

CONF-840872--14-Rev.

NOTICE

PORTIONS OF THIS REPORT ARE ILLEGIBLE. It has been reproduced from the best available copy to permit the broadest possible availability.

Los Alamos National Laboratory is operated by the University of California for the United States Department of Energy under contract W-7405-ENG-36

LA-UR--84-2336-Rev.

DE84 016488

TITLE: ALGORITHM FOR SINGULAR DECOMPOSITION

AUTHOR(S): D. C. ROSS

SUBMITTED TO REAL-TIME SIGNAL PROCESSING VII, SPIE 28th Annual International Technical Symposium and Instrument Exhibition, San Diego, CA, August 1984.

MASTER

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

By acceptance of this article the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution or to allow others to do so, for U.S. Government purposes.

The Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy.

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED

Los Alamos Los Alamos National Laboratory Los Alamos, New Mexico 87545

Algorithm for Singular Value Decomposition

D. C. Ross

Electronics Division
Los Alamos National Laboratory
P.O. Box 1663, Los Alamos, New Mexico 87545

Abstract

An iterative algorithm for the singular value decomposition (SVD) of a non-zero $m \times n$ matrix M is described and illustrated numerically. Derivations of the algorithm and sufficient conditions for convergence are outlined.

SVD is one of the most important procedures in digital processing of signals and images, and in applied mathematics generally. SVD provides an effective way to find the rank of a matrix, to compress data, to find the pseudo-inverse of a matrix and, in general, to calculate with rectangular and square asymmetric matrices almost as easily as with square symmetric matrices. The eigensystem of covariance matrices and other symmetric matrices of the form $A \bar{A}$ may be found accurately from the SVD of A . The theory of the SVD is well presented in a famous 1958 paper by Cornelius Lanczos in American Mathematical Monthly. The basic facts are that any non-zero matrix M of rank r may be written as the product of three factors: (1) $m \times r$ partial isometry \bar{U} , (2) positive-definite $r \times r$ diagonal matrix D , (3) $r \times n$ partial isometry \bar{V} .

$$M = \bar{U} D \bar{V} \quad \bar{U} \bar{U}^T = \bar{V} \bar{V}^T = I \quad D = \bar{U} M \bar{V}^T$$

The pseudo-inverse of M is then: $M^+ = \bar{V} D^{-1} \bar{U}^T$.

Each pass of the algorithm produces one set of corresponding singular elements; that is, one diagonal element of D along with the corresponding rows of \bar{U} and \bar{V} . Each pass starts with a trial row of \bar{U} and a trial row of \bar{V} . A trial singular value corresponding to the starting rows is then found along with measures of error. The algorithm then finds a new pair of trial singular rows to start the next iteration. The matrix M is progressively deflated and the deflated matrix is used to start each pass, but the original matrix is used for all computations within each pass to avoid unnecessary accumulation of round-off error.

1. An important procedure that is often used in the digital processing of signals and images, and in applied mathematics generally, is the singular value decomposition (SVD) of a non-zero $m \times n$ matrix M of rank r .

$$M = \bar{U} D \bar{V} = \sum_{k=1}^r \bar{U}_k \gamma_k \langle \bar{V}_k \quad \bar{U} \bar{U}^T = I = \bar{V} \bar{V}^T \tag{1}$$

$$D = \bar{U} M \bar{V}^T \quad \gamma_k = \langle \bar{U}_k M \bar{V}_k \rangle \quad \bar{U} M = D \bar{V} \quad M \bar{V}^T = \bar{U} D \tag{2}$$

D is a positive-definite $r \times r$ diagonal matrix and I the $r \times r$ identity matrix (Lanczos 1958). Each set

$$\{ \bar{U}_k \}, \gamma_k, \langle \bar{V}_k \rangle \quad k = 1, 2, \dots, r \tag{3}$$

is called a set of corresponding singular elements of the matrix M and their product is an $m \times n$ matrix called the corresponding dyad. The $n \times m$ matrix $M^+ = \bar{V} D^{-1} \bar{U}^T$ is the pseudo-inverse of M . The $m \times m$ matrix $M M^+ = \bar{U} \bar{U}^T$ and the $n \times n$ matrix $M^+ M = \bar{V} \bar{V}^T$ represent perpendicular projectors on the column space and on the row space, respectively, of the matrix M .

The main purpose of this paper is to present an iterative algorithm for the SVD. The Euclidean norm is used throughout, along with the standard inner product and Dirac notation (Huggins 1963). A vector $\langle X \rangle$ in a primal space (a space of signals for example) is represented by a row $\langle X$. A vector $|P\rangle$ in the dual space (for example, the set of linear measurements derived appropriate for the given signals) is represented by a column P . Bases are orthonormal unless otherwise indicated. Though our results may be extended easily to complex vector spaces, we confine attention in this paper to real spaces so that \bar{m} means transposition of rows and columns.

2. It is convenient to scale \underline{M} initially so that the magnitude of its largest element is near one and then to scale \underline{D} appropriately at the end of the algorithm. A unit trial m -row $\langle \underline{U}$ and a unit trial n -row $\langle \underline{V}$ are needed to start the algorithm. The normalized transpose of the largest column of \underline{M} is a reasonable choice for $\langle \underline{U}$ along with the normalized largest row of \underline{M} for $\langle \underline{V}$.

