

# Lecture 1: Matrix Decompositions (Chapter 4 of Textbook A)

Jinwoo Shin

AI503: Mathematics for AI

This lecture slide is based upon https://yung-web.github.io/home/courses/mathml.html (made by Prof. Yung Yi, KAIST EE)

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

## Summary



- How to summarize matrices: determinants and eigenvalues
- How matrices can be decomposed: Cholesky decomposition, diagonalization, singular value decomposition
- How these decompositions can be used for matrix approximation

# Roadmap



#### (1) Determinant and Trace

- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

# Determinant: Motivation (1)



• For 
$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$
,  $\mathbf{A}^{-1} = \frac{1}{a_{11}a_{22}-a_{12}a_{21}} \begin{pmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{pmatrix}$ .

- **A** is invertible iff  $a_{11}a_{22} a_{12}a_{21} \neq 0$
- Let's define  $det(\mathbf{A}) = a_{11}a_{22} a_{12}a_{21}$ .
- Notation: det(**A**) or |whole matrix|
- What about  $3 \times 3$  matrix? By doing some algebra (e.g., Gaussian elimination),

 $\begin{array}{cccc} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{array} = a_{11}a_{22}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} \\ \end{array}$ 

# Determinant: Motivation (2)

KAIST A

• Try to find some pattern ...

 $\begin{aligned} &a_{11}a_{22}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} \\ &- a_{31}a_{22}a_{13} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33} = \\ &a_{11}(-1)^{1+1}\det(\boldsymbol{A}_{1,1}) + a_{12}(-1)^{1+2}\det(\boldsymbol{A}_{1,2}) \\ &+ a_{13}(-1)^{1+3}\det(\boldsymbol{A}_{1,3}) \end{aligned}$ 

-  $A_{k,j}$  is the submatrix of A that we obtain when deleting row k and column j.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$
 gives the term  $a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix}$ 
$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$
 gives the term  $a_{12} \begin{pmatrix} - \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{32} \end{vmatrix}$ 
$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$
 gives the term  $a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$ 

source: www.cliffsnotes.com

- This is called Laplace expansion.
- Now, we can generalize this and provide the formal definition of determinant.

## **Determinant: Formal Definition**

#### Determinant

For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , for all  $j = 1, \ldots, n$ ,

- 1. Expansion along column *j*: det( $\mathbf{A}$ ) =  $\sum_{k=1}^{n} (-1)^{k+j} a_{kj} \det(\mathbf{A}_{k,j})$
- 2. Expansion along row j:  $det(\mathbf{A}) = \sum_{k=1}^{n} (-1)^{k+j} a_{jk} det(\mathbf{A}_{j,k})$
- All expansion are equal, so no problem with the definition.
- Theorem. det( $\mathbf{A}$ )  $\neq$  0  $\iff$  rk( $\mathbf{A}$ ) =  $n \iff \mathbf{A}$  is invertible.

## **Determinant:** Properties



- (1)  $det(\boldsymbol{AB}) = det(\boldsymbol{A}) det(\boldsymbol{B})$
- (2)  $det(\boldsymbol{A}) = det(\boldsymbol{A}^{\mathsf{T}})$
- (3) For a regular  $\boldsymbol{A}$ , det $(\boldsymbol{A}^{-1}) = 1/\det(\boldsymbol{A})$
- (4) For two similar matrices  $\boldsymbol{A}, \boldsymbol{A}'$  (i.e.,  $\boldsymbol{A}' = \boldsymbol{S}^{-1} \boldsymbol{A} \boldsymbol{S}$  for some  $\boldsymbol{S}$ ), det $(\boldsymbol{A}) = \det(\boldsymbol{A}')$
- (5) For a triangular matrix<sup>1</sup>  $\boldsymbol{T}$ , det( $\boldsymbol{T}$ ) =  $\prod_{i=1}^{n} T_{ii}$
- (6) Adding a multiple of a column/row to another one does not change det(A)
- (7) Multiplication of a column/row with  $\lambda$  scales det(**A**): det( $\lambda$ **A**) =  $\lambda$ <sup>n</sup>**A**
- (8) Swapping two rows/columns changes the sign of det(A)
  - Using (5)-(8), Gaussian elimination (reaching a triangular matrix) enables to compute the determinant.

