogonal) it follows that

$$\|(\hat{A}_T^+ A - I) x_0\|^2 = \tilde{x}_0^T \left[\begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} + N \right] \tilde{x}_0 = \|(\hat{A}_S^+ A - I) x_0\|^2 + \tilde{x}_0^T N \tilde{x}_0.$$

The theorem follows from this equation, (6) and (13). \square

Note that for a *UTV* factorization chosen so that $M \neq 0$ and $N \neq 0$ then it follows from (9) and (10) that M and N are symmetric indefinite matrices. Therefore if $M, N \neq 0$ in the *UTV* factorization of any matrix A , by (11) there exists solution vectors x_0 and noise vectors δb such that the truncated *UTV* solution is closer to x_0 than is the truncated *SVD* solution. We will see in our numerical experiments in Section 6 that it is frequently true that x_T is closer to x_0 than x_S is (and, conversely, that x_S is frequently closer to x_0 than x_T is). In Section 4 and 5 we will use Theorem 1 to explore reasons why this is true.

A result from [21] that we will need later relates the singular values of A , \widehat{T} , E and G . If $\|E\| < \sigma_k(\widehat{T})$ and if T is lower triangular then

$$\sigma_j(\widehat{T}) \leq \sigma_j(A) \leq \sigma_j(\widehat{T}) / \left[1 - \frac{\|H\|^2}{\sigma_k^2(\widehat{T}) - \|E\|^2} \right]^{1/2}, \text{ for } 1 \leq j \leq k \quad (14)$$

and

$$\sigma_{k+j}(A) \leq \sigma_j(E) \leq \sigma_{k+j}(A) / \left[1 - \frac{\|H\|^2}{\sigma_k^2(\widehat{T}) - \|E\|^2} \right]^{1/2}, \text{ for } 1 \leq j \leq n - k. \quad (15)$$

If T is upper triangular and if $\|G\| < \sigma_k(\widehat{T})$ then equations (14) and (15) are also true with H and E replaced, respectively, by F and G .

In the later sections we will also use some of the results [9] which we collect here. These results bound $\sin \theta$, the sine of the angle between the subspaces spanned by \widehat{V} and \widehat{V}_S , and $\sin \phi$, the sine of the angle between the subspaces spanned by \widehat{U} and \widehat{U}_S . Let \widetilde{U} and \widetilde{V} be defined by (7) and (8). Assume that $\|E\| < \sigma_k(\widehat{T})$ and $\|G\| < \sigma_k(\widehat{T})$. If T is lower triangular then

$$\sin \phi = \|\widetilde{U}_{12}\| = \|\widetilde{U}_{21}\| \leq \frac{\sigma_k(\widehat{T})\|H\|}{\sigma_k^2(\widehat{T}) - \|E\|^2} \quad \text{and} \quad (16)$$

$$\sin \theta = \|\widetilde{V}_{12}\| = \|\widetilde{V}_{21}\| \leq \frac{\|H\|\|E\|}{\sigma_k^2(\widehat{T}) - \|E\|^2}. \quad (17)$$

If T is upper triangular then

$$\sin \phi = \|\widetilde{U}_{12}\| = \|\widetilde{U}_{21}\| \leq \frac{\|F\|\|G\|}{\sigma_k^2(\widehat{T}) - \|G\|^2} \quad \text{and} \quad (18)$$

$$\sin \theta = \|\widetilde{V}_{12}\| = \|\widetilde{V}_{21}\| \leq \frac{\sigma_k(\widehat{T})\|F\|}{\sigma_k^2(\widehat{T}) - \|G\|^2}. \quad (19)$$

Note that in most case of interest to us, we will have $\|H\| \leq \|E\|$ and $\|F\| \leq \|G\|$. Assuming this, $\sin \theta$ and $\sin \phi$ can be small for either of two reasons: (1) $\|E\| \ll \sigma_k(\widehat{T})$ and $\|G\| \ll \sigma_k(\widehat{T})$ or (2) $\|H\| \ll \|E\|$ and $\|F\| \ll \|G\|$. The first condition will be true if there is a sufficiently large gap, at singular value k , in the singular values of A and if the *UTV* factorization is rank-revealing (as defined in the next section). The second condition can be achieved by some of the algorithms for calculating *UTV* factorizations, even when there is not a gap in the singular values.

3 Calculating UTV factorizations

There are a number of algorithms for calculating UTV factorizations [11, 21, 25]. We will discuss the algorithm in [21] and a variation of this algorithm. One nice feature of this algorithm is that if the algorithm is stopped after one step it produces a UTV factorization which uses a single QR factorization and, as the algorithm continues with more steps, it approaches the singular value decomposition [21]. The algorithm in [21] does not include interchanges in the columns of A . We will consider a variation that can include column interchanges in the algorithm. At step i the algorithm produces the factorization $A = U_i T_i V_i^T$.

Algorithm:

For $i = 1$ let $A = U_1 T_1 V_1^T$ where U_1, T_1 and V_1 are formed either by $A = Q_1 R_1 \Pi_1^T$ with $U_1 = Q_1, T_1 = R_1$ and $V_1 = \Pi_1$ or $A^T = Q_1 R_1 \Pi_1^T$ with $U_1 = \Pi_1, T_1 = R_1^T$ and $V_1 = Q_1$. In this second case we will also use L_1 to indicate T_1 since T_1 is lower triangular.

For $i \geq 2$, if T_{i-1} is upper triangular form $T_{i-1}^T = Q_i R_i \Pi_i^T$ and let

$$T_i \equiv L_i = R_i^T, U_i = U_{i-1} \Pi_i, V_i = V_{i-1} Q_i.$$

For $i \geq 2$, if T_{i-1} is lower triangular form $T_{i-1} = Q_i R_i \Pi_i^T$ and let

$$T_i = R_i, U_i = U_{i-1} Q_i, V_i = V_{i-1} \Pi_i.$$

To determine when to stop this algorithm one can use the bounds on $\|x_T - x_S\|$ in Theorem 3.3 of [10]. We will select the initial permutation Π_1 using the standard pivoting technique [1, 4, 20] for QR factorizations and let $\Pi_i = I$ for $i \geq 2$. We will also consider a variation where at each step Π_i is chosen by the standard pivoting technique. We will see shortly that often the two alternatives produce identical low-rank solutions. In later section when we use “the algorithm” or “the algorithm of Section 3” we will refer to the first, simpler alternative ($\Pi_i = I, i \geq 2$).

We will find it useful to introduce notation to describe the first few steps of the algorithm. When T_1 is upper triangular we use QRP to indicate the first step of the algorithm, QRLP the next step, QRLRP the next step, etc. When T_1 is lower triangular we use QLP for the first step, QLRLP for the next step, QLRLP for the next step, etc. Here the Q indicates that we are using QR factorizations, P indicates that we use pivoting at the first step and the middle letters indicate the history of the steps of the algorithm. When we are calculating x_T using one of these factorizations we will use TQRP, TQRLP, TQLP, TQLRP, etc. to indicate that we are using a truncated factorization – we only need to calculate a portion of U, V and T .

