LECTURE 7: THE SINGULAR VALUE DECOMPOSITION

1. Mappings as data

1.1. Vector spaces of mappings and matrix representations

A vector space \mathscr{L} can be formed from all linear mappings from the vector space $\mathscr{U} = (U, S, +, \cdot)$ to another vector space $\mathscr{V} = (V, S, +, \cdot)$

$$\mathscr{L} = \{L, S, +, \cdot\}, \ L = \{f | f \colon U \to V, f(a\mathbf{u} + b\mathbf{v}) = af(\mathbf{u}) + bf(\mathbf{v})\},\$$

with addition and scaling of linear mappings defined by (f + g)(u) = f(u) + g(u) and (af)(u) = af(u). Let $B = \{u_1, u_2, ...\}$ denote a basis for the domain U of linear mappings within \mathcal{L} , such that the linear mapping $f \in \mathcal{L}$ is represented by the matrix

$$\boldsymbol{A} = [\boldsymbol{f}(\boldsymbol{u}_1) \ \boldsymbol{f}(\boldsymbol{u}_2) \ \dots].$$

When the domain and codomain are the real vector spaces $U = \mathbb{R}^n$, $V = \mathbb{R}^m$, the above is a standard matrix of real numbers, $A \in \mathbb{R}^{m \times n}$. For linear mappings between infinite dimensional vector spaces, the matrix is understood in a generalized sense to contain an infinite number of columns that are elements of the codomain *V*. For example, the indefinite integral is a linear mapping between the vector space of functions that allow differentiation to any order,

$$\int : \mathscr{C}^{\infty} \to \mathscr{C}^{\infty} \ v(t) = \int u(t) \, \mathrm{d}t$$

and for the monomial basis $B = \{1, t, t^2, ...\}$, is represented by the generalized matrix

$$\boldsymbol{A} = \left[\begin{array}{ccc} t & \frac{1}{2}t^2 & \frac{1}{3}t^3 & \dots \end{array} \right].$$

Truncation of the MacLaurin series $u(t) = \sum_{j=1}^{\infty} u_j t^j$, with $u_j = u^{(j)}(0) / j! \in \mathbb{R}$ to *n* terms, and sampling of $u \in \mathscr{C}^{\infty}$ at points t_1, \ldots, t_m , forms a standard matrix of real numbers

$$\boldsymbol{A} = \left[\boldsymbol{t} \quad \frac{1}{2} \boldsymbol{t}^2 \quad \frac{1}{3} \boldsymbol{t}^3 \quad \dots \right] \in \mathbb{R}^{m \times n}, \ \boldsymbol{t}^j = \left[\begin{array}{c} \boldsymbol{t}_1^j \\ \vdots \\ \boldsymbol{t}_m^j \end{array} \right].$$

Values of function $u \in \mathscr{C}^{\infty}$ at t_1, \ldots, t_m are approximated by

$$\boldsymbol{u} = \boldsymbol{B}\boldsymbol{x} = [u(t_1) \dots u(t_m)]^T,$$

with x denoting the coordinates of u in basis B. The above argument states that the coordinates y of v, the primitive of u are given by

$$y = A x$$
,

as can be indeed verified through term-by-term integration of the MacLaurin series.

As to be expected, matrices can also be organized as vector space \mathcal{M} , which is essentially the representation of the associated vector space of linear mappings,

$$\mathcal{M} = (M, S, +, \cdot) \quad M = \{A \mid A = [f(u_1) \mid f(u_2) \mid \dots \}\}$$

The addition C = A + B and scaling S = aR of matrices is given in terms of the matrix components by

$$c_{ij} = a_{ij} + b_{ij}, s_{ij} = ar_{ij}$$

1.2. Measurement of mappings

From the above it is apparent that linear mappings and matrices can also be considered as data, and a first step in analysis of such data is definition of functionals that would attach a single scalar label to each linear mapping of matrix. Of particular interest is the definition of a norm functional that characterizes in an appropriate sense the size of a linear mapping.

