

The equation  $Av = \sigma u$  is close but different. Now we have **two vectors  $v$  and  $u$** . Our matrix  $A$  is probably rectangular, and full of data. What part of that data matrix is important? The *Singular Value Decomposition* (SVD) finds its simplest pieces  $\sigma uv^T$ . Those pieces are matrices (column  $u$  times row  $v^T$ ). Every matrix is built from these orthogonal pieces. **Data science meets linear algebra in the SVD.**

Finding those pieces  $\sigma uv^T$  is the object of Principal Component Analysis (PCA).

The big factorization for data science is the “SVD” of  $A$ —when the first factor  $C$  has  $r$  *orthogonal* columns and the second factor  $R$  has  $r$  *orthogonal* rows.

Here are five important factorizations, with the standard choice of letters (usually  $A$ ) for the original product matrix and then for its factors. This book will explain all five.

$$A = LU \quad A = QR \quad S = Q\Lambda Q^T \quad A = X\Lambda X^{-1} \quad A = U\Sigma V^T$$

**5**  $A = U\Sigma V^T$  is the **Singular Value Decomposition** of any matrix  $A$  (square or not).

**Singular values**  $\sigma_1, \dots, \sigma_r$  in  $\Sigma$ . Orthonormal **singular vectors** in  $U$  and  $V$ .

**3. Orthogonal subspaces.** Equation (1) looked at  $Ax = 0$ . Every row of  $A$  is multiplying that nullspace vector  $x$ . So each row (and all combinations of the rows) will be orthogonal to  $x$  in  $N(A)$ . **The row space of  $A$  is orthogonal to the nullspace of  $A$ .**

$$Ax = \begin{bmatrix} \text{row 1} \\ \vdots \\ \text{row } m \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad A^T y = \begin{bmatrix} (\text{column 1})^T \\ \vdots \\ (\text{column } n)^T \end{bmatrix} \begin{bmatrix} y \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad (7)$$

From  $A^T y = 0$ , the columns of  $A$  are all orthogonal to  $y$ . Their combinations (the whole column space) will also be orthogonal to  $y$ . **The column space of  $A$  is orthogonal to the nullspace of  $A^T$ .** This produces the “Big Picture of Linear Algebra” in Figure I.6.

Notice the dimensions  $r$  and  $n - r$  adding to  $n$ . The whole space  $\mathbf{R}^n$  is accounted for. Every vector  $v$  in  $\mathbf{R}^n$  has a row space component  $v_r$  and a nullspace component  $v_n$  with  $v = v_r + v_n$ . A row space basis ( $r$  vectors) together with a nullspace basis ( $n - r$  vectors) produces a basis for all of  $\mathbf{R}^n$  ( $n$  vectors).

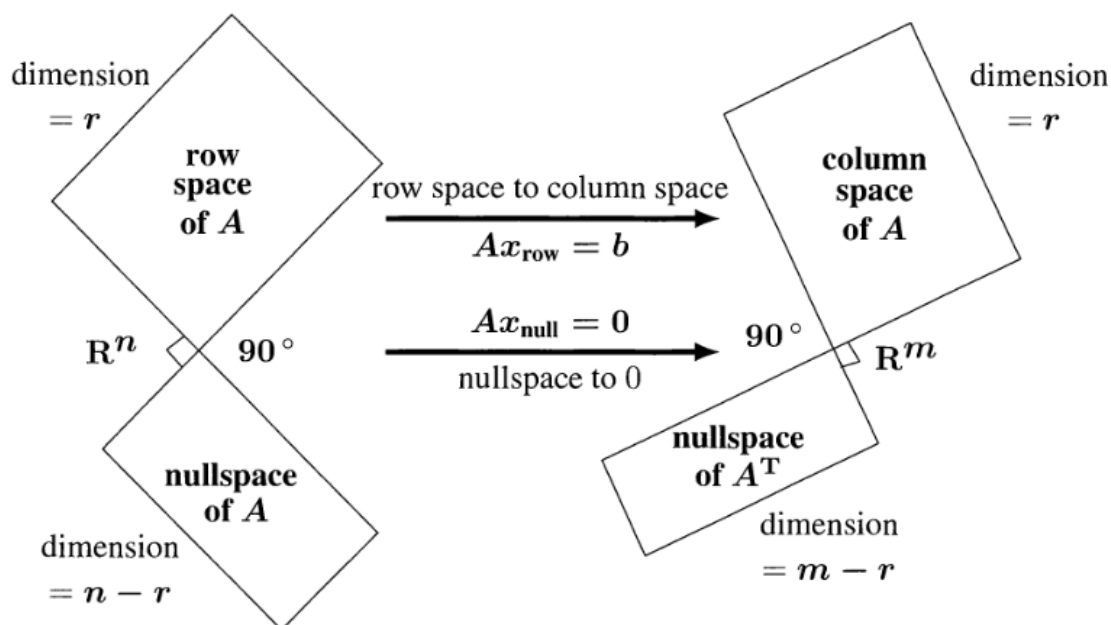


Figure I.6: Two pairs of orthogonal subspaces. The dimensions add to  $n$  and add to  $m$ . **This is the Big Picture**—two subspaces in  $\mathbf{R}^n$  and two subspaces in  $\mathbf{R}^m$ .

I will mention a big improvement. It comes from the Singular Value Decomposition. The SVD is the most important theorem in data science. It finds orthonormal bases  $v_1, \dots, v_r$  for the row space of  $A$  and  $u_1, \dots, u_r$  for the column space of  $A$ . Well, Gram-Schmidt can do that. The special bases from the SVD have the extra property that each pair ( $v$  and  $u$ ) is connected by  $A$ :

$$\text{Singular vectors} \quad Av_1 = \sigma_1 u_1 \quad Av_2 = \sigma_2 u_2 \quad \cdots \quad Av_r = \sigma_r u_r. \quad (8)$$

In Figure I.6, imagine the  $v$ 's on the left and the  $u$ 's on the right. For the bases from the SVD, multiplying by  $A$  takes an orthogonal basis of  $v$ 's to an orthogonal basis of  $u$ 's.