The above algorithm, without column interchanges, was used in [21], to calculate the singular value decomposition. The paper [21] focuses a block implementation of the above algorithm. Results concerning the convergence of the above algorithm as a tool to estimate singular values and singular vectors is discussed in [7]. Stewart [26, 27] discusses QRLP, with pivoting at both steps, as a tool for estimating singular values and for constructing low-rank approximations. We should note that Stewart uses the designation QLP to refer to what we have called QRLP. In [19] the TQLRP algorithm is described and an example is presented where it works well for regularization. Also we should note that the block-row low-rank approximation produced by TQRP is

the same as the approximate solution to (1) produced by LAPACK's xGELSY [1] or by the algorithm HFTI in [20]. Mathematically these algorithms are identical. Our comparison in this paper of TQRP with TSVD provides a comparison of the accuracy of xGELSY and xGELSD, the two recommended tools in LAPACK 3.0 for solving rank-deficient problems.

We now show that there is a close relationship between solutions produced by block-row/column low-rank approximations and by corner low-rank approximations. We also show that often identical low-rank solutions to (1) are produced by the algorithm if $\Pi_i = I$, $i \geq 2$, or if Π_i is chosen by standard column pivoting.

Theorem 2 *Let $U_i T_i V_i^T$ be the decomposition of A at step $i \geq 1$ of the algorithm. Using the notation of (2) and (3) with subscripts added to indicate the step number in the algorithm, define $\rho_i = \|E_i\|/\sigma_k(\hat{T}_i)$ if T_i is lower triangular and $\rho_i = \|G_i\|/\sigma_k(\hat{T}_i)$ if T_i is upper triangular.*

(a) *Assume \hat{T}_{i-1} is nonsingular and that Π_i in the algorithm has the form $\Pi_i = \begin{pmatrix} \Pi_a & 0 \\ 0 & \Pi_b \end{pmatrix}$ where Π_a and Π_b , respectively, are $k \times k$ and $(n-k) \times (n-k)$ permutation matrices. Then the block-row/column rank k solution x_B calculated from $U_{i-1} T_{i-1} V_{i-1}^T$ is the same as the corner rank k solution x_C calculated from $U_i T_i V_i^T$.*

(b) *If the initial factorization has the property that $\rho_1 < 1$ and if Π_i , $i \geq 2$, is chosen by the standard pivoting algorithm of [4] then for all $i \geq 2$, Π_i is of the form required in part (a).*

(c) *Assume that $\rho_1 < 1$. Then the corner low-rank solution produced at each step of the algorithm using standard column pivoting is identical to the corner low-rank solution produced at the corresponding step of the algorithm where standard column pivoting is used at the first step and no pivoting is used at each following step.*

Proof. Assume that T_{i-1} is lower triangular. We will use the notation of the algorithm, equation (2) and equation (3) except that we will add subscripts to E , \hat{V} , \hat{L} , H , \hat{R} to indicate the step number of the algorithm.

To prove part (a), note that at step $i-1$ the rank k block-row/column approximation to A is $\hat{A}_{i-1} = U_{i-1}(\hat{L}_{i-1}^T \ H_{i-1}^T \ 0)^T \hat{V}_{i-1}^T$ and by properties of pseudoinverses [3, 20], $x_B = \hat{A}_{i-1}^+ b = \hat{V}_{i-1} \begin{pmatrix} \hat{L}_{i-1} \\ H_{i-1} \\ 0 \end{pmatrix}^+ U_{i-1}^T b$. However by our assumption on Π_i and by the constructions of the algorithm,

$$\begin{pmatrix} \hat{L}_{i-1} \\ H_{i-1} \\ 0 \end{pmatrix} = Q_i \begin{pmatrix} \hat{R}_i \\ 0 \end{pmatrix} \Pi_a^T \text{ and so } \begin{pmatrix} \hat{L}_{i-1} \\ H_{i-1} \\ 0 \end{pmatrix}^+ = \Pi_a (\hat{R}_i^{-1} \ 0) Q_i^T.$$

It now follows that, $x_B = \hat{V}_{i-1} \Pi_a (\hat{R}_i^{-1} \ 0) Q_i^T U_{i-1}^T b = \hat{V}_i (\hat{R}_i^{-1} \ 0) U_i^T b = x_C$. The proof for the case that T_{i-1} is upper triangular is similar.

To show part (b) we again assume that T_{i-1} is lower triangular. We will use induction. Assume that after the first step of the algorithm and prior to step i , that the permutation matrices in the algorithm have the form of part (a). It then follows easily from Theorem 2.1 of [21] and its proof, that $\sigma_k(\hat{L}_{i-1}) \geq \sigma_k(\hat{T}_1)$ and, if T_1 is lower triangular, $\|E_1\| \geq \|E_{i-1}\|$ or, if T_1 is upper triangular, $\|G_1\| \geq \|E_{i-1}\|$. Therefore, by the assumption of part (b), $\sigma_k(\hat{L}_{i-1}) > \|E_{i-1}\|$. If Π_i is of the form of part (a), then for $1 \leq j \leq k$, it follows that $|r_{jj}| \geq \sigma_k(\hat{L}_{i-1}^T, H_{i-1}^T) \geq \sigma_k(\hat{L}_{i-1})$. Now

suppose, on the other hand, that at step i the column interchanges in the standard pivoted QR factorization move a column of L_{i-1} with column index larger than k into column j , where $1 \leq j \leq k$. The diagonal entry r_{jj} in the QR factorization of L_{i-1} will satisfy $|r_{jj}| \leq \|E_{i-1}\|$. It follows that this last type interchange is not possible since $\sigma_k(\widehat{L}_{i-1}) > \|E_{i-1}\|$ and since standard column pivoting will move to column j , the column of the remaining unprocessed columns that will make $|r_{jj}|$ as large as possible. The proof when T_{i-1} is upper triangular is similar.

Part (c) follows from part (b) and the block structure of Π_i for $i \geq 2$. Our proof is somewhat tedious and we omit it here. Contact the author for the details. \square

In [27] Stewart notes that in the QRLP factorization if there is a substantial gap in diagonal entries of R_1 “it is unlikely that the pivoting process will interchange columns” across column k in constructing L_2 . Theorem 2 proves under a mild condition on R_1 ($\rho_1 < 1$) that Stewart’s observation is true. Note that $\rho_1 < 1$ will be true if there is a modest gap in the singular values of A and if the initial QR factorization is rank-revealing [3, p. 22]. Part (c) shows that if $\rho_1 < 1$ pivoting is not necessary after the first step in the algorithm in the sense that the corner solutions are the same with pivoting or without pivoting. This is also true for block-row/column solutions by part (a) of the theorem. Our numerical experience indicates that even when $\rho_1 > 1$ pivoting at steps after the first step usually makes little difference in the quality of the solution. However, pivoting at the first step is often critical.

The theorem also shows that there is a close connection between solutions produced by rank-revealing QR factorizations and rank-revealing UTV factorizations. A rank-revealing QR factorization of A [3, p. 22] has the properties $\|G\| = O(s_{k+1})$ and $\sigma_k(\widehat{R}) = O(s_k)$. A rank-revealing ULV factorization of A [10, p. 456] has the properties $\|(H \ E)\| = O(s_{k+1})$ and $\sigma_k(\widehat{L}) = O(s_k)$. Assume that, at step 1 of the algorithm, $A = U_1 T_1 V_1^T = QR\Pi^T$ is a rank-revealing QR factorization of A and that Π_2 has the form of part (a) of Theorem 1, for example if no pivoting is done at step 2. It is not hard to show that the UTV factorization of step 2 of the algorithm will be a rank-revealing ULV factorization of A . It follows from Theorem 2 that the regularized solution produced from a block-row low-rank approximation using a rank-revealing QR factorization is identical to the corner low-rank solution using a related rank-revealing UTV factorization.