Consider first the case of finite matrices with real components $A \in \mathbb{R}^{m \times n}$ that represent linear mappings between real vector spaces $f : \mathbb{R}^m \to \mathbb{R}^n$. The columns a_1, \ldots, a_n of $A \in \mathbb{R}^{m \times n}$ could be placed into a single column vector c with mn components

$$\boldsymbol{c} = \begin{bmatrix} \boldsymbol{a}_1 \\ \vdots \\ \boldsymbol{a}_n \end{bmatrix}.$$

Subsequently the norm of the matrix A could be defined as the norm of the vector c. An example of this approach is the Frobenius norm

$$\|\boldsymbol{A}\|_{F} = \|\boldsymbol{c}\|_{2} = \left(\sum_{i=1}^{m} \sum_{j=1}^{n} |a_{ij}|^{2}\right)^{1/2}.$$

A drawback of the above approach is that the structure of the matrix and its close relationship to a linear mapping is lost. A more useful characterization of the size of a mapping is to consider the amplification behavior of linear mapping. The motivation is readily understood starting from linear mappings between the reals $f: \mathbb{R} \to \mathbb{R}$, that are of the form f(x) = ax. When given an argument of unit magnitude |x| = 1, the mapping returns a real number with magnitude |a|. For mappings $f: \mathbb{R}^2 \to \mathbb{R}^2$ within the plane, arguments that satisfy $||\mathbf{x}||_2 = 1$ are on the unit circle with components $\mathbf{x} = [\cos \theta \sin \theta]$ have images through f given analytically by

$$f(\mathbf{x}) = A\mathbf{x} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 \end{bmatrix} \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} = \cos \theta \mathbf{a}_1 + \sin \theta \mathbf{a}_2$$

and correspond to ellipses.



Figure 1. Mapping of unit circle by f(x) = Ax, $A = \begin{bmatrix} 2 & -1 \\ 3 & 1 \end{bmatrix}$.

: $n=250$; $h=2.0*\pi/n$; $\theta=h*(1:n)$; $c=cos.(\theta)$; $s=sin.(\theta)$;	
∴ a1=[2; 3]; a2=[-1; 1]; A=[a1 a2]	
$\left[\begin{array}{cc} 2 & -1 \\ 3 & 1 \end{array}\right]$	(1)
: fx = c.*a1[1]+s.*a2[1]; fy = c.*a1[2]+s.*a2[2];	
<pre> clf(); grid("on"); plot(c,s); axis("equal");</pre>	

∴ plot(fx, fy, "r");
∴ F=svd(A); U=F.U; Σ=Diagonal(F.S); Vt=F.Vt; V=Vt';
:. $\sigma 1 = \Sigma [1, 1]; \sigma 2 = \Sigma [2, 2];$
∴ z=[0; 0]; u1=o1*[z U[:,1]]; u2=o2*[z U[:,2]];
∴ v1=[z V[:,1]]; v2=[z V[:,2]];
<pre> cd(homedir()*"/courses/MATH661/images");</pre>
∴ plot(u1[1,:],u1[2,:],"r");
∴ plot(u2[1,:],u2[2,:],"r");
∴ plot(v1[1,:],v1[2,:],"b");
∴ plot(v2[1,:],v2[2,:],"b");
∴ savefig("L08Fig01.eps")

From the above the mapping associated A amplifies some directions more than others. This suggests a definition of the size of a matrix or a mapping by the maximal amplification unit norm vectors within the domain.

DEFINITION. For vector spaces U, V with norms $\| \|_U : U \to \mathbb{R}_+, \| \|_V : V \to \mathbb{R}_+$, the induced norm of $f: U \to V$ is

$$\|f\| = \sup_{\|x\|_U=1} \|f(x)\|_V.$$

DEFINITION. For vector spaces \mathbb{R}^n , \mathbb{R}^m with norms $\|\|^{(n)}$: $U \to \mathbb{R}_+$, $\|\|^{(m)}$: $V \to \mathbb{R}_+$, the induced norm of matrix $A \in \mathbb{R}^{m \times n}$ is

$$\|A\| = \sup_{\|x\|^{(n)}=1} \|Ax\|^{(m)}.$$

In the above, any vector norm can be used within the domain and codomain.

2. The Singular Value Decomposition (SVD)

The fundamental theorem of linear algebra partitions the domain and codomain of a linear mapping $f: U \to V$. 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

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

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.