## I.8 Singular Values and Singular Vectors in the SVD

The best matrices (real symmetric matrices  $S$ ) have real eigenvalues and orthogonal eigenvectors. But for other matrices, the eigenvalues are complex or the eigenvectors are not orthogonal. If  $A$  is not square then  $Ax = \lambda x$  is impossible and eigenvectors fail (left side in  $\mathbf{R}^m$ , right side in  $\mathbf{R}^n$ ). We need an idea that succeeds for every matrix.

The Singular Value Decomposition fills this gap in a perfect way. In our applications,  $A$  is often a matrix of data. The rows could tell us the age and height of 1000 children. Then  $A$  is 2 by 1000: definitely rectangular. Unless height is exactly proportional to age, the rank is  $r = 2$  and that matrix  $A$  has two positive singular values  $\sigma_1$  and  $\sigma_2$ .

The key point is that we need **two sets of singular vectors**, the  $u$ 's and the  $v$ 's. For a real  $m$  by  $n$  matrix, the  $n$  right singular vectors  $v_1, \dots, v_n$  are orthogonal in  $\mathbf{R}^n$ . The  $m$  left singular vectors  $u_1, \dots, u_m$  are perpendicular to each other in  $\mathbf{R}^m$ . The connection between  $n$   $v$ 's and  $m$   $u$ 's is not  $Ax = \lambda x$ . That is for eigenvectors. **For singular vectors, each  $Av$  equals  $\sigma u$ :**

$$\boxed{Av_1 = \sigma_1 u_1 \quad \dots \quad Av_r = \sigma_r u_r} \quad \boxed{Av_{r+1} = 0 \quad \dots \quad Av_n = 0} \quad (1)$$

I have separated the first  $r$   $v$ 's and  $u$ 's from the rest. That number  $r$  is the *rank of  $A$* , the number of independent columns (and rows). Then  $r$  is the dimension of the column space and the row space. **We will have  $r$  positive singular values in descending order  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ .** The last  $n - r$   $v$ 's are in the nullspace of  $A$ , and the last  $m - r$   $u$ 's are in the nullspace of  $A^T$ .

Our first step is to write equation (1) in matrix form. All of the right singular vectors  $v_1$  to  $v_n$  go in the columns of  $V$ . The left singular vectors  $u_1$  to  $u_m$  go in the columns of  $U$ . Those are **square orthogonal matrices** ( $V^T = V^{-1}$  and  $U^T = U^{-1}$ ) because their columns are orthogonal unit vectors. Then equation (1) becomes the full SVD, with square matrices  $V$  and  $U$ :

$$\boxed{AV = U\Sigma \quad A \begin{bmatrix} v_1 & \dots & v_r & \dots & v_n \end{bmatrix} = \begin{bmatrix} u_1 & \dots & u_r & \dots & u_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & 0 \\ & \ddots & & & \\ & & \sigma_r & & \\ \hline & & & 0 & 0 \end{bmatrix}} \quad (2)$$

You see  $Av_k = \sigma_k u_k$  in the first  $r$  columns above. That is the important part of the SVD. It shows the basis of  $v$ 's for the row space of  $A$  and then  $u$ 's for the column space. After the positive numbers  $\sigma_1, \dots, \sigma_r$  on the main diagonal of  $\Sigma$ , the rest of that matrix is all zero from the nullspaces of  $A$  and  $A^T$ .

The eigenvectors give  $AX = X\Lambda$ . But  $AV = U\Sigma$  needs **two sets of singular vectors**.

## 1.8. Singular Values and Singular Vectors in the SVD

57

$$\text{Example 1} \quad \boxed{AV = U\Sigma} \quad \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 3\sqrt{5} & \\ & \sqrt{5} \end{bmatrix}$$

The matrix  $A$  is not symmetric, so  $V$  is different from  $U$ . The rank is 2, so there are two singular values  $\sigma_1 = 3\sqrt{5}$  and  $\sigma_2 = \sqrt{5}$ . Their product  $3 \cdot 5 = 15$  is the determinant of  $A$  (in this respect singular values are like eigenvalues). The columns of  $V$  are orthogonal and the columns of  $U$  are orthogonal. Those columns are unit vectors after the divisions by  $\sqrt{2}$  and  $\sqrt{10}$ , so  **$V$  and  $U$  are orthogonal matrices**:  $V^T = V^{-1}$  and  $U^T = U^{-1}$ .

That orthogonality allows us to go from  $AV = U\Sigma$  to the usual and famous expression of the SVD: Multiply both sides of  $AV = U\Sigma$  by  $V^{-1} = V^T$ .

$$\boxed{\text{The Singular Value Decomposition of } A \text{ is } A = U\Sigma V^T.} \quad (3)$$

Then column-row multiplication of  $U\Sigma$  times  $V^T$  separates  $A$  into  $r$  pieces of rank 1:

**Pieces of the SVD**

$$A = U\Sigma V^T = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T. \quad (4)$$

In the 2 by 2 example, the first piece is more important than the second piece because  $\sigma_1 = 3\sqrt{5}$  is greater than  $\sigma_2 = \sqrt{5}$ . To recover  $A$ , add the pieces  $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T$ :

$$\frac{3\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} + \frac{\sqrt{5}}{\sqrt{10}\sqrt{2}} \begin{bmatrix} -3 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \end{bmatrix} = \frac{3}{2} \begin{bmatrix} 1 & 1 \\ 3 & 3 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 3 & -3 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix}$$

This simplified because  $\sqrt{5}/\sqrt{10}\sqrt{2}$  equals  $1/2$ . Notice that the right singular vectors  $(1, 1)$  and  $(-1, 1)$  in  $V$  are transposed to rows  $\mathbf{v}_1^T, \mathbf{v}_2^T$  of  $V^T$ . We have not yet explained how  $V$  and  $U$  and  $\Sigma$  were computed!