<sup>&</sup>lt;sup>1</sup>This includes diagonal matrices.



• Definition. The trace of a square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is defined as

$$\mathsf{tr}(\boldsymbol{A}) := \sum_{i=1}^n a_{ii}$$

• 
$$tr(\boldsymbol{A} + \boldsymbol{B}) = tr(\boldsymbol{A}) + tr(\boldsymbol{B})$$

- $tr(\alpha A) = \alpha tr(A)$
- $tr(I_n) = n$



- tr(AB) = tr(BA) for  $A \in \mathbb{R}^{n \times k}$  and  $B \in \mathbb{R}^{k \times n}$
- tr(AKL) = tr(KLA), for  $A \in \mathbb{R}^{a \times k}$ ,  $K \in \mathbb{R}^{k \times l}$ ,  $L \in \mathbb{R}^{l \times a}$
- $tr(xy^{\mathsf{T}}) = tr(y^{\mathsf{T}}x) = y^{\mathsf{T}}x \in \mathbb{R}$
- A linear mapping Φ : V → V, represented by a matrix A and another matrix B.
   A and B use different bases, where B = S<sup>-1</sup>AS

$$\operatorname{tr}(\boldsymbol{B}) = \operatorname{tr}(\boldsymbol{S}^{-1}\boldsymbol{A}\boldsymbol{S}) = \operatorname{tr}(\boldsymbol{A}\boldsymbol{S}\boldsymbol{S}^{-1}) = \operatorname{tr}(\boldsymbol{A})$$

 Message. While matrix representations of linear mappings are basis dependent, but their traces are not.



• Definition. For  $\lambda \in \mathbb{R}$  and a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the characteristic polynomial of  $\mathbf{A}$  is defined as:

,

$$p_{\boldsymbol{A}}(\lambda) := \det(\boldsymbol{A} - \lambda \boldsymbol{I})$$

$$= c_0 + c_1 \lambda + c_2 \lambda^2 + \dots + c_{n-1} \lambda^{n-1} + (-1)^n \lambda^n,$$
where  $c_0 = \det(\boldsymbol{A})$  and  $c_{n-1} = (-1)^{n-1} \operatorname{tr}(\boldsymbol{A}).$ 
• Example. For  $\boldsymbol{A} = \begin{pmatrix} 4 & 2 \\ 1 & 3 \end{pmatrix},$ 

$$p_{\boldsymbol{A}}(\lambda) = \begin{vmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{vmatrix} = (4 - \lambda)(3 - \lambda) - 2 \cdot 1$$

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny



Definition. Consider a square matrix *A* ∈ ℝ<sup>n×n</sup>. Then, λ ∈ ℝ is an eigenvalue of *A* and *x* ∈ ℝ<sup>n</sup> \ {0} is the corresponding eigenvector of *A* if

$$Ax = \lambda x$$

- Equivalent statements
  - $\circ \lambda$  is an eigenvalue.
  - $(\mathbf{A} \lambda \mathbf{I}_n)\mathbf{x} = 0$  can be solved non-trivially, i.e.,  $\mathbf{x} \neq 0$ .
  - $\mathsf{rk}(\boldsymbol{A} \lambda \boldsymbol{I}_n) < n.$
  - $\det(\mathbf{A} \lambda \mathbf{I}_n) = 0 \iff$  The characteristic polynomial  $p_{\mathbf{A}}(\lambda) = 0$ .

#### Example



• For 
$$\mathbf{A} = \begin{pmatrix} 4 & 2 \\ 1 & 3 \end{pmatrix}$$
,  $p_{\mathbf{A}}(\lambda) = \begin{vmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{vmatrix} = (4 - \lambda)(3 - \lambda) - 2 \cdot 1 = \lambda^2 - 7\lambda + 10$ 

- Eigenvalues  $\lambda = 2$  or  $\lambda = 5$ .
- Eigenvector  $E_5$  for  $\lambda = 5$  $\begin{pmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{pmatrix} \mathbf{x} = 0 \implies \begin{pmatrix} -1 & 2 \\ 1 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0 \implies E_5 = \operatorname{span} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$
- Eigenvector  $E_2$  for  $\lambda = 2$ . Similarly, we get  $E_2 = \text{span}[\begin{pmatrix} 1 \\ -1 \end{pmatrix}]$
- Message. Eigenvectors are not unique.