3. A trial singular value corresponding to the trial rows is then found by

$$\gamma = \langle \underline{U} \underline{M} \underline{V} \rangle = \langle \underline{V} \underline{M} \underline{U} \rangle \quad (4)$$

which must be positive, by reversing the sign of one of the trial rows if necessary. The squared norms of the residual rows

$$\langle \underline{R} \underline{R} \rangle = \langle \underline{V} \underline{M} - \gamma \underline{U} \cdot \quad \quad \quad \langle \underline{S} \underline{S} \rangle = \langle \underline{U} \underline{M} - \gamma \langle \underline{V}$$

are then found by

$$\begin{aligned} \langle \underline{R} \underline{R} \rangle &= v^2 - \gamma^2 & v^2 &\triangleq \langle \underline{V} \underline{M} \underline{M} \underline{V} \rangle \\ \langle \underline{S} \underline{S} \rangle &= \mu^2 - \gamma^2 & \mu^2 &\triangleq \langle \underline{U} \underline{M} \underline{M} \underline{U} \rangle \end{aligned} \quad \tau \triangleq \langle \underline{R} \underline{R} \rangle + \langle \underline{S} \underline{S} \rangle \quad (5)$$

If both residuals are sufficiently small as measured by τ , then $\langle \underline{U} \rangle$, γ , $\langle \underline{V} \rangle$ closely equals one of the sets of corresponding singular elements of \underline{M} as can be seen from the following theorems.

Theorem 1: $\langle \underline{R} \rangle \perp \langle \underline{U} \rangle$
 Proof: $\langle \underline{R} | \underline{U} \rangle = \langle \underline{R} \underline{U} \rangle = \langle \underline{V} \underline{M} \underline{U} \rangle - \langle \underline{U} \underline{M} \underline{V} \rangle \langle \underline{U} \underline{U} \rangle = \langle \underline{U} \underline{M} \underline{V} \rangle - \langle \underline{U} \underline{M} \underline{V} \rangle = 0$ □

Theorem 2: $\langle \underline{S} \rangle \perp \langle \underline{V} \rangle$
 Proof: $\langle \underline{S} | \underline{V} \rangle = \langle \underline{S} \underline{V} \rangle = \langle \underline{U} \underline{M} \underline{V} \rangle - \langle \underline{U} \underline{M} \underline{V} \rangle \langle \underline{V} \underline{V} \rangle = \langle \underline{U} \underline{M} \underline{V} \rangle - \langle \underline{U} \underline{M} \underline{V} \rangle = 0$ □

Theorem 3: $\gamma \langle \underline{U} \rangle$ is the \perp projection of $\langle \underline{V} \underline{M} |$ onto $\langle \underline{U} \rangle$.
 Proof: $\gamma \langle \underline{U} \rangle$ is collinear with $\langle \underline{U} | \underline{M} \underline{V} \rangle$. $\langle \underline{U} \rangle$ is self-adjoint and idempotent.
 $\gamma \langle \underline{U} \rangle = \langle \underline{U} \underline{M} \underline{V} \rangle \langle \underline{U} \rangle = \langle \underline{V} \underline{M} | \underline{U} \rangle \langle \underline{U} \rangle = \left[\langle \underline{V} \underline{M} | \right] \left[\underline{U} \right] \langle \underline{U} \rangle$ □

Theorem 4: $\gamma \langle \underline{V} \rangle$ is the \perp projection of $\langle \underline{U} \underline{M} |$ onto $\langle \underline{V} \rangle$.
 Proof: Similar to Theorem 3. □

Theorem 5: $\left\{ \langle \underline{U} \rangle = \langle \underline{U}_k \rangle, \langle \underline{V} \rangle = \langle \underline{V}_k \rangle, \gamma = \gamma_k \right\} \Rightarrow \left\{ \langle \underline{R} \rangle = \langle \underline{0} \rangle, \langle \underline{S} \rangle = \langle \underline{0} \rangle \right\}$
 Proof: Substitute premises into Equation (1). □

Theorem 6: $\left\{ \langle \underline{U} \rangle, \gamma, \langle \underline{V} \rangle \text{ is one of the sets of corresponding singular elements of } \underline{M} \right\} \Rightarrow \left\{ \begin{array}{l} \langle \underline{U} \rangle \text{ is an eigenrow of } \underline{M} \underline{M} \\ \text{belonging to eigenvalue } \gamma^2; \\ \langle \underline{V} \rangle \text{ is an eigenrow of } \underline{M} \underline{M} \\ \text{belonging to eigenvalue } \gamma^2. \end{array} \right.$

Theorem 7: $\left\{ \langle \underline{R} \rangle = \langle \underline{0} \rangle, \langle \underline{S} \rangle = \langle \underline{0} \rangle \right\} \Rightarrow \left\{ \langle \underline{U} \rangle = \langle \underline{U}_k \rangle, \langle \underline{V} \rangle = \langle \underline{V}_k \rangle, \gamma = \gamma_k \right\}$