## The Reduced Form of the SVD

The full form  $AV = U\Sigma$  in equation (2) can have a lot of zeros in  $\Sigma$  when the rank of  $A$  is small and its nullspace is large. Those zeros contribute nothing to matrix multiplication. The heart of the SVD is in the first  $r$   $\mathbf{v}$ 's and  $\mathbf{u}$ 's and  $\sigma$ 's. We can reduce  $AV = U\Sigma$  to  $AV_r = U_r \Sigma_r$  by removing the parts that are sure to produce zeros. This leaves the **reduced SVD where  $\Sigma_r$  is now square**:

$$AV_r = U_r \Sigma_r \quad A \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_r \\ \text{row space} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_r \\ \text{column space} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \quad (5)$$

We still have  $V_r^T V_r = I_r$  and  $U_r^T U_r = I_r$  from those orthogonal unit vectors  $\mathbf{v}$ 's and  $\mathbf{u}$ 's. But when  $V_r$  and  $U_r$  are not square, we can no longer have two-sided inverses:  $V_r V_r^T \neq I$  and  $U_r U_r^T \neq I$ .

**Example**  $V_r = \begin{bmatrix} 1/3 \\ 2/3 \\ 2/3 \end{bmatrix}$  and  $V_r^T V_r = [1]$  but  $V_r V_r^T = \frac{1}{9} \begin{bmatrix} 1 & 2 & 2 \\ 2 & 4 & 4 \\ 2 & 4 & 4 \end{bmatrix} = \text{rank } 1$ .

Problem 21 shows that **we still have**  $A = U_r \Sigma_r V_r^T$ . The rest of  $U\Sigma V^T$  contributes nothing to  $A$ , because of those blocks of zeros in  $\Sigma$ . The key formula is still  $A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$ . The SVD sees only the  $r$  nonzeros in the diagonal matrix  $\Sigma$ .

## The Important Fact for Data Science

Why is the SVD so important for this subject and this book ? Like the other factorizations  $A = LU$  and  $A = QR$  and  $S = Q\Lambda Q^T$ , it separates the matrix into rank one pieces. A special property of the SVD is that **those pieces come in order of importance**. The first piece  $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$  is the closest rank one matrix to  $A$ . More than that is true: *The sum of the first  $k$  pieces is best possible for rank  $k$ .*

$A_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \dots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$  is the best rank  $k$  approximation to  $A$ :

**Eckart-Young**

$$\text{If } B \text{ has rank } k \text{ then } \|A - A_k\| \leq \|A - B\|. \quad (6)$$

To interpret that statement you need to know the meaning of the symbol  $\|A - B\|$ . This is the “**norm**” of the matrix  $A - B$ , a measure of its size (like the absolute value of a number). The Eckart-Young theorem is proved in Section I.9.

Our first job is to find the  $\mathbf{v}$ 's and  $\mathbf{u}$ 's for equation (1), to reach the SVD.

### First Proof of the SVD

Our goal is  $A = U\Sigma V^T$ . We want to identify the two sets of singular vectors, the  $\mathbf{u}$ 's and the  $\mathbf{v}$ 's. One way to find those vectors is to form the symmetric matrices  $A^T A$  and  $AA^T$ :

$$A^T A = (V\Sigma^T U^T) (U\Sigma V^T) = V\Sigma^T \Sigma V^T \quad (7)$$

$$AA^T = (U\Sigma V^T) (V\Sigma^T U^T) = U\Sigma \Sigma^T U^T \quad (8)$$

Both (7) and (8) produced symmetric matrices. Usually  $A^T A$  and  $AA^T$  are different. Both right hand sides have the special form  $Q\Lambda Q^T$ . Eigenvalues are in  $\Lambda = \Sigma^T \Sigma$  or  $\Sigma \Sigma^T$ . **Eigenvectors are in  $Q = V$  or  $Q = U$ .** So we know from (7) and (8) how  $V$  and  $U$  and  $\Sigma$  connect to the symmetric matrices  $A^T A$  and  $AA^T$ .

**$V$  contains orthonormal eigenvectors of  $A^T A$**

**$U$  contains orthonormal eigenvectors of  $AA^T$**

**$\sigma_1^2$  to  $\sigma_r^2$  are the nonzero eigenvalues of both  $A^T A$  and  $AA^T$**

We are not quite finished, for this reason. **The SVD requires that  $A\mathbf{v}_k = \sigma_k \mathbf{u}_k$ .** It connects each right singular vector  $\mathbf{v}_k$  to a left singular vector  $\mathbf{u}_k$ , for  $k = 1, \dots, r$ . When I choose the  $\mathbf{v}$ 's, that choice will decide the signs of the  $\mathbf{u}$ 's. If  $S\mathbf{u} = \lambda\mathbf{u}$  then also  $S(-\mathbf{u}) = \lambda(-\mathbf{u})$  and I have to know the correct sign. More than that, there is a whole plane of eigenvectors when  $\lambda$  is a double eigenvalue. When I choose two  $\mathbf{v}$ 's in that plane, then  $A\mathbf{v} = \sigma\mathbf{u}$  will tell me both  $\mathbf{u}$ 's. This is in equation (9).