2.1. Orthogonal matrices

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

$$\boldsymbol{U} = [\boldsymbol{u}_1 \ \boldsymbol{u}_2 \ \dots \ \boldsymbol{u}_m] \in \mathbb{R}^{m \times m}, \boldsymbol{V} = [\boldsymbol{v}_1 \ \boldsymbol{v}_2 \ \dots \ \boldsymbol{v}_n] \in \mathbb{R}^{n \times n}$$

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 $u_i^T u_j = 1$. Such behavior is concisely expressed as a matrix product

$$\boldsymbol{U}^T \boldsymbol{U} = \boldsymbol{I}_m,$$

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

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

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

$$\boldsymbol{U}\boldsymbol{X} = \boldsymbol{U} \left[\begin{array}{cccc} \boldsymbol{x}_1 & \boldsymbol{x}_2 & \dots & \boldsymbol{x}_m \end{array} \right] = \left[\begin{array}{ccccc} \boldsymbol{e}_1 & \boldsymbol{e}_2 & \dots & \boldsymbol{e}_m \end{array} \right].$$

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

$$\boldsymbol{U}\boldsymbol{x}_{j} = \boldsymbol{e}_{j} \Rightarrow \boldsymbol{U}^{T} \boldsymbol{U}\boldsymbol{x}_{j} = \boldsymbol{U}^{T} \boldsymbol{e}_{j} \Rightarrow \boldsymbol{I}_{m} \boldsymbol{x}_{j} = \begin{bmatrix} \boldsymbol{u}_{1}^{T} \\ \boldsymbol{u}_{2}^{T} \\ \vdots \\ \boldsymbol{u}_{m}^{T} \end{bmatrix} \boldsymbol{e}_{j} \Rightarrow \boldsymbol{x}_{j} = (\boldsymbol{U}^{T})_{j},$$

where $(U^T)_i$ is the j^{th} column of U^T , hence $X = U^T$, leading to

$$\boldsymbol{U}^T \boldsymbol{U} = \boldsymbol{I} = \boldsymbol{U} \boldsymbol{U}^T$$

Note that the second equality

$$\begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix} \begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix}^T = \begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix} \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} = u_1 u_1^T + u_2 u_2^T + \dots + u_m u_m^T = I$$

acts as normalization condition on the matrices $U_i = u_i u_i^T$.

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

2.2. Intrinsic basis of a linear mapping

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

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

Imposing the condition that U, V are orthogonal leads to

$$Us = y \Rightarrow s = U^T y, Vr = x \Rightarrow r = V^T x,$$

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

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

From the above the orthogonal bases U, V and scaling coefficients Σ that are sought must satisfy $A = U \Sigma V^T$.

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

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

with properties:

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

Proof. The proof of the SVD makes use of properties of the norm, concepts from analysis and complete induction. Adopting the 2-norm set $\sigma_1 = ||A||_2$,

$$\sigma_1 = \sup_{\|\bm{x}\|_2 = 1} \|\bm{A}\,\bm{x}\|_2$$

The domain $\|\mathbf{x}\|_2 = 1$ is compact (closed and bounded), and the extreme value theorem implies that $f(\mathbf{x}) = A\mathbf{x}$ attains its maxima and minima, hence there must exist some vectors $\mathbf{u}_1, \mathbf{v}_1$ of unit norm such that $\sigma_1 \mathbf{u}_1 = A \mathbf{v}_1 \Rightarrow \sigma_1 = \mathbf{u}_1^T A \mathbf{v}_1$. Introduce orthogonal bases $\mathbf{U}_1, \mathbf{V}_1$ for $\mathbb{R}^m, \mathbb{R}^n$ whose first column vectors are $\mathbf{u}_1, \mathbf{v}_1$, and compute

$$\boldsymbol{U}_1^T \boldsymbol{A} \boldsymbol{V}_1 = \begin{bmatrix} \boldsymbol{u}_1^T \\ \vdots \\ \boldsymbol{u}_m^T \end{bmatrix} \begin{bmatrix} \boldsymbol{A} \boldsymbol{v}_1 & \dots & \boldsymbol{A} \boldsymbol{v}_n \end{bmatrix} = \begin{bmatrix} \sigma_1 & \boldsymbol{w}^T \\ \boldsymbol{0} & \boldsymbol{B} \end{bmatrix} = \boldsymbol{C}.$$