# Properties (1)



- If x is an eigenvector of A, so are all vectors that are collinear<sup>2</sup>.
- *E<sub>λ</sub>*: the set of all eigenvectors for eigenvalue λ, spanning a subspace of ℝ<sup>n</sup>. We call this eigensapce of *A* for λ.
- $E_{\lambda}$  is the solution space of  $(\mathbf{A} \lambda \mathbf{I})\mathbf{x} = 0$ , thus  $E_{\lambda} = \ker(\mathbf{A} \lambda \mathbf{I})$
- Geometric interpretation
  - The eigenvector corresponding to a nonzero eigenvalue points in a direction stretched by the linear mapping.
  - The eigenvalue is the factor of stretching.
- Identity matrix I: one eigenvalue  $\lambda = 1$  and all vectors  $x \neq 0$  are eigenvectors.

 $<sup>^{2}</sup>$ Two vectors are collinear if they point in the same or the opposite direction.

# Properties (2)



- **A** and  $\mathbf{A}^{\mathsf{T}}$  share the eigenvalues, but not necessarily eigenvectors.
- For two similar matrices A, A' (i.e., A' = S<sup>-1</sup>AS for some S), they possess the same eigenvalues.
  - Meaning: A linear mapping Φ has eigenvalues that are independent of the choice of basis of its transformation matrix.
  - Symmetric, positive definite matrices always have positive, real eigenvalues.

determinant, trace, eigenvalues: all invariant under basis change

# Examples for Geometric Interpretation (1)

- 1.  $A = ( \begin{smallmatrix} \frac{1}{2} & 0 \\ 0 & 2 \end{smallmatrix} )$ , det(A) = 1
  - $\circ \ \lambda_1 = \tfrac{1}{2}, \lambda_2 = 2$
  - eigenvectors: canonical basis vectors
  - area preserving, just vertical horizontal) stretching.

2. 
$$\boldsymbol{A} = (\begin{smallmatrix} 1 & \frac{1}{2} \\ 0 & 1 \end{smallmatrix})$$
,  $\det(\boldsymbol{A}) = 1$   
 $\circ \ \lambda_1 = \lambda_2 = 1$ 

- eigenvectors: colinear over the horiontal line
- area preserving, shearing

3. 
$$\boldsymbol{A} = \begin{pmatrix} \cos(\frac{\pi}{6}) - \sin(\frac{\pi}{6}) \\ \sin(\frac{\pi}{6}) & \cos(\frac{\pi}{6}) \end{pmatrix}$$
, det $(\boldsymbol{A}) = 1$ 

- $\circ~$  Rotation by  $\pi/6$  counter-clockwise
- only complex eigenvalues (no eigenvectors)
- area preserving



# Examples for Geometric Interpretation (2)



4. 
$$\mathbf{A} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$
,  $\det(\mathbf{A}) = 0$   
 $\circ \lambda_1 = 0, \lambda_2 = 2$ 

• Mapping that collapses a 2D onto 1D

• area collapses

5. 
$$\boldsymbol{A} = ( \begin{smallmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{smallmatrix} )$$
,  $\det(\boldsymbol{A}) = 3/4$   
 $\circ \lambda_1 = 0.5, \lambda_2 = 1.5$ 

• area scales by 75%, shearing and stretching



# Properties (3)



- For *A* ∈ ℝ<sup>n×n</sup>, *n* distinct eigenvalues ⇒ eigenvectors are linearly independent, which form a basis of ℝ<sup>n</sup>.
  - Converse is not true.
  - Example of *n* linearly independent eigenvectors for less than *n* eigenvalues???
- Determinant. For (possibly repeated) eigenvalues  $\lambda_i$  of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i$$

• Trace. For (possibly repeated) eigenvalues  $\lambda_i$  of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\mathsf{tr}(oldsymbol{A}) = \sum_{i=1}^n \lambda_i$$

• Message. det(A) is the area scaling and tr(A) is the circumference scaling

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

# LU Decomposition





- The Gaussian elimination is the processing of reaching an upper triangular matrix
- Gaussian elimination: multiplying the matrices corresponding to two elementary operations ((i) row multiplication by *a* and (ii) adding two rows downward)
- The above elementary operations are the low triangular matrices (LTM), and their inverses and their product are all LTMs.