The plan is to start with the  $v$ 's. **Choose orthonormal eigenvectors  $v_1, \dots, v_r$  of  $A^T A$ . Then choose  $\sigma_k = \sqrt{\lambda_k}$ . To determine the  $u$ 's we require  $Av = \sigma u$ :**

$$\boxed{\text{v's then u's} \quad A^T A v_k = \sigma_k^2 v_k \quad \text{and then} \quad u_k = \frac{A v_k}{\sigma_k} \quad \text{for } k = 1, \dots, r} \quad (9)$$

This is the proof of the SVD! Let me check that those  $u$ 's are eigenvectors of  $AA^T$ :

$$AA^T u_k = AA^T \left( \frac{A v_k}{\sigma_k} \right) = A \left( \frac{A^T A v_k}{\sigma_k} \right) = A \frac{\sigma_k^2 v_k}{\sigma_k} = \sigma_k^2 u_k \quad (10)$$

The  $v$ 's were chosen to be orthonormal. I must check that the  $u$ 's are also orthonormal:

$$u_j^T u_k = \left( \frac{A v_j}{\sigma_j} \right)^T \left( \frac{A v_k}{\sigma_k} \right) = \frac{v_j^T (A^T A v_k)}{\sigma_j \sigma_k} = \frac{\sigma_k}{\sigma_j} v_j^T v_k = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases} \quad (11)$$

Notice that  $(AA^T)A = A(A^T A)$  was the key to equation (10). The law  $(AB)C = A(BC)$  is the key to a great many proofs in linear algebra. Moving the parentheses is a powerful idea. This is the *associative law*.

Finally we have to choose the last  $n - r$  vectors  $v_{r+1}$  to  $v_n$  and the last  $m - r$  vectors  $u_{r+1}$  to  $u_m$ . This is easy. **These  $v$ 's and  $u$ 's are in the nullspaces of  $A$  and  $A^T$ .** We can choose any orthonormal bases for those nullspaces. They will automatically be orthogonal to the first  $v$ 's in the row space of  $A$  and the first  $u$ 's in the column space. This is because the whole spaces are orthogonal:  $N(A) \perp C(A^T)$  and  $N(A^T) \perp C(A)$ . *The proof of the SVD is complete.*

Now we have  $U$  and  $V$  and  $\Sigma$  in the full size SVD of equation (1). You may have noticed that the eigenvalues of  $A^T A$  are in  $\Sigma^T \Sigma$ , and *the same numbers  $\sigma_1^2$  to  $\sigma_r^2$  are also eigenvalues of  $AA^T$  in  $\Sigma \Sigma^T$* . An amazing fact:  **$BA$  always has the same nonzero eigenvalues as  $AB$ : 5 pages ahead.**

**Example 1** (completed) Find the matrices  $U, \Sigma, V$  for  $A = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix}$ .

With rank 2, this  $A$  has two positive singular values  $\sigma_1$  and  $\sigma_2$ . We will see that  $\sigma_1$  is larger than  $\lambda_{\max} = 5$ , and  $\sigma_2$  is smaller than  $\lambda_{\min} = 3$ . Begin with  $A^T A$  and  $AA^T$ :

$$A^T A = \begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix} \quad AA^T = \begin{bmatrix} 9 & 12 \\ 12 & 41 \end{bmatrix}$$

Those have the same trace (50) and the same eigenvalues  $\sigma_1^2 = 45$  and  $\sigma_2^2 = 5$ . The square roots are  $\sigma_1 = \sqrt{45}$  and  $\sigma_2 = \sqrt{5}$ . Then  $\sigma_1 \sigma_2 = 15$  and this is the determinant of  $A$ .

A key step is to find the eigenvectors of  $A^T A$  (with eigenvalues 45 and 5):

$$\begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 45 \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 25 & 20 \\ 20 & 25 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = 5 \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Then  $v_1$  and  $v_2$  are those orthogonal eigenvectors rescaled to length 1. Divide by  $\sqrt{2}$ .

**Right singular vectors**  $v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$   $v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$  **Left singular vectors**  $u_i = \frac{Av_i}{\sigma_i}$

Now compute  $Av_1$  and  $Av_2$  which will be  $\sigma_1 u_1 = \sqrt{45} u_1$  and  $\sigma_2 u_2 = \sqrt{5} u_2$ :

$$Av_1 = \frac{3}{\sqrt{2}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \sqrt{45} \frac{1}{\sqrt{10}} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \sigma_1 u_1$$

$$Av_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \sqrt{5} \frac{1}{\sqrt{10}} \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \sigma_2 u_2$$

The division by  $\sqrt{10}$  makes  $u_1$  and  $u_2$  orthonormal. Then  $\sigma_1 = \sqrt{45}$  and  $\sigma_2 = \sqrt{5}$  as expected. The Singular Value Decomposition of  $A$  is  $U$  times  $\Sigma$  times  $V^T$ .

$$\boxed{U = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sqrt{45} & \\ & \sqrt{5} \end{bmatrix} \quad V = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}} \quad (12)$$

$U$  and  $V$  contain orthonormal bases for the column space and the row space of  $A$  (both spaces are just  $\mathbf{R}^2$ ). The real achievement is that those two bases diagonalize  $A$ :  $AV$  equals  $U\Sigma$ . The matrix  $A = U\Sigma V^T$  splits into two rank-one matrices, columns times rows, with  $\sqrt{2}\sqrt{10} = \sqrt{20}$ .

$$\sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T = \frac{\sqrt{45}}{\sqrt{20}} \begin{bmatrix} 1 & 1 \\ 3 & 3 \end{bmatrix} + \frac{\sqrt{5}}{\sqrt{20}} \begin{bmatrix} 3 & -3 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 4 & 5 \end{bmatrix} = A.$$

Every matrix is a sum of rank one matrices with orthogonal  $u$ 's and orthogonal  $v$ 's.

