

## THE SINGULAR VALUE DECOMPOSITION

SYNOPSIS. The linear algebra framework in which matrix-vector multiplication represents a linear combination and also corresponds to evaluation of a linear mapping has been shown to be complete through the FTLA. In order to effectively solve the main problems of linear algebra within this framework, bases must be constructed for the fundamental subspaces of a matrix. In order for the bases to be computationally efficient they should be orthonormal. The existence of such basis sets is guaranteed by the singular value decomposition theorem.

### 1. Orthogonal matrices

The fundamental theorem of linear algebra partitions the domain and codomain of a linear mapping  $f: U \rightarrow V$ . For real vectors spaces  $U = \mathbb{R}^n$ ,  $V = \mathbb{R}^m$  the partition properties are stated in terms of spaces of the associated matrix  $A$  as

$$C(A) \oplus N(A^T) = \mathbb{R}^m \quad C(A) \perp N(A^T) \quad C(A^T) \oplus N(A) = \mathbb{R}^n \quad C(A^T) \perp N(A).$$

The dimension of the column and row spaces  $r = \dim C(A) = \dim C(A^T)$  is the rank of the matrix,  $n - r$  is the nullity of  $A$ , and  $m - r$  is the nullity of  $A^T$ . A infinite number of bases could be defined for the domain and codomain. It is of great theoretical and practical interest to define bases with properties that facilitate insight or computation.

The above partitions of the domain and codomain are orthogonal, and suggest searching for orthogonal bases within these subspaces. Introduce a matrix representation for the bases

$$U = [u_1 \ u_2 \ \dots \ u_m] \in \mathbb{R}^{m \times m}, \quad V = [v_1 \ v_2 \ \dots \ v_n] \in \mathbb{R}^{n \times n},$$

with  $C(U) = \mathbb{R}^m$  and  $C(V) = \mathbb{R}^n$ . Orthogonality between columns  $u_i, u_j$  for  $i \neq j$  is expressed as  $u_i^T u_j = 0$ . For  $i = j$ , the inner product is positive  $u_i^T u_i > 0$ , and since scaling of the columns of  $U$  preserves the spanning property  $C(U) = \mathbb{R}^m$ , it is convenient to impose  $u_i^T u_i = 1$ . Such behavior is concisely expressed as a matrix product

$$U^T U = I_m,$$

with  $I_m$  the identity matrix in  $\mathbb{R}^m$ . Expanded in terms of the column vectors of  $U$  the first equality is

$$[u_1 \ u_2 \ \dots \ u_m]^T [u_1 \ u_2 \ \dots \ u_m] = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} [u_1 \ u_2 \ \dots \ u_m] = \begin{bmatrix} u_1^T u_1 & u_1^T u_2 & \dots & u_1^T u_m \\ u_2^T u_1 & u_2^T u_2 & \dots & u_2^T u_m \\ \vdots & \vdots & \ddots & \vdots \\ u_m^T u_1 & u_m^T u_2 & \dots & u_m^T u_m \end{bmatrix} = I_m.$$

It is useful to determine if a matrix  $X$  exists such that  $UX = I_m$ , or

$$UX = U [x_1 \ x_2 \ \dots \ x_m] = [e_1 \ e_2 \ \dots \ e_m].$$

The columns of  $X$  are the coordinates of the column vectors of  $I_m$  in the basis  $U$ , and can readily be determined

$$Ux_j = e_j \Rightarrow U^T Ux_j = U^T e_j \Rightarrow I_m x_j = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} e_j \Rightarrow x_j = (U^T)_j,$$

where  $(U^T)_j$  is the  $j^{\text{th}}$  column of  $U^T$ , hence  $\mathbf{x} = U^T$ , leading to

$$U^T U = I = U U^T.$$

Note that the second equality

$$[ \mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m ] [ \mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m ]^T = [ \mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m ] \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_m^T \end{bmatrix} = \mathbf{u}_1 \mathbf{u}_1^T + \mathbf{u}_2 \mathbf{u}_2^T + \dots + \mathbf{u}_m \mathbf{u}_m^T = I$$

acts as normalization condition on the matrices  $U_j = \mathbf{u}_j \mathbf{u}_j^T$ .

DEFINITION. A square matrix  $U$  is said to be orthogonal if  $U^T U = U U^T = I$ .