Another useful consequence of part (a) is that the results that we develop for corner low-rank solutions to (1) lead directly to results for block-row/column low-rank solutions to (1).

For some of our later results we will need to assume that the UTV factorization is rank-revealing. The first factorization in the algorithm uses the QR factorization with standard column pivoting. There are contrived examples [3, p. 105] where standard column pivoting is not rank-revealing. To overcome this potential problem the algorithm could be started with a QR factorization that guarantees to be reveal rank [3, pp. 22, 108] or one could include pivoting at the second step [27]. Our numerical experience suggests that this is not necessary for examples that are not contrived.

Finally, for later use we would like to present some of the results of [21] that concern the convergence of the algorithm. We define ρ_i as in Theorem 2.

$$\text{if } T_i \text{ is lower triangular then } \|H_i\| \leq \rho_1 \rho_2 \dots \rho_{i-1} \sigma_k(\widehat{T}_1) \leq \rho_1^{i-1} \sigma_k(\widehat{T}_1) \quad (20)$$

$$\text{if } T_i \text{ is upper triangular then } \|F_i\| \leq \rho_1 \rho_2 \dots \rho_{i-1} \sigma_k(\widehat{T}_1) \leq \rho_1^{i-1} \sigma_k(\widehat{T}_1) \quad (21)$$

The inequalities (20) and (21) indicate that, if $\rho_1 < 1$, then the off-diagonal blocks F_i and G_i are forced to zero as the algorithm proceeds. These results combined with (14) and (15) show that the singular values of \hat{T}_i , E_i and G_i converge to singular values of A . Inequalities (20) and (21) combined with (16) - (19) show that if $\rho_1 < 1$ then $\sin \phi$ and $\sin \theta$ approach zero as the algorithm proceeds.

4 Perturbation errors

We now compare the perturbation error terms in (4), (5) and (6) for corner low-rank approximations calculated using a truncated UTV decomposition with the corresponding error terms when using a truncated SVD decomposition. We will assume in (4) that the regularization error term $(CA - I)x_0$ is sufficiently small so that

$$x - x_0 = C(\delta b) \quad (22)$$

is a good approximation. We will also assume initially that the UTV factorization has T lower triangular. As we noted following Theorem 1, equation (11) implies that in some cases $\|x_T - x_0\|$ will be smaller than $\|x_S - x_0\|$. The following theorem concerns the probability that this occurs, if we assume (22).

Theorem 3 *Let x_T be calculated using a corner low-rank UTV approximation to A with T lower triangular. Assume that (22) is true, that $\|E\| < \sigma_k(\hat{T})$ and that the components of δb come from uncorrelated zero mean Gaussian random variables with common variance (Gaussian white noise). Then as $\sin \phi$ approaches 0 the probability that $\|x_T - x_0\|$ is less than $\|x_S - x_0\|$ approaches one-half.*

Proof. Due to (22) in equation (11) we can assume that N is 0. By (9) we can write M as $M = \begin{pmatrix} -M_{11} & M_{12} \\ M_{12}^T & M_{22} \end{pmatrix}$ with $M_{11} = \hat{D}^{-1} \tilde{V}_{21}^T \tilde{V}_{21} \hat{D}^{-1}$, $M_{12} = \tilde{U}_{11}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12}$ and $M_{22} = \tilde{U}_{12}^T \hat{T}^{-T} \hat{T}^{-1} \tilde{U}_{12}$. Then it follows that

$$\|M_{11}\| \leq \frac{s_{k+1}}{s_k} (\tan \theta) \|M_{12}\| \leq (\tan \phi) \|M_{12}\| \quad \text{and} \quad \|M_{22}\| \leq (\tan \phi) \|M_{12}\| \quad (23)$$

where $\tan \phi = \sin \phi / \sqrt{1 - \sin^2 \phi}$ and $\tan \theta = \sin \theta / \sqrt{1 - \sin^2 \theta}$. These inequalities follow from (12), $\|E\| < \sigma_k(\hat{T})$ and the identities $M_{22} = \tilde{U}_{12}^T \tilde{U}_{11}^{-T} M_{12}$ and $M_{11} = M_{12} (D_0 0)^T \tilde{V}_{22}^{-1} \tilde{V}_{21} \hat{D}^{-1}$ which are consequences of (12) and properties of orthogonal matrices. For $\sin \phi$ small (23) implies that the diagonal blocks of M are small relative to the off-diagonal blocks. Consider the matrix $\tilde{M} = \begin{pmatrix} 0 & M_{12} \\ M_{12}^T & 0 \end{pmatrix}$ formed by the off-diagonal blocks. Note that the eigenvalues of \tilde{M} come in plus and minus pairs of equal magnitude. Since we are assuming that δb is governed by Gaussian white noise it follows from Theorem 4.4.8 and Corollary 5.4.2 of [24] that the distribution governing $\tilde{b}^T \tilde{M} \tilde{b}$ is symmetric and that the probability that $\tilde{b}^T \tilde{M} \tilde{b}$ is negative is one-half. Due to (23) the theorem follows from a continuity argument. \square

By the comments following (19) $\sin \phi$ will be small when $\|H\|$ is sufficiently small or when there is a sufficiently large gap in the singular values of A and the UTV factorization is rank-revealing. It follows under the conditions of the theorem that if $\sin \phi$ is small then x_T will be closer to x_0 than x_S is approximately half the time and, conversely, x_S will be closer approximately half the time. Our numerical experiments

support this. They also suggest that in some cases even when $\sin \phi$ is not small x_T is still frequently as close or closer to x_0 than x_S is.

We assumed in this theorem that the noise is Gaussian white noise. According to [28] “Gaussian white noise is a common occurrence in many signal processing systems.”

It is also of interest to look at the expected value of $\|x_T - x_0\|^2$ relative to the expected value of $\|x_S - x_0\|^2$ which we do in Theorem 4. The following lemma is used in the proof of Theorem 4.

Lemma 1 *Assume that $u \in R^m$ has components that come from uncorrelated zero mean random variables with common variance (white noise) and that the expected value of $\|u\|^2$, indicated by $E(\|u\|^2)$, is Δ^2 . Then for an m by m matrix A , $E(u^T A u) = \Delta^2 \text{trace}(A)/m$. Also for an n by m matrix A , $E(\|A u\|^2) = \Delta^2 \|A\|_F^2/m$ where $\|A\|_F$ indicates the Frobenius norm.*

Proof. Since $\Delta^2 = E(\sum_{i=1}^m u_i^2) = \sum_{i=1}^m E(u_i^2)$ then $E(u_i^2) = \Delta^2/m$. Now for an m by m matrix A , $E(u^T A u) = \sum \sum a_{ij} E(u_i u_j) = \sum a_{ii} E(u_i^2) = \text{trace}(A) \Delta^2/m$. Also for an n by m matrix A , $E(\|A u\|^2) = E(u^T A^T A u) = \Delta^2 \text{trace}(A^T A)/m = \Delta^2 \|A\|_F^2/m$. \square

Theorem 4 *Let x_T be calculated using a corner low-rank UTV approximation to A . Define $\sin \phi$ as in (16). Assume that T is lower triangular, that (22) is true and that the components of δb correspond to white noise. Then*

$$0 \leq \frac{E(\|x_T - x_0\|^2) - E(\|x_S - x_0\|^2)}{E(\|x_T - x_0\|^2)} \leq \sin^2 \phi. \quad (24)$$