4. If either residual differs appreciably from zero, then the algorithm is repeated with the trial rows perturbed in a direction intended to reduce the residuals. Let $\langle \underline{U} \rangle$ and $\langle \underline{V} \rangle$ be replaced with normalized versions of $\langle \underline{U} \rangle + \langle \underline{P} \rangle$ and $\langle \underline{V} \rangle + \langle \underline{Q} \rangle$ where $\langle \underline{P} \rangle$ and $\langle \underline{Q} \rangle$ are small perturbations. The residual norms squared then become

$$\begin{aligned} \langle \underline{R} \underline{R} \rangle &= \frac{\langle \underline{V} + \langle \underline{Q} \rangle \underline{M} \underline{M} [\underline{V} + \langle \underline{Q} \rangle] + \langle \underline{Q} \rangle}{[\langle \underline{V} + \langle \underline{Q} \rangle] [\underline{V} + \langle \underline{Q} \rangle]} - \frac{\left\{ [\langle \underline{U} + \langle \underline{P} \rangle] \underline{M} [\underline{V} + \langle \underline{Q} \rangle] \right\}^2}{[\langle \underline{U} + \langle \underline{P} \rangle] [\underline{U} + \langle \underline{P} \rangle] [\langle \underline{V} + \langle \underline{Q} \rangle] [\underline{V} + \langle \underline{Q} \rangle]} \quad (6) \\ \langle \underline{S} \underline{S} \rangle &= \frac{\langle \underline{U} + \langle \underline{P} \rangle \underline{M} \underline{M} [\underline{U} + \langle \underline{P} \rangle] + \langle \underline{P} \rangle}{[\langle \underline{U} + \langle \underline{P} \rangle] [\underline{U} + \langle \underline{P} \rangle]} - \frac{\left\{ [\langle \underline{U} + \langle \underline{P} \rangle] \underline{M} [\underline{V} + \langle \underline{Q} \rangle] \right\}^2}{[\langle \underline{U} + \langle \underline{P} \rangle] [\underline{U} + \langle \underline{P} \rangle] [\langle \underline{V} + \langle \underline{Q} \rangle] [\underline{V} + \langle \underline{Q} \rangle]} \end{aligned}$$

Setting both residuals to zero and dropping terms in $\langle \underline{P} \rangle$ and $\langle \underline{Q} \rangle$ of second and higher order yields

$$\begin{cases} [v^2 + 2 \langle \underline{Q} \underline{M} \underline{M} \underline{V} \rangle + 2v^2 \langle \underline{P} \underline{U} \rangle] - [\gamma^2 + 2\gamma \langle \underline{P} \underline{M} \underline{V} \rangle + 2\gamma \langle \underline{Q} \underline{M} \underline{U} \rangle] = 0 \\ [u^2 + 2 \langle \underline{P} \underline{M} \underline{M} \underline{U} \rangle + 2u^2 \langle \underline{Q} \underline{V} \rangle] - [\gamma^2 + 2\gamma \langle \underline{P} \underline{M} \underline{V} \rangle + 2\gamma \langle \underline{Q} \underline{M} \underline{U} \rangle] = 0 \end{cases} \quad (7)$$

which can be rewritten in partitioned matrix form.

$$\begin{bmatrix} \langle \underline{P} \rangle & \langle \underline{Q} \rangle \end{bmatrix} \begin{bmatrix} \underline{M} \underline{V} \rangle \gamma - \langle \underline{U} \rangle v^2 & \underline{M} \underline{V} \rangle \gamma - \underline{M} \underline{M} \underline{U} \rangle \\ \underline{M} \underline{U} \rangle \gamma - \underline{M} \underline{M} \underline{V} \rangle & \underline{M} \underline{U} \rangle \gamma - \langle \underline{V} \rangle u^2 \end{bmatrix} \approx \frac{1}{2} \begin{bmatrix} v^2 - \gamma^2 & 0 \\ 0 & u^2 - \gamma^2 \end{bmatrix} \quad (8)$$

$$\begin{array}{|c|c|} \hline \begin{array}{c} m \\ \gamma \langle \underline{v} \tilde{M} - \nu^2 \langle \underline{u} \end{array} & \begin{array}{c} n \\ \gamma \langle \underline{u} \tilde{M} - \langle \underline{v} \tilde{M} \tilde{M} \end{array} \\ \hline \gamma \langle \underline{v} \tilde{M} - \langle \underline{u} \tilde{M} \tilde{M} & \gamma \langle \underline{u} \tilde{M} - \mu^2 \langle \underline{v} \end{array} \\ \hline \end{array} \triangleq \underline{W} \quad (9)$$

$$\begin{array}{|c|c|} \hline \nu^2 - \gamma^2 & \mu^2 - \gamma^2 \\ \hline \end{array} \triangleq \underline{T} \quad (10)$$

With these definitions, our problem reduces to solving Equation (11) for $\langle \underline{z}$ and then the direct sum of $\langle \underline{p}$ and $\langle \underline{q}$ is given approximately by $\langle \underline{z}$.

$$\langle \underline{z} \tilde{W} = \frac{1}{2} \langle \underline{T} \quad (11)$$

$$\begin{array}{|c|c|} \hline \begin{array}{c} m \\ \langle \underline{p} \end{array} & \begin{array}{c} n \\ \langle \underline{q} \end{array} \\ \hline \end{array} \approx \langle \underline{z} \quad (12)$$

The general solution for $\langle \underline{z}$ in Equation (11) is $(1/2) \langle \underline{T} [\underline{W} \tilde{W}]^{-1} \underline{W} + \phi \langle \underline{A}$ where $\langle \underline{A}$ is any unit row orthogonal to $\langle \underline{T} [\underline{W} \tilde{W}]^{-1} \underline{W}$. Clearly, the solution of least norm is the one with $\phi = 0$.