## 2. Intrinsic basis of a linear mapping

Given a linear mapping  $f: U \rightarrow V$ , expressed as  $\mathbf{y} = f(\mathbf{x}) = A\mathbf{x}$ , the simplest description of the action of  $A$  would be a simple scaling, as exemplified by  $\mathbf{g}(\mathbf{x}) = a\mathbf{x}$  that has as its associated matrix  $aI$ . Recall that specification of a vector is typically done in terms of the identity matrix  $\mathbf{b} = I\mathbf{b}$ , but may be more insightfully given in some other basis  $A\mathbf{x} = I\mathbf{b}$ . This suggests that especially useful bases for the domain and codomain would reduce the action of a linear mapping to scaling along orthogonal directions, and evaluate  $\mathbf{y} = A\mathbf{x}$  by first re-expressing  $\mathbf{y}$  in another basis  $U$ ,  $U\mathbf{s} = I\mathbf{y}$  and re-expressing  $\mathbf{x}$  in another basis  $V$ ,  $V\mathbf{r} = I\mathbf{x}$ . The condition that the linear operator reduces to simple scaling in these new bases is expressed as  $s_i = \sigma_i r_i$  for  $i = 1, \dots, \min(m, n)$ , with  $\sigma_i$  the scaling coefficients along each direction which can be expressed as a matrix vector product  $\mathbf{s} = \Sigma \mathbf{r}$ , where  $\Sigma \in \mathbb{R}^{m \times n}$  is of the same dimensions as  $A$  and given by

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix}.$$

Imposing the condition that  $U, V$  are orthogonal leads to

$$U\mathbf{s} = \mathbf{y} \Rightarrow \mathbf{s} = U^T \mathbf{y}, \quad V\mathbf{r} = \mathbf{x} \Rightarrow \mathbf{r} = V^T \mathbf{x},$$

which can be replaced into  $\mathbf{s} = \Sigma \mathbf{r}$  to obtain

$$U^T \mathbf{y} = \Sigma V^T \mathbf{x} \Rightarrow \mathbf{y} = U \Sigma V^T \mathbf{x}.$$

From the above the orthogonal bases  $U, V$  and scaling coefficients  $\Sigma$  that are sought must satisfy  $A = U \Sigma V^T$ . The SVD theorem states that the matrix factors  $U, \Sigma, V$  do indeed exist.

THEOREM. Every matrix  $A \in \mathbb{R}^{m \times n}$  has a *singular value decomposition* (SVD)

$$A = U \Sigma V^T,$$

with properties:

1.  $\mathbf{U} \in \mathbb{R}^{m \times m}$  is an orthogonal matrix,  $\mathbf{U}^T \mathbf{U} = \mathbf{I}_m$ ;
2.  $\mathbf{V} \in \mathbb{R}^{n \times n}$  is an orthogonal matrix,  $\mathbf{V}^T \mathbf{V} = \mathbf{I}_n$ ;
3.  $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$  is diagonal,  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p)$ ,  $p = \min(m, n)$ , and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ .

The scaling coefficients  $\sigma_j$  are called the *singular values* of  $\mathbf{A}$ . The columns of  $\mathbf{U}$  are called the *left singular vectors*, and those of  $\mathbf{V}$  are called the *right singular vectors*. Carrying out computation of the matrix products

$$\mathbf{A} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_r^T \\ \mathbf{v}_{r+1}^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} \Rightarrow$$

$$\mathbf{A} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 \mathbf{v}_1^T \\ \sigma_2 \mathbf{v}_2^T \\ \vdots \\ \sigma_r \mathbf{v}_r^T \\ \vdots \\ 0 \end{bmatrix}$$

leads to a representation of  $\mathbf{A}$  as a sum

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T,$$

with  $r \leq \min(m, n)$ . Written out in full, the above sum is

$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T.$$

Each product  $\mathbf{u}_i \mathbf{v}_i^T$  is a matrix of rank one, and is called a rank-one update. Truncation of the above sum to  $p$  terms leads to an approximation of  $\mathbf{A}$

$$\mathbf{A} \cong \mathbf{A}_p = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^T.$$

In very many cases the singular values exhibit rapid decay,  $\sigma_1 \gg \sigma_2 \gg \dots$ , such that the approximation above is an accurate representation of the matrix  $\mathbf{A}$  for  $p \ll r$ .

