

EECS 275 Matrix Computation

Ming-Hsuan Yang

Electrical Engineering and Computer Science
University of California at Merced
Merced, CA 95344
<http://faculty.ucmerced.edu/mhyang>



Lecture 5

Overview

- Matrix properties via singular value decomposition (SVD)
- Geometric interpretation of SVD
- Applications

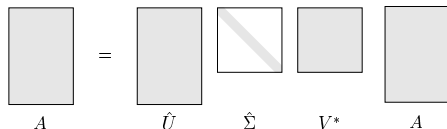
Reading

- Chapter 5 of *Numerical Linear Algebra* by Lloyd Trefethen and David Bau
- Chapter 3 of *Matrix Computations* by Gene Golub and Charles Van Loan
- Chapter 3 of *Mathematical Modeling of Continuous Systems* by Carlo Tomasi
- Chapter 5 of *Matrix Analysis and Applied Linear Algebra* by Carl Meyer

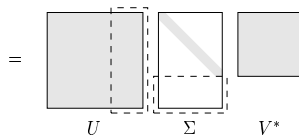
Full and reduced SVD

- Let $A \in \mathbb{R}^{m \times n}$
- Reduced SVD: $A = \hat{U} \hat{\Sigma} \hat{V}^\top$, $\hat{U} \in \mathbb{R}^{m \times n}$, $\hat{\Sigma} \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times n}$
- Full SVD: $A = U \Sigma V^\top$, $U \in \mathbb{R}^{m \times m}$, $\Sigma \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times n}$

Reduced SVD ($m \geq n$)



Full SVD ($m \geq n$)



Uniqueness

- First note that σ_1 and \mathbf{v}_1 can be uniquely determined by $\|A\|_2$
- Suppose in addition to \mathbf{v}_1 , there is another linearly independent vector \mathbf{w} with $\|\mathbf{w}\|_2 = 1$ and $\|A\mathbf{w}\|_2 = \sigma_1$
- Define a unit vector \mathbf{v}_2 , orthogonal to \mathbf{v}_1 as a linear combination of \mathbf{v}_1 and \mathbf{w}

$$\mathbf{v}_2 = \frac{\mathbf{w} - (\mathbf{v}_1^\top \mathbf{w})\mathbf{v}_1}{\|\mathbf{w} - (\mathbf{v}_1^\top \mathbf{w})\mathbf{v}_1\|_2}$$

- Since $\|A\|_2 = \sigma_1$, $\|A\mathbf{v}_2\|_2 \leq \sigma_1$, but this must be an equality, for otherwise $\mathbf{w} = c\mathbf{v}_1 + s\mathbf{v}_2$ for some constants c and s with $|c|^2 + |s|^2 = 1$, we would have $\|A\mathbf{w}\| < \sigma_1$
- \mathbf{v}_2 is a second right singular vector of A corresponding to σ_1
- Once σ_1 , \mathbf{v}_1 , and \mathbf{v}_2 are determined, the remainder of SVD is determined by the action of A on the space orthogonal to \mathbf{v}_1
- Since \mathbf{v}_1 is unique up to a sign, the orthogonal space is uniquely defined and so are the remaining singular values

Matrix properties via SVD

Theorem

The rank of A is r , the number of nonzero singular values.

Proof.

The rank of a diagonal matrix is equal to the number of its nonzero entries, and in SVD, $A = U\Sigma V^T$ where U and V are of full rank. Thus, $\text{rank}(A) = \text{rank}(\Sigma) = r$ □

Theorem

$$\|A\|_2 = \sigma_1, \text{ and } \|A\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_r^2}$$

Proof.

As U and V are orthogonal, $A = U\Sigma V^T$, $\|A\|_2 = \|\Sigma\|_2$. By definition, $\|\Sigma\|_2 = \max_{\|x\|=1} \|\Sigma x\|_2 = \max\{|\sigma_i|\} = \sigma_1$. Likewise, $\|A\|_F = \|\Sigma\|_F$, and by definition $\|\Sigma\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_r^2}$ □

Eigenvalue decomposition

- From linear algebra, $A\mathbf{x} = \lambda\mathbf{x}$, λ is an eigenvalue, and \mathbf{x} is an eigenvector
- For m eigenvectors,

$$A[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_m \end{bmatrix}$$

and

$$AX = X\Lambda$$

where Λ is an $m \times m$ diagonal matrix whose entries are the eigenvalues of A , and $X \in \mathbb{R}^{m \times m}$ contains linearly independent eigenvector of A

- The eigenvalue decomposition of A

$$A = X\Lambda X^{-1}$$

SVD and eigenvalue decomposition

- SVD uses two different bases (the sets of left and right singular vectors), whereas the eigenvalue decomposition uses just one (eigenvectors)
- SVD uses orthonormal bases, whereas the eigenvalue decomposition uses a basis that generically is not orthogonal
- Not all matrices have an eigenvalue decomposition, but all matrices have a SVD

Matrix properties via SVD (cont'd)

Theorem

The nonzero singular values of A are the square roots of the nonzero eigenvalues of AA^\top or $A^\top A$ (they have the same nonzero eigenvalues).