The method used here to locate a minimum of τ is an example of the Newton-Raphson algorithm which is analogous to the use of Newton's Method to locate the vertex of an erect parabola starting from any point on the parabola. As shown in Figure 1, doubling the computed perturbation greatly accelerates convergence, assuming in our case that $\langle \underline{p}$ and $\langle \underline{q}$ are small. Applying this principle to our SVD algorithm, we redefine $\langle \underline{z}$ as follows.

$$\begin{array}{|c|c|} \hline \begin{array}{c} m \\ \langle \underline{p} \end{array} & \begin{array}{c} n \\ \langle \underline{q} \end{array} \\ \hline \end{array} \approx \langle \underline{z} = \langle \underline{T} [\underline{W} \tilde{W}]^{-1} \underline{W} \quad (13)$$

After approximate perturbation rows $\langle \underline{p}$ and $\langle \underline{q}$ are found from Equation (13) and added respectively to $\langle \underline{u}$ and $\langle \underline{v}$, the resulting sums are normalized and used as new trial rows for the next iteration of the inner loop of the algorithm.

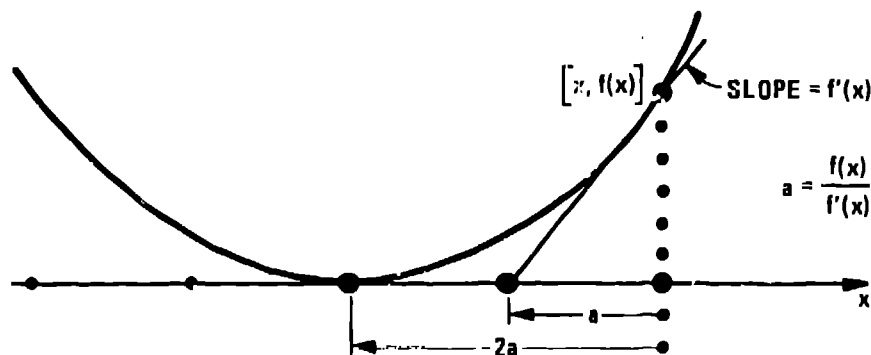


Figure 1. Acceleration of convergence.

5. The inner loop of the algorithm, defined in Sections 3 and 4, is repeated until both residual norms are sufficiently small as measured by their sum τ . The current pair of trial rows and corresponding trial singular value is then selected as a close approximation to one of the sets of corresponding singular elements of \underline{M} . This step ends the pass that started with the selection of initial trial singular rows and returns control to the outer loop of the algorithm for the selection of a new pair of trial singular rows to use in restarting the inner loop. To accomplish this objective, the dyad $\langle \underline{u} \rangle \gamma \langle \underline{v} \rangle$ is found from the results of the previous pass, subtracted from \underline{M} to form a deflated matrix \underline{N} which is then used to determine a new pair of starting rows $\langle \underline{u} \rangle$ and $\langle \underline{v} \rangle$. \underline{N} is progressively deflated at the end of each pass. The original matrix \underline{M} is used for all calculations in the inner loop of the algorithm in order to avoid unnecessary accumulation of round-off error.

A reasonable stopping rule for the outer loop of the algorithm can be based on the computation of the Euclidean norm of the deflated matrix \underline{N} at the end of each pass and stopping the algorithm when the norm of \underline{N} is sufficiently small. Additional measures of error in the SVD representation of \underline{M} may be obtained from the Euclidean norm of any of the following matrices: $\underline{M} - \underline{U} \underline{D} \underline{V}$, $\underline{U} \underline{U} - \underline{I}$, $\underline{V} \underline{V} - \underline{I}$. These measures are relative if applied to the scaled version of \underline{M} and absolute if applied to \underline{M} .

6. An outline of the inner loop of the SVD algorithm is given in Table 1 along with a measure of complexity of each step. For the sake of simplicity, only multiplications are counted. The total number of multiplications in each iteration in which all four tests are passed is $4mn + 12(m + n) + 10$. Tests 2-4 will be discussed later in this paper. Entry to the inner loop from the outer loop is at Step 0 and return to the outer loop follows Step 5.

Step	Computations	Multiplications
0	Given: $m \times n$ matrix \underline{M} , starting m -row $\langle \underline{U}$, starting n -row $\langle \underline{V}$	
1	$\langle \underline{U} \underline{M}$ and $\langle \underline{V} \underline{M}$	$2mn$
2	μ^2 and ν^2	$m + n$
3	γ , including check	$m + n$
4	$\langle \underline{R} \underline{\tilde{R}} \rangle$, $\langle \underline{S} \underline{\tilde{S}} \rangle$, τ	1
5	Test 1: Compare τ with stopping limit of inner loop.	
6	Test 2: Compare with previous iteration to ensure substantial reduction in τ	
7	$\gamma \langle \underline{V} \underline{\tilde{M}}$, $\nu^2 \langle \underline{U}$, $\langle \underline{U} \underline{M} \underline{\tilde{M}}$	$2m + mn$
8	$\gamma \langle \underline{U} \underline{M}$, $\mu^2 \langle \underline{V}$, $\langle \underline{V} \underline{\tilde{M}} \underline{M}$	$2n + mn$
9	Assemble $\langle \underline{T}$ and \underline{W} .	
10	$\underline{W} \underline{\tilde{W}}$ and Det $[\underline{W} \underline{\tilde{W}}]$	$3(m + n) +$
11	Test 3: Compare determinant with lower limit.	
12	$[\underline{W} \underline{\tilde{W}}]^{-1}$	3
13	$\langle \underline{T} [\underline{W} \underline{\tilde{W}}]^{-1}$	4
14	$\begin{bmatrix} \underline{P} & \underline{Q} \end{bmatrix} = \langle \underline{T} [\underline{W} \underline{\tilde{W}}]^{-1} \underline{W}$	$2(m + n)$
15	Test 4: Compare norm of perturbation row with upper limit to prevent excessive overshoot.	$m + n$
16	$\langle \underline{U} + \underline{P}$ and $\langle \underline{V} + \underline{Q}$	
17	Normalize to find new $\langle \underline{U}$ and $\langle \underline{V}$.	$2(m + n)$
18	Return to Step 1.	