• 
$$(\boldsymbol{E}_{k}\boldsymbol{E}_{k-1}\cdot\boldsymbol{E}_{1})\boldsymbol{A} = \boldsymbol{U} \implies \boldsymbol{A} = \underbrace{(\boldsymbol{E}_{1}^{-1}\cdots\boldsymbol{E}_{k-1}^{-1}\boldsymbol{E}_{k}^{-1})}_{\boldsymbol{L}}\boldsymbol{U}$$



- A real number: decomposition of two identical numbers, e.g.,  $9 = 3 \times 3$
- Theorem. For a symmetric, positive definite matrix  $\mathbf{A}$ ,  $\mathbf{A} = \mathbf{L}\mathbf{L}^{\mathsf{T}}$ , where
  - $\circ~\textbf{\textit{L}}$  is a lower-triangular matrix with positive diagonals
  - Such a *L* is unique, called Cholesky factor of *A*.
- Applications
  - (a) factorization of covariance matrix of a multivariate Gaussian variable
  - (b) linear transformation of random variables
  - (c) fast determinant computation:  $det(\mathbf{A}) = det(\mathbf{L}) det(\mathbf{L}^{\mathsf{T}}) = det(\mathbf{L})^2$ , where  $det(\mathbf{L}) = \prod_i l_{ii}$ . Thus,  $det(\mathbf{A}) = \prod_i l_{ii}^2$ .

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

### **Diagonal Matrix and Diagonalization**



• Diagonal matrix. zero on all off-diagonal elements,  $\boldsymbol{D} = \begin{pmatrix} \boldsymbol{u}_1 & \cdots & \boldsymbol{v} \\ \vdots & & \vdots \\ 0 & \cdots & d_n \end{pmatrix}$ 

$$oldsymbol{D}^k = egin{pmatrix} d_1^k & \cdots & 0 \ dots & & dots \ 0 & \cdots & d_n^k \end{pmatrix}, \quad oldsymbol{D}^{-1} = egin{pmatrix} 1/d_1 & \cdots & 0 \ dots & & dots \ 0 & \cdots & 1/d_n \end{pmatrix}, \quad \det(oldsymbol{D}) = d_1 d_2 \cdots d_n$$

- Definition.  $A \in \mathbb{R}^{n \times n}$  is diagonalizable if it is similar to a diagonal matrix D, i.e.,  $\exists$  an invertible  $P \in \mathbb{R}^{n \times n}$ , such that  $D = P^{-1}AP$ .
- Definition.  $A \in \mathbb{R}^{n \times n}$  is orthogonally diagonalizable if it is similar to a diagonal matrix D, i.e.,  $\exists$  an orthogonal  $P \in \mathbb{R}^{n \times n}$ , such that  $D = P^{-1}AP = P^{\mathsf{T}}AP$ .



- $\boldsymbol{A}^k = \boldsymbol{P} \boldsymbol{D}^k \boldsymbol{P}^{-1}$
- $det(\boldsymbol{A}) = det(\boldsymbol{P}) det(\boldsymbol{D}) det(\boldsymbol{P}^{-1}) = det(\boldsymbol{D}) = \prod_i d_{ii}$
- Many other things ...
- Question. Under what condition is **A** diagonalizable (or orthogonally diagonalizable) and how can we find **P** (thus **D**)?

# Diagonalizablity, Algebraic/Geometric Multiplicity



- Definition. For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with an eigenvalue  $\lambda_i$ ,
  - the algebraic multiplicity  $\alpha_i$  of  $\lambda_i$  is the number of times the root appears in the characteristic polynomial.
  - the geometric multiplicity  $\zeta_i$  of  $\lambda_i$  is the number of linearly independent eigenvectors associated with  $\lambda_i$  (i.e., the dimension of the eigenspace spanned by the eigenvectors of  $\lambda_i$ )
- Example. The matrix  $\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 0 & 2 \end{pmatrix}$  has two repeated eigenvalues  $\lambda_1 = \lambda_2 = 2$ , thus  $\alpha_1 = 2$ . However, it has only one distinct unit eigenvector  $\mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ , thus  $\zeta_1 = 1$ .
- Theorem.  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is diagonalizable  $\iff \sum_{i} \alpha_{i} = \sum_{i} \zeta_{i} = n$ .

