

Lecture 05: A Refresher in Linear Algebra

Nikola Zlatanov

Innopolis University
Advanced Statistics

20-th of Feb to 27-th of Feb, 2023

Spectral Theorem

- Thm (Spectral Thm): \forall symmetric $n \times n$ matrix \mathbf{A} , the i -th eigenvalue and eigenvector are defined as the unit vector, \mathbf{u}_i , and scalar λ_i that satisfy the following equality

$$\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i \quad (1)$$

Moreover, the following holds

- all eigenvalues, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ are real,
- all unit eigenvectors, $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$, corresponding to eigenvalues, $\lambda_1, \lambda_2, \dots, \lambda_n$, respectively, are orthonormal, and therefore they form the orthonormal basis of \mathbb{R}^n

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

Spectral Decomposition (SD)

- Cor (Spectral Decomposition (SD)): \forall symmetric $n \times n$ matrix \mathbf{A} , can be written as

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

where $r = \text{rank}(\mathbf{A})$ and $\lambda_i = 0$, for $r < i \leq n$

Proof: $\forall \mathbf{x}$, can be written as $\mathbf{x} = \sum_{i=1}^n (\mathbf{u}_i^T \mathbf{x}) \mathbf{u}_i$, which is the basis expansion of \mathbf{x} . Multiplying both sides by \mathbf{A} , we get

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \sum_{i=1}^n (\mathbf{u}_i^T \mathbf{x}) \mathbf{A}\mathbf{u}_i \stackrel{(a)}{=} \sum_{i=1}^n (\mathbf{u}_i^T \mathbf{x}) \lambda_i \mathbf{u}_i = \sum_{i=1}^n \lambda_i \mathbf{u}_i (\mathbf{u}_i^T \mathbf{x}) \\ &= \left(\sum_{i=1}^n (\lambda_i \mathbf{u}_i \mathbf{u}_i^T) \right) \mathbf{x} \stackrel{(a)}{=} \mathbf{A}\mathbf{x} \end{aligned} \quad (3)$$

where (a) comes from (1) and (b) holds only if (2) holds. Q.E.D.

Matrix Form of the SD

- Thm (Matrix Form of the SD): Let $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. Then,

$$\mathbf{A} = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{bmatrix}$$

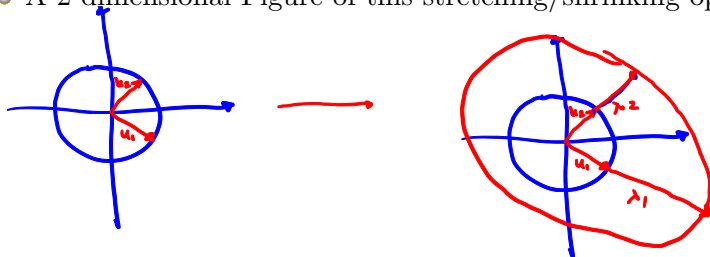
$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

Proof: Prove at home.

- The matrix \mathbf{U} is called an orthogonal matrix.
- A matrix \mathbf{U} is orthogonal if it satisfies any of the following equivalent properties
 - \mathbf{U} is invertible and $\mathbf{U}^{-1} = \mathbf{U}^T$
 - $\mathbf{U}\mathbf{U}^T = \mathbf{I}$
 - $\mathbf{U}^T\mathbf{U} = \mathbf{I}$
 - The columns of \mathbf{U} are mutually orthonormal
 - The rows of \mathbf{U} are mutually orthonormal

Geometric Form of the SD Thm

- Geometrically, we can view \mathbf{A} as a operator that stretches or shrinks the orthonormal basis $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ by $\lambda_1, \lambda_2, \dots, \lambda_n$, respectively.
- Hence, it transforms a n -dimensional unit sphere into an n -dimensional ellipsoid.
- A 2-dimensional Figure of this stretching/shrinking operation:



Optimization Form of the SD Thm:

- We can also view the SD as a series of optimization problems as

$$\lambda_1 = \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \mathbf{x}^T A \mathbf{x}, \text{ where: } \mathbf{u}_1 = \arg \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \mathbf{x}^T A \mathbf{x} \quad (4)$$

$$\lambda_2 = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{u}_1}} \mathbf{x}^T A \mathbf{x}, \text{ where: } \mathbf{u}_2 = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{u}_1}} \mathbf{x}^T A \mathbf{x} \quad (5)$$

\vdots

$$\lambda_i = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{u}_1, \dots, \mathbf{u}_{i-1}\}}} \mathbf{x}^T A \mathbf{x}, \text{ where: } \mathbf{u}_i = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{u}_1, \dots, \mathbf{u}_{i-1}\}}} \mathbf{x}^T A \mathbf{x} \quad (6)$$

\vdots

$$\lambda_r = \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \mathbf{x}^T A \mathbf{x}, \text{ where: } \mathbf{u}_r = \arg \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \mathbf{x}^T A \mathbf{x} \quad (7)$$

Singular Value Decomposition (SVD)

- All of the above holds for symmetric matrices. We need to extend this to non-symmetric matrices.
- Thm (Singular Value Decomposition (SVD)): \forall general $m \times n$ matrix \mathbf{A} , can be written as

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T, \text{ where } r = \text{rank}(\mathbf{A})$$

$$\text{and } \sigma_i = 0, \text{ for } r < i \leq \min\{m, n\} \quad (8)$$

where $\sigma_1 \geq \sigma_2 \geq \dots \sigma_r \geq 0$. The vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ are mutually orthonormal and thereby form an orthonormal basis on \mathbb{R}^m , and the vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are mutually orthonormal and thereby form an orthonormal basis on \mathbb{R}^n .