Proof. Let us assume that $E(\|\delta b\|^2) = \Delta^2$. Since $\|x_S - x_0\|^2 = \|\widehat{A}_S^+ \delta b\|^2 = \|\widehat{V}_S \widehat{D}^{-1} \widehat{U}_S^T \delta b\|^2$ and $\|x_T - x_0\|^2 = \|\widehat{A}_T^+ \delta b\|^2 = \|\widehat{V} \widehat{T}^{-1} \widehat{U}^T \delta b\|^2$ it follows from the lemma that $E(\|x_S - x_0\|^2) = \Delta^2 \|\widehat{D}^{-1}\|_F^2/m$ and $E(\|x_T - x_0\|^2) = \Delta^2 \|\widehat{T}^{-1}\|_F^2/m$. The left inequality in (24) is true since $E(\|x_S - x_0\|^2) = \Delta^2 \|\widehat{D}^{-1}\|_F^2/m$, $E(\|x_T - x_0\|^2) = \Delta^2 \|\widehat{T}^{-1}\|_F^2/m$, since the Frobenius norm squared is the sum of the square of the singular values and due to the left inequality in (14). By (11), (22) and Lemma 1, $E(\|x_T - x_0\|^2) - E(\|x_S - x_0\|^2) = E(\delta b^T M \delta b) = \Delta^2 \text{trace}(M)/m = \Delta^2 [\text{trace}(M_{22}) - \text{trace}(M_{11})]/m \leq \Delta^2 \text{trace}(M_{22})/m = \Delta^2 \|\widehat{T}^{-1} \widehat{U}_{12}\|_F^2/m \leq \Delta^2 \|\widehat{U}_{12}\|_F^2 \|\widehat{T}^{-1}\|_F^2/m = \Delta^2 (\sin^2 \phi) \|\widehat{T}^{-1}\|_F^2/m$. The theorem now follows. \square

The left hand inequality in (24) implies under the conditions of the theorem that the truncated SVD solutions will, on average, be better than truncated UTV solutions. However, the right hand term suggests, as we will see in our numerical experiments, that often the difference, on average, will not be large and the truncated UTV and SVD will be similar in accuracy. Note that by the comments following equation (19) the size of $\sin \phi$ is related to the size of a gap in the singular values of A and to the size of H . Also note that $\sin^2 \phi$ in (24) can be small even for modest $\sin \phi$.

Theorems 3 and 4 are applicable to corner low-rank UTV approximations when T is lower triangular. When T is upper triangular one can prove, although we will not do so here, that equation (24) is valid except that $\sin^2 \phi$ must be replaced by $\sin^2 \theta$. Our numerical experiments suggest that there are results similar to Theorem 3 for the case that T is upper triangular.

5 Regularization errors

We now compare the regularization error terms in (4), (5) and (6) for corner low-rank approximations calculated using a truncated UTV decomposition with the corresponding error terms when using a truncated SVD decomposition. We will assume in (4) that the perturbation error term $C(\delta b)$ is sufficiently small so that

$$x - x_0 = (CA - I)x_0 \quad (25)$$

is a good approximation. We will also assume that the UTV factorization has T lower triangular.

Some of our results in this section will involve the values of components of $U_S^T b_0$. The discrete Picard condition [14, p. 507] is that these components decay to zero somewhat faster than the singular values. The condition is required for regularization to produce useful solutions [13, 14, 15]. We will call these components of $U_S^T b_0$ the ‘‘Picard coefficients’’ to indicate their connection to the Picard condition (the term Fourier coefficients is sometimes used). If we let \tilde{D} be the $n \times n$ diagonal matrix consisting of the first n rows of D in the singular value decomposition $A = U_S D V_S^T$, we will model the rate of decrease of the Picard coefficients by assuming that the first n components of $U_S^T b_0$ are equal to the components of $\tilde{D}^{p+1} w$ where \tilde{D}^{p+1} indicates the $(p+1)^{st}$ power of \tilde{D} , $p \geq 0$ and w is a vector whose components do not depend on p or the singular values of A . Following Hansen [13, 14, 15] who defines a similar parameter, we will call p the relative decay rate of the Picard or Fourier coefficients.

It will be useful to assume a particular form for the underlying noiseless solution x_0 . We assume that

$$x_0 = V_S \tilde{D}^p w \quad (26)$$

where \tilde{D} is first n rows of D . We have two motivations for this choice. First, with this x_0 , $b_0 = Ax_0 = U_S \begin{pmatrix} \tilde{D}^{p+1} & 0 \end{pmatrix}^T w$ so that the first n Picard coefficients are $\tilde{D}^{p+1} w$. Therefore p is the relative decay rate of the Picard coefficients. If $p > 0$ the Picard condition will be satisfied. Also note that by (26) x_0 is a linear combination of the singular vectors of A . Due to the factor \tilde{D}^p , if the singular values decrease sufficiently rapidly or if p is sufficiently large, the contribution of higher index singular vectors will be small. It is often the case that the lower index singular vectors correspond to smoothly varying functions [17, 18]. If these assumptions are true, as is frequently the case, x_0 will be smoothly varying. We also note that (26) is equivalent to the model [22, p. 640] for characterizing smooth solutions x_0 . We conclude that (26) provides a method, that has been used by others, to generate a class of smoothly varying solutions x_0 that satisfy the Picard condition.

Our results will involve the decay rate p of the Picard coefficients for smaller values of p since these values of p appear to be useful in many practical applications. For example we looked at 16 examples from Hansen’s Regularization Tools [16]. Most of the problems in [16] come from the literature and all share characteristic features of ill-posed problems. For each example we made a rough estimate of p by estimating the slope of a graph of the log of the Picard coefficients versus the log of the singular values (for values not dominated by errors). In 14 of the 16 cases the rough estimate was 1 or less. Our theorems in this section will assume $0 \leq p \leq 1$ (Theorem 5) and $0 \leq p \leq 2$ (Theorem 6).

As we noted following Theorem 1 and in Section 4, equation (11) implies that in some cases $\|x_T - x_0\|$ will be smaller than $\|x_S - x_0\|$. The following theorem concerns the probability that this occurs, if we assume (25).

Theorem 5 Let x_T be calculated using a corner low-rank UTV approximation to A with T lower triangular. Assume that x_0 satisfies (26) with $0 \leq p \leq 1$, that (25) is true, that $\|E\| < \sigma_k(\hat{T})$ and that w follows Gaussian white noise. Then as $\sin \phi$ approaches 0 the probability that $\|x_T - x_0\|$ is less than $\|x_S - x_0\|$ approaches one-half.