# Orthogonally Diagonaliable and Symmetric Matrix

KAIST AI

Theorem.  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is orthogonally diagonalizable  $\iff \mathbf{A}$  is symmetric.

- Question. . How to find **P** (thus **D**)?
- Spectral Theorem. If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric,
  - (a) the eigenvalues are all real
  - (b) the eigenvectors to different eigenvalues are perpendicular.
  - (c) there exists an orthogonal eigenbasis
- For (c), from each set of eigenvectors, say {x<sub>1</sub>,..., x<sub>k</sub>} associated with a particular eigenvalue, say λ<sub>j</sub>, we can construct another set of eigenvectors {x'<sub>1</sub>,..., x'<sub>k</sub>} that are orthonormal, using the Gram-Schmidt process.
- Then, all eigenvectors can form an orthornormal basis.

L4(4)

### Example



• Example.  $\mathbf{A} = \begin{pmatrix} 3 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 3 \end{pmatrix}$ .  $p_{\mathbf{A}}(\lambda) = -(\lambda - 1)^2(\lambda - 7)$ , thus  $\lambda_1 = 1, \lambda_2 = 7$  $E_1 = \operatorname{span}[\begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}]$ ,  $E_7 = \operatorname{span}[\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}]$  $\circ (111)^{\mathsf{T}}$  is perpendicular to  $(-110)^{\mathsf{T}}$  and  $(-101)^{\mathsf{T}}$ 

• 
$$\begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}$$
 and  $\begin{pmatrix} -1/2 \\ -1/2 \\ 1 \end{pmatrix}$  (for  $\lambda = 1$ ) and  $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$  (for  $\lambda = 7$ ) are the orthogonal basis in  $\mathbb{R}^3$ .

• After normalization, we can make the orthonormal basis.

# Eigendecomposition



- Theorem. The following is equivalent.
  - (a) A square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  can be factorized into  $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ , where  $\mathbf{P} \in \mathbb{R}^{n \times n}$  and  $\mathbf{D}$  is the diagonal matrix whose diagonal entries are eigenvalues of  $\mathbf{A}$ .
  - (b) The eigenvectors of **A** form a basis of  $\mathbb{R}^n$  (i.e., The *n* eigenvectors of **A** are linearly independent)
- The above implies the columns of *P* are the *n* eigenvectors of *A* (because *AP* = *PD*)
- $\boldsymbol{P}$  is an orthogonal matrix, so  $\boldsymbol{P}^{\mathsf{T}} = \boldsymbol{P}^{-1}$
- **A** is symmetric, then (b) holds (Spectral Theorem).

## Example of Orthogonal Diagonalization (1)



- Eigendecomposition for  $\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$
- Eigenvalues:  $\lambda_1 = 1, \lambda_2 = 3$
- (normalized) eigenvectors:  $\boldsymbol{p}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \ \boldsymbol{p}_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$
- **p**<sub>1</sub> and **p**<sub>2</sub> linearly independent, so A is diagonalizable.

• 
$$\boldsymbol{P} = (\boldsymbol{p}_1 \ \boldsymbol{p}_2) = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$
  
•  $\boldsymbol{D} = \boldsymbol{P}^{-1} \boldsymbol{A} \boldsymbol{P} = \begin{pmatrix} 1 & 0 \\ 0 & 3 \end{pmatrix}$ . Finally, we get  $\boldsymbol{A} = \boldsymbol{P} \boldsymbol{D} \boldsymbol{P}^{-1}$ 