7. Several numerical examples of the application of the SVD algorithm to some relatively small matrices (3×5) have been developed and studied. In several of these examples, the starting rules normally employed with the algorithm were ignored and intentionally poor choices were made for the starting rows. In all of these examples, one iteration of the inner loop of the algorithm resulted in substantial reduction of the residual norms and of the sum of the angles between the trial vectors and the singular pair to which the trial vectors were converging. In view of the tedium of the required calculations even for small matrices, only one numerical example is presented here. In this example, the initial trial rows are determined by reasonable starting rules.

Given that $\underline{M} =$

.640	-.640	1.000	.104	.640
.480	-.480	.816	.288	.480
-.300	.300	.240	.820	-.300

We obtain starting rows by normalizing the largest row and largest column of \underline{M} .

$$\begin{aligned} \langle \underline{u} = & \begin{bmatrix} 1.088 & .816 & .240 \end{bmatrix} \frac{1}{\sqrt{1.9072}} \\ \langle \underline{v} = & \begin{bmatrix} .640 & -.640 & 1.088 & .384 & .640 \end{bmatrix} \frac{1}{\sqrt{2.56}} \\ \langle \underline{u} \underline{M} = & \begin{bmatrix} .735691 & -.735691 & 1.781014 & .615200 & .735691 \end{bmatrix} \\ \langle \underline{v} \underline{M} = & \begin{bmatrix} 1.6 & 1.2 & 0.0 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \gamma &= \langle \underline{v} \underline{M} \underline{u} \rangle = 1.969567 & \text{Check: } \langle \underline{u} \underline{M} \underline{v} \rangle &= 1.969567 \\ \gamma^2 &= 3.879194 & \mu^2 &= 3.909394 & \nu^2 &= 4.000000 \\ \langle \underline{r} \underline{r} \rangle &= \nu^2 - \gamma^2 = 0.120806 & \langle \underline{s} \underline{s} \rangle &= \mu^2 - \gamma^2 = 0.030200 \\ \tau &= \langle \underline{r} \underline{r} \rangle \div \langle \underline{s} \underline{s} \rangle = 0.151006 & \langle \underline{t} = & \begin{bmatrix} .120806 & .030200 \end{bmatrix} \end{aligned}$$

$$\underline{W} = \begin{bmatrix} .00000 & .00000 & -.69514 & -.15101 & .15101 & .00000 & .25168 & -.15101 \\ .00000 & .00000 & -.17378 & -.11476 & .11476 & .06161 & .27342 & -.11476 \end{bmatrix}$$

$$\underline{W} \underline{W}^{-1} = \begin{bmatrix} .614974 & .241605 \\ .241605 & .148263 \end{bmatrix} \quad [\underline{W} \underline{W}]^{-1} = \frac{1}{.032805} \begin{bmatrix} .148263 & -.241605 \\ -.241605 & .614974 \end{bmatrix}$$

$$\begin{aligned} \langle \underline{t} [\underline{W} \underline{W}]^{-1} \underline{w} = & \begin{bmatrix} \langle \underline{p} & \langle \underline{q} \\ -.00000 & .00000 & -.16869 & -.01173 & .01173 & -.01994 & -.00706 & -.01173 \end{bmatrix} \end{aligned}$$

After adding these perturbation rows to the initial trial rows and normalizing, we have a new pair of trial rows with which to start the next iteration.

$$\begin{aligned} \langle \underline{u} = & \begin{bmatrix} .79998 & .59998 & .00517 \end{bmatrix} \\ \langle \underline{v} = & \begin{bmatrix} .40000 & -.40000 & .68000 & .23998 & .40000 \end{bmatrix} \end{aligned}$$

In this example, the SVD of \underline{M} is known exactly and we can deduce that the algorithm is converging to the following set of corresponding singular elements.

$$\begin{aligned} \langle \underline{u}_1 = & \begin{bmatrix} .00 & .60 & .00 \end{bmatrix} & \gamma_1 &= 2.00 \\ \langle \underline{v}_1 = & \begin{bmatrix} .40 & -.40 & .68 & .24 & .40 \end{bmatrix} \end{aligned}$$

The angle between the trial row $\langle \underline{u}$ and the singular row $\langle \underline{u}_1$ is reduced by one iteration of the inner loop of the algorithm from $\arccos(.984784) = 10.008^\circ$ to $\arccos(.999972) = 0.429^\circ$, that is, by a factor of more than 20. The initial trial row $\langle \underline{v}$ was chosen equal to $\langle \underline{v}_1$ and the effect of one iteration is to change the angle between $\langle \underline{v}$ and $\langle \underline{v}_1$ from 0° to $\arccos(.999995) = 0.177^\circ$. Thus, even when the algorithm is converging normally, one of the two angles may increase by a small amount and this effect is more than cancelled by a large reduction in the other angle. This behavior is related to the zig-zag approach to solution called "hemstitching" which is common with algorithms of the Newton-Raphson or steepest-descent type.