Proof. Due to (25) it follows that M is 0 in equation (11). Due to (26) and (10) we can write $\tilde{x}_0^T N \tilde{x}_0 = w^T N_p w$ where

$$N_p = \tilde{D}^p N \tilde{D}^p = \begin{pmatrix} \hat{D}^p \tilde{V}_{21}^T \tilde{V}_{21} \hat{D}^p & \hat{D}^p \tilde{V}_{21}^T \tilde{V}_{22} D_0^p \\ D_0^p \tilde{V}_{22}^T \tilde{V}_{21} \hat{D}^p & -D_0^p \tilde{V}_{12}^T \tilde{V}_{12} D_0^p \end{pmatrix} = \begin{pmatrix} N_{11} & N_{12} \\ N_{12}^T & -N_{22} \end{pmatrix}. \quad (27)$$

By equation (12) and properties of orthogonal matrices it follows that $N_{11} = N_{12}(D_0^{1-p} 0) \tilde{U}_{22}^{-1} \tilde{U}_{21} \hat{D}^{p-1}$ and $N_{22} = D_0^p \tilde{V}_{12}^T \tilde{V}_{11}^{-T} \hat{D}^{-p} N_{12}$. From these identities, $0 \leq p \leq 1$ and $\|E\| < \sigma_k(\hat{T})$ it follows that $\|N_{11}\| \leq (s_{k+1}/s_k)^{1-p} (\tan \phi) \|N_{12}\| \leq (\tan \phi) \|N_{12}\|$ and that $\|N_{22}\| \leq (s_{k+1}/s_k)^p (\tan \theta) \|N_{12}\| \leq (\tan \phi) \|N_{12}\|$. The rest of the proof follows in manner very similar to the proof of Theorem 3. \square

It follows under the conditions of the theorem that if $\sin \phi$ is small then $\|x_T - x_0\|$ will be smaller than $\|x_S - x_0\|$ approximately half the time and, conversely, $\|x_S - x_0\|$ will be smaller approximately half the time. One condition of the theorem is that $0 \leq p \leq 1$. Numerical experiments suggest that the conclusion of the theorem is also true for $1 \leq p < 2$. They also suggest that in some cases even when $\sin \phi$ is not small x_T is still frequently as close or closer to x_0 than x_S is.

For regularization errors we can also prove a useful result about the expected value of the errors.

Theorem 6 Let x_T be calculated using a corner low-rank UTV approximation to A . Assume that T is lower triangular, that (25) is true, that x_0 satisfies (26) with $0 \leq p \leq 2$ and that the components of w correspond to white noise. If

$$\alpha = \left(\frac{s_k}{s_{k+1}} \right)^p [\|H\| + (\sin \theta) \|E\|] \|E\|_F / s_k^2 \quad \text{then} \quad (28)$$

$$-\sin^2 \theta \leq \frac{E(\|x_T - x_0\|^2) - E(\|x_S - x_0\|^2)}{E(\|x_S - x_0\|^2)} \leq \alpha^2. \quad (29)$$

Proof. Lemma 1 and (25) imply for $p \geq 0$ that $E(\|x_S - x_0\|^2) = \tau^2 \|D_0^p\|_F^2 / n$, where $\tau^2 \equiv E(\|w\|^2)$. Lemma 1, (11), (25) and (27) imply $E(\|x_T - x_0\|^2) - E(\|x_S - x_0\|^2) = \tau^2 \text{trace}(N_p) / n$. By (27) it follows that $-(\sin^2 \theta) \|D_0^p\|_F^2 \leq -\|\tilde{V}_{12} D_0^p\|_F^2 \leq \text{trace}(N_p) \leq \|\tilde{V}_{21} \hat{D}^p\|_F^2$. For T lower triangular (12) implies that $\tilde{V}_{21} \hat{D}^p = E^T H \tilde{V}_{11} + E^T E \tilde{V}_{21}$. Therefore for $0 \leq p \leq 2$ we have $\text{trace}(N_p) \leq \|\tilde{V}_{21} \hat{D}^p\|_F^2 \leq \|E^T H \tilde{V}_{11} + E^T E \tilde{V}_{21}\|_F^2 \leq s_k^{2(p-2)} (\|H\| + (\sin \theta) \|E\|)^2 \|E\|_F^2$. Since $s_{k+1}^{2p} \leq \|D_0^p\|_F^2$ the theorem follows. \square

For TQRLP or at any subsequent step of the algorithm with T lower triangular, it follows from (16), (17), (20) and (28) that, for $0 \leq p < 2$, $\sin \phi$, $\sin \theta$ and α will be small either if $\|H\|$ is sufficiently small or if the UTV factorization is rank-revealing and there is a sufficiently large gap in the singular values. Note that $\sin^2 \theta$ and α^2 may be small even for modest $\sin \theta$ and α .

The right hand bound in (29) increases in magnitude as p increases. This suggests that the solutions produced by a truncated UTV factorization will be best, relative to those produced by the SVD, for smaller values of p . As mentioned earlier, in practice

values of p one or less appear to be common. For larger values of p , accuracy close to that of the SVD can be achieved by using additional steps in the algorithm. As seen in (20) if $\rho_1 < 1$ these steps will force $\|H_i\|$ and $\sin\theta_i$ to become small. Note the [22, p. 644] discusses the effect of p on classical Tikhonov regularization.

Theorems 5 and 6 assume equation (26), $x_0 = V_S \tilde{D}^p w$, applies where w is governed by Gaussian white noise (Theorem 5) or white noise (Theorem 6). These may be only rough models of solutions x_0 as they appear in practical applications. However note that Neumaier [22, p. 641] comments that a model equivalent to (26) is “a frequently used assumption” and, in addition, the model has been used with similar statistical assumptions about the components of w [2, 22]. A conclusion from Theorem 5 is, under the conditions of the theorem, that $\|x_T - x_0\|$ is frequently smaller than $\|x_S - x_0\|$. This is consistent with our experiments using examples from [16] where x_0 is not chosen randomly (see Table 2).

6 Implementation and numerical experiments

Before discussing our numerical experiments we will discuss some implementation issues for the algorithm of Section 3 and the efficiency of the algorithm. For a point of comparison we note that there are three classical methods for solving least squares problems – the QR factorization without column interchanges, the QR factorizations with column interchanges and the SVD. The first algorithm is not reliable for solving rank-deficient problems but we include it for comparison of the efficiency of the algorithms. These algorithms are implemented in LAPACK as xGELS, xGELSY and xGELSD. For large n and $m \geq n$ (but not too much bigger than n) the approximate flop counts are $2mn^2 - 2/3n^3$, $2mn^2 - 2/3n^3$ and $4mn^2 - 4/3n^3$ [23], respectively, for xGELS, xGELSY, and xGELSD for the full rank case. For the case that $m = n$ these counts predict run times in the ratio 1:1:2 for the three algorithms. However in practice xGELS makes more effective use of the potential speedup in BLAS-3 calculations and the actual run time ratios depend on the computer architecture and matrix size. As illustrations of potential actual run-time ratios note that LAPACK [1, p. 72] reports ratios of 1:1:4 for 900 by 900 matrices run on an Compaq AlphaServer DS-20, Ren [23, p. 94] reports ratios of 1:1.3:3.4 for 1600 by 1600 matrices run on an IBM RS 6000/590 and our numerical experiments indicate ratios of 1:2.1:4.7 for 1600 by 1600 matrices run on a 700 MHertz Pentium computer.