### Example of Orthogonal Diagonalization (2)

# KAIST A

•  $A = \begin{pmatrix} 1 & 2 & 2 \\ 2 & 1 & 2 \\ 2 & 2 & 1 \end{pmatrix}$ • Eigenvalues:  $\lambda_1 = -1, \lambda_2 = 5$   $(\alpha_1 = 2, \alpha_2 = 1)$ •  $E_{-1} = \operatorname{span} \begin{bmatrix} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} \end{bmatrix} \xrightarrow{\text{Gram-Schmidt}}$  $\operatorname{span} [\frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \frac{1}{\sqrt{6}} \begin{pmatrix} -1 \\ 1 \\ 2 \end{pmatrix} ]$ 

• 
$$E_5 = \operatorname{span}\left[\frac{1}{\sqrt{3}}\begin{pmatrix}1\\1\\1\\1\end{pmatrix}
ight]$$
  
•  $P = \begin{pmatrix}-1/\sqrt{2} & -1/\sqrt{6} & 1/\sqrt{3}\\1/\sqrt{2} & -1/\sqrt{6} & 1/\sqrt{3}\\0 & 2/\sqrt{6} & 1/\sqrt{3}\end{pmatrix}$   
•  $D = P^{\mathsf{T}}AP = \begin{pmatrix}-1 & 0 & 0\\0 & -1 & 0\\0 & 0 & 5\end{pmatrix}$ 

# Eigendecomposition: Geometric Interpretation





Question. Can we generalize this beautiful result to a general matrix  $A \in \mathbb{R}^{m \times n}$ ?

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

# Storyline



- Eigendecomposition (also called EVD: EigenValue Decomposition): (Orthogoanl) Diagonalization for symmetric matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .
- Extensions: Singular Value Decomposition (SVD)
  - 1. First extension: diagonalization for non-symmetric, but still square matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$
  - 2. Second extension: diagonalization for non-symmetric, and non-square matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$
- Background. For *A* ∈ ℝ<sup>m×n</sup>, a matrix *S* := *A*<sup>T</sup>*A* ∈ ℝ<sup>n×n</sup> is always symmetric, positive semidefinite.
  - Symmetric, because  $\boldsymbol{S}^{\mathsf{T}} = (\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A})^{\mathsf{T}} = \boldsymbol{A}^{\mathsf{T}}\boldsymbol{A} = \boldsymbol{S}.$
  - Positive semidefinite, because  $\mathbf{x}^{\mathsf{T}} \mathbf{S} \mathbf{x} = \mathbf{x}^{\mathsf{T}} \mathbf{A}^{\mathsf{T}} \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^{\mathsf{T}} (\mathbf{A} \mathbf{x}) \ge 0.$
  - If  $rk(\mathbf{A}) = n$ , then symmetric and positive definite.

## Singular Value Decomposition

Theorem. A ∈ ℝ<sup>m×n</sup> with rank r ∈ [0, min(m, n)]. The SVD of A is a decomposition of the form

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathsf{T}}, \qquad \qquad \boldsymbol{\varepsilon} \boldsymbol{A} = \boldsymbol{\varepsilon} \boldsymbol{U} \boldsymbol{\varepsilon} \boldsymbol{\Sigma} \boldsymbol{v}^{\mathsf{T}} \boldsymbol{\varepsilon}^{\mathsf{T}}$$

with an orthogonal matrix  $\boldsymbol{U} = (\boldsymbol{u}_1 \cdots \boldsymbol{u}_m) \in \mathbb{R}^{m \times m}$  and an orthogonal matrix  $\boldsymbol{V} = (\boldsymbol{v}_1 \cdots \boldsymbol{v}_n) \in \mathbb{R}^{n \times n}$ . Moreoever,  $\boldsymbol{\Sigma}$  s an  $m \times n$  matrix with  $\boldsymbol{\Sigma}_{ii} = \sigma_i \geq 0$  and  $\boldsymbol{\Sigma}_{ij} = 0, i \neq j$ , which is uniquely determined for  $\boldsymbol{A}$ .

- Note
  - The diagonal entries  $\sigma_i$ , i = 1, ..., r are called singular values.
  - $\boldsymbol{u}_i$  and  $\boldsymbol{v}_i$  are called left and right singular vectors, respectively.