Proof.

From definition,

$$AA^\top = (U\Sigma V^\top)(U\Sigma V^\top)^\top = U\Sigma V^\top V\Sigma U^\top = U \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2) U^\top \quad \square$$

Theorem

For $A \in \mathbb{R}^{m \times m}$, $|\det(A)| = \prod_{i=1}^m \sigma_i$

Proof.

$$|\det(A)| = |\det(U\Sigma V^\top)| = |\det(U)| |\det(\Sigma)| |\det(V^\top)| = |\det(\Sigma)| = \prod_{i=1}^m \sigma_i \quad \square$$

Low-rank approximation

Theorem

(Eckart-Young 1936) Let $A = U\Sigma V^\top = U \operatorname{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) V^\top$. For any ν with $0 \leq \nu \leq r$, $A_\nu = \sum_{i=1}^\nu \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$,

$$\|A - A_\nu\|_2 = \min_{\operatorname{rank}(B) \leq \nu} \|A - B\|_2 = \sigma_{\nu+1}$$

Proof.

Suppose there is some B with $\operatorname{rank}(B) \leq \nu$ such that $\|A - B\|_2 < \|A - A_\nu\|_2 = \sigma_{\nu+1}$. Then there exists an $(n - \nu)$ -dimensional subspace $W \in \mathbb{R}^n$ such that $\mathbf{w} \in W \Rightarrow B\mathbf{w} = 0$. Then

$$\|A\mathbf{w}\|_2 = \|(A - B)\mathbf{w}\|_2 \leq \|A - B\|_2 \|\mathbf{w}\|_2 < \sigma_{\nu+1} \|\mathbf{w}\|_2$$

Thus W is a $(n - \nu)$ -dimensional subspace where $\|A\mathbf{w}\| < \sigma_{\nu+1} \|\mathbf{w}\|$. But there is a $(\nu + 1)$ -dimensional subspace where $\|A\mathbf{w}\| \geq \sigma_{\nu+1} \|\mathbf{w}\|$, namely the space spanned by the first $\nu + 1$ right singular vector of A . Since the sum of the dimensions of these two spaces exceeds n , there must be a nonzero vector lying in both, and this is a contradiction. □

Low-rank approximation

Theorem

A is the sum of r rank one matrices: $A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$

Theorem

(Eckart-Young 1936) Let $A = U \Sigma V^\top = U \operatorname{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) V^\top$. For any ν with $0 \leq \nu \leq r$, $A_\nu = \sum_{i=1}^\nu \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$,

$$\|A - A_\nu\|_2 = \min_{\operatorname{rank}(B) \leq \nu} \|A - B\|_2 = \sigma_{\nu+1}$$

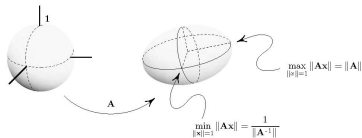
Proof.

Let $\Sigma_\nu = U(A - A_\nu)V^\top$, then

$$\begin{aligned}\Sigma_\nu &= U (\operatorname{diag}(\sigma_1, \dots, \sigma_\nu, \sigma_{\nu+1}, \dots, \sigma_p) - \operatorname{diag}(\sigma_1, \dots, \sigma_\nu, 0, \dots, 0)) V^\top \\ &= U \operatorname{diag}(0, \dots, 0, \sigma_{\nu+1}, \dots, \sigma_p) V^\top\end{aligned}$$

consequently $\|A - A_\nu\|_2 = \|\Sigma_\nu\|_2 = \sigma_{\nu+1}$. □

Geometric interpretation of Eckart-Young theorem



- What is the best approximation of a hyperellipsoid by a line segment?
 - ▶ Take the line segment to be the longest axis
- Next, what is the best approximation by a two-dimensional ellipsoid?
 - ▶ Take the ellipsoid spanned by the longest and the second longest axis
- Continue and improve the approximation by adding into our approximation the largest axis of the hyperellipsoid not yet included
- Reminiscent of techniques used in image compression, machine learning, and functional analysis (e.g., matching pursuit)