In normal use of the algorithm, we must rely on τ , the total residual norm squared, to measure distance between our trial solution and the final solution. Let us carry the calculation of our numerical example through the second iteration as far as finding τ .

$$\begin{aligned} \langle \underline{u} \underline{M} = & \begin{bmatrix} .79843 & -.79843 & 1.36120 & .48423 & .79843 \end{bmatrix} \\ \langle \underline{v} \underline{M} = & \begin{bmatrix} 1.59994 & 1.19999 & -.00002 \end{bmatrix} \end{aligned}$$

$$\gamma = \langle \underline{U} \underline{M} \tilde{\underline{V}} \rangle = 1.99994 \quad \text{Check: } \langle \underline{V} \tilde{\underline{M}} \underline{\underline{U}} \rangle = 1.99993$$

$$\mu^2 = 3.99982 \quad \nu^2 = 3.99994 \quad \tau = .00020 + .00008 = .00028$$

Note that one iteration of the inner loop has reduced both $\langle \underline{R} \tilde{\underline{R}} \rangle$ and $\langle \underline{S} \tilde{\underline{S}} \rangle$, and that it reduced $\sqrt{\tau}$ by a factor of more than 20.

8. There are several special cases to consider. First, we note that the null matrix $\underline{0}$ has no SVD. If the source of matrices to be factored by the algorithm can possibly produce $\underline{0}$, then the outer loop of the algorithm must first test the given matrix to ensure that it is not $\underline{0}$ before proceeding. Second, it is important that the γ found in Step 1 of the inner loop be positive, so we need to consider situations where γ could be zero. If $\langle \underline{U}$ is orthogonal to all of the singular rows $\langle \underline{U}_k$, then $\langle \underline{U} \underline{M} \rangle = \langle \underline{0}$ and $\gamma = 0$. Similarly, if $\langle \underline{V}$ is orthogonal to all of the singular rows $\langle \underline{V}_k$, then $\langle \underline{V} \underline{M} \rangle = \langle \underline{0}$ and $\gamma = 0$. Another way that γ can be zero is: $\langle \underline{U} = \langle \underline{U}_i$ and $\langle \underline{V} = \langle \underline{V}_k$ where neither i nor k exceed r and $i \neq k$. Thus, the check in Step 3 must ensure that the two ways of finding γ agree within allowable round-off error and that γ exceed zero by more than allowable round-off error. If $\gamma \approx 0$, then either $\langle \underline{U}$ or $\langle \underline{V}$ must be replaced before proceeding. If $\mu^2 < \nu^2$, replace $\langle \underline{U}$; if $\mu^2 \geq \nu^2$, replace $\langle \underline{V}$. The replacement ought to be orthogonal to all previous trial rows.

Let us consider cases in which $\text{Det}[\underline{W} \tilde{\underline{W}}] = 0$. The cases described above that lead to $\gamma = 0$ also lead to zero for the determinant. The cure is the same as that described in the preceding paragraph. Another way that the determinant can be zero is: $\langle \underline{U} = \langle \underline{U}_k$, $\langle \underline{V} = \langle \underline{V}_k$, $k \leq r$. This case corresponds to normal convergence behavior and causes no problem, because any iteration that would lead to $\text{Det}[\underline{W} \tilde{\underline{W}}] = 0$ with $\gamma \neq 0$ would be terminated at Step 5 and the determinant calculation in Step 10 would not be reached. Further study may show that the check in Step 3 may eliminate the need for Test 3.

9. In order to study the behavior of the algorithm and sufficient conditions for convergence, let the trial rows be described in terms of deviations from a pair of corresponding singular rows $\langle \underline{U}_k$ and $\langle \underline{V}_k$ of the given $m \times n$ matrix \underline{M} of rank r . Let \underline{N} be the deflated matrix defined by: $\underline{N} = \underline{M} - \underline{U}_k \underline{V}_k$.

$$\begin{cases} \langle \underline{U} = \sqrt{1-x^2} \langle \underline{U}_k + x \langle \underline{F} & 0 \leq x^2 < 1 \\ \langle \underline{V} = \sqrt{1-y^2} \langle \underline{V}_k + y \langle \underline{G} & 0 \leq y^2 < 1 \end{cases}$$

The rows $\langle \underline{F}$ and $\langle \underline{G}$ may be any unit rows such that $\langle \underline{F} \perp \langle \underline{U}_k$ and $\langle \underline{G} \perp \langle \underline{V}_k$ and such that $\langle \underline{F} \underline{N} \underline{G} \rangle \geq 0$ for (x,y) in the first or third quadrants or $\langle \underline{F} \underline{N} \underline{G} \rangle \leq 0$ for (x,y) in the second or fourth quadrants.