We will consider the construction of a block-row/column low-rank approximate solution to (1) using the algorithm in Section 3. In the algorithm if k , the effective numerical rank, is less than n it is not necessary to do complete QR factorizations. One can start the algorithm with a QR factorization with the usual pivoting scheme [4] and stop the initial factorization when, for example, norms of the columns of E_1 (or G_1) are small. The matrices $E_i, i \geq 1$ and $G_i, i \geq 1$, need not be factored in order to calculate the solutions, x_T , at subsequent steps of the algorithm. An efficient way to implement the algorithm is to begin with the initial partial factorization just described. At subsequent steps one can construct orthogonal factorizations that successively update the block triangular structure of T while keeping the structure (lower or upper triangular) of \hat{T} fixed. The solution x_T produced by this implementation of the algorithm is identical to the solution produced by the implementation described in Section 3. With this implementation the overall flop count for i steps of the algorithm is approximately $4k^2(n - 2/3k) + 2k(n - k)(2n - k)i$ where, for simplicity, the count is for the case where $m = n$. For any $m \geq n$ and for $i = 1$ the flop count for

the algorithm is approximately $4mnk - 2k^2m - 2k^3/3$. For $k < n$ this is less than the flop count for xGELSY. Also we can show that for $m \geq n$ and $i \leq 2$ the flop count of the algorithm is less than the theoretical flop count for xGELSD. The advantage of the algorithm is most striking in the low-rank case where $k \ll n$. In this case, for $m \geq n$, the leading order term in the flop count is $4kmni$ which is substantially less than the theoretical counts for xGELSY and xGELSD. Alternative algorithms for the low-rank case are discussed in [5, 11]. The smallest flop count of the algorithms in [11] is $12mnk$. The flop count for the algorithm of Section 3 for $i = 1$ will also be smaller than the count for the algorithm in [5].

The actual time required by an algorithm depends on details of its implementation and the computer architecture as well as flop counts. As discussed earlier the block-row low-rank solution produced by TQRP (or by Theorem 2 the corner low-rank solution produced by TQRLP) will be the same as the solution produced by LAPACK's xGELSY. However LAPACK does a complete factorization of A not a partial factorization as discussed in the last paragraph. The routine xGELSY can be modified to incorporate the partial factorization. The tests in xGELSY for the determination of the effective rank can be moved into LAPACK's factorization routine xGEPQ3. If these tests are inserted in xGEPQ3 immediately after the call to LAPACK's xLAQPS, the efficient BLAS-3 calls in xGEPQ3 will not be affected. The solution x_T produced by this modification is identical to that of the unmodified LAPACK but the modified algorithm will run much more quickly for low-rank problems. This is illustrated in Figure 1. The sample problems in Figure 1 were generated by LAPACK's xLATMS and have a gap in the singular values at the indicated numerical rank. They were run on a 700 MHertz Pentium computer using BLAS routines supplied by Intel. From the graph it is clear that for low-rank problems the modification of xGELSY is much more efficient than the existing implementations of LAPACK's routines xGELSY and xGELSD. For $k = 25$ the run-time ratios are approximately 1:14:31. Due to this substantial speedup it also clear that in the low-rank case the algorithm of Section 3 will remain more efficient than the LAPACK routines if the algorithm is continued with some additional steps.

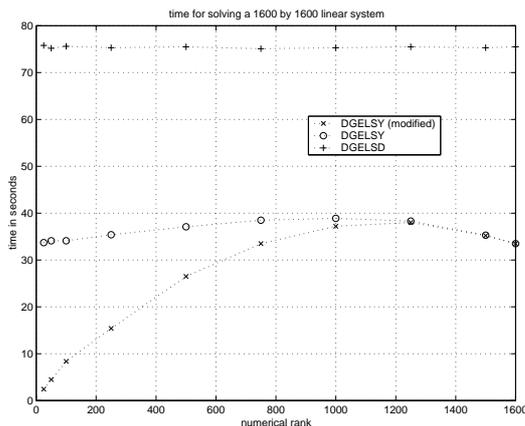


Figure 1: Timings for the modification of DGELSY, DGELSY and DGELSD.

We should add here that another issue that can be important in choosing an algorithm to solve (1) is the ability to easily do updates and downdates. This is more

easily done with a UTV factorization than the SVD [25]. Also note that [26, 27] discuss implementation issues for the QRLP algorithm including the observation that in the low-rank case the savings in stopping the reduction are substantial. Finally, we should note that it is well known [12, p. 250] that truncating the QR factorization reduces the flop count in the factorization to approximately $4mnk$ for small k .

We now present experiments that focus on the accuracy of the algorithm. For our first test results we generated random 64×64 matrices A using REGUTM from [16]. We chose the singular values of A in three manners. To describe the first type of selection let us define quantities which we call the “gap” and the “spread” where the gap is s_{16}/s_{17} and the spread is $s_1/s_{16} = s_{17}/s_{64}$. Singular values s_2 to s_{15} were selected from a log-uniform distribution over s_1 to s_{16} and singular values s_{18} to s_{63} from a log-uniform over s_{17} to s_{64} . For this selection of singular values we fixed the effective numerical rank at 16. In the second selection of singular values we selected 10 singular values from a log-uniform distribution from 1 to .001, 10 singular values equal to .001, and 44 singular values from a log-uniform distribution from .001 to .000001. We again fixed the effective numerical rank at 16. This selection is designed to test the algorithm in an extreme case, forcing the algorithm to select the numerical rank in the middle of a cluster of identical singular values. To describe the third selection of singular values we will use a quantity which we call the “mean gap.” We let $s_1 = 1$ and $s_{64} = (\text{meangap})^{63}$ and chose s_2 through s_{63} from a log-uniform distribution from s_1 to s_{64} . For this selection of singular values the geometric mean of s_j/s_{j+1} , $j = 1, 2, \dots, 63$ equals the mean gap. In this case the singular values decay gradually and there is not an obvious gap to help select the effective rank k . A variety of approaches have been suggested in the literature [18] for selecting k . For simplicity and to focus on the approximation scheme, not the technique for selecting k , for the third choice of singular values and for each regularization algorithm we selected the effective rank by calculating a regularized solution, x , for each numerical rank $k < n$. Among all lower rank approximations, the approximation that minimizes $\|x - x_0\|$ was selected.

For each matrix A we chose the underlying noiseless solution $x_0 = V\tilde{D}^p w$, $b_0 = Ax_0$ and noise vectors $\delta b = \Delta \|b_0\| v$. We used seven different noise to signal ratios, $\Delta = .3, .1, .01, .001, 10^{-4}, 10^{-6}$ and 10^{-10} . This wide range of noise levels should produce cases where the regularization error dominates (Δ small), where the perturbation error dominates (Δ large) and cases in between these extremes. We selected 100 random matrices as described above. For each matrix and for each of the seven noise levels we selected 100 random values of $x_0 = \tilde{D}^p w$, with w selected from white noise and for each x_0 we selected a random noise vector δb , with v selected from white noise for a total of 70,000 ($= 100 \times 7 \times 100$) samples. For each sample we calculated x_S as well as five different solutions x_T . To calculate x_T we used block-row/column low-rank approximations for the TQLP, TQLRP, TQRP, TQRLP and TQRLRP factorizations. To summarize the results in a concise manner, for each low-rank approximation we calculated the mean value of $\|x_T - x_0\|/\|x_S - x_0\| - 1$ over all 70,000 samples. These results are in Table 1. Also in the last column of Table 1 we indicate the percent of the cases where $\|x_T - x_0\|$ is smaller than $\|x_S - x_0\|$ for the block-row low-rank solutions using the TQRP factorization. These solutions can be produced by LAPACK’s xGELSY.