Theorem

For any ν with $0 \leq \nu \leq r$, $A_\nu = \sum_{i=1}^{\nu} \sigma_i \mathbf{u}_i \mathbf{v}_i^T$,

$$\|A - A_\nu\|_F = \min_{\text{rank}(B) \leq \nu} = \sqrt{\sigma_{\nu+1}^2 + \cdots + \sigma_r^2}$$

Sensitivity of square systems

- If

$$A = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = U \Sigma V^\top$$

is the SVD of A , then

$$\mathbf{x} = A^{-1} \mathbf{b} = (U \Sigma V^\top)^{-1} \mathbf{b} = \sum_{i=1}^n \frac{\mathbf{u}_i^\top \mathbf{b}}{\sigma_i} \mathbf{v}_i$$

- Small changes in A or \mathbf{b} can induce relatively large changes in \mathbf{x} if σ_n is small
- The magnitude of σ_n has bearing on the sensitivity of the $A\mathbf{x} = \mathbf{b}$ problem
- The solution \mathbf{x} is increasingly sensitive to perturbations

Condition

- Consider the parameterized system

$$(A + \varepsilon F)\mathbf{x}(\varepsilon) = \mathbf{b} + \varepsilon \mathbf{f} \quad \mathbf{x}(0) = \mathbf{x}$$

where $F \in \mathbb{R}^{n \times n}$ and $\mathbf{f} \in \mathbb{R}^n$

- If A is nonsingular, then $\mathbf{x}(\varepsilon)$ is differentiable in a neighborhood of zero
- Moreover, $\dot{\mathbf{x}} = A^{-1}(\mathbf{f} - F\mathbf{x})$ and the Taylor series expansion

$$\mathbf{x}(\varepsilon) = \mathbf{x} + \varepsilon \dot{\mathbf{x}}(0) + O(\varepsilon^2)$$

- Using any vector norm

$$\frac{\|\mathbf{x}(\varepsilon) - \mathbf{x}\|}{\|\mathbf{x}\|} \leq |\varepsilon| \|A^{-1}\| \left\{ \frac{\|\mathbf{f}\|}{\|\mathbf{x}\|} + \|F\| \right\} + O(\varepsilon^2)$$

Condition number

- For square matrices A , define the condition number by

$$\kappa(A) = \|A\| \|A^{-1}\|$$

with the convention that $\kappa(A) = \infty$ for singular A

- Using the inequality $\|\mathbf{b}\| \leq \|A\| \|\mathbf{x}\|$ it follows that

$$\frac{\|\mathbf{x}(\varepsilon) - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \kappa(A)(\rho_A + \rho_b) + O(\varepsilon^2)$$

where

$$\rho_A = |\varepsilon| \frac{\|F\|}{\|A\|} \quad \text{and,} \quad \rho_b = |\varepsilon| \frac{\|f\|}{\|\mathbf{b}\|}$$

represent the relative errors in A and \mathbf{b}

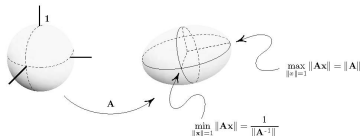
- The relative error in \mathbf{x} is $\kappa(A)$ times the relative error in A and \mathbf{b}
- The condition $\kappa(A)$ quantifies the sensitivity of the $A\mathbf{x} = \mathbf{b}$ problem

Condition number (cont'd)

- Note that $\kappa(\cdot)$ depends on the underlying norm

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_1(A)}{\sigma_n(A)}$$

- $\kappa_2(A)$ measures the elongation of the hyperellipsoid $\{A\mathbf{x} : \|\mathbf{x}\|_2 = 1\}$



- If $\kappa(A)$ is large, then A is said to be an ill-conditioned matrix
- $\kappa_\alpha(\cdot)$ and $\kappa_\beta(\cdot)$ on $\mathbb{R}^{n \times n}$ are equivalent if constants c_1 and c_2 can be found such that $c_1 \kappa_\alpha(A) \leq \kappa_\beta(A) \leq c_2 \kappa_\alpha(A)$, e.g.,

$$\begin{aligned} \frac{1}{n} \kappa_2(A) &\leq \kappa_1(A) \leq n \kappa_2(A) \\ \frac{1}{n} \kappa_\infty(A) &\leq \kappa_2(A) \leq n \kappa_\infty(A) \\ \frac{1}{n^2} \kappa_1(A) &\leq \kappa_\infty(A) \leq n^2 \kappa_1(A) \end{aligned}$$