**Question:** If  $S = Q\Lambda Q^T$  is symmetric positive definite, what is its SVD?

**Answer:** The SVD is exactly  $U\Sigma V^T = Q\Lambda Q^T$ . The matrix  $U = V = Q$  is orthogonal. And the eigenvalue matrix  $\Lambda$  becomes the singular value matrix  $\Sigma$ .

**Question:** If  $S = Q\Lambda Q^T$  has a negative eigenvalue ( $Sx = -\alpha x$ ), what is the singular value and what are the vectors  $v$  and  $u$ ?

**Answer:** The singular value will be  $\sigma = +\alpha$  (positive). One singular vector (either  $u$  or  $v$ ) must be  $-x$  (reverse the sign). Then  $Sx = -\alpha x$  is the same as  $Sv = \sigma u$ . The two sign changes cancel.

**Question:** If  $A = Q$  is an orthogonal matrix, why does every singular value equal 1?

**Answer:** All singular values are  $\sigma = 1$  because  $A^T A = Q^T Q = I$ . Then  $\Sigma = I$ . But  $U = Q$  and  $V = I$  is only one choice for the singular vectors  $u$  and  $v$ :

$$Q = U\Sigma V^T \text{ can be } Q = QII^T \text{ or any } Q = (QQ_1)IQ_1^T.$$

**Question :** Why are all eigenvalues of a square matrix  $A$  less than or equal to  $\sigma_1$  ?

*Answer :* Multiplying by orthogonal matrices  $U$  and  $V^T$  does not change vector lengths :

$$\|Ax\| = \|U\Sigma V^T x\| = \|\Sigma V^T x\| \leq \sigma_1 \|V^T x\| = \sigma_1 \|x\| \text{ for all } x. \quad (13)$$

An eigenvector has  $\|Ax\| = |\lambda| \|x\|$ . Then (13) gives  $|\lambda| \|x\| \leq \sigma_1 \|x\|$  and  $|\lambda| \leq \sigma_1$ .

**Question :** If  $A = xy^T$  has rank 1, what are  $u_1$  and  $v_1$  and  $\sigma_1$  ? **Check that  $|\lambda_1| \leq \sigma_1$ .**

*Answer :* The singular vectors  $u_1 = x/\|x\|$  and  $v_1 = y/\|y\|$  have length 1. Then  $\sigma_1 = \|x\| \|y\|$  is the only nonzero number in the singular value matrix  $\Sigma$ . Here is the SVD :

$$\text{Rank 1 matrix} \quad xy^T = \frac{x}{\|x\|} (\|x\| \|y\|) \frac{y^T}{\|y\|} = u_1 \sigma_1 v_1^T.$$

*Observation* The only nonzero eigenvalue of  $A = xy^T$  is  $\lambda = y^T x$ . The eigenvector is  $x$  because  $(xy^T)x = x(y^T x) = \lambda x$ . Then  $|\lambda_1| = |y^T x| \leq \sigma_1 = \|y\| \|x\|$ .

The key inequality  $|\lambda_1| \leq \sigma_1$  becomes exactly the Schwarz inequality.

**Question :** What is the Karhunen-Loève transform and its connection to the SVD ?

*Answer :* KL begins with a covariance matrix  $V$  of a zero-mean random process.  $V$  is symmetric and positive definite or semidefinite. In general  $V$  could be an infinite matrix or a covariance function. Then the KL expansion will be an infinite series.

The eigenvectors of  $V$ , in order of decreasing eigenvalues  $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq 0$ , are the basis functions  $u_i$  for the KL transform. The expansion of any vector  $v$  in an orthonormal basis  $u_1, u_2, \dots$  is  $v = \sum (u_i^T v) u_i$ .

In this stochastic case, that transform decorrelates the random process: the  $u_i$  are independent. More than that, the ordering of the eigenvalues means that the first  $k$  terms, stopping at  $(u_k^T v) u_k$ , minimize the expected square error. This fact corresponds to the Eckart-Young Theorem in the next section I.9.

The KL transform is a stochastic (random) form of Principal Component Analysis.

## The Geometry of the SVD

The SVD separates a matrix into  $A = U\Sigma V^T$  : **(orthogonal)  $\times$  (diagonal)  $\times$  (orthogonal)**. In two dimensions we can draw those steps. The orthogonal matrices  $U$  and  $V$  rotate the plane. The diagonal matrix  $\Sigma$  stretches it along the axes. Figure I.11 shows **rotation** times **stretching** times **rotation**. **Vectors  $x$  on the unit circle go to  $Ax$  on an ellipse.**