- The vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ are called “left singular value vectors” and the vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ are called “right singular value vectors”.

Informal Proof of SVD

- Since (8) holds, let's see what are $\mathbf{A}^T \mathbf{A}$ and $\mathbf{A} \mathbf{A}^T$

$$\begin{aligned} \mathbf{A} \mathbf{A}^T &= \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \left(\sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right)^T = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \sum_{j=1}^r \sigma_j \mathbf{v}_j \mathbf{u}_j^T \\ &= \sum_{i=1}^r \sum_{j=1, j \neq i}^r \sigma_i \mathbf{u}_i (\mathbf{v}_i^T \mathbf{v}_j) \mathbf{u}_j^T + \sum_{i=1}^r \sigma_i^2 \mathbf{u}_i (\mathbf{v}_i^T \mathbf{v}_i) \mathbf{u}_i^T = \sum_{i=1}^r \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^T \end{aligned}$$

Similairy,

$$\mathbf{A}^T \mathbf{A} = \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \rightarrow \sigma_i^2 = \lambda(\mathbf{A}^T \mathbf{A})$$

Hence, $\sigma_i^2 = \lambda_i(\mathbf{A} \mathbf{A}^T) = \lambda_i(\mathbf{A}^T \mathbf{A})$. Vectors \mathbf{u}_i 's and \mathbf{v}_i 's are the eigenvectors of $\mathbf{A} \mathbf{A}^T$ and $\mathbf{A}^T \mathbf{A}$, respectively. Now finish the proof.

Matrix form of the SVD

- Thm (Matrix form of the SVD): Let $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]$, $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]$ and $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{\min\{m,n\}})$, where $\mathbf{\Sigma}$ is a diagonal $m \times n$ matrix. Then,

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

Proof: Prove at home.

Geometric Form of the SVD Thm

- Geometric Form of the SVD Thm: We can see \mathbf{A} as a operator that first rotates $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ by $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r$, respectively, and then stretches/shrinks the corresponding rotated vectors by $\sigma_1, \sigma_2, \dots, \sigma_n$, respectively.
- Hence, it rotates a n -dimensional unit sphere to another sphere and then transforms it into an n -dimensional ellipsoid.
- A 2-dimensional Figure of this stretching/shrinking operation:



Optimization Form of the SVD Thm:

• Optimization Form of the SVD Thm:

$$\sigma_1 = \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2, \text{ where: } \mathbf{v}_1 = \arg \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 \quad (9)$$

$$\sigma_2 = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{v}_1}} \|\mathbf{A}\mathbf{x}\|_2, \text{ where: } \mathbf{v}_2 = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{v}_1}} \|\mathbf{A}\mathbf{x}\|_2 \quad (10)$$

\vdots

$$\sigma_i = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{v}_1, \dots, \mathbf{v}_{i-1}\}}} \|\mathbf{A}\mathbf{x}\|_2, \text{ where: } \mathbf{v}_i = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{v}_1, \dots, \mathbf{v}_{i-1}\}}} \|\mathbf{A}\mathbf{x}\|_2 \quad (11)$$

\vdots

$$\sigma_r = \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2, \text{ where: } \mathbf{v}_r = \arg \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 \quad (12)$$

Optimization Form of the SVD Thm:

Moreover

$$\sigma_1 = \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{x}^T \mathbf{A}\|_2, \text{ where: } \mathbf{u}_1 = \arg \max_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{x}^T \mathbf{A}\|_2 \quad (13)$$

$$\sigma_2 = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{v}_1}} \|\mathbf{x}^T \mathbf{A}\|_2, \text{ where: } \mathbf{u}_2 = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \mathbf{u}_1}} \|\mathbf{x}^T \mathbf{A}\|_2 \quad (14)$$

\vdots

$$\sigma_i = \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{u}_1, \dots, \mathbf{u}_{i-1}\}}} \|\mathbf{x}^T \mathbf{A}\|_2, \text{ where: } \mathbf{u}_i = \arg \max_{\substack{\mathbf{x}: \|\mathbf{x}\|_2=1 \\ \mathbf{x} \perp \{\mathbf{u}_1, \dots, \mathbf{u}_{i-1}\}}} \|\mathbf{x}^T \mathbf{A}\|_2 \quad (15)$$