### SVD: How It Works (for $\mathbf{A} \in \mathbb{R}^{n \times n}$ )

KAIST A

- $\mathbf{A} \in \mathbb{R}^{n \times n}$  with rank  $r \leq n$ . Then,  $\mathbf{A}^{\mathsf{T}}\mathbf{A}$  is symmetric.
- Orthogonal diagonalization of **A**<sup>T</sup>**A**:

$$\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A} = \boldsymbol{V}\boldsymbol{D}\boldsymbol{V}^{\mathsf{T}}$$

- $\boldsymbol{D} = \begin{pmatrix} \lambda_1 \\ \ddots \\ \lambda_n \end{pmatrix}$  and an orthogonal matrix  $\boldsymbol{V} = (\boldsymbol{v}_1 \cdots \boldsymbol{v}_n)$ , where  $\lambda_1 \geq \cdots \geq \lambda_r \geq \lambda_{r+1} = \cdots \lambda_n = 0$  are the eigenvalues of  $\boldsymbol{A}^T \boldsymbol{A}$  and  $\{\boldsymbol{v}_i\}$  are orthonormal.
- All  $\lambda_i$  are positive

$$\forall \boldsymbol{x} \in \mathbb{R}^{n}, \left\|\boldsymbol{A}\boldsymbol{x}\right\|^{2} = \boldsymbol{A}\boldsymbol{x}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x} = \boldsymbol{x}^{\mathsf{T}}\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x} = \lambda_{i}\left\|\boldsymbol{x}\right\|^{2}$$

- $\mathsf{rk}(\mathbf{A}) = \mathsf{rk}(\mathbf{A}^{\mathsf{T}}\mathbf{A}) = \mathsf{rk}(D) = \mathsf{r}$
- Choose  $\boldsymbol{U}' = \left( \boldsymbol{u}_1 \ \cdots \ \boldsymbol{u}_r \right)$ , where

$$oldsymbol{u}_i = rac{oldsymbol{A}oldsymbol{v}_i}{\sqrt{\lambda_i}}, \ 1 \leq i \leq r.$$

We can construct {u<sub>i</sub>}, i = r + 1, ..., n, so that U = (u<sub>1</sub> ... u<sub>n</sub>) is an orthonormal basis of ℝ<sup>n</sup>.

• Define 
$$\Sigma = \begin{pmatrix} \sqrt{\lambda_1} & \\ & \ddots & \\ & & \sqrt{\lambda_n} \end{pmatrix}$$

- Then, we can check that  $\boldsymbol{U}\boldsymbol{\Sigma}=\boldsymbol{A}\boldsymbol{V}.$
- Similar arguments for a general *A* ∈ ℝ<sup>m×n</sup> (see pp. 104)

# Example



• 
$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{pmatrix}$$
  
•  $\mathbf{A}^{\mathsf{T}} \mathbf{A} = \begin{pmatrix} 5 & -2 & 1 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} = \mathbf{V} \mathbf{D} \mathbf{V}^{\mathsf{T}},$   
 $\mathbf{D} = \begin{pmatrix} 6 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} \frac{5}{\sqrt{30}} & \frac{-2}{\sqrt{30}} & \frac{1}{\sqrt{30}} \\ 0 & \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ \frac{-1}{\sqrt{6}} & \frac{-2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \end{pmatrix}$ 

•  $rk(\mathbf{A}) = 2$  because we have two singular values  $\sigma_1 = \sqrt{6}$  and  $\sigma_2 = 1$ 

• 
$$\Sigma = \begin{pmatrix} \sqrt{6} & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$

• 
$$\boldsymbol{u}_1 = \boldsymbol{A}\boldsymbol{v}_1/\sigma_1 = \begin{pmatrix} \frac{1}{\sqrt{5}} \\ \frac{-2}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} \end{pmatrix}$$
  
•  $\boldsymbol{u}_2 = \boldsymbol{A}\boldsymbol{v}_2/\sigma_2 = \begin{pmatrix} \frac{2}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} \end{pmatrix}$   
•  $\boldsymbol{U} = \begin{pmatrix} \boldsymbol{u}_1 & \boldsymbol{u}_2 \end{pmatrix} = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 & 2 \\ -2 & 1 \end{pmatrix}$ 

• Then, we can see that 
$$\mathbf{A} = \mathbf{U} \Sigma V^{\mathsf{T}}$$
.