Most of the entries in the table are positive, which indicates that on average the truncated SVD solutions are better. However except for TQLP, if $p < 2$, the truncated UTV solutions are on average not far from the truncated SVD solutions. For example,

Problem Properties			Method					%
g	s	p	TQLP	TQLRP	TQRP	TQRLP	TQRLRP	
100	100	.5	.0015	2.1×10^{-7}	2.8×10^{-5}	7.7×10^{-8}	-7.8×10^{-10}	50
100	100	1	.23	1.2×10^{-7}	6.8×10^{-5}	1.0×10^{-5}	8.4×10^{-10}	50
100	100	1.5	5.3	1.8×10^{-6}	4.1×10^{-4}	3.2×10^{-4}	2.2×10^{-10}	49
100	100	2	817	1.3×10^{-5}	.14	.14	-2.3×10^{-8}	44
100	100	3	6882	3.3×10^{-6}	.65	.65	-2.0×10^{-8}	43
100	10^4	1	.14	3.5×10^{-7}	1.7×10^{-4}	6.5×10^{-6}	2.1×10^{-9}	50
100	1	1	.29	-3.8×10^{-8}	7.3×10^{-5}	3.5×10^{-5}	-3.4×10^{-10}	49
10	100	1	.32	3.1×10^{-5}	5.9×10^{-3}	9.1×10^{-4}	9.9×10^{-7}	47
4	100	1	.40	.0016	.040	.0091	6.5×10^{-5}	42
1	100	1	.36	.082	.13	.099	.063	39
cluster	1		.11	.044	.067	.046	.037	45
m.g. 10	1		.28	.0081	.024	.015	.0044	49
m.g. 10	2		81	.019	.18	.16	.0076	44
m.g. 10	3		1163	.11	11	11	.016	23
m.g. 4	1		.29	.015	.038	.024	.0094	48
m.g. 1.2	1		.34	.041	.091	.065	.024	32

Table 1. Mean value of $\frac{\|x_T - x_0\| - \|x_S - x_0\|}{\|x_S - x_0\|}$ for block-row/column low-rank approximations and, in the last column, the percent of the runs where, for TQRP, x_T is closer to x_0 than is x_S . In the table g stands for gap, s for spread and m. g. for mean gap. These terms and the term cluster are defined in the text. p is the decay rate in the Picard coefficients. Each entry summarizes 70,000 samples.

for $p < 2$, $\|x_T - x_0\|$ was within 15% of $\|x_S - x_0\|$ on average for all the methods except TQLP. Remarkably this is true for runs with a small gap or no gap in the singular values and when the numerical rank is selected in the middle of a cluster of singular values. As p increases more steps of the algorithm are required to match the accuracy of the SVD. As mentioned earlier, smaller values of p appear to be more common in practice. Note that we arrive at these same general conclusions by looking at the cases where the rank of the low-rank approximation is fixed at 16 or the “mean gap” cases where the rank is chosen dynamically. Also note that the last column of the table indicates, as suggested by Theorems 3 and 5, that if there is a sufficiently large gap in the singular values and if p is not large then $\|x_T - x_0\|$ is smaller than $\|x_S - x_0\|$ close to 50% of the time for block-row TQRP solutions. For the problems with a small or no gap and $p = 1$ the percent of the cases where the block-row TQRP solution is closer to x_0 than is the TSVD solution varied between 45% for the cluster example to 32% for the runs with a mean gap of 1.2.

The examples so far have been artificial. To test examples from practice or used elsewhere in the literature we looked at problems from Hansen’s Regularization Tools [16]. Our sample consists of Hansen’s baart, deriv2 (with 3 different solutions), fox-good, heat (with 3 parameter values), ilaplace (with 4 different solutions), phillips, shaw, spikes, and wing, for a total of 16 different examples. This is all the ill-conditioned examples in Regularization Tools, except for parallax and ursell for which

x_0 is not supplied and blur which is parameterized differently from the other examples. Most of these example do not have a clear gap in the singular value spectrum and so we need a technique to choose the numerical rank k . For simplicity and to focus on the approximation scheme, not the technique for selecting k , for each regularization algorithm we selected the effective rank by calculating a regularized solution, x , for each numerical rank $k < n$ and then among all lower rank approximations, selecting the approximation that minimizes $\|x - x_0\|$.

For each of the 16 examples we looked at the seven noise levels used in Table 1, for a total of 112 cases. For each case we chose 100 random noise vectors, applied a variety of regularization methods and calculated the mean values of $(\|x_T - x_0\| - \|x_S - x_0\|)$ and of $\|x_S - x_0\|$. In each of the 112 cases we used the x_0 supplied by Regularization Tools, not a randomly chosen x_0 . Each mean value is a mean over 100 different random noise vectors δb for a fixed x_0 . In Table 2 we summarize the results for the block-row/column low-rank solutions produced by TQLP, TQLRP, TQRP and TQRLRP factorizations. Each entry counts the number of the 112 cases where $\frac{\text{mean}(\|x_T - x_0\| - \|x_S - x_0\|)}{\text{mean}(\|x_S - x_0\|)}$ is in the indicated range.

Range	less than	-50%	-10%	-5%	-1%	1%	5%	10%	50%
Method	-50%	-10%	-5%	-1%	1%	5%	10%	50%	or more
TQRP	0	13	8	17	50	11	8	5	0
TQRLP	0	11	12	16	53	7	8	5	0
TQRLRP	0	7	7	7	77	12	0	2	0
TQLP	9	16	4	10	37	8	5	17	6
TQLRP	0	5	8	8	65	15	5	6	0

Table 2. Counts for examples with characteristic features of ill-posed problems from [16] of the number of cases, out of 112 cases, where $\frac{\text{mean}(\|x_T - x_0\| - \|x_S - x_0\|)}{\text{mean}(\|x_S - x_0\|)}$ is in the range indicated in the first row of the table. Rows two through six correspond to block-row/column solutions for different truncated factorizations.

In this table in some cases the truncated SVD solutions are closer to x_0 and in others the truncated UTV solutions are closer. However, overall for this set of problems the truncated UTV algorithm, even when stopped at the first step, appears to work as well as the truncated SVD. The table also indicates that additional steps of the algorithm bring values of $\text{mean}(\|x_T - x_0\|)$ closer to values of $\text{mean}(\|x_S - x_0\|)$. We also kept track of the percent of the time that $\|x_T - x_0\|$ was less than $\|x_S - x_0\|$ over the 11,200 ($= 112 \times 100$) samples. These percents were 51%, 54%, 51%, 57% and 57%, respectively, for the block-row/column solutions corresponding to the TQRP, TQRLP, TQRLRP, TQLP, and TQLRP factorizations. It is interesting to note that the results for these test problems, where A and x_0 are not random, seem to favor the truncated UTV solutions more than do the results for test problems involving randomly generated examples. The reason for this merits further investigation.

In order to understand Table 2 better it is useful to look at a specific case, for example, the Phillips example of [16] when the noise to signal ratio, $\|\delta b\|/\|b\|$, equals 0.1. We can illustrate the results in the table by looking at $\|x_T - x_0\|$ for TQRP and $\|x_S - x_0\|$ for a few typical values of δb . For the Phillips example the underlying true solution x_0 provided by [16] has $\|x_0\| = 2.99$. Six typical values of $\|x_T - x_0\|$ are .25, .29, .17, .30, .25 and .29 and the corresponding values of $\|x_S - x_0\|$ are .21, .34,

.14, .34, .23 and .33. Overall the magnitude of these values are quite similar and the two methods have approximately the same accuracy. The differences between these values are, respectively, .04, $-.05$, .03, $-.04$, .02 and $-.04$ and the corresponding values of $\|x_T - x_S\|$ are .09, .09, .07, .09, .04 and .41, respectively. In the table the sample size was 100, not 6. For these 100 values $\frac{\text{mean}(\|x_T - x_0\| - \|x_S - x_0\|)}{\text{mean}(\|x_S - x_0\|)}$ was $-.042$ and this example is one of the 17 entries in the table for the TQR method with $\frac{\text{mean}(\|x_T - x_0\| - \|x_S - x_0\|)}{\text{mean}(\|x_S - x_0\|)}$ between -5% and -1% .