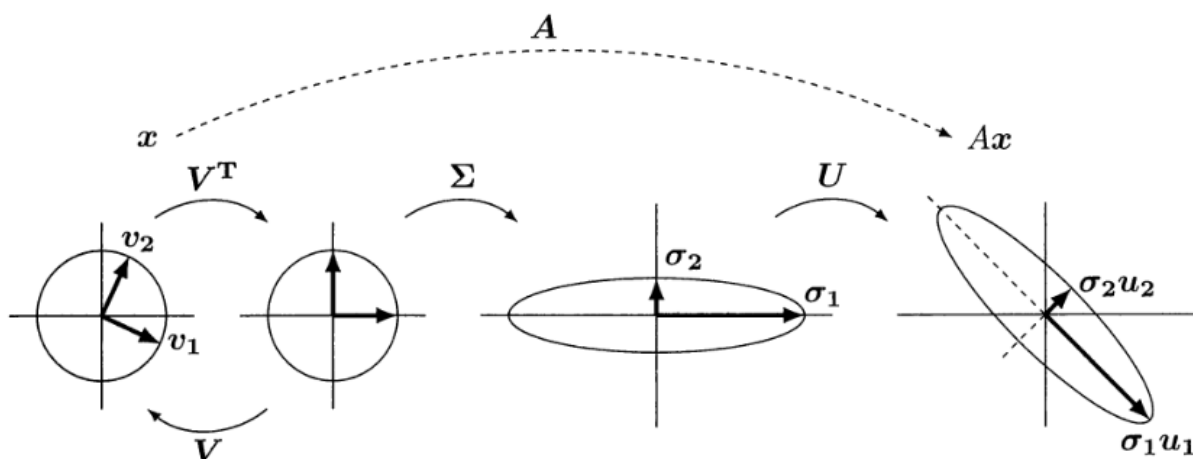


Figure I.10:  $U$  and  $V$  are rotations and possible reflections.  $\Sigma$  stretches circle to ellipse.

This picture applies to a 2 by 2 invertible matrix (because  $\sigma_1 > 0$  and  $\sigma_2 > 0$ ). First is a rotation of any  $x$  to  $V^T x$ . Then  $\Sigma$  stretches that vector to  $\Sigma V^T x$ . Then  $U$  rotates to  $Ax = U \Sigma V^T x$ . We kept all determinants positive to avoid reflections. The four numbers  $a, b, c, d$  in the matrix connect to *two angles*  $\theta$  and  $\phi$  and *two numbers*  $\sigma_1$  and  $\sigma_2$ .

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \sigma_1 & \\ & \sigma_2 \end{bmatrix} \begin{bmatrix} \cos \phi & \sin \phi \\ -\sin \phi & \cos \phi \end{bmatrix}. \quad (14)$$

*Question.* If the matrix is symmetric then  $b = c$  and  $A$  has only 3 (not 4) parameters. How do the 4 numbers  $\theta, \phi, \sigma_1, \sigma_2$  reduce to 3 numbers for a symmetric matrix  $S$ ?

### The First Singular Vector $v_1$

The next page will establish a new way to look at  $v_1$ . The previous pages chose the  $v$ 's as eigenvectors of  $A^T A$ . Certainly that remains true. But there is a valuable way to understand these singular vectors **one at a time instead of all at once**. We start with  $v_1$  and the singular value  $\sigma_1$ .

**Maximize the ratio  $\frac{\|Ax\|}{\|x\|}$ . The maximum is  $\sigma_1$  at the vector  $x = v_1$ .**

(15)

The ellipse in Figure I.10 showed why the maximizing  $x$  is  $v_1$ . When you follow  $v_1$  across the page, it ends at  $Av_1 = \sigma_1 u_1$  (the longest axis of the ellipse). Its length started at  $\|v_1\| = 1$  and ended at  $\|Av_1\| = \sigma_1$ .

But we aim for an independent approach to the SVD! We are not assuming that we already know  $U$  or  $\Sigma$  or  $V$ . How do we recognize that the ratio  $\|Ax\|/\|x\|$  is a maximum when  $x = v_1$ ? Calculus tells us that the first derivatives must be zero. The derivatives will be easier if we square our function:

**Problem : Find the maximum value  $\lambda$  of** 
$$\frac{\|Ax\|^2}{\|x\|^2} = \frac{x^T A^T A x}{x^T x} = \frac{x^T S x}{x^T x}. \quad (16)$$

This “Rayleigh quotient” depends on  $x_1, \dots, x_n$ . Calculus uses the quotient rule, so we need

$$\frac{\partial}{\partial x_i} (x^T x) = \frac{\partial}{\partial x_i} (x_1^2 + \dots + x_i^2 + \dots + x_n^2) = 2(x)_i \quad (17)$$

$$\frac{\partial}{\partial x_i} (x^T S x) = \frac{\partial}{\partial x_i} \left( \sum_i \sum_j S_{ij} x_i x_j \right) = 2 \sum_j S_{ij} x_j = 2 (Sx)_i \quad (18)$$

The quotient rule finds  $\partial/\partial x_i (x^T S x / x^T x)$ . Set those  $n$  partial derivatives of (16) to zero:

$$(x^T x) 2(Sx)_i - (x^T S x) 2(x)_i = 0 \text{ for } i = 1, \dots, n \quad (19)$$

**Equation (19) says that the best  $x$  is an eigenvector of  $S = A^T A$ !**

**$2Sx = 2\lambda x$  and the maximum value of  $\frac{x^T S x}{x^T x} = \frac{\|Ax\|^2}{\|x\|^2}$  is an eigenvalue  $\lambda$  of  $S$ .**

The search is narrowed to eigenvectors of  $S = A^T A$ . The eigenvector that maximizes is  $x = v_1$ . The eigenvalue is  $\lambda_1 = \sigma_1^2$ . Calculus has confirmed the solution (15) of the maximum problem—the first piece of the SVD.

For the full SVD, we need *all* the singular vectors and singular values. To find  $v_2$  and  $\sigma_2$ , we adjust the maximum problem so it looks only at vectors  $x$  orthogonal to  $v_1$ .

**Maximize  $\frac{\|Ax\|}{\|x\|}$  under the condition  $v_1^T x = 0$ . The maximum is  $\sigma_2$  at  $x = v_2$ .**

“Lagrange multipliers” were invented to deal with constraints on  $x$  like  $v_1^T x = 0$ . And Problem 3 gives a simple direct way to work with this condition  $v_1^T x = 0$ .

In the same way, every singular vector  $v_{k+1}$  gives the maximum ratio over all vectors  $x$  that are perpendicular to the first  $v_1, \dots, v_k$ . The left singular vectors would come from maximizing  $\|A^T y\|/\|y\|$ . We are always finding the axes of an ellipsoid and the eigenvectors of symmetric matrices  $A^T A$  or  $AA^T$ .