# KAIST A

# $\mathsf{EVD} \ (\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}) \text{ vs. SVD } (\mathbf{A} = \mathbf{U}\Sigma \mathbf{V}^{\mathsf{T}})$

- SVD: always exists, EVD: square matrix and exists if we can find a basis of eigenvectors (such as symmetric matrices)
- **P** in EVD is not necessarily orthogonal (only true for symmetric **A**), but **U** and **V** are orthogonal (so representing rotations)
- Both EVD and SVD: (i) basis change in the domain, (ii) independent scaling of each new basis vector and mapping from domain to codomain, (iii) basis change in the codomain. The difference: for SVD, different vector spaces of domain and codomain.
- SVD and EVD are closely related through their projections
  - The left-singular (resp. right-singular) vectors of **A** are eigenvectors of  $AA^{T}$  (resp.  $A^{T}A$ )
  - The singular values of **A** are the square roots of eigenvalues of  $AA^T$  and  $A^TA$
  - When **A** is symmetric, EVD = SVD (from spectral theorem)

# Different Forms of SVD



When rk(A) = r, we can construct SVD as the following with only non-zero diagonal entries in Σ:



 We can even truncate the decomposed matrices, which can be an approximation of *A*: for *k* < *r*

$$\boldsymbol{A} \approx \overbrace{\boldsymbol{U}}^{m \times k} \overbrace{\boldsymbol{\Sigma}}^{k \times k} \overbrace{\boldsymbol{V}^{\mathsf{T}}}^{k \times n}$$

We will cover this in the next slides.

# Matrix Approximation via SVD











(d)  $A_3, \sigma_3 \approx 26, 125.$ 

(e)  $A_4$ ,  $\sigma_4 \approx 20,232$ . (f)  $A_5$ ,  $\sigma_5 \approx 15,436$ .

- $\mathbf{A} = \sum_{i=1}^{r} \sigma_i \, \mathbf{u}_i \mathbf{v}_i^{\mathsf{T}}$ , where  $\mathbf{A}_i$  is the outer product<sup>3</sup> of  $\mathbf{u}_i$  and  $\mathbf{v}_i$
- Rank k-approximation:  $\hat{A}(k) = \sum_{i=1}^{k} \sigma_i A_i, \ k < r$

<sup>&</sup>lt;sup>3</sup>If  $\boldsymbol{u}$  and  $\boldsymbol{v}$  are both nonzero, then the outer product matrix  $\boldsymbol{u}\boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}$  always has matrix rank 1. Indeed, the columns of the outer product are all proportional to the first column.

# How Close $\hat{A}(k)$ is to A?

- Definition. Spectral Norm of a Matrix. For  $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ ,  $\|\boldsymbol{A}\|_2 := \max_{\boldsymbol{x}} \frac{\|\boldsymbol{A}\boldsymbol{x}\|_2}{\|\boldsymbol{x}\|_2}$ 
  - As a concept of length of A, it measures how long any vector x can at most become, when multiplied by A
- Theorem. Eckart-Young. For  $A \in \mathbb{R}^{m \times n}$  of rank r and  $B \in \mathbb{R}^{m \times n}$  of rank k, for any  $k \leq r$ , we have:

$$\hat{\boldsymbol{A}}(k) = \arg\min_{\mathsf{rk}(\boldsymbol{B})=k} \|\boldsymbol{A} - \boldsymbol{B}\|_2, \text{ and } \|\boldsymbol{A} - \hat{\boldsymbol{A}}(k)\|_2 = \sigma_{k+1}$$

- Quantifies how much error is introduced by the SVD-based approximation
- $\hat{A}(k)$  is optimal in the sense that such SVD-based approximation is the best one among all rank-k approximations.
- In other words, it is a projection of the full-rank matrix *A* onto a lower-dimensional space of rank-at-most-*k* matrices.

# Roadmap



- (1) Determinant and Trace
- (2) Eigenvalues and Eigenvectors
- (3) Cholesky Decomposition
- (4) Eigendecomposition and Diagonalization
- (5) Singular Value Decomposition
- (6) Matrix Approximation
- (7) Matrix Phylogeny

# Phylogenetic Tree of Matrices





L4(7)



# Questions?