\vdots

$$\sigma_r = \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{x}^T \mathbf{A}\|_2, \text{ where: } \mathbf{u}_r = \arg \min_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{x}^T \mathbf{A}\|_2 \quad (16)$$

Optimization Form of the SVD Thm:

- In all previous optimization problems, we can remove the restriction of searching only over unit vectors, i.e., $\|\mathbf{x}\|_2 = 1$, by optimizing $\frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}$ instead of $\|\mathbf{Ax}\|_2$ and by optimizing $\frac{\|\mathbf{x}^T \mathbf{A}\|_2}{\|\mathbf{x}\|_2}$ instead of $\|\mathbf{x}^T \mathbf{A}\|_2$
- Hence,

$$\sigma_1 = \max_{\mathbf{x}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2},$$

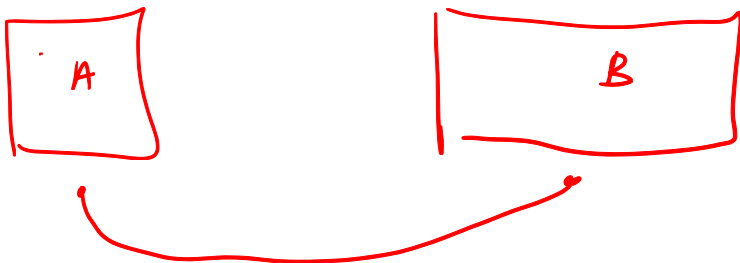
$$\sigma_r = \min_{\mathbf{x}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2},$$

- And this can be looked at as the maximum and minimum possible distortion that the transformation \mathbf{A} can impose on any vector \mathbf{x} .
 Example:

$$\sigma_r \|\mathbf{x}\|_2 \leq \|\mathbf{Ax}\|_2 \leq \sigma_1 \|\mathbf{x}\|_2, \quad \forall \mathbf{x}$$

Matrix Norms

- Later on we would like to know how close is matrix A to B
- Figure:



- Hence, we need some metric that measures distances between two matrices
- There are many such natural metrics. The most easy ones are the Frobenius norm and the Operator norm

Frobenius Norm

- Def: Frobenius (a.k.a. Hilbert Schmidt) norm of matrix \mathbf{A} , denoted by matrix $\|\mathbf{A}\|_F^2$ is

$$\|\mathbf{A}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^m a_{ij}^2 = \text{tr}(\mathbf{A}^T \mathbf{A}) \quad (17)$$

- Hence, Frobenius norm of matrix \mathbf{A} is simply a sum of all of its squared elements.
- Note that Frobenius norm does not distinguish between rows and columns. It pretends that the matrix is a vector.
- Def: Frobenius inner (dot) product between matrix \mathbf{A} and \mathbf{B} , denoted by $\langle \mathbf{A}, \mathbf{B} \rangle_F$, is

$$\langle \mathbf{A}, \mathbf{B} \rangle_F = \sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{ij} = \text{tr}(\mathbf{A}^T \mathbf{B}) \quad (18)$$

Hence,

$$\langle \mathbf{A}, \mathbf{A} \rangle_F = \sum_{i=1}^n \sum_{j=1}^m a_{ij}^2 = \text{tr}(\mathbf{A}^T \mathbf{A}) = \|\mathbf{A}\|_F^2 \quad (19)$$

Frobenius Norm

Properties:

- Orthogonal Invariance: For any orthogonal matrices U and V

$$\|UA\|_F = \|AV\|_F = \|A\|_F^2$$

Proof:

$$\text{tr}(A) = \text{tr}(A^T)$$

$$\|UA\|_F^2 = \langle UA, UA \rangle_F = \text{tr}(AU^T UA) = \text{tr}(AA) = \|UA\|_F^2$$

- Euclidean norm of singular values: $\|A\|_F^2 = \sum_{i=1}^r \sigma_i^2$

Proof:

$$\text{tr}(AA^T)$$

$$\begin{aligned} \|A\|_F^2 &= \|U^T \Sigma V\|_F^2 = \text{tr}(V^T \Sigma^T U U^T \Sigma V) = \text{tr}(V^T \Sigma^T \Sigma V) \\ &= \text{tr}(\Sigma^T V V^T \Sigma) = \text{tr}(\Sigma^T \Sigma) = \sum_{i=1}^r \sigma_i^2 \end{aligned} \quad (20)$$

Operator Norm

- Def: Operator norm is the largest singular value i.e.,

$$\|\mathbf{A}\|_{op} = \sigma_1 = \max_{\mathbf{x}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}, \quad (21)$$

- Properties:

$$\|\mathbf{A}\|_{op} \leq \|\mathbf{A}\|_F \leq \sqrt{n} \|\mathbf{A}\|_{op} \quad (22)$$

Proof:

$$\sigma_1^2 \leq \sum_{i=1}^r \sigma_i^2 \leq n \sigma_1^2 \quad (23)$$

taking square root, we obtain the desired result.