- For any p -norm, we have $\kappa(A) \geq 1$, and matrices with small conditional number are said to be well-conditioned

Minimum norm least square solution

Theorem

The minimum norm least squares solution to a linear system $A\mathbf{x} = \mathbf{b}$, that is, the shortest vector \mathbf{x} that achieves $\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|$ is unique, and is given by

$$\hat{\mathbf{x}} = V\Sigma^\dagger U^\top \mathbf{b}$$

where

$$\Sigma^\dagger = \begin{bmatrix} 1/\sigma_1 & & & 0 & \cdots & 0 \\ & \ddots & & & & \\ & & 1/\sigma_r & & & \\ & & & 0 & & \\ & & & & \ddots & \\ & & & & & 0 & 0 & \cdots & 0 \end{bmatrix}$$

- The matrix $A^\dagger = V\Sigma^\dagger U^\top$ is the pseudoinverse of A

Minimum norm solution

- The minimum norm solution to $A\mathbf{x} = \mathbf{b}$ is the vector that minimizes $\|A\mathbf{x} - \mathbf{b}\|$,

$$\|U\Sigma V^T \mathbf{x} - \mathbf{b}\| = \|U(\Sigma V^T \mathbf{x} - U^T \mathbf{b})\|$$

and so it is equivalent to solve $\|\Sigma V^T \mathbf{x} - U^T \mathbf{b}\|$

- Let $\mathbf{y} = V^T \mathbf{x}$ and $\mathbf{c} = U^T \mathbf{b}$, it becomes

$$\|\Sigma \mathbf{y} - \mathbf{c}\|$$

$$\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \ddots & \cdots & 0 \\ & & \sigma_r & \\ \vdots & & & 0 \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_r \\ y_{r+1} \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} c_1 \\ \vdots \\ c_r \\ c_{r+1} \\ \vdots \\ c_m \end{bmatrix}$$

Minimum norm solution (cont'd)

- The optimal \mathbf{y} has the following components

$$\begin{aligned} y_i &= \frac{c_i}{\sigma_i} & \text{for } i = 1, \dots, r \\ y_i &= 0 & \text{for } i = r + 1, \dots, n \end{aligned}$$

- In vector form

$$\mathbf{y} = \Sigma^\dagger \mathbf{c}$$

- Notice there is no other choice for \mathbf{y} , which is therefore unique: minimum residual forces the choice of y_1, \dots, y_r , and minimum norm solution forces the other entries of \mathbf{y}
- The minimum norm least squares solution is

$$\hat{\mathbf{x}} = V\mathbf{y} = V\Sigma^\dagger \mathbf{c} = V\Sigma^\dagger U^\top \mathbf{b}$$

- The residual is

$$\|A\mathbf{x} - \mathbf{b}\|^2 = \|\Sigma\mathbf{y} - \mathbf{c}\|^2 = \sum_{i=r+1}^m c_i^2 = \sum_{i=r+1}^m (\mathbf{u}_i^\top \mathbf{b})^2$$

which is the projection of \mathbf{b} onto the complement of the range of A

Least squares solution of homogeneous linear systems

Theorem

For $A\mathbf{x} = \mathbf{0}$ or $\min_{\|\mathbf{x}\|=1} \|A\mathbf{x}\|$. Let $A = U\Sigma V^T$, the solution is

$$\mathbf{x} = \alpha_1 \mathbf{v}_{n-k+1} + \dots + \alpha_k \mathbf{v}_n$$

where k is the largest integer such that

$$\sigma_{n-k+1} = \dots = \sigma_n, \text{ and } \alpha_1^2 + \dots + \alpha_k^2 = 1$$

Proof.

Consider the unit-norm least square solution

$$\|A\mathbf{x}\| = \|U\Sigma V^T \mathbf{x}\| \equiv \|\Sigma V^T \mathbf{x}\| = \|\Sigma \mathbf{y}\|$$

where $\mathbf{y} = V^T \mathbf{x}$. Thus the unit norm vector \mathbf{y} that minimizes the norm

$$\sigma_1^2 y_1^2 + \dots + \sigma_n^2 y_n^2$$

which is achieved by concentrating all the mass of \mathbf{y} w.r.t smallest σ

$$y_1 = \dots = y_{n-k} = 0$$

and thus $\mathbf{x} = V\mathbf{y} = y_1 \mathbf{v}_1 + \dots + y_{n-k+1} \mathbf{v}_{n-k+1} + \dots + y_n \mathbf{v}_n$ and

$$\alpha_1 = y_{n-k+1}, \dots, \alpha_k = y_n$$