The singular vector matrices  $\mathbf{U}, \mathbf{V}$  specify the intrinsic directions within  $\mathbb{R}^m, \mathbb{R}^n$  along which the matrix  $\mathbf{A}$  acts as a simple scaling transformation. For example, applying the linear mapping to the  $\mathbf{v}_1$  vector,  $\mathbf{f}(\mathbf{v}_1) = \mathbf{A} \mathbf{v}_1$ , leads to

$$\mathbf{A} \mathbf{v}_1 = \left( \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^T \right) \mathbf{v}_1 = \sum_{i=1}^p \sigma_i \mathbf{u}_i (\mathbf{v}_i^T \mathbf{v}_1) = \sigma_1 \mathbf{u}_1.$$

Since  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ , the above states that the input direction most amplified by the  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$  mapping is  $\mathbf{v}_1$  and the result is the vector  $\sigma_1 \mathbf{u}_1$ . The two-norm of  $\mathbf{v}_1$  is equal to one and that of  $\sigma_1 \mathbf{u}_1$  is  $\sigma_1$ . The conclusion is that  $\sigma_1$  is the maximal amplification factor in the two-norm

$$\sigma_1 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2,$$

and the above satisfies the properties of a norm over matrices leading to the definition

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2.$$

The largest singular value is thus the two-norm of a matrix.

### 3. SVD solution of linear algebra problems

The SVD can be used to solve common problems within linear algebra.

**Linear systems.** To change from vector coordinates  $\mathbf{b}$  in the canonical basis  $\mathbf{I} \in \mathbb{R}^{m \times m}$  to coordinates  $\mathbf{x}$  in some other basis  $\mathbf{A} \in \mathbb{R}^{m \times m}$ , a solution to the equation  $\mathbf{I}\mathbf{b} = \mathbf{A}\mathbf{x}$  can be found by the following steps.

1. Compute the SVD,  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{A}$ ;
2. Find the coordinates of  $\mathbf{b}$  in the orthogonal basis  $\mathbf{U}$ ,  $\mathbf{c} = \mathbf{U}^T \mathbf{b}$ ;
3. Scale the coordinates of  $\mathbf{c}$  by the inverse of the singular values  $y_i = c_i / \sigma_i$ ,  $i = 1, \dots, m$ , such that  $\mathbf{\Sigma}\mathbf{y} = \mathbf{c}$  is satisfied;
4. Find the coordinates of  $\mathbf{y}$  in basis  $\mathbf{V}^T$ ,  $\mathbf{x} = \mathbf{V}\mathbf{y}$ .

**Least squares.** In the above  $\mathbf{A}$  was assumed to be a basis, hence  $r = \text{rank}(\mathbf{A}) = m$ . If columns of  $\mathbf{A}$  do not form a basis,  $r < m$ , then  $\mathbf{b} \in \mathbb{R}^m$  might not be reachable by linear combinations within  $C(\mathbf{A})$ . The closest vector to  $\mathbf{b}$  in the 2-norm is however found by the same steps, with the simple modification that in Step 3, the scaling is carried out only for non-zero singular values,  $y_i = c_i / \sigma_i$ ,  $i = 1, \dots, r$ .

**The pseudo-inverse.** From the above, finding either the solution of  $\mathbf{A}\mathbf{x} = \mathbf{I}\mathbf{b}$  or the best approximation possible if  $\mathbf{A}$  is not of full rank, can be written as a sequence of matrix multiplications using the SVD

$$(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)\mathbf{x} = \mathbf{b} \Rightarrow \mathbf{U}(\mathbf{\Sigma}\mathbf{V}^T\mathbf{x}) = \mathbf{b} \Rightarrow (\mathbf{\Sigma}\mathbf{V}^T\mathbf{x}) = \mathbf{U}^T\mathbf{b} \Rightarrow \mathbf{V}^T\mathbf{x} = \mathbf{\Sigma}^+ \mathbf{U}^T\mathbf{b} \Rightarrow \mathbf{x} = \mathbf{V}\mathbf{\Sigma}^+ \mathbf{U}^T\mathbf{b},$$

where the matrix  $\mathbf{\Sigma}^+ \in \mathbb{R}^{n \times m}$  (notice the inversion of dimensions) is defined as a matrix with elements  $\sigma_i^{-1}$  on the diagonal, and is called the pseudo-inverse of  $\mathbf{\Sigma}$ . Similarly the matrix

$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+ \mathbf{U}^T$$

that allows stating the solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$  simply as  $\mathbf{x} = \mathbf{A}^+ \mathbf{b}$  is called the *pseudo-inverse* of  $\mathbf{A}$ . Note that in practice  $\mathbf{A}^+$  is not explicitly formed. Rather the notation  $\mathbf{A}^+$  is simply a concise reference to carrying out steps 1-4 above.