## The Singular Vectors of $A^T$

The SVD connects  $v$ 's in the row space to  $u$ 's in the column space. When we transpose  $A = U\Sigma V^T$ , we see that  $A^T = V\Sigma^T U^T$  goes the opposite way, from  $u$ 's to  $v$ 's:

$$A^T u_k = \sigma_k v_k \text{ for } k = 1, \dots, r \quad A^T u_k = 0 \text{ for } k = r + 1, \dots, m \quad (20)$$

Multiply  $Av_k = \sigma_k u_k$  by  $A^T$ . Remember  $A^T Av_k = \sigma_k^2 v_k$  in equation (9). Divide by  $\sigma_k$ .

## A Different Symmetric Matrix Also Produces the SVD

We created the SVD from two symmetric matrices  $A^T A$  and  $AA^T$ . Another good way uses one symmetric block matrix  $S$ . *This matrix has  $r$  pairs of plus and minus eigenvalues.* The nonzero eigenvalues of this matrix  $S$  are  $\sigma_k$  and  $-\sigma_k$ , and its size is  $m + n$ :

$$S = \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} \text{ has eigenvectors } \begin{bmatrix} u_k \\ v_k \end{bmatrix} \text{ and } \begin{bmatrix} -u_k \\ v_k \end{bmatrix}.$$

We can check those eigenvectors directly, remembering  $Av_k = \sigma_k u_k$  and  $A^T u_k = \sigma_k v_k$ :

$$\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} \pm u_k \\ v_k \end{bmatrix} = \begin{bmatrix} Av_k \\ \pm A^T u_k \end{bmatrix} = \sigma_k \begin{bmatrix} u_k \\ v_k \end{bmatrix} \text{ and } -\sigma_k \begin{bmatrix} -u_k \\ v_k \end{bmatrix}. \quad (21)$$

That gives  $2r$  eigenvalues. The eigenvectors are orthogonal:  $-u_k^T u_k + v_k^T v_k = -1 + 1$ . Can you see the other  $(m - r) + (n - r)$  eigenvectors with  $\lambda = 0$  for that block matrix? They must involve the remaining  $u$ 's and  $v$ 's in the nullspaces of  $A^T$  and  $A$ .

## Submatrices Have Smaller Singular Values

The approach to  $\|A\| = \sigma_1$  by maximizing  $\|Ax\|/\|x\|$  makes it easy to prove this useful fact. The norm of a submatrix cannot be larger than the norm of the whole matrix:  $\sigma_1(B) \leq \sigma_1(A)$ .

**If  $B$  keeps  $M \leq m$  rows and  $N \leq n$  columns of  $A$ , then  $\|B\| \leq \|A\|$ .** (22)

*Proof* Look at vectors  $y$  with nonzeros only in the  $N$  positions that correspond to columns in  $B$ . Certainly maximum of  $\|By\|/\|y\| \leq$  maximum of  $\|Ax\|/\|x\|$ .

Reduce  $\|By\|$  further by looking only at the  $M$  components that correspond to rows of  $B$ . So removing columns and rows cannot increase the norm  $\sigma_1$ , and  $\|B\| \leq \|A\|$ .

## The SVD for Derivatives and Integrals

This may be the clearest example of the SVD. It does not start with a matrix (but we will go there). Historically, the first SVD was not for vectors but for *functions*. Then  $A$  is not a matrix but an *operator*. One example is the operator that integrates every function. Another example is the (unbounded) operator  $D$  that takes the derivative :

$$\begin{array}{ll} \textbf{Operators on functions} & \\ \textbf{Integral and derivative} & Ax(s) = \int_0^s x(t) dt \quad \text{and} \quad Dx(t) = \frac{dx}{dt}. \end{array} \quad (23)$$

Those operators are linear (or calculus would be a lot more difficult than it is). In some way  $D$  is the inverse of  $A$ , by the Fundamental Theorem of Calculus. More exactly  $D$  is a left inverse with  $DA = I$ : derivative of integral equals original function.

But  $AD \neq I$  because the derivative of a constant function is zero. Then  $D$  has a nullspace, like a matrix with dependent columns.  **$D$  is the pseudoinverse of  $A$ !** Sines and cosines are the  $u$ 's and  $v$ 's for  $A = \text{integral}$  and  $D = \text{derivative}$ :

$$Av = \sigma u \text{ is } A(\cos kt) = \frac{1}{k}(\sin kt) \quad \text{Then } D(\sin kt) = k(\cos kt). \quad (24)$$

The simplicity of those equations is our reason for including them in the book. We are working with *periodic functions*:  $x(t + 2\pi) = x(t)$ . The input space to  $A$  contains the *even functions* like  $\cos t = \cos(-t)$ . The outputs from  $A$  (and the inputs to  $D$ ) are the *odd functions* like  $\sin t = -\sin(-t)$ . Those input and output spaces are like  $\mathbf{R}^n$  and  $\mathbf{R}^m$  for an  $m$  by  $n$  matrix.

The special property of the SVD is that the  $v$ 's are orthogonal, and so are the  $u$ 's. Here those singular vectors have become very nice functions—the *cosines are orthogonal to each other and so are the sines*. Their inner products are integrals equal to zero:

$$v_k^T v_j = \int_0^{2\pi} (\cos kt) (\cos jt) dt = 0 \quad \text{and} \quad u_k^T u_j = \int_0^{2\pi} (\sin kt) (\sin jt) dt = 0.$$

Notice that the inner product of functions  $x_1$  and  $x_2$  is the integral of  $x_1(t)x_2(t)$ . This copies into function space (*Hilbert space*) the dot product that adds  $y \cdot z = \sum y_i z_i$ . In fact the symbol  $\int$  was somehow created from  $\Sigma$  (and integrals are the limits of sums).