In the above \mathbf{w}^T is a row vector with n-1 components $\mathbf{u}_1^T \mathbf{A} \mathbf{v}_j$, j = 2, ..., n, and $\mathbf{u}_i^T \mathbf{A} \mathbf{v}_1$ must be zero for \mathbf{u}_1 to be the direction along which the maximum norm $\|\mathbf{A}\mathbf{v}_1\|$ is obtained. Introduce vectors

$$\mathbf{y} = \begin{bmatrix} \sigma_1 \\ \mathbf{w} \end{bmatrix}, \mathbf{z} = \mathbf{C} \, \mathbf{y} = \begin{bmatrix} \sigma_1^2 + \mathbf{w}^T \mathbf{w} \\ \mathbf{B} \mathbf{w} \end{bmatrix}$$

and $\|C\mathbf{y}\|_2 = \|\mathbf{z}\|_2 \ge \sigma_1^2 + \mathbf{w}^T \mathbf{w} + \|\mathbf{B}\mathbf{w}\|_1 \ge \sigma_1^2 + \mathbf{w}^T \mathbf{w} = \|\mathbf{y}\|_2^2 = \sqrt{\sigma_1^2 + \mathbf{w}^T \mathbf{w}} \|\mathbf{y}\|_2$. From $\|U_1^T A V_1\| = \|A\| = \sigma_1 = \|C\| \ge \sigma_1^2 + \mathbf{w}^T \mathbf{w}$ it results that $\mathbf{w} = \mathbf{0}$. By induction, assume that \mathbf{B} has a singular value decomposition, $\mathbf{B} = U_2 \Sigma_2 V_2^T$, such that

$$\boldsymbol{U}_1^T \boldsymbol{A} \boldsymbol{V}_1 = \begin{bmatrix} \boldsymbol{\sigma}_1 & \boldsymbol{0}^T \\ \boldsymbol{0} & \boldsymbol{U}_2 \boldsymbol{\Sigma}_2 \boldsymbol{V}_2^T \end{bmatrix} = \begin{bmatrix} 1 & \boldsymbol{0}^T \\ \boldsymbol{0} & \boldsymbol{U}_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\sigma}_1 & \boldsymbol{0}^T \\ \boldsymbol{0} & \boldsymbol{\Sigma}_2 \end{bmatrix} \begin{bmatrix} 1 & \boldsymbol{0}^T \\ \boldsymbol{0} & \boldsymbol{\Sigma}_2^T \end{bmatrix},$$

and the orthogonal matrices arising in the singular value decomposition of A are

$$\boldsymbol{U} = \boldsymbol{U}_1 \begin{bmatrix} \mathbf{1} & \mathbf{0}^T \\ \mathbf{0} & \boldsymbol{U}_2 \end{bmatrix}, \boldsymbol{V}^T = \begin{bmatrix} \mathbf{1} & \mathbf{0}^T \\ \mathbf{0} & \boldsymbol{V}_2^T \end{bmatrix} \boldsymbol{V}_1^T.$$

The scaling coefficients σ_j are called the *singular values* of A. The columns of U are called the *left singular vectors*, and those of V are called the *right singular vectors*.

The fact that the scaling coefficients are norms of A and submatrices of A, $\sigma_1 = ||A||$, is crucial importance in applications. Carrying out computation of the matrix products

leads to a representation of A as a sum

$$\boldsymbol{A} = \sum_{i=1}^{r} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^T, r \leq \min(m, n).$$
$$\boldsymbol{A} = \sigma_1 \boldsymbol{u}_1 \boldsymbol{v}_1^T + \sigma_2 \boldsymbol{u}_2 \boldsymbol{v}_2^T + \dots + \sigma_r \boldsymbol{u}_r \boldsymbol{v}_r^T$$

Each product $u_i v_i^T$ 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 A

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

In very many cases the singular values exhibit rapid, exponential decay, $\sigma_1 \gg \sigma_2 \gg \cdots$, such that the approximation above is an accurate representation of the matrix A.



Figure 2. Successive SVD approximations of Andy Warhol's painting, Marilyn Diptych (~1960), with k = 10, 20, 40 rank-one updates.

3. SVD solution of linear algebra problems

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

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

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

Best 2-norm approximation. In the above A was assumed to be a basis, hence $r = \operatorname{rank}(A) = m$. If columns of A do not form a basis, r < m, then $b \in \mathbb{R}^m$ might not be reachable by linear combinations within C(A). The closest vector to b in the 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, ..., r.

The pseudo-inverse. From the above, finding either the solution of Ax = Ib or the best approximation possible if A is not of full rank, can be written as a sequence of matrix multiplications using the SVD

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

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

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

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