$$\begin{aligned} \langle \underline{U} \underline{M} \rangle &= \gamma_k \sqrt{1-x^2} \langle \underline{V}_k + x \langle \underline{F} \underline{N} \\ \langle \underline{V} \underline{M} \rangle &= \gamma_k \sqrt{1-y^2} \langle \underline{U}_k + y \langle \underline{G} \underline{N} \\ \mu^2 &= \langle \underline{U} \underline{M} \tilde{\underline{M}} \underline{\underline{U}} \rangle = \gamma_k^2 (1-x^2) + x^2 \langle \underline{F} \underline{N} \tilde{\underline{N}} \tilde{\underline{F}} \rangle \\ \nu^2 &= \langle \underline{V} \underline{M} \tilde{\underline{M}} \underline{\underline{V}} \rangle = \gamma_k^2 (1-y^2) + y^2 \langle \underline{G} \underline{N} \tilde{\underline{N}} \tilde{\underline{G}} \rangle \\ \gamma &= \langle \underline{U} \underline{M} \tilde{\underline{V}} \rangle = \langle \underline{V} \tilde{\underline{M}} \underline{\underline{U}} \rangle = \gamma_k \sqrt{(1-x^2)(1-y^2)} + xy \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle \end{aligned}$$

At this point, we can see that the condition on the sign of $\langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle$ is necessary to ensure that γ is positive for all (x,y) .

$$\begin{aligned} \langle \underline{R} \tilde{\underline{R}} \rangle &= \nu^2 - \gamma^2 = \gamma_k^2 (1-y^2)x^2 + y^2 \langle \underline{G} \underline{N} \tilde{\underline{N}} \tilde{\underline{G}} \rangle - x^2 y^2 \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle^2 \\ &\quad - 2xy \sqrt{(1-x^2)(1-y^2)} \gamma_k \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle \\ \langle \underline{S} \tilde{\underline{S}} \rangle &= \mu^2 - \gamma^2 = \gamma_k^2 (1-x^2)y^2 + x^2 \langle \underline{F} \underline{N} \tilde{\underline{N}} \tilde{\underline{F}} \rangle - x^2 y^2 \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle^2 \\ &\quad - 2xy \sqrt{(1-x^2)(1-y^2)} \gamma_k \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle \\ \tau &= \langle \underline{R} \tilde{\underline{R}} \rangle + \langle \underline{S} \tilde{\underline{S}} \rangle = \gamma_k^2 (x^2 + y^2 - 2x^2 y^2) + x^2 \langle \underline{F} \underline{N} \tilde{\underline{N}} \tilde{\underline{F}} \rangle + y^2 \langle \underline{G} \underline{N} \tilde{\underline{N}} \tilde{\underline{G}} \rangle \\ &\quad - 2x^2 y^2 \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle^2 - 4xy \sqrt{(1-x^2)(1-y^2)} \gamma_k \langle \underline{F} \underline{N} \tilde{\underline{G}} \rangle \quad (14) \end{aligned}$$

Note that $\tau(x,y) \geq 0$ and $\tau(0,0) = 0$.

In order to develop a geometric picture of the $\tau(x,y)$ surface, let us consider a few of its cross-sections: $x = 0$; $y = 0$; $x = y$.

$$x = 0 \Rightarrow \tau = \left[\gamma_R^2 + \langle \underline{G} \underline{N} \underline{N} \underline{G} \rangle \right] y^2 \tag{15}$$

$$y = 0 \Rightarrow \tau = \left[\gamma_R^2 + \langle \underline{F} \underline{N} \underline{N} \underline{F} \rangle \right] x^2 \tag{16}$$

$$\begin{aligned} x = y \Rightarrow \tau &= 2\gamma_R^2(x^2 - x^4) + x^2 \frac{\langle \underline{F} \underline{N} \underline{N} \underline{F} \rangle}{4x^2(1 - x^2)\gamma_k} - 2x^4 \frac{\langle \underline{F} \underline{N} \underline{G} \rangle^2}{\langle \underline{F} \underline{N} \underline{G} \rangle} + x^2 \langle \underline{G} \underline{N} \underline{N} \underline{G} \rangle \\ &= x^2(2C^2 + D^2) - x^4(2C^2) \quad \text{where } C^2 = (\gamma_k - \langle \underline{F} \underline{N} \underline{G} \rangle)^2 \\ &\quad D^2 = \langle \underline{F} \underline{N} (\underline{I} - \underline{G}) \rangle \langle \underline{G} \rangle^2 \langle \underline{N} \underline{F} \rangle + \langle \underline{G} \underline{N} (\underline{I} - \underline{F}) \rangle \langle \underline{F} \rangle^2 \langle \underline{N} \underline{G} \rangle \end{aligned} \tag{17}$$

Note that \underline{F} , \underline{G} , $\underline{I} - \underline{F}$, and $\underline{I} - \underline{G}$ represent perpendicular projectors.

We see that Equations (15,16) represent parabolas and Equation (17) represents a quartic that approaches a parabola for small x^2 as shown in Figure 2. Our algorithm was designed to work well with paraboloids, so we expect convergence for almost any y^2 when x^2 is small and for almost any x^2 when y^2 is small. To obtain some initial ideas of conditions for convergence when neither x^2 nor y^2 are small, let us consider Equation (17) and its first derivative.

$$d\tau/dx = 2x[(2C^2 + D^2) - 2x^2(2C^2)] \tag{18}$$

The stationary points of $\tau(x,x)$ are at $x = 0$ and at $x = \pm x_M$ where

$$x_M^2 = \frac{2C^2 + D^2}{4C^2} \geq \frac{1}{2} = (\sin 45^\circ)^2 \tag{19}$$

and $\tau(x_M, x_M) = (2C^2 + D^2)^2/8C^2$ while $\tau(1,1) = D^2$. The stationary point at $x = 0$ is a global minimum where $\tau(0,0) = 0$. The other stationary points are local maxima of $\tau(x,x)$ which do not occur if $2C^2 \leq D^2$. Note that the smallest x_M^2 is 1/2 corresponding to the case where both trial vectors are 45° from one of the vectors in a pair of singular vectors.