## Finite Differences

The discrete form of a derivative is a **finite difference**. The discrete form of an integral is a **sum**. Here we choose a 4 by 3 matrix  $D$  that corresponds to the backward difference  $f(x) - f(x - \Delta x)$ :

$$D = \begin{bmatrix} 1 & & & \\ -1 & 1 & & \\ & -1 & 1 & \\ & & -1 & 1 \end{bmatrix} \quad \text{with} \quad D^T = \begin{bmatrix} 1 & -1 & 0 & 0 \\ & 1 & -1 & 0 \\ & & 1 & -1 \\ & & & 1 \end{bmatrix}. \quad (25)$$

To find singular values and singular vectors, compute  $D^T D$  (3 by 3) and  $DD^T$  (4 by 4):

$$D^T D = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \quad \text{and} \quad DD^T = \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}. \quad (26)$$

Their nonzero eigenvalues are always the same!  $DD^T$  also has a zero eigenvalue with eigenvector  $\mathbf{u}_4 = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$ . This is the discrete equivalent of the function  $f(x) = \frac{1}{2}$  with  $df/dx = 0$ .

The nonzero eigenvalues of both symmetric matrices  $D^T D$  and  $DD^T$  are

$$\lambda_1 = \sigma_1^2(D) = 2 + \sqrt{2} \quad \lambda_2 = \sigma_2^2(D) = 2 \quad \lambda_3 = \sigma_3^2(D) = 2 - \sqrt{2} \quad (27)$$

The eigenvectors  $\mathbf{v}$  of  $D^T D$  are the right singular vectors of  $D$ . They are **discrete sines**. The eigenvectors  $\mathbf{u}$  of  $DD^T$  are the left singular vectors of  $D$ . They are **discrete cosines**:

$$\sqrt{2} \mathbf{V} = \begin{bmatrix} \sin \frac{\pi}{4} & \sin \frac{2\pi}{4} & \sin \frac{3\pi}{4} \\ \sin \frac{2\pi}{4} & \sin \frac{4\pi}{4} & \sin \frac{6\pi}{4} \\ \sin \frac{3\pi}{4} & \sin \frac{6\pi}{4} & \sin \frac{9\pi}{4} \end{bmatrix} \quad \sqrt{2} \mathbf{U} = \begin{bmatrix} \cos \frac{1}{2} \frac{\pi}{4} & \cos \frac{1}{2} \frac{2\pi}{4} & \cos \frac{1}{2} \frac{3\pi}{4} & 1 \\ \cos \frac{3}{2} \frac{\pi}{4} & \cos \frac{3}{2} \frac{2\pi}{4} & \cos \frac{3}{2} \frac{3\pi}{4} & 1 \\ \cos \frac{5}{2} \frac{\pi}{4} & \cos \frac{5}{2} \frac{2\pi}{4} & \cos \frac{5}{2} \frac{3\pi}{4} & 1 \\ \cos \frac{7}{2} \frac{\pi}{4} & \cos \frac{7}{2} \frac{2\pi}{4} & \cos \frac{7}{2} \frac{3\pi}{4} & 1 \end{bmatrix}.$$

These are the famous DST and DCT matrices—**Discrete Sine Transform** and **Discrete Cosine Transform**. The DCT matrix has been the backbone of JPEG image compression. Actually JPEG increases  $U$  to 8 by 8, which reduces the “blockiness” of the image. 8 by 8 blocks of pixels are transformed by a two-dimensional DCT—then compressed and transmitted. Orthogonality of these matrices is the key in Section IV.4.

Our goal was to show the discrete form of the beautiful Singular Value Decomposition  $D(\sin kt) = k(\cos kt)$ . You could correctly say that this is only one example. But Fourier is always present for linear equations with constant coefficients—and always important.

In signal processing the key letters are LTI: **Linear Time Invariance**.

## The Polar Decomposition $A = QS$

**Every complex number  $x + iy$  has the polar form  $re^{i\theta}$ .** A number  $r \geq 0$  multiplies a number  $e^{i\theta}$  on the unit circle. We have  $x + iy = r \cos \theta + ir \sin \theta = re^{i\theta}$ . Think of these numbers as 1 by 1 matrices. Then  $e^{i\theta}$  is an *orthogonal matrix*  $Q$  and  $r \geq 0$  is a *positive semidefinite matrix* (call it  $S$ ). The **polar decomposition** extends the same idea to  $n$  by  $n$  matrices: orthogonal times positive semidefinite,  $A = QS$ .

Every real square matrix can be factored into  $A = QS$ , where  $Q$  is *orthogonal* and  $S$  is *symmetric positive semidefinite*. If  $A$  is invertible,  $S$  is positive definite.

<b>Polar decomposition</b>	$A = U\Sigma V^T = (UV^T)(V\Sigma V^T) = (Q)(S).$	$(28)$
----------------------------	---------------------------------------------------	--------

The first factor  $UV^T$  is  $Q$ . The product of orthogonal matrices is orthogonal. The second factor  $V\Sigma V^T$  is  $S$ . It is positive semidefinite because its eigenvalues are in  $\Sigma$ .

If  $A$  is invertible then  $\Sigma$  and  $S$  are also invertible.  **$S$  is the symmetric positive definite square root of  $A^T A$ ,** because  $S^2 = V\Sigma^2 V^T = A^T A$ . So the eigenvalues of  $S$  are the singular values of  $A$ . The eigenvectors of  $S$  are the singular vectors  $v$  of  $A$ .

---