7 Conclusions

We have discussed the application of the algorithm of [21] to solving ill-posed and rank-deficient problems. The algorithm constructs a UTV factorization of A by using one or more QR factorizations. The following are some of our results.

- The block-row solution produced by a rank-revealing QR factorization is identical to the corner solution produced by a related rank-revealing UTV factorization. (See Theorem 2 and the comments following the theorem.)
- If there is a modest gap in the singular values so that $\rho_1 < 1$, pivoting is not needed after the first step in the algorithm of Section 3. (See Theorem 2.)
- We have presented an implementation of LAPACK's xGELSY that, in the low-rank case, is substantially faster than the implementation of xGELSY currently in LAPACK. (See Figure 1 and the discussion prior to Figure 1.)

Some of our results concern the accuracy, relative to truncated SVD solutions, of the solutions to (1) produced by truncated UTV factorizations. The results suggest the following recommendations about the appropriate choice of a method to use to construct regularized solutions to the system (1).

- If one can identify and evaluate the accuracy of typical examples then we recommend that a variety of methods of regularization be compared for these examples. Our results indicate that although in some examples a relatively expensive method such as the truncated SVD will produce the best solution in other examples cheaper methods will calculate solutions as close or closer to the underlying desired solution. (See Theorems 3 and 5 and Tables 1 and 2.)
- If the initial QR factorization is rank-revealing, if the desired regularized solution corresponds to a sufficiently large gap in the singular values, and if p , the decay rate in the Picard coefficients, is not too large, as is often true in practice, then we recommend using the block-row truncated QRP solution. On average this truncated QRP solution will be very close to the accuracy of the truncated SVD solution and it can be calculated more quickly, dramatically so for low-rank problems. (See Theorem 4, Theorem 6, Figure 1, Table 1 and Table 2.)
- If the desired solution does not correspond to a gap in the singular values our experimental results with random examples suggest that truncated SVD solutions are, on average, somewhat better than truncated UTV solutions but, for $p < 2$, the difference may be modest (see Table 1). The case where there is not a gap in the singular values merits further investigation. For this case Stewart [27] conjectures for the QRLP decomposition that “the analysis of this decomposition will not be simple.”

We also did test runs for the set of problems of [16], which have characteristic features of ill-posed problems. In some cases the truncated SVD solutions were closer to the desired solution and in others the truncated UTV solutions were closer. Overall for this set of problems the truncated UTV algorithm, even when stopped at the first step, appeared to work as well as the truncated SVD algorithm (see Table 2).

Acknowledgments. The helpful comments and suggestions from the anonymous referees are gratefully acknowledged.

References

- [1] E. Anderson, Z. Bai, C. Bischof, , S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide, Third Edition*. SIAM, Philadelphia, 1999.
- [2] M. Bertero, C. De Mol, and G. A. Viano. The stability of inverse problems. In H. P. Baltes, editor, *Scattering in Optics*, pages 161–214. Springer-Verlag, New York, 1980.
- [3] A. Bjorck. *Numerical Methods for Least Squared Problems*. SIAM, Philadelphia, 1996.
- [4] P. Businger and G. H. Golub. Linear least squares solutions by Householder transformations. *Numerische Mathematik*, 7:269–276, 1965.
- [5] T. F. Chan and P. C. Hansen. Low-rank revealing QR factorizations. *Numerical Linear Algebra with Applications*, 1:33–44, 1991.
- [6] T. F. Chan and P. C. Hansen. Some applications of the rank revealing QR factorization. *SIAM Journal on Scientific and Statistical Computing*, 13:727–741, 1992.
- [7] S. Chandrasekaran and I.C.F. Ipsen. Analysis of a QR algorithm for computing singular values. *SIAM J. Matrix Anal. App.*, 16:520–535, 1995.
- [8] R. D. Fierro. Perturbation analysis for two-sided (or complete) orthogonal decompositions. *SIAM Journal on Matrix Analysis and Applications*, 17:383–400, 1996.
- [9] R. D. Fierro and J. R. Bunch. Bounding the subspaces from rank revealing two-sided orthogonal decompositions. *SIAM J. Matrix Anal. App.*, 16:743–759, 1995.
- [10] R. D. Fierro and P. C. Hansen. Accuracy of TSVD solutions computed from rank-revealing decompositions. *Numer. Math.*, 70:453–471, 1995.
- [11] R. D. Fierro and P. C. Hansen. Low-rank revealing UTV decompositions. *Numerical Algorithms*, 15:37–55, 1997.
- [12] G. Golub and C. F. Van Loan. *Matrix Computations*. John Hopkins, Baltimore, 1996.
- [13] P. C. Hansen. The discrete Picard condition for discrete ill-posed problems. *BIT*, 30:658–672, 1990.

- [14] P. C. Hansen. Truncated singular value decomposition solutions to discrete ill-posed problems with ill-determined numerical rank. *SIAM Journal on Scientific and Statistical Computing*, 11:503–518, 1990.
- [15] P. C. Hansen. Analysis of discrete ill-posed problems by means of the L-curve. *SIAM Review*, 34:561–580, 1992.
- [16] P. C. Hansen. Regularization tools: A matlab package for analysis and solution of discrete ill-posed problems. *Numer. Algorithms*, 6:1–35, 1994.
- [17] P. C. Hansen. Test matrices for regularization methods. *SIAM J. Sci. Comput.*, 16:506–512, 1995.
- [18] P. C. Hansen. *Rank-Deficient and Discrete Ill-Posed Problems*. SIAM, Philadelphia, 1998.
- [19] Y. Hosoda. Truncated least-squares least-norm solutions by applying the QR decomposition twice. *Transactions of the Information Processing Society of Japan*, 40:1051–1055, 1999.
- [20] C. L. Lawson and R. J. Hanson. *Solving Least Squares Problems*. Prentice-Hall, Englewood Cliffs, N. J., 1974.
- [21] R. Mathias and G. W. Stewart. A block QR algorithm and the singular value decomposition. *Linear Algebra and Its Applications*, 182:91–100, 1993.
- [22] A. Neumaier. Solving ill-conditioned and singular linear systems: A tutorial on regularization. *SIAM Review*, 40:636–666, 1999.
- [23] H. Ren. On the error analysis and implementation of some eigenvalue decomposition and singular value decomposition algorithms. UT, CS-96-336, LAPACK Working Note 115, 1996.
- [24] V. K. Rohatgi. *An Introduction to Probability Theory and Mathematical Statistics*. John Wiley and Sons, NY, 1976.
- [25] G. W. Stewart. An updating algorithm for subspace tracking. *IEEE Trans. Signal Proc.*, 40:1535–1541, 1992.
- [26] G. W. Stewart. *Matrix Algorithms Volume 1: Basic Decompositions*. SIAM, Philadelphia, 1998.
- [27] G. W. Stewart. The QLP approximation to the singular value decomposition. *SIAM J. Sci. Comput.*, 20:1336–1348, 1999.
- [28] C. W. Therrien. *Discrete Random Signals and Statistical Signal Processing*. Prentice Hall, Englewood Cliffs, NJ, 1992.