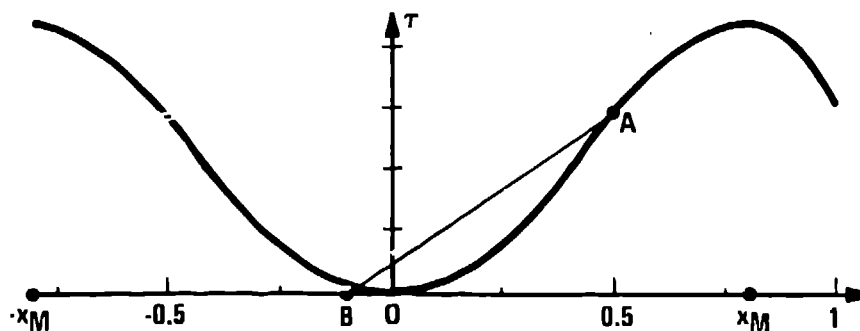


Figure 2. $\tau(x,x)$ for $D^2 < 2C^2$.

A somewhat oversimplified view of the behavior of the algorithm can be seen by studying Figure 2 where point A represents the initial choice of trial rows. One iteration of the inner loop of the algorithm would yield a new trial represented by the point B which would lie at or near the origin if $\tau(x,y)$ were paraboloidal or nearly paraboloidal. Because of the quartic nature of the $\tau(x,y)$ surface, the point B will lie beyond the point nearest to the origin along the direction of the perturbation vector. That is, the algorithm will overshoot. The amount of overshoot will be small for small x_A^2 but will become very large for x_A^2 near x_M^2 . However, as long as $x_A^2 < x_M^2$, the algorithm will make changes in the correct direction. This observation suggests that the length of each perturbation vector obtained, at least in the first iteration in each pass, be checked and reduced, if necessary, to some preset maximum allowable length. Thus, Test 4 is included at Step 15 of the inner loop just before the vector addition and renormalization steps. Calculations have been made on the example of Figure 2 and on several other examples for two values of the maximum length of the perturbation vector: $1/\sqrt{2}$ and $1/2$. Best results have been obtained with this limit set at $1/\sqrt{2}$.

If x_A^2 exceeds x_M^2 , the algorithm may make changes in the wrong direction. Apparently, many, if not all, of such cases may be detected by noting that the new pair of trial rows

does not yield a substantial decrease in \mathcal{T} . It appears that such cases may be treated by reversing the direction of the calculated perturbation via Test 2 at Step 6 of the inner loop. Further study is needed to set the limiting reduction ratio used in Test 2 and to determine whether Test 2 may be eliminated by selection of good starting rules.

The choice of the limit on the norm of the perturbation vector also requires further study and computer testing. In order to speed convergence in most cases, it appears to be desirable to set this limit at about $1/\sqrt{2}$, that is, $\langle \underline{p} \ \underline{p} \rangle + \langle \underline{q} \ \underline{q} \rangle \leq 1/2$.

One result of the analysis of the algorithm to date is Theorem 8 which leads to the following conjecture: The algorithm converges for all starting points inside the square S of side $\sqrt{2}$ centered on the origin of the (x,y) plane, that is, for

$$x^2 < 1/2 > y^2$$

An alternative statement of the conjecture is: The algorithm converges if each trial vector is less than 45° from the corresponding vector in a singular pair.

Theorem 8: The only stationary point of $\mathcal{T}(x,y)$ inside the square S is the global minimum at the origin.

Proof: Equate to zero the partial derivatives of \mathcal{T} with respect to x and y. Transfer radicals to the opposite side of the equals signs. Multiply the resulting equations together to form one equation: $f(x,y) = 0$. Note that $f(x,y) = 0$ at the origin and at no other points inside S. □

Study of the effect of the ratio $D^2/2C^2$ on the $\mathcal{T}(x,y)$ surface leads to the further conjecture that the algorithm converges for all (x,y) if the following conditions are satisfied.

$$\langle \underline{F} \ \underline{N} \ \underline{\tilde{N}} \ \underline{\tilde{F}} \rangle > [\delta_k^2 + \langle \underline{F} \ \underline{N} \ \underline{\tilde{G}} \rangle^2] < \langle \underline{G} \ \underline{\tilde{N}} \ \underline{N} \ \underline{\tilde{G}} \rangle \quad (20)$$

These conditions are unlikely to be satisfied on the first pass of the algorithm but are usually satisfied on subsequent passes associated with calculation of the smaller singular values.

The algorithm appears to be particularly applicable to the updating of the SVD of a matrix of observations as each new set of observations is obtained and to the checking and refinement of results obtained from other SVD algorithms. Further study and computer testing is needed to put the algorithm in final form and to determine its rate of convergence under various conditions. In the limited number of numerical examples developed to date, the convergence of the algorithm is surprisingly fast.

References

1. Cornelius Lanczos
"Linear Systems in Self-Adjoint Form"
American Mathematical Monthly
v 65 (1958), pp 665-679
2. William H. Muggins
"An Algebra for Signal Representation"
Yearbook of the Society for General Systems Research
v VIII (1963